

Comprehending Queries

Torsten Grust

**Dissertation der Universität Konstanz
Fakultät für Mathematik und Informatik
1999**

1. Referent: Prof. Dr. Marc H. Scholl, Universität Konstanz
2. Referent: Prof. dr. Peter M.G. Apers, Universität Twente (Niederlande)

Tag der mündlichen Prüfung: 30. September 1999

Abstract

There are no compelling reasons why database-internal query representations have to be designated by *operators*. This text describes a world in which *datatypes* determine the comprehension of queries. In this world, a datatype is characterized by its algebra of value constructors. These algebras are principal. Query operators are secondary in the sense that they simply box (recursive) programs that describe how to form a query result by application of datatype constructors. Often, operators will be unboxed to inspect and possibly rewrite these programs. Query optimization then means to deal with the transformation of programs.

The predominant role of the constructor algebras suggests that this model understands queries as mappings between such algebras. The key observation that makes the whole approach viable is that (a) *homomorphic* mappings are expressive enough to cover declarative user query languages like OQL or recent SQL dialects, and, at the same time, (b) a single program form suffices to express homomorphisms between constructor algebras. Reliance on a single combining form, *catamorphisms*, renders the query programs susceptible to *Constructive Algorithmics*, an effective and extensive algebraic theory of program transformations.

The text then takes a step from catamorphisms towards a higher-level query representation based on the categorical notion of *monads*. In a nutshell, monads are algebras exhibiting exactly the structure that is needed to support the interpretation of a query calculus, the *monad comprehension calculus*. Built on top of the abstract monad notion, the calculus maps a variety of query constructs (*e. g.*, bulk operations, aggregates, and quantifiers) to few syntactic forms. The uniformity of the calculus facilitates the analysis and transformation, especially the normalization, of its expressions. Few but generic calculus rewriting rules suffice to implement query transformations that would otherwise require extensive rule sets.

The text rediscovers well-known query optimization knowledge on sometimes unusual paths that are more practicable to follow for an optimizer, though. Solutions previously proposed by others can be simplified and generalized mainly due to the clear account of the structure of queries that the monad comprehension calculus—thanks to its density—provides. The calculus effectively supports query optimization in the presence of grouping, various forms of nesting, aggregates, and quantifiers. Although built on top of abstract concepts like homomorphisms and monads, this query model is specific enough to grasp implementation issues, such as the generation of stream-based (pipelined) query execution plans, whose treatment has traditionally been delayed until query runtime.

It is the main objective of this thesis to show that catamorphisms and monad comprehensions enable a comprehension of queries that is *effective* and easily *exploitable* inside a query optimizer.

Zusammenfassung (Summary in German)

Obwohl es dafür keine zwingenden Gründe gibt, setzen Anfrageoptimierer fast ausschließlich interne Anfragerepräsentationen ein, die durch *Operatoren* charakterisiert sind. In diesem Text beschreiben wir dagegen eine Welt, in der *Datentypen* das Verständnis von Anfragen bestimmen. Innerhalb dieser Welt wird ein Datentyp durch die Algebra seiner Konstruktoren beschrieben. Diese Konstruktoralgebren sind zentral. Operatoren hingegen sind zweitrangig: ein Operator ist lediglich eine Box in der ein (rekursives) Programm eingekapselt ist, welches die Berechnung von Anfrageergebnissen durch die Anwendung von Konstruktoren beschreibt. Immer wieder werden wir diese Box aufbrechen, um das Verhalten des Operators genauer zu studieren oder sein internes Programm zu transformieren. Anfrageoptimierung bedeutet dann vor allem, sich mit Programmtransformationen auseinanderzusetzen.

In einer Welt, die durch Datentypen und ihre Konstruktoralgebren charakterisiert ist, sind Anfragen Abbildungen zwischen diesen Algebren. Zwei grundlegende Beobachtungen machen diese Sicht auf Anfragen praktikabel: (a) schon die *Homomorphismen* zwischen Konstruktoralgebren sind bereits ausdrucksstark genug, um deklarative Anfragesprachen wie OQL oder neuere SQL-Dialekte zu begreifen und (b) gleichzeitig ist *eine* Programmform ausreichend, um diese Homomorphismen zu implementieren. Der disziplinierte Einsatz dieser Programmform, die *Catamorphismen*, macht die Programme zugänglich für die *Constructive Algorithmics*, eine umfassende und zugleich effektive algebraische Theorie der Programmtransformation.

Der nächste Schritt führt uns dann von Catamorphismen zu einer Anfragedarstellung, die auf dem kategoriellen Begriff der *Monade* basiert. Monaden besitzen exakt die algebraische Struktur, die für die Interpretation eines generischen Anfragekalküls, dem *Monad Comprehension Kalkül*, notwendig ist. Der Monad Comprehension Kalkül erlaubt lediglich eine geringe Anzahl syntaktischer Formen, die jedoch ein breites Spektrum von Anfragekonstrukten (mengenorientierte Operationen, Aggregation, Quantifikation) uniform abbilden können. Die Analyse und Transformation, besonders die Normalisierung, der Ausdrücke dieses Kalküls profitiert signifikant von dieser Uniformität. Wenige generische Transformationsregeln im Kalkül können Anfragetransformationen beschreiben, die sonst durch umfassende Regelmengen implementiert werden.

Der Text entdeckt bekannte Optimierungstechniken teilweise erneut und bewegt sich dabei auf ungewöhnlichen Pfaden, die von einem Optimierer jedoch effektiv nachvollzogen werden können. Mehrmals können wir etablierte

Strategien vereinfachen und generalisieren, was vor allem durch die effektiv zu analysierende Struktur ermöglicht wird, die der Monad Comprehension Kalkül Anfragen aufprägt. Der Kalkül unterstützt die Optimierung von Anfragen, die Gruppierung, verschiedene Formen von Nestung, Aggregate und Quantoren einsetzen. Obwohl abstrakte Konzepte wie Homomorphismen und Monaden die Grundlage liefern, ist diese Anfragerepräsentation spezifisch genug, um auch die eigentliche Anfrageimplementation zu unterstützen (wie etwa die Generierung von Anfrageplänen, die die Erzeugung von Zwischenresultaten vermeiden). Damit können Probleme statisch angegangen werden, die traditionell dynamisch zur Anfrageausführungszeit betrachtet wurden.

Catamorphismen und Monad Comprehension Kalkül ermöglichen eine Anfragerepräsentation, die innerhalb eines Optimierers *effektiv analysierbar und manipulierbar* ist. Dies zu zeigen, ist das eigentliche Ziel dieser Arbeit.

Acknowledgments

I have no doubt that you can accomplish nothing without a challenging and lively setting. During the last five years I was lucky enough to live and work in such an environment and—in some sense—this thesis is my humble try to acknowledge the scientific and personal support I have received throughout.

Marc Scholl, my advisor, somehow managed to step up onto the stage in just the right moment and I lost little time to accept his invitation to join his newly founded research group at the University of Konstanz. To me, Marc has been the ideal advisor. I have never faced, literally, closed doors when I was seeking for his guidance or help. At the same time he gave me all the freedom you need to get your bearing as a young researcher. There have been periods where my work has led me away from the core database research field into the world of functional programming, but Marc never hesitated to encourage me to follow this alien path. Maybe that's what I value the most about him.

Throughout the five years I had many, often amazing, encounters with the members of our research community. During a visit at the University of Twente I got to know Peter Apers and his group and I'm particularly grateful that Peter accepted to review this thesis. I have found the discussions with him, Annita Wilschut, and David Spelt especially exciting. Without the help of Maarten Fokkinga there would be more errors than you see here now.

Sometimes even short exchanges of thoughts can change the way you see things. I had such experiences when I met Peter Buneman (and the participants of the FDM workshop in London, July 1997), Erik Meijer, Doaitse Swierstra, and Phil Wadler in person, and during e-mail conversation with Mitch Cherniack and Phil Trinder.

I often felt that there was exactly the right mixture of people around in the Konstanz Computer Science department. The students (especially Andres Löh) and my colleagues proved to be constant sources of useful feedback and good humor. Ulrik Brandes and Dieter Gluche, my fellow contemporaries in the department, never failed to be sympathetic towards my matters. (Probably that's only possible if you share the same goals as well as uncounted rounds of beer.)

Finally, I owe thanks to Brita. Somehow she reliably managed to distract me when there was not much more on my mind than banana $\not\perp$ brackets.

Contents

1	Prelude	1
2	Categorical Datatypes and Monads	5
2.1	A Datatype-based Query Representation	5
2.2	Basic Category Theory	7
2.3	Categorical Collection Type Constructors	12
2.4	Catamorphisms	16
2.5	Type Functors	24
2.6	Fusion	28
2.7	Datatype Equations	32
2.8	Monads	38
2.9	Monad Comprehensions	46
2.10	Union Representation	52
2.11	Monad Comprehensions in Functional Programming	58
3	Query Compilation and Normalization	63
3.1	A Catamorphic Interpretation of OQL	65
3.2	Monad Comprehension Normalization	85
4	Combinators	95
4.1	Defining and Detecting Combinators	98
4.2	Combinators From Monad Comprehensions	102
5	Comprehending Queries	113
5.1	Parallelizing Group-By Queries	114
5.2	Aggregation by Decomposition	117
5.3	Normalization of KOLA Combinator Queries	122
6	Query Deforestation	127
6.1	Streaming	128
6.2	Cheap Deforestation	131

7 Finalization	137
References	141
Index	155

Chapter 1

Prelude

1 How would you go about and try to comprehend database queries?

Open the cover and start to dismantle a database system. Then unbox the query engine and disassemble its parts. You end up with a set of query operators that may be combined in various but well-defined ways, the engine's operator algebra. The algebra tells you how the operators fit together and you start to play and combine operators to form queries. But as you reach for a *join* operator you hear clatter. In fact, all operators clatter as you shake them. Apparently you cannot develop a complete comprehension of the query engine if you do not further unbox the inner workings of the operators.—This text is an exploration of what you will discover as soon as you break an operator's case.

2 In the course of this exploration we will soon realize that there is a single principle action that is pervasive inside all operators: the *construction* of values, which we will represent by the function symbol *cons*. Once we unfold the operators and inspect their definitions, we will find these operators to merely provide structure—a program—that controls application(s) of *cons*.

We will encounter construction in various instantiations, *e. g.*, as insertion of an element into a collection or incremental computation of an aggregate value, but these instantiations share so many properties that we will often do without telling them apart. Unlike other query models, we do not let collection types play an exceptional role. This sets the scene for a comprehension of queries which acknowledges the presence of query constructs other than bulk operations.

The predominance of *cons* motivates the starting point of our exploration. From the start, we let the construction of values dominate our understanding of queries and then work our way bottom-up. To effectively reason about construction we will exploit algebras of value constructors. In fact, these

are the primary algebras we will work with and it is the programs that build terms over these algebras that cause the clatter you have heard. As operators merely box such terms we can study the action of operators by actually examining how they construct values. The resulting operator algebra, however, is secondary in this text.

3 Unboxing the operators gives us a rather fine control over value construction and we will see how a query optimizer can benefit from this control. At the same time, unboxing also implies that we have to deal with what we find inside: programs. Query analysis and transformation in this model amounts to analyze and transform programs. Throughout the entire text we will make good use of techniques native to the program transformation domain and establish well-known as well as invent novel query transformations this way.

To base a query optimizer on these techniques means that we have to be restrictive about the program forms we may admit. Only then we can assure that the optimizer can operate as a program transformation system free of the need for external guidance or *Eureka steps*.

Here, the algebras of value constructors provide a point of reference. The only program forms we will admit are those that mimic the structure (of the recursive type) of the values they analyze and construct, *i. e.*, those that perform structural recursion. This restrictive discipline will render programs as *homomorphic mappings* between algebras of value constructors. It is worthwhile to dwell on this thought a little longer.

Let F denote a function that, given a type A as argument, encodes the structure of terms (expressions) that construct values of type A . A constructor algebra α then is nothing but a mapping from FA to A , in symbols $\alpha:FA \rightarrow A$. Consider a second algebra of identical structure $\beta:FB \rightarrow B$. If h denotes a query, in our model, it is a mapping between values, say $h:A \rightarrow B$. The above restriction, however, says that the program for h may not be arbitrary but has to respect the term structure F . Such a structure-preserving computation will be denoted Fh in this text. To get a bearing on this process, the overall picture of the involved mappings is

$$\begin{array}{ccc}
 FA & \xrightarrow{\alpha} & A \\
 \downarrow Fh & & \downarrow h \\
 FB & \xrightarrow{\beta} & B
 \end{array}$$

Intuitively, h meets the restrictions, if both paths from FA to B denote the

same function. We will trade this intuition for precise statements about query programs using the language of category theory. Basic categorical vocabulary suffices, however. We perceive category theory as the vehicle not the cargo of this text.

4 The restrictions we impose on query programs may seem rather rigid at first sight but actually they are not. The expressive power of these programs is sufficient to cover orthogonal query languages for complex value databases, like OQL or newer dialects of SQL. Everything can be reduced to a single recursive program form, the *catamorphism*, which provides all the control structure we need. At the same time, this restriction can lead to new insights into compositions of programs—and thus complex queries—due to the properties catamorphisms exhibit.

To narrow the gap between user level query syntax and catamorphisms, we will exploit a calculus, the *monad comprehension calculus*, as a mediator between the two worlds. For now, think of monads as algebras that offer just the right measure of structure to interpret the constructs of a query calculus. Monad comprehensions will be our query representation of choice throughout major portions of this text as they are accessible to the human eye as well as an effective way to manipulate queries inside an optimizer. Once we remove the calculus’ syntactic sugar, however, we realize that we are still operating with catamorphisms.

5 This work draws ideas and methods on a variety of sources, some of which are somewhat alien to the query optimization domain. The text continuously walks the fine line—if there is any—between query optimization, category theory, program transformation, type theory, and functional programming.

(a) We let *category theory* play the role that set theory has in the world of relational databases. The categorical view provides a measure of abstraction that enables important generalizations and elegant reasoning at the same time. While set-theoretic accounts of query optimization dominate the field of research by far, others have paved the way for a categorical model of queries [20, 114, 124].

(b) There exists an extensive theory, the *Constructive Algorithmics* [7, 8, 9, 42, 57, 67, 84], on the transformation of programs that were built from a small set of combining forms. This theory understands programs as objects that are subject to calculation just like numbers in arithmetic. We will establish core query transformations through calculation with programs.

(c) Type theory, especially *parametric polymorphism* and the laws it justifies for free [116, 122], constitutes another field we will benefit from.

(d) Last but not least, we perceive query transformation and optimization as a *functional programming activity*. Superficially, this concerns a number of notational conventions we adopted. More deeply, note that we generate query results solely through the side-effect free construction of values from simpler constituents. In fact, we find an approach to query optimization that does otherwise hard to imagine: referential transparency is the key to painless transformational programming and equational reasoning. Functional composition will be the predominant way of forming complex queries from simpler ones. Finally, when it comes to the generation of query execution plans, we will establish connections to implementation techniques for lazily evaluated functional programming languages [18, 51, 92].

6 The following exploits this toolbox to develop a comprehension of queries that is geared to be *effective* and easily *exploitable* inside a database query optimizer.

At places, we will rediscover well-known query optimization knowledge using unusual paths that are more practicable to follow for an optimizer, however (Chapter 4). At places, we can generalize and at the same time simplify solutions that have been proposed by others (Chapter 5). The query model is abstract enough to stress the common ground of a diversity of query constructs from bulk operations to quantifiers. This makes the model an ideal target for the translation of declarative OQL-like user query languages, including recent feature additions to these languages. The monad comprehension calculus effectively supports optimization in presence of, *e. g.*, grouping, nesting, quantifiers, and aggregates in queries (Chapter 3). The query model is specific enough to provide the necessary hints and handles to serve as an effectively manipulable representation of query execution plans—this provides a static account (*i. e.*, at query optimization time) of query runtime issues that have been traditionally tackled on the implementation level only. Finally, the model has already shown its suitability as a platform on which the rapid prototypical development of a query engine for complex value databases is viable (Chapter 6).

To get us going, Chapter 2 lays the categorical foundation upon everything else will rest.¹ You feel better if you know where that clatter actually comes from.

¹To quote Philip Wadler, “*No knowledge of category theory is assumed.*” [117]

Chapter 2

Categorical Datatypes and Monads

2.1 A Datatype-based Query Representation

7 *Datatypes* lie at the very heart of the query representation used throughout this text. The language we are going to develop in this and the upcoming chapter is designed around types, *not* operations. All datatypes which we will encounter in the sequel themselves induce a minimal set of functions which allow for the construction of arbitrary values of that type. Exactly these *constructors* (and their functional composition) will form the initial query language we will propose here. Once we have thus marked the start of our path, we will then pave our way towards a query representation that will be more suitable for particular stages of the query compilation process as well as the human mind and eye. We will often see, however, that query optimization can benefit from working with this rather fine granulated view of datatypes.

This does *not* mean that an actual implementation of our representation is forced to always operate at this fine level of granularity. It rather provides us with handles to manipulate small granules during query translation and optimization. Before a query is shipped to the underlying query engine for execution it is then our responsibility to identify those parts of the query expression which may be executed in a *bulk-oriented* fashion.

Note that there are *no mutators* defined for the different types. Query evaluation in this framework will be solely concerned with the construction of instances of a type from simpler constituents. The language will thus be *referentially transparent* or, put another way, *purely functional*. This feature, in turn, will smooth the way for query optimization based on equational reasoning (“*replace equals by equals.*”)

As query representation (not implementation) is the focus of this chapter

we view datatypes as being *abstract*. We will devise an interface—which will consist of the list of constructors and, in some cases, certain *equations* that have to be obeyed—for the types and leave their implementation for later. Any datatype providing a compatible interface may be incorporated into language without significant effort.

8 The instances of a type and its constructors will be modeled as an *algebra* for that type. There may be more than one algebra for a particular type. Our presentation of the material will rely on the tools of *category theory*. Category theory provides a succinct notation for the concepts involved and seems to supply just the right measure of abstraction that we strive for in this chapter. Moreover, the query language’s core operators will turn out to be familiar categorical concepts. At a first sight the formal machinery might seem awkward, but important questions like the guaranteed existence or uniqueness of the operators so defined are answered for free by the underlying theory. Categorical insights will prove useful in all query translation phases, notably during the generation of query execution plans (see Chapter 6). Following ideas of Fokkinga [42, 43], we will not only express the structural aspects of the constructor algebras using category theory but also the equations which certain types are expected to fulfill. This will save us from introducing the otherwise necessary machinery of signatures for algebras, syntax for terms, as well as variables and their (scoped) bindings.

The discussion in this chapter proceeds as follows. Right after we have provided the categorical view of algebras we will examine *catamorphisms*, functions whose recursion pattern mimics that of its source type. Despite their simplicity, catamorphisms provide the core concept upon which almost all of the query language will be built. One can prove, given few prerequisites, that compositions of catamorphisms collapse into a single catamorphism, which will later justify efficient optimizations. We will then generalize our observations to cover polymorphic type constructors (like `bag α` , α denoting an arbitrary type) as these are basic building blocks of type systems for query languages for complex collection types. Being so far, the notion of *equation* is formalized which will make level for the categorical representation of datatypes with commutative and/or idempotent constructor functions like `bag` and `set`. We already include a step towards a more high-level calculus-like query representation, *monad comprehensions*, which is shown to be quite easily definable on top of the algebraic datatypes presented so far.

The material presented here is not original but the arranging is unique in the sense that it is geared towards the problem of the formalization of datatypes with equations for database query languages. This is especially

true for the derivation of multi-monad comprehensions. We refer to [10, 42, 43, 84, 117] for an extensive categorical treatment of the ideas developed in the following paragraphs.

9 To let the reader gain familiarity with the involved notation and style of reasoning, we will provide more detailed lines of arguments in this chapter than in the rest of the thesis. Equational proofs (also called *calculations* or *rewrites* in the sequel) use the style adopted by the *Squiggol* [7, 11] and *Constructive Algorithmics* [7, 42, 57, 67, 84] communities because it allows for the convenient annotation of a sequence of replacements of equals by equals or similar connected chains of reasoning. In [3], [Backhouse](#) convincingly discusses the merits and benefits of this notation.

An equation display like (let f, g, h denote functions)

$$\begin{aligned} & (f \cdot g) = h \\ \equiv & \quad \{ f \cdot g = h \} \\ & h = h \\ \equiv & \quad \{ \text{reflexivity of } = \} \\ & \text{true} , \end{aligned}$$

represents the derivation of the last term from the first by step-wise replacement of equals by equals. The intermediate results are given. Justifications for the replacements are shown explicitly in curly braces $\{\cdot\}$. You may have noted that the above calculation is carried out at the function level and that no function is ever applied to a particular argument. Such calculations are called *point-free* (in contrast to *point-wise*). We shall mainly use point-free reasoning throughout this text as the introduction of a dummy function argument (x say) is not needed due to the principle of extensionality and merely clutters the display of the calculation.

2.2 Basic Category Theory

10 Category theory. A few basic concepts of category theory are sufficient to master the theory of datatypes we will use as the basis to build our query language on. These are: *category* (and the inherently related *object*, *morphism*, *target*, *source*), *initiality*, *isomorphism*, *functor*, *product*, and *sum*. Later on we will meet *natural transformations* and *monads*. We give a short introduction to all of these concepts before we delve into categorical datatypes. This is not an introduction to category theory in its own right. Our perception of the material is heavily biased towards the goal we try to

reach: a concise and formal yet sufficiently abstract framework in which we can talk about polymorphic (collection) type constructors. The theory does even better: it will also provide us with fundamental query operators as these are induced by the types and their algebras themselves.

Introductory texts on category theory may be found in [2, 10, 42, 63, 64, 79, 98, 103, 104]. Among these texts, Bird and de Moor [10] put particular emphasis on the exploration of categorical datatypes.

11 Category. A *category* is a system of *objects* and *morphisms* (also called *arrows*). Fix a category \mathbf{C} (categories will be written using a **bold font**). Every morphism f in \mathbf{C} uniquely determines two \mathbf{C} -objects, its *source* ($\text{src } f$) and *target* ($\text{tgt } f$). If $\text{src } f = A$ and $\text{tgt } f = B$ we write $f: A \rightarrow_{\mathbf{C}} B$, $A \rightarrow B$ is called the *type* of f in \mathbf{C} . If the category \mathbf{C} is clear from the context we declare it to be the default category and write $f: A \rightarrow B$ instead.

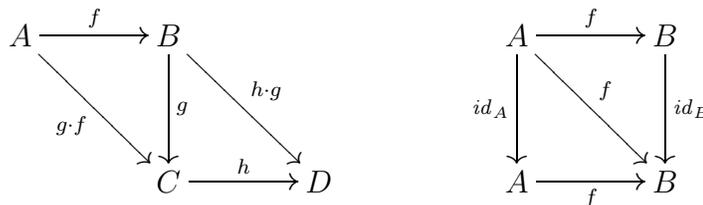
Two \mathbf{C} -morphisms $f: A \rightarrow B$ and $g: B \rightarrow C$ may be *composed* to obtain the morphism $g \cdot f$ of type $A \rightarrow C$ in \mathbf{C} (*closure*). Composition is associative, *i. e.*, for morphisms $f: A \rightarrow B$, $g: B \rightarrow C$, and $h: C \rightarrow D$ we have

$$h \cdot (g \cdot f) = (h \cdot g) \cdot f .$$

Each object A comes equipped with an *identity morphism* $\text{id}_A: A \rightarrow A$ which acts as a neutral element with respect to composition (let $f: A \rightarrow B$):

$$f \cdot \text{id}_A = f = \text{id}_B \cdot f .$$

It is sometimes insightful to display objects and morphisms as *diagrams*, *i. e.*, directed graphs having their nodes and edges labeled by objects and morphisms, respectively. The diagram *commutes*, if any two paths (compositions of morphisms) connecting the same pair of objects denote the same morphism. Associativity of composition and neutrality of identity morphisms may thus be depicted as two commutative diagrams:



While we will use diagrams to illustrate the typing of morphisms, we shall rarely employ them to conduct proofs. Proofs are carried out in the calculational style introduced in Paragraph 9.

In what follows we fix the default category to be the category **Set**. In **Set**, objects denote sets while a morphism f is to be interpreted as a total function with domain $\text{src } f$ and range (codomain) $\text{tgt } f$. If we interpret composition of morphisms as function composition and let id_A denote the identity function $\lambda x.x$ on A , we obey the category axioms. The type of a morphism f may then be understood as the type of a function f in the sense of typed functional programming languages. This is exactly the interpretation we will need during our discussion of datatypes.

12 Initiality. An object 1 in category **C** is *the initial object*, if **C** contains exactly one morphism of type $1 \rightarrow A$ for any object A of **C**. 1 is uniquely determined *up to isomorphism*: let $1'$ be another initial object. Then there are morphisms $f:1 \rightarrow 1'$ (by initiality of 1) and $g:1' \rightarrow 1$ (by initiality of $1'$). Since the identity morphisms $\text{id}_1:1 \rightarrow 1$ and $\text{id}_{1'}:1' \rightarrow 1'$ are the only arrows of these types, we can conclude that $g \cdot f = \text{id}_1$ and $f \cdot g = \text{id}_{1'}$. These two equations exactly establish the categorical notion of 1 and $1'$ being isomorphic. The initial object in **Set** is \emptyset .

In due course, the initiality of an algebra for a datatype will assert the existence as well as uniqueness of the fundamental operators for that datatype.

13 Functor. A *functor* $F:\mathbf{C} \rightarrow \mathbf{D}$ is a mapping from category **C** to category **D**. This implies that F specifies an object as well as a morphism mapping. Each object A of **C** is mapped to its image FA in **D**. The morphism mapping “preserves” arrows: if $f:A \rightarrow_{\mathbf{C}} B$ then $Ff:FA \rightarrow_{\mathbf{D}} FB$. Additionally, F is required to preserve identities and to distribute through composition:

$$F\text{id}_A = \text{id}_{FA} \quad \text{and} \quad F(g \cdot f) = Fg \cdot Ff. \quad \text{FUNCTOR}$$

Note that if $G:\mathbf{D} \rightarrow \mathbf{E}$ is another functor, then the composition $GF:\mathbf{C} \rightarrow \mathbf{E}$, with $GF A$ defined to mean $G(FA)$ (likewise for the morphism mapping), is a functor as is the identity functor $\text{ld}:\mathbf{C} \rightarrow \mathbf{C}$ with $\text{ld} A = A$ and $\text{ld} f = f$. (Categories and functors are the objects and morphisms in category **Cat** respectively.)

Which role do functors play in the context of categorical datatypes? A functor $F:\mathbf{Set} \rightarrow \mathbf{Set}$ (also called an *endofunctor* in **Set**) can perfectly model a *type constructor* or, alternatively, *type former*. Functional programming languages like Haskell or ML provide a type constructor `list` that lifts any type A to the type of lists of A , namely `list A`. This is the object mapping part. If we take `list f = map f` (`map` being the standard higher-order function

that applies a function to all elements of a list) we obtain a valid morphism mapping, because

$$f : A \rightarrow B \quad \Rightarrow \quad \text{map } f : \text{list } A \rightarrow \text{list } B$$

and

$$\text{map } id_A = id_{\text{list } A} \quad \text{and} \quad \text{map } (f \cdot g) = \text{map } f \cdot \text{map } g .$$

This indeed turns the list type constructor with *map* into an endofunctor $\text{list} : \mathbf{Set} \rightarrow \mathbf{Set}$. The important thing to observe here is the close connection between the type constructor and its associated morphism mapping. We will formally justify this intuition soon.

14 Product and sum. Given two objects A and B of \mathbf{C} we can construct their *product object* $A \times B$ with accompanying left and right *projections* $outl : A \times B \rightarrow A$ and $outr : A \times B \rightarrow B$ if we can find for all morphisms $f : C \rightarrow A$ and $g : C \rightarrow B$ a morphism $f \Delta g$ (“*f split g*”) in \mathbf{C} so that the following diagram commutes:

$$\begin{array}{ccccc} A & \xleftarrow{outl} & A \times B & \xrightarrow{outr} & B \\ & \searrow f & \uparrow f \Delta g & \nearrow g & \\ & & C & & \end{array}$$

The diagram asserts the *cancellation* properties for products $outl \cdot (f \Delta g) = f$ and $outr \cdot (f \Delta g) = g$.

The notion of categorical *sum* or *coproduct* is completely dual. The sum of A and B , denoted $A + B$, and its *injections* $inl : A \rightarrow A + B$ and $inr : B \rightarrow A + B$ can be constructed if we are able to find for all morphisms $f : A \rightarrow C$ and $g : B \rightarrow C$ a morphism $f \nabla g$ (“*f junc g*”) in \mathbf{C} so that the diagram

$$\begin{array}{ccccc} A & \xrightarrow{inl} & A + B & \xleftarrow{inr} & B \\ & \searrow f & \downarrow f \nabla g & \swarrow g & \\ & & C & & \end{array}$$

commutes. It displays the cancellation for sums: $(f \nabla g) \cdot inl = f$ and $(f \nabla g) \cdot inr = g$.

It may be helpful to perceive Δ and ∇ as mnemonics for the arrow heads of *split* and *junc* in the corresponding diagrams above.

If we can build the product (sum) for any two objects, \mathbf{C} is said to *have products (sums)*. Category **Set** has products and sums: define $A \times B$ to be the cartesian product of pairs (x, y) with $x \in A, y \in B$ and $A + B$ to be the disjoint sum of sets A and B . The apparent projections and injections then are $outl(x, y) = x$, $outr(x, y) = y$ and $inl x = Lx$, $inr y = Ry$ respectively (L and R being tags to distinguish the origin set of the injected elements). To complete our definition of product and sum in **Set** finally define

$$f \Delta g = \lambda x.(f x, g x) \quad \text{and} \quad f \nabla g = \lambda x.\text{case } x \text{ of } \begin{array}{l} Ly \rightarrow f y \\ Ry \rightarrow g y \end{array},$$

so that ∇ essentially implements case selection.

15 Bifunctors and polynomial functors. Fix \mathbf{C} to be category with products. In \mathbf{C} , we can lift the product former \times to be applicable to morphisms, too: for morphisms $f:A \rightarrow B$ and $g:C \rightarrow D$, define

$$f \times g = (f \cdot outl) \Delta (g \cdot outr): A \times C \rightarrow B \times D.$$

With the definition of Δ from the previous paragraph, we derive the meaning of $f \times g$ in **Set** as follows:

$$\begin{aligned} & f \times g \\ = & \quad \{ \text{definition of } \times \text{ on morphisms} \} \\ & (f \cdot outl) \Delta (g \cdot outr) \\ = & \quad \{ \text{definition of } \Delta \text{ in } \mathbf{Set} \} \\ & \lambda z.((f \cdot outl) z, (g \cdot outr) z) \\ = & \quad \{ \text{set } z = (x, y) \} \\ & \lambda(x, y).((f \cdot outl)(x, y), (g \cdot outr)(x, y)) \\ = & \quad \{ outl \text{ and } outr \text{ in } \mathbf{Set} \} \\ & \lambda(x, y).(f x, g y). \end{aligned}$$

For sums we can argue dually and define $+$ on morphisms $f:A \rightarrow B$, $g:C \rightarrow D$ to mean

$$f + g = (inl \cdot f) \nabla (inr \cdot g): A + C \rightarrow B + D$$

A similar calculation like the one exercised above reveals the meaning of the thus lifted $+$ in **Set**:

$$f + g = \lambda x.\text{case } x \text{ of } \begin{array}{l} Ly \rightarrow L(f y) \\ Ry \rightarrow R(g y) \end{array},$$

i. e., $f + g$ uses the tags to scrutinize its argument for application of f or g and finally injects its argument into the result by restoring the original tag.

Now that \times and $+$ act on objects *and* morphisms one might suspect (or wish) that they behave like functors. Indeed, defined like above, \times and $+$ are *bifunctors*, *i. e.*, functors of two arguments: given the morphisms $f:A \rightarrow_{\mathbf{C}} B$, $g:B \rightarrow_{\mathbf{C}} C$ and $h:D \rightarrow_{\mathbf{D}} E$, $j:E \rightarrow_{\mathbf{D}} F$, \odot is a bifunctor if $f \odot h:A \odot D \rightarrow_{\mathbf{E}} B \odot E$ and

$$id_A \odot id_D = id_{A \odot D} \quad \text{and} \quad (g \cdot f) \odot (j \cdot h) = (g \odot j) \cdot (f \odot h).$$

We can alternatively view \odot as a regular functor $\odot:\mathbf{C} \times \mathbf{D} \rightarrow \mathbf{E}$, where $\mathbf{C} \times \mathbf{D}$ denotes the *product category* of categories \mathbf{C} and \mathbf{D} . The objects of $\mathbf{C} \times \mathbf{D}$ are all pairs (C, D) with C being an object of \mathbf{C} and D being an object of \mathbf{D} . Morphisms (and likewise their composition) are defined coordinate-wise in $\mathbf{C} \times \mathbf{D}$: if $f:C_1 \rightarrow_{\mathbf{C}} C_2$ and $g:D_1 \rightarrow_{\mathbf{D}} D_2$ then $(f, g):(C_1, D_1) \rightarrow_{\mathbf{C} \times \mathbf{D}} (C_2, D_2)$.

A bifunctor can be used to combine two given functors F and G to give a functor $F \odot G$ defined to mean $(F \odot G)A = FA \odot GA$ (likewise for morphisms).

We will solely use \times and $+$ in place of \odot . As usual, \times binds stronger than $+$ so that $F + G \times H$ parses as $F + (G \times H)$. Functor composition binds weakest.

Functors constructed from the identity functor Id , *constant functors* K_A (with $K_A B = A$ and $K_A f = id_A$), bifunctors \times and $+$, and functor composition only are the so-called *polynomial functors*. Polynomial endofunctors in **Set** form the basis of our system of type constructors: they provide *tuples* (built by \times), *unions* (built by $+$) and simple or *atomic* types (built by Id and K respectively, *i. e.*, objects in **Set**).

2.3 Categorical Collection Type Constructors

Now that we have got the necessary categorical machinery at our disposal, we are ready to complete our type system by the introduction of collection type constructors.

16 We have already mentioned in Paragraph 7 that our view of datatypes will be non-monolithic, *i. e.*, the actual construction of a complex value of a particular type will be explicitly visible. To represent any such complex value, we need to assemble it from its constituents (which might be atomic or complex themselves) by means of its type's constructor functions.

For collection types, we will adopt the *insert representation* [111] throughout this text. In such a representation, it is sufficient for a collection type

constructor \mathbb{T}^1 to come equipped with just two constructor functions²

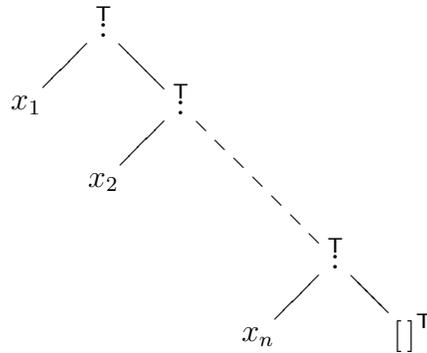
- $nil^{\mathbb{T}}$ is a constant function representing an empty \mathbb{T} -collection, while
- $cons^{\mathbb{T}}(x, xs)$ inserts x into the \mathbb{T} -collection xs .

Issues of polymorphism will be addressed soon. For now, you may assume x to be of some fixed *element type* E . The object in **Set** representing (or: “the type of”) \mathbb{T} -collections of E -typed elements will be $\mathbb{T}E$.

The economical use of constructors is one of the major advantages of the insert representation. See Section 2.10 for a discussion of alternative representations. Any finite \mathbb{T} -collection value may be built by a finite iterative application of the insertion constructor $cons^{\mathbb{T}}$ to the seed $nil^{\mathbb{T}}$. Put differently: the type of \mathbb{T} -collections is *generated* by $nil^{\mathbb{T}}$ and $cons^{\mathbb{T}}$ (here and in what follows, we will abbreviate $nil^{\mathbb{T}}$ by $[]^{\mathbb{T}}$ and $cons^{\mathbb{T}}$ by the right-associative infix operator \mathbb{T} whenever it renders expressions more readable):

$$[x_1, x_2, \dots, x_n]^{\mathbb{T}} = x_1^{\mathbb{T}}(x_2^{\mathbb{T}}(\dots(x_n^{\mathbb{T}}[]^{\mathbb{T}})\dots)) = x_1^{\mathbb{T}}x_2^{\mathbb{T}}\dots^{\mathbb{T}}[]^{\mathbb{T}},$$

The application of this composition of $nil^{\mathbb{T}}$ and $cons^{\mathbb{T}}$ to elements x_i is sometimes called the *spine* of $[x_1, \dots, x_n]^{\mathbb{T}}$, mnemonically depicted as



The above uses $[x_1, x_2, \dots, x_n]^{\mathbb{T}}$ as a handy shorthand notation for a \mathbb{T} -collection, but we should keep in mind that we are actually constructing the value by forming its spine.

In the sequel we will mainly use **set**, **bag**, and **list** as collection type constructors and the following paragraphs will introduce the respective functors

¹Here we use the **sans-serif** font that has been reserved to denote functors which will be justified in Paragraph 3.4 where we establish \mathbb{T} 's functoriality.

²We realize that the names *nil* (“nothing in list”) and *cons* (“construct”) have already been taken by the functional programming community to denote the constructor functions of the **list** type constructor. To emphasize various useful analogies, however, we refrained from inventing new constructor names for collection formers other than **list**.

for these. For **set** as well as **bag**, the same value may be denoted by different spines. As the obvious interpretations for $cons^{\text{set}}(x, xs)$ and nil^{set} are $\{x\} \cup xs$ and \emptyset , we have that $cons^{\text{set}}$ is *left-commutative* and *left-idempotent*, i. e.,

$$y^{\text{set}}(x^{\text{set}}xs) = x^{\text{set}}(y^{\text{set}}xs) \quad \text{and} \quad x^{\text{set}}(x^{\text{set}}xs) = x^{\text{set}}xs.$$

A **bag** can contain an element x multiple times but does not preserve the order of insertions so that $cons^{\text{bag}}$ is left-commutative only. Constructor $cons^{\text{list}}$ is neither. Section 2.7 will formalize the equalities of spines by the notion of *datatype equations*.

17 Fix an element type E and recall that we started out with the aim to model a collection type constructor \mathbb{T} as an algebra. The *carrier* set of this algebra will clearly be $\mathbb{T}E$ as already noted in the previous paragraph. The only operations in this algebra will be the *constructor* functions $nil^{\mathbb{T}}$ and $cons^{\mathbb{T}}$. In **Set**, these are nothing else but the ordinary arrows³ $nil^{\mathbb{T}}: 1 \rightarrow \mathbb{T}E$ and $cons^{\mathbb{T}}: E \times \mathbb{T}E \rightarrow \mathbb{T}E$. As both constructor functions share carrier $\mathbb{T}E$ as their target object, we can alternatively specify their types by stating:

$$nil^{\mathbb{T}} \nabla cons^{\mathbb{T}}: 1 + E \times \mathbb{T}E \rightarrow \mathbb{T}E.$$

The categorical representation of the thus induced algebra is remarkably concise.

18 Definition. Let F denote an endofunctor in category \mathbf{C} . An **F-algebra** is an arrow $\alpha = FA \rightarrow A$. The object A of \mathbf{C} is called the **carrier** of α . \parallel

19 $\mathbb{T}E$ collections in insert representation are modeled by the algebra $nil^{\mathbb{T}} \nabla cons^{\mathbb{T}}: F(\mathbb{T}E) \rightarrow \mathbb{T}E$, with F being the polynomial endofunctor $K_1 + K_E \times \text{Id}$ in **Set**. Note that

$$\begin{aligned} & F(\mathbb{T}E) \\ = & \{ \text{definition of } F \} \\ & (K_1 + K_E \times \text{Id})(\mathbb{T}E) \\ = & \{ \text{bifunctors } + \text{ and } \times \} \\ & K_1(\mathbb{T}E) + K_E(\mathbb{T}E) \times \text{Id}(\mathbb{T}E) \end{aligned}$$

³An arrow $f: 1 \rightarrow A$ is also called an *element* of object A . In **Set**, such arrows represent constant functions and for convenience we shall sometimes identify the arrow and its target object. In type theory, 1 is often called the *unit* type consisting of the only element $()$ (also referred to as *unit*).

$$= \{ \text{functors } \mathbf{K} \text{ and } \text{Id} \} \\ 1 + E \times \mathbb{T}E ,$$

and analogously $\mathbf{F}f = id_1 + id_E \times f$.

20 Definition. Pin down two \mathbf{F} -algebras $\alpha = \mathbf{F}A \rightarrow A$ and $\beta = \mathbf{F}B \rightarrow B$. An **F-homomorphism** from α to β is an arrow $h: A \rightarrow B$ such that

$$h \cdot \alpha = \beta \cdot \mathbf{F}h , \quad \text{HOM}$$

which we can equivalently render as the commuting diagram

$$\begin{array}{ccc} \mathbf{F}A & \xrightarrow{\alpha} & A \\ \mathbf{F}h \downarrow & & \downarrow h \\ \mathbf{F}B & \xrightarrow{\beta} & B \end{array}$$

As an important note, observe that for any object A its identity arrow $id: A \rightarrow A$ is an \mathbf{F} -homomorphism. For any two \mathbf{F} -homomorphisms f and g , we have that $f \cdot g$ is an \mathbf{F} -homomorphism as well. This already provides everything that is needed to form a category (*cf.* Paragraph 11): we define $\mathbf{Alg}(\mathbf{F})$ to denote the category whose objects are \mathbf{F} -algebras. The arrows of $\mathbf{Alg}(\mathbf{F})$ are the \mathbf{F} -homomorphisms. \parallel

21 A simple calculation reveals that property **HOM** of h indeed coincides with the well-known concept of h being a homomorphism between algebras: let $\mathbf{F} = \mathbf{K}_1 + \mathbf{K}_E \times \text{Id}$ once more denote the functor encoding the constructors of a type in insert representation introduced in Paragraph 19. Furthermore fix the \mathbf{F} -algebras $\alpha = e \nabla \otimes: \mathbf{F}A \rightarrow A$ and $\beta = z \nabla \oplus: \mathbf{F}B \rightarrow B$ (see [10, p. 47] for a proof that any algebra for functor \mathbf{F} with carrier A has the form $e \nabla \otimes$ with arrows $e: 1 \rightarrow A$ and $\otimes: E \times A \rightarrow A$). The calculation is as follows:

$$\begin{aligned} & h \cdot \alpha = \beta \cdot \mathbf{F}h \\ \equiv & \{ \alpha = e \nabla \otimes, \beta = z \nabla \oplus \} \\ & h \cdot (e \nabla \otimes) = (z \nabla \oplus) \cdot \mathbf{F}h \\ \equiv & \{ \text{sum: } h \cdot (f \nabla g) = (h \cdot f) \nabla (h \cdot g) \} \\ & (h \cdot e) \nabla (h \cdot \otimes) = (z \nabla \oplus) \cdot \mathbf{F}h \\ \equiv & \{ \text{unfold } \mathbf{F}, \text{ polynomial functor} \} \\ & (h \cdot e) \nabla (h \cdot \otimes) = (z \nabla \oplus) \cdot (id_1 + id_E \times h) \end{aligned}$$

$$\begin{aligned}
&\equiv \{ \text{sum: } (h \nabla j) \cdot (f + g) = (h \cdot f) \nabla (j \cdot g) \} \\
&\quad (h \cdot e) \nabla (h \cdot \otimes) = (z \cdot id_1) \nabla (\oplus \cdot (id_E \times h)) \\
&\equiv \{ \text{identity arrow } \} \\
&\quad (h \cdot e) \nabla (h \cdot \otimes) = z \nabla (\oplus \cdot (id_E \times h)) \\
&\equiv \{ \text{cancellation } \} \\
&\quad h \cdot e = z \\
&\quad \wedge h \cdot \otimes = \oplus \cdot (id_E \times h) \\
&\equiv \{ \text{point-wise reasoning } \} \\
&\quad h e = z \\
&\quad \wedge h(x \otimes xs) = x \oplus (h xs)
\end{aligned}$$

The two last equations constitute the familiar statement of h being a homomorphism from α to β .

2.4 Catamorphisms

22 Up to now we have informally called nil^\top and $cons^\top$ the *constructors* for a datatype former \top . The intuition behind the constructor notion is that

- every element of the carrier $\top E$ of $\alpha = nil \nabla cons : F(\top E) \rightarrow \top E$ can be constructed using applications of nil^\top and $cons^\top$ to values in E only. No other elements are in $\top E$, *i. e.*, there is no “junk.”
- Secondly, given that there are no equations yet, two elements of $\top E$ are identical only if their spines are exactly the same. Jumping ahead to Section 2.7 in which we introduce equations, this means that two elements are identical only if this can be derived from the equations of the datatype. This prevents us from accidentally “confusing” two elements of $\top E$ with each other.

Let us shed more categorical light on this here.

Assume for now that $\mathbf{Alg}(\mathbf{F})$ has an initial object, say $\tau : FT \rightarrow T$. Actually this is no strong assumption at all in our setting as for polynomial functors \mathbf{F} , the existence of an initial object in $\mathbf{Alg}(\mathbf{F})$ is guaranteed (see [81] for a proof of this statement).

23 Definition. Let $\tau : FT \rightarrow T$ be the initial and $\alpha : FA \rightarrow A$ an arbitrary object of $\mathbf{Alg}(\mathbf{F})$. An \mathbf{F} -homomorphism of type $T \rightarrow A$ (which is, due to the initiality of τ , guaranteed to exist and to be unique, *cf.* Paragraph 12) from

τ to α is a so-called **F-catamorphism**. These special homomorphisms will accompany us through the entire text, so we will adopt the succinct notation $\langle \alpha \rangle_{\mathbb{F}}$ for the catamorphism from τ to α .⁴ If the functor is evident from the context we take the freedom to omit the subscript.

The uniqueness of $\langle \alpha \rangle$ indeed is its prominent property. Any other F-homomorphism h from τ to α clearly has to be identical with $\langle \alpha \rangle$:

$$h \cdot \tau = \alpha \cdot \mathbb{F}h \quad \equiv \quad h = \langle \alpha \rangle . \quad \text{CATA}$$

The morphisms involved in law **CATA** are shown in the commuting diagram below. A (unique) homomorphic translation between algebras like the one depicted here will form the basis of our query language.

$$\begin{array}{ccc} \mathbb{F}T & \xrightarrow{\tau} & T \\ \mathbb{F}\langle \alpha \rangle \downarrow & & \downarrow \langle \alpha \rangle \\ \mathbb{F}A & \xrightarrow{\alpha} & A \end{array}$$

//

24 Initial algebras. What does the initiality of τ buy us? First, Paragraph 12 introduced initiality as a concept that determined the initial object in a category up to isomorphism. Of course the same is true in category $\mathbf{Alg}(\mathbb{F})$: suppose we can identify another initial algebra τ' , we know that τ and τ' are isomorphic and therefore *abstractly the same*. But this notion of abstraction perfectly coincides with the independence of representation of an abstract datatype. Abstraction is inherent to initiality. In what follows, we will use the term *abstract datatype* to refer to the isomorphism class of the initial algebra τ .

Second, as the next paragraph will show, τ is an isomorphism from $\mathbb{F}T$ (the type of expressions built from operators in τ and elements of the carrier T) to T . As τ is a morphism in **Set**, we thus know that τ is a bijective mapping from $\mathbb{F}T$ to T :

- τ is surjective, so any element in T is represented by at least one expression in $\mathbb{F}T$. There is no “*junk*.”
- τ is injective, so we may conclude that the elements of T have a unique spine representation in $\mathbb{F}T$. Consequently there is no “*confusion*.”

⁴The “*banana brackets*” $\langle \cdot \rangle$ —presumably introduced by Lambert Meertens—are the established catamorphism symbol in the Squiggol community.

Most importantly however, initiality justifies the principle of defining functions on T by *structural recursion*. As each element e of T is constructed by its unique spine in FT , *i. e.*, a composition of the constructors of τ , we can define a function h by case selection on the constructor whose application is outermost in the spine of e . Each case branch may then descend further down the spine and recursively apply h .

We can also see this as follows: as $\tau:FT \rightarrow T$ is an isomorphism, we are assured of the existence of an “inverse” isomorphism $\bar{\tau}:T \rightarrow FT$ with $\tau \cdot \bar{\tau} = id_T$ and $\bar{\tau} \cdot \tau = id_{FT}$. Observe that the type of $\bar{\tau}$ suggests to interpret this arrow as a *destructor* or *parser* mapping values in T to the expression (spine) in FT that constructs this value. This gives us more insight into the nature of the action of a catamorphism $h = \langle \alpha \rangle$ for some F -algebra α :

$$\begin{aligned} h &= \langle \alpha \rangle \\ &\equiv \{ \text{CATA} \} \\ &\quad \langle \alpha \rangle \cdot \tau = \alpha \cdot F\langle \alpha \rangle \\ &\equiv \{ \tau \cdot \bar{\tau} = id_T \} \\ &\quad \langle \alpha \rangle = \alpha \cdot F\langle \alpha \rangle \cdot \bar{\tau}, \end{aligned}$$

i. e., h forms the spine of its argument, recursively applies h by walking down the spine, and finally builds the result by interpreting the spine in algebra α .

The following lemma guarantees the existence of $\bar{\tau}$ and justifies the above.

25 Lambek’s Lemma [71]. Let $\tau:FT \rightarrow T$ be the initial algebra of endofunctor F . Then τ is an isomorphism, *i. e.*, there exists an arrow $\bar{\tau}:T \rightarrow FT$ such that $\tau \cdot \bar{\tau} = \bar{\tau} \cdot \tau = id$.

Proof. Applying the endofunctor F to the algebra τ yields the F -algebra $F\tau:F(FT) \rightarrow FT$. The unique mapping from τ to $F\tau$ —which is guaranteed to exist due to the initiality of τ —is $\langle F\tau \rangle:T \rightarrow FT$. We now argue that $\bar{\tau} = \langle F\tau \rangle$.

In the diagram below, the left square merely states that the catamorphism $\langle F\tau \rangle$ is an F -homomorphism which is evident. Both paths in the right square denote the very same arrow composition so that it trivially commutes. “Pasting” commutative diagrams yields a commutative diagram.

$$\begin{array}{ccccc} FT & \xrightarrow{F\langle F\tau \rangle} & F(FT) & \xrightarrow{F\tau} & FT \\ \tau \downarrow & & \downarrow F\tau & & \downarrow \tau \\ T & \xrightarrow{\langle F\tau \rangle} & FT & \xrightarrow{\tau} & T \end{array}$$

The diagram reveals $\tau \cdot \langle \mathbf{F}\tau \rangle$ to be a mapping from the algebra τ to τ itself, so we have $\tau \cdot \langle \mathbf{F}\tau \rangle = id$ by initiality. All that remains to be done now is to carefully look at the left square from which we get

$$\begin{aligned}
& \langle \mathbf{F}\tau \rangle \cdot \tau = \mathbf{F}\tau \cdot \mathbf{F}\langle \mathbf{F}\tau \rangle \\
\equiv & \quad \{ \text{functoriality of } \mathbf{F} \} \\
& \langle \mathbf{F}\tau \rangle \cdot \tau = \mathbf{F}(\tau \cdot \langle \mathbf{F}\tau \rangle) \\
\equiv & \quad \{ \tau \cdot \langle \mathbf{F}\tau \rangle = id \text{ established above} \} \\
& \langle \mathbf{F}\tau \rangle \cdot \tau = \mathbf{F}id \\
\equiv & \quad \{ \text{functoriality of } \mathbf{F} \} \\
& \langle \mathbf{F}\tau \rangle \cdot \tau = id ,
\end{aligned}$$

which completes the proof. \square

26 Lambek's Lemma suggests an alternative view of datatype semantics [72]. Observe that the initial object $\tau: FT \rightarrow T$ in $\mathbf{Alg}(\mathbf{F})$ may be understood as a *fixpoint* of endofunctor \mathbf{F} in the underlying category \mathbf{C} since FT and T are isomorphic and therefore abstractly the same ($\mathbf{C} = \mathbf{Set}$ in our case). Since τ is initial it is indeed the *least fixpoint* of \mathbf{F} if we interpret the arrows in \mathbf{C} as a pre-order. This observation has led to the so-called *least fixpoint semantics* which defines the algebra for a datatype as $\tau = \mu\mathbf{F}$ (symbol μ denotes the least fixpoint operator).

27 Let us take up functor $\mathbf{F} = 1 + \mathbf{K}_E \times \mathbf{Id}$ of datatypes in insert representation again (E fixed). Algebra $\tau = nil^\top \nabla cons^\top: \mathbf{F}(TE) \rightarrow TE$ as defined in Paragraph 19 is initial in $\mathbf{Alg}(\mathbf{F})$.

The proof of this claim necessarily involves to show that for any homomorphism h from τ to an arbitrary \mathbf{F} -algebra $\alpha = e \nabla \otimes: \mathbf{F}A \rightarrow A$ we have that $h = \langle \alpha \rangle = \langle e \nabla \otimes \rangle$ for a suitable definition of $\langle \cdot \rangle$, *i. e.*, we have to establish the validity of law **CATA**. We can see this as follows:

$$\begin{aligned}
& h \cdot \tau = \alpha \cdot \mathbf{F}h \\
\equiv & \quad \{ \tau = nil^\top \nabla cons^\top, \alpha = e \nabla \otimes \} \\
& h \cdot (nil^\top \nabla cons^\top) = (e \nabla \otimes) \cdot \mathbf{F}h \\
\equiv & \quad \{ \text{sum} \} \\
& (h \cdot nil^\top) \nabla (h \cdot cons^\top) = (e \nabla \otimes) \cdot \mathbf{F}h \\
\equiv & \quad \{ \mathbf{F} = 1 + \mathbf{K}_E \times \mathbf{Id}, \text{ functor} \} \\
& (h \cdot nil^\top) \nabla (h \cdot cons^\top) = (e \nabla \otimes) \cdot (id_1 + id_E \times h)
\end{aligned}$$

This observation immediately yields the so-called *reflection law* for catamorphisms:

$$(\tau) = id . \quad \text{CATA-REFLECT}$$

CATA-INS-REP obviously adheres to the general recursion scheme of *primitive recursive functions* [89]. The simplicity of the action of $(\cdot)_F$ is all the more remarkable as it is already powerful enough to express the core operators of relational algebra, nested relational algebra [13, 20], and—as we will discuss to some depth in this text—modern orthogonal query languages for complex collection types like ODMG’s OQL [23]. We can be more specific here: as the catamorphism combinator $(\cdot)_F$ is a generic variant of the *sri* structural recursion operator on collections in insert representation it is capable of expressing exactly the PTIME queries, *i. e.*, those computable in polynomial time [13, 110, 111].

28 A functional inclined programmer has surely noticed that the action **CATA-INS-REP** of catamorphism $(e \nabla \otimes)$ on a spine is exactly that of the well-known higher-order function *foldr* on a list:

$$\begin{aligned} \text{foldr} & : (E \times A \rightarrow A) \rightarrow A \rightarrow \text{list}E \rightarrow A \\ \text{foldr} \otimes e \text{ nil}^{\text{list}} & = e \\ \text{foldr} \otimes e (\text{cons}^{\text{list}}(x, xs)) & = x \otimes (\text{foldr} \otimes e xs) . \end{aligned}$$

As it happens, *foldr* is a catamorphism for lists in insert representation ($F = 1 + K_E \times \text{Id}, T = \text{list}$), *i. e.*, we have $\text{foldr} \otimes e = (e \nabla \otimes)$:

$$\begin{array}{ccc} 1 + E \times \text{list}E & \xrightarrow{\text{nil}^{\text{list}} \nabla \text{cons}^{\text{list}}} & \text{list}E \\ \downarrow \text{id}_1 + \text{id}_E \times (\text{foldr} \otimes e) & & \downarrow \text{foldr} \otimes e \\ 1 + E \times A & \xrightarrow{e \nabla \otimes} & A \end{array}$$

F-catamorphisms—with F being the endofunctor describing the insert representation of a datatype—indeed define a generalized *foldr* for any datatype former T.

More generally even, executing a calculation like the one in Paragraph 27 automatically provides us with a *fold* for any endofunctor F and datatype constructor T.

29 Example. A wide range of list-processing functions turn out to be catamorphisms from the list algebra $\tau = nil^{list} \nabla cons^{list} : F(list\ E) \rightarrow list\ E$ to other F-algebras, *i. e.*, these functions compute their result by spine transformations alone.

Let Num denote the object in **Set** representing the set of natural numbers \mathbb{N} . The function $length : list\ E \rightarrow Num$ which computes the length of a list is a catamorphism from τ to the F-algebra $0 \nabla (+ \cdot (K_1 \times id)) : F\ Num \rightarrow Num$ with $K_x\ y = x$. Note that $+$ denotes the addition on Num in this example. We have $length = \llbracket 0 \nabla (+ \cdot (K_1 \times id)) \rrbracket$ as the following calculation in which we apply the $length$ catamorphism to the argument $[0, 1, 2]^{list}$ exemplifies:

$$\begin{aligned}
& \llbracket 0 \nabla (+ \cdot (K_1 \times id)) \rrbracket [0, 1, 2]^{list} \\
= & \quad \{ \text{list insert representation} \} \\
& \llbracket 0 \nabla (+ \cdot (K_1 \times id)) \rrbracket (0^{list} : 1^{list} : 2^{list} : []^{list}) \\
= & \quad \{ \text{CATA-INS-REP} \} \\
& (+ \cdot (K_1 \times id)) (0, \llbracket 0 \nabla (K_1 \times id) \rrbracket (1^{list} : 2^{list} : []^{list})) \\
= & \quad \{ \text{calculation} \} \\
& 1 + (\llbracket 0 \nabla (K_1 \times id) \rrbracket (1^{list} : 2^{list} : []^{list})) \\
= & \quad \{ \text{CATA-INS-REP} \} \\
& 1 + ((+ \cdot (K_1 \times id)) (1, \llbracket 0 \nabla (K_1 \times id) \rrbracket (2^{list} : []^{list}))) \\
= & \quad \{ \text{calculation} \} \\
& 1 + 1 + (\llbracket 0 \nabla (K_1 \times id) \rrbracket (2^{list} : []^{list})) \\
= & \quad \{ \text{CATA-INS-REP} \} \\
& 1 + 1 + ((+ \cdot (K_1 \times id)) (2, \llbracket 0 \nabla (K_1 \times id) \rrbracket ([]^{list}))) \\
= & \quad \{ \text{calculation} \} \\
& 1 + 1 + 1 + (\llbracket 0 \nabla (K_1 \times id) \rrbracket ([]^{list})) \\
= & \quad \{ \text{CATA-INS-REP} \} \\
& 1 + 1 + 1 + 0 \\
= & \quad \{ \text{arithmetic} \} \\
& 3.
\end{aligned}$$

◇

30 Example. A spine resulting from a catamorphism application may perfectly have a greater or smaller height than the argument spine. The selection function $filter\ p : list\ E \rightarrow list\ E$, with $p : E \rightarrow Bool$ being some predicate on the list elements of type E , is a catamorphism from τ to the F-algebra

$nil^{\text{list}} \nabla ((cons^{\text{list}} \nabla outr) \cdot (p \cdot outl)?): F(\text{list}E) \rightarrow \text{list}E$. Here, $?$ denotes the *infix conditional* defined by

$$\begin{aligned} ? & : (A \rightarrow Bool) \rightarrow A \rightarrow A + A \\ p?x & = \text{if } px \text{ then } inl\ x \text{ else } inr\ x \end{aligned}$$

The morphism $(cons^{\text{list}} \nabla outr) \cdot (p \cdot outl)?: E \times \text{list}E \rightarrow \text{list}E$ implements a *conditional* $cons^{\text{list}}$ in **Set** as a straightforward calculation shows:

$$\begin{aligned} & (cons^{\text{list}} \nabla outr) \cdot (p \cdot outl)? \\ \equiv & \{ \text{unfold } ? \} \\ & (cons^{\text{list}} \nabla outr) \cdot (\lambda z. \text{if } (p \cdot outl)\ z \text{ then } inl\ z \text{ else } inr\ z) \\ \equiv & \{ \text{conditional, } (cons^{\text{list}} \nabla outr) \text{ strict} \} \\ & \lambda z. \text{if } (p \cdot outl)\ z \text{ then } ((cons^{\text{list}} \nabla outr) \cdot inl)\ z \\ & \quad \text{else } ((cons^{\text{list}} \nabla outr) \cdot inr)\ z \\ \equiv & \{ \text{sum cancellation} \} \\ & \lambda z. \text{if } (p \cdot outl)\ z \text{ then } cons^{\text{list}}\ z \text{ else } outr\ z \\ \equiv & \{ \text{set } z = (x, xs), \text{ outl, outr in } \mathbf{Set} \} \\ & \lambda(x, xs). \text{if } px \text{ then } cons^{\text{list}}(x, xs) \text{ else } xs . \end{aligned}$$

(As the default category we are working in is **Set**, we will sometimes alternatively denote morphisms in this category by λ -expressions directly. The context of the current calculation will determine which notation suits best to improve readability.)

We thus have $filter\ p = \langle nil^{\text{list}} \nabla ((cons^{\text{list}} \nabla outr) \cdot (p \cdot outl)?) \rangle$. To illustrate the action of $filter\ p$ on a spine let us instantiate E as Num and the filter predicate as $p = \lambda x. x \neq 0$:

$$filter(\lambda x. x \neq 0) \left(\begin{array}{c} \text{list} \\ \vdots \\ 0 \quad \text{list} \\ \quad \vdots \\ \quad 1 \quad \text{list} \\ \quad \quad \vdots \\ \quad \quad 2 \quad \text{list} \\ \quad \quad \quad \vdots \\ \quad \quad \quad []^{\text{list}} \end{array} \right) = \begin{array}{c} \text{list} \\ \vdots \\ 1 \quad \text{list} \\ \quad \vdots \\ \quad 2 \quad \text{list} \\ \quad \quad \vdots \\ \quad \quad []^{\text{list}} \end{array}$$

◇

2.5 Type Functors

31 The collection type constructors \mathbb{T} that we have met up to now have been monomorphic in the sense that were assuming a *fixed* element type E . To remedy this shortcoming, this section introduces real polymorphic type formers.

In the case of algebras for polymorphic type constructors \mathbb{T} , there is the need to additionally parameterize the algebra's endofunctor by the object that represents the element type. This observation naturally leads to the use of bifunctors to model such algebras, *i. e.*, we enhance our definition of F -algebras as follows:

32 Definition. Let F denote a bi-endofunctor in a category \mathbf{C} and let E, A be objects in \mathbf{C} . A **polymorphic F -algebra** is an arrow $\alpha : F(E, A) \rightarrow A$. The object A is called the **carrier** of α .

Using the lingo of type theory, polymorphic algebras defined like above have a second-order type [22]. Type parameter E is *universally quantified* so that the proper type for α is $\forall E. F(E, A) \rightarrow A$. Once we have fixed a specific element type E' , we can revert to the situation in which the monofunctor $F_{E'}$ (the left *section* of the bifunctor F) is sufficient to model the resulting, now monomorphic, algebra:

$$\alpha_{E'} : F_{E'} A \rightarrow A \quad \text{with} \quad F_{E'} A = F(E', A) \quad \text{and} \quad F_{E'} f = F(id_{E'}, f) .$$

//

33 Polymorphic algebras already provide everything we need to define the algebra of a polymorphic collection type constructor \mathbb{T} in insert representation. As before (see Paragraph 19), $\mathbb{T}E$ will be its carrier but the polymorphic algebra differs in that we pass the element type E as an additional argument.

$$\alpha = cons^{\mathbb{T}} \nabla nil^{\mathbb{T}} : F(E, \mathbb{T}E) \rightarrow \mathbb{T}E$$

with

$$\begin{aligned} F(E, A) &= 1 + E \times A \\ \wedge \quad F(f, g) &= id_1 + f \times g \quad . \end{aligned}$$

Fixing a specific element type E' leads to the algebra functor we were used to until now so that our earlier discussion of monomorphic algebras applies without revision:

$$\begin{aligned} F_{E'}(\mathbb{T}E') &= F(E', \mathbb{T}E') = 1 + E' \times \mathbb{T}E' \\ \wedge \quad F_{E'} f &= F(id_{E'}, f) = id_1 + id_{E'} \times f \quad , \end{aligned}$$

i. e., $F_{E'} = K_1 + K_{E'} \times Id$.

$$\top f = (\tau \cdot F(f, id)) . \qquad \text{TYPE-FUNCTOR}$$

The pair (τ, \top) is called the **initial type** [10] induced by the bifunctor F .

We omit the details of formally establishing that this indeed turns \top into a functor. The proof of $\top id = id$ is a trivial consequence of F 's bifunctoriality and law **CATA-REFLECT**. The distribution property follows immediately with the help of law **CATA-MAP-FUSION** (see Paragraph 39 below and [10, 42, 57]). //

36 Example. For the type constructor $\top = \text{list}$ we have $\text{list}f = \text{map}f$ where map denotes the well-known higher-order function which applies function f to every element of its list argument. For the sake of uniformity we write $\text{list}f$ instead of $\text{map}f$ throughout this text as we will also encounter the type functors **set** and **bag** among others.

The type functor concept is rather generic and gives us a tool to derive map -like functions for any initial datatype.

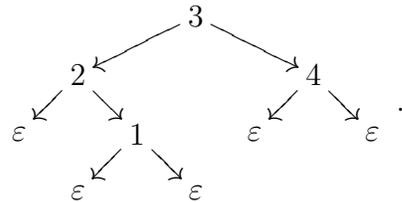
Consider the abstract datatype **bintree** of binary trees with inner nodes labeled with a tag of type E . Such trees are built by means of the two constructor functions

$$\begin{aligned} \text{empty} & : 1 \rightarrow \text{bintree}E \\ \text{node} & : E \times \text{bintree}E \times \text{bintree}E \rightarrow \text{bintree}E \end{aligned}$$

where empty denotes the empty tree ε and $\text{node}(n, t_1, t_2)$ constructs a tree with root n and left and right subtrees t_1 and t_2 respectively. The constructor expression

$$\text{node}(3, \text{node}(2, \text{empty}, \text{node}(1, \text{empty}, \text{empty})), \text{node}(4, \text{empty}, \text{empty}))$$

thus builds the binary tree



The initial algebra for this datatype is $\tau = \text{empty} \nabla \text{node} : F(E, \top E) \rightarrow \top E$ where $F(E, \top E) = 1 + E \times \top E \times \top E$, $F(f, g) = id_1 + f \times g \times g$, and of course $\top = \text{bintree}$.

How does the associated type functor look like? By instantiating the ingredients of **TYPE-FUNCTOR** and conducting the simple calculation below, we see that `bintree`'s map functor indeed is a *tree map*: the catamorphism for the functor F defined above (whose derivation we leave to the kind reader) recursively descends into the two subtrees of an inner node until an empty tree ε is encountered which is left untouched by the tree map:

$$\begin{aligned}
& \text{bintree } f \\
= & \quad \{ \text{definition of type functor} \} \\
& \quad (\tau \cdot F(f, id)) \\
= & \quad \{ \text{unfold } \tau, \text{unfold } F \} \\
& \quad ((\text{empty} \nabla \text{node}) \cdot (id_1 + f \times id \times id)) \\
= & \quad \{ \text{sum} \} \\
& \quad (\text{empty} \nabla (\text{node} \cdot (f \times id \times id))) .
\end{aligned}$$

As another and final example let us investigate a non-collection type former. The type `maybeE` is the same as E but it additionally comes with a distinguished element *nothing* which may serve as, *e. g.*, an error value. Type former `maybe` has constructors

$$\begin{aligned}
\text{nothing} & : 1 \rightarrow \text{maybe}E \\
\text{just} & : E \rightarrow \text{maybe}E .
\end{aligned}$$

The accompanying initial algebra is $\tau = \text{nothing} \nabla \text{just} : F(E, \text{maybe}E) \rightarrow \text{maybe}E$ with $F(E, \text{maybe}E) = 1 + E$, $F(f, g) = id_1 + f$. Its type functor turns out to be a specialized function application operator that propagates the error value *nothing* and otherwise simply applies the given function:

$$\begin{aligned}
& \text{maybe } f \\
= & \quad \{ \text{definition of type functor} \} \\
& \quad (\tau \cdot F(f, id)) \\
= & \quad \{ \text{unfold } \tau, \text{unfold } F \} \\
& \quad ((\text{nothing} \nabla \text{just}) \cdot (id_1 + f)) \\
= & \quad \{ \text{sum} \} \\
& \quad (\text{nothing} \nabla (\text{just} \cdot f)) \\
\equiv & \quad \{ \text{point-wise reasoning} \} \\
& \quad \text{maybe } f \text{ nothing} = \text{nothing} \\
& \quad \wedge \text{ maybe } f (\text{just } x) = \text{just}(f x) .
\end{aligned}$$

Among other uses, `maybe f` is helpful during the implementation of and reasoning about functional exception handling mechanisms. \diamond

2.6 Fusion

37 Now that it has now become more clear that spine transformations will be the basic building blocks from which we will construct queries, we can imagine that a query will typically be built by the *functional composition* of such transformations.

In other words, if Q denotes a query expression then it will exhibit a structure that may be roughly sketched as

$$Q \equiv \dots h \cdot g \cdot f \dots$$

where the functions f , g , and h represent spine transformers, in particular catamorphisms and type functor applications (*maps*). Building query expressions like this leads to a clear and modular way of thinking about query construction itself: complex queries may be formulated by incrementally adding a basic transformer to an already existing sequence of function compositions. This *compositional* flavor of internal query representation facilitates to offer the same flexibility at the user level (where this taste of query formulation is often referred to as *orthogonality* if there is nothing but the language's type system that stops the user from freely combining query expressions).

We will not go into the details here as Chapter 6 will investigate the matter more closely, but let us note that the compositional style does not come for free:

- Each function generates an *intermediate result* that has to be taken up later by the next function in the composition sequence. The management of intermediate results may be quite costly.
- Catamorphisms essentially traverse the spine of their arguments while performing their transformation task and each catamorphism in the sequence does so from the spine's root downwards. It may be desirable to combine two or more of these transformations in such a way that a *single* walk of the spine is sufficient.

What we are calling for here is the *fusion* of two adjacent spine transformers into one [37, 44].

Two well-known theorems provide the handles we need to fuse two catamorphisms as well as a type functor application followed by a catamorphism in the way just described. In Chapter 6, we will accompany these two laws with a third theorem—also known as the *acid rain theorem* [112]—which justifies an effective fusion technique that will prove especially useful during the derivation of a stream-based (or pipelined) program from a given query.

38 Promotion Theorem [3, 56, 57]. Let F be an endofunctor in category \mathbf{C} so that $\mathbf{Alg}(F)$ has an initial object τ . For any two F -algebras $f:FA \rightarrow A$ and $g:FB \rightarrow B$ and an F -homomorphism $h:A \rightarrow B$ from f to g we have that

$$h \cdot \langle f \rangle = \langle g \rangle . \quad \text{CATA-FUSION}$$

Proof. A single calculation suffices to establish the claim:

$$\begin{aligned} & h \cdot \langle f \rangle \cdot \tau \\ = & \quad \{ \text{CATA: } \langle f \rangle \cdot \tau = f \cdot F\langle f \rangle \} \\ & h \cdot f \cdot F\langle f \rangle \\ = & \quad \{ \text{HOM: } h \cdot f = g \cdot Fh \} \\ & g \cdot Fh \cdot F\langle f \rangle \\ = & \quad \{ F \text{ functor} \} \\ & g \cdot F(h \cdot \langle f \rangle) \\ \Rightarrow & \\ & h \cdot \langle f \rangle \cdot \tau = g \cdot F(h \cdot \langle f \rangle) \\ \equiv & \quad \{ \text{CATA, uniqueness of } \langle g \rangle \} \\ & h \cdot \langle f \rangle = \langle g \rangle . \end{aligned}$$

□

As any F -catamorphism is a special F -homomorphism, law **CATA-FUSION** indeed is the catamorphism fusion tool we have appealed for above.

39 Theorem. If (τ, \mathbb{T}) denotes the initial type (Definition 35) for the endofunctor F , then we can fuse a composition of an F -catamorphism with its type functor as follows (under typings $f:A \rightarrow B$, $g:FB \rightarrow B$):

$$\langle g \rangle \cdot \mathbb{T}f = \langle g \cdot F(f, id) \rangle . \quad \text{CATA-MAP-FUSION}$$

Proof. The proof invokes law **CATA-FUSION** as a subroutine which is hardly surprising because the arrow mapping of the type functor \mathbb{T} has been defined as a specific F -catamorphism (see Definition 35). We calculate as follows (*cf.* [10]):

$$\begin{aligned} & \langle g \rangle \cdot \mathbb{T}f = \langle g \cdot F(f, id) \rangle \\ \equiv & \quad \{ \text{TYPE-FUNCTOR for } \mathbb{T} \} \\ & \langle g \rangle \cdot \langle \tau \cdot F(f, id) \rangle = \langle g \cdot F(f, id) \rangle \end{aligned}$$

$$\begin{aligned}
&\Leftarrow \{ \text{CATA-FUSION} \} \\
&\quad (g) \cdot \tau \cdot F(f, id) = g \cdot F(f, id) \cdot F(id, (g)) \\
&\equiv \{ \text{CATA: } (g) \cdot \tau = g \cdot F(id, (g)) \} \\
&\quad g \cdot F(id, (g)) \cdot F(f, id) = g \cdot F(f, id) \cdot F(id, (g)) \\
&\equiv \{ F \text{ bifunctor} \} \\
&\quad g \cdot F(id, (g)) \cdot F(f, id) = g \cdot F(id, (g)) \cdot F(f, id) \\
&\equiv \{ \text{reflexivity of } = \} \\
&\quad \text{true} ,
\end{aligned}$$

which establishes the law as a logical consequence of the trivial premise *true*. \square

40 Example. To become acquainted with the flavor of program transformations which now have been justified by the fusion laws, let us take a closer look at a simple spine transformer composition and how to derive its fused equivalent.

The expression $filter\ p \cdot map\ f : listA \rightarrow listB$ traverses a given argument spine twice. During the first traversal, $map\ f$ generates a new spine by replacing each list element x by $f\ x$. The second traversal then takes up this intermediate spine to build the final result from all elements that satisfy predicate p .

We can clearly accomplish the same by walking down the spine only once. In order not to alter the meaning of the expression we only have to take care to apply f to the list elements *before* predicate p checks for their inclusion in the final result. Fusion allows us to automatically derive an expression that does just this. No intermediate result is produced.

As before, we are working with collections in insert representation, *i. e.*, the underlying algebra functor is $F(f, g) = id_1 + f \times g$. Although the example expression operates on lists only this is not principle to the method, so let us generalize the argument to some arbitrary type constructor \mathbb{T} . We have $map\ f = \mathbb{T}f$, making the expression subject to law **CATA-MAP-FUSION**. Once we have unfolded *filter*'s definition as an F-catamorphism (see Example 30), we immediately spot the possibility for fusion:

$$\begin{aligned}
&filter\ p \cdot \mathbb{T}f \\
&= \{ \text{unfold } filter \} \\
&\quad ((nil^{\mathbb{T}} \nabla ((cons^{\mathbb{T}} \nabla outr) \cdot (p \cdot outl)?)) \cdot \mathbb{T}f \\
&= \{ \text{CATA-MAP-FUSION} \} \\
&\quad ((nil^{\mathbb{T}} \nabla ((cons^{\mathbb{T}} \nabla outr) \cdot (p \cdot outl)?)) \cdot F(f, id))
\end{aligned}$$

$$\begin{aligned}
&= \{ \text{definition of } \mathbf{F} \} \\
&\quad \langle (nil^\top \nabla ((cons^\top \nabla outr) \cdot (p \cdot outl)?)) \cdot (id_1 + f \times id) \rangle \\
&= \{ \text{sum} \} \\
&\quad \langle nil^\top \nabla ((cons^\top \nabla outr) \cdot (p \cdot outl)? \cdot (f \times id)) \rangle,
\end{aligned}$$

which already completes the fusion step. We can rewrite this expression a bit further so that it reveals its nature of actually being a specialized *filter*:

$$\begin{aligned}
&\langle nil^\top \nabla ((cons^\top \nabla outr) \cdot (p \cdot outl)? \cdot (f \times id)) \rangle \\
&= \{ p? \cdot f = (f + f) \cdot (p \cdot f)? \} \\
&\quad \langle nil^\top \nabla ((cons^\top \nabla outr) \cdot (f \times id + f \times id) \cdot (p \cdot outl \cdot (f \times id)?)) \rangle \\
&= \{ \text{product, sum} \} \\
&\quad \langle nil^\top \nabla (((cons^\top \cdot (f \times id)) \nabla (outr \cdot (f \times id))) \cdot (p \cdot f \cdot outl?)) \rangle \\
&= \{ \text{product} \} \\
&\quad \langle nil^\top \nabla (((cons^\top \cdot (f \times id)) \nabla outr) \cdot (p \cdot f \cdot outl?)) \rangle.
\end{aligned}$$

We finally ended up with a single spine transformer instead of two. Subsequent transformation steps could now proceed with fusion given that the newly derived catamorphism is adjacent to other transformers in a larger composition chain. While we do not formulate this as our primary optimization goal, “excessive” fusion—possibly hand in hand with helper transformations that rely on commutativity, distributivity, or similar algebraic properties—could thus be exploited to minimize the spine transformer count in a given expression. \diamond

41 The development of strategies for the removal of intermediate data structures, for which [Wadler \[118\]](#) has coined the term *deforestation*, has always been a major driving force of the program transformation community.⁵ The impact of this research on our work will be the focus of [Chapter 6](#).

Taken as a whole, the transformation exercised in the last example may look rather involved. This is especially true for the step in which the infix conditional is rewritten: $p? \cdot f = (f + f) \cdot (p \cdot f)?$. While this is no deep insight at all, the step must be regarded as a so-called *Eureka step* in this context, *i. e.*, a rewriting that is not immediately motivated by the goal which the current calculation strives for. Even worse, asking if a particular morphism is a homomorphism which, *e. g.*, is necessary to check the applicability of law [CATA-FUSION](#), is undecidable in general [[14](#), [20](#)].

⁵Deforestation removes intermediate data structures which [Wadler](#) collectively refers to as trees, hence the name *deforestation*.

Such obstacles hinder the construction of a completely unguided program transformation system, an instance of which—a database query optimizer—we are discussing in this text.

As it will turn out, however, we do *not* need the full power of fusion at our disposal to generate efficient programs from query expressions. While still retaining the expressiveness needed to cover modern orthogonal query languages, the query compiler (Chapter 3) will emit expressions in a form that is amenable to fusion using a simple yet powerful one-step transformation, *cheap deforestation*. This process will be fully automatic (see Chapter 6 and [50, 51, 59, 74, 112]).

2.7 Datatype Equations

42 The algebras we have met up to now have all been *anarchic* or *free* in the sense that they do not obey any *equations*. An equation for a datatype—actually for its algebra to be precise—is imposed by the designer of that type at will, or, put more positive, equations provide the instrument to introduce a controlled form of *confusion*: two distinct algebra expressions may denote the very same object if the equality of these expressions is derivable from a set of specified equations. It is important not to confuse equations with the notion of *law* which represents provable statements like the **CATA-FUSION** theorem.

We have chosen `list` as the predominant example of a type former until now because its algebra is indeed free. The equational theory of the list algebra in insert representation is empty: there is no confusion between a list spine (in $F(\text{list}E)$) and the actual list value (in $\text{list}E$) it represents. Note that this is just the isomorphism statement of Lambek’s Lemma:

$$F(\text{list}E) \begin{array}{c} \xrightarrow{\tau} \\ \xleftarrow{\bar{\tau}} \end{array} \text{list}E \quad ,$$

which, in this context, can be read as: syntactical equality means semantical equality.

The need for datatypes that obey equations arises with the problem of modeling non-free algebras like those for the type formers `bag` and `set`. No isomorphism like the above exists between $F(\text{bag}E)$ or $F(\text{set}E)$ and their respective carriers: the left-commutativity of bag and set (the latter additionally being left-idempotent or left-absorptive) renders these constructors *non-injective* (see Paragraph 16). The `set` spines $[1, 2, 4]^{\text{set}}$ and $[2, 4, 4, 1, 2]^{\text{set}}$ both represent the set value $\{1, 2, 4\}$, for example.

43 How does the presence of confusion affect query processing in our catamorphism-based framework? To retain the determinism of the catamorphism action on a spine we need the action to be indifferent towards the non-injectivity of the underlying algebra's constructors. Otherwise we would in general have $(\llbracket e \nabla \otimes \rrbracket [1, 2, 4]^{\text{set}} \neq \llbracket e \nabla \otimes \rrbracket [2, 4, 4, 1, 2]^{\text{set}}$ which clearly is not what we aim for.

To be more precise, assume a type \mathbb{T} in insert representation with a left-commutative constructor \oplus . A catamorphism $(\llbracket e \nabla \otimes \rrbracket)$ over values of \mathbb{T} does not care about the actual order of constructor applications, *i. e.*, we have

$$\begin{aligned} & (\llbracket e \nabla \otimes \rrbracket) \cdot \oplus \cdot (\text{outl } \Delta (\oplus \cdot \text{outr})) \\ &= \otimes \cdot ((\text{outl} \cdot \text{outr}) \Delta (\otimes \cdot (\text{id} \times \text{outr}))) \cdot (\text{id} \times (\text{id} \times (\llbracket e \nabla \otimes \rrbracket))) . \end{aligned} \quad \text{CATA-COMM}$$

Read point-wise and slightly simplified, this is equivalent to the equality

$$(\llbracket e \nabla \otimes \rrbracket) (y \oplus x \oplus xs) = x \otimes ((\llbracket e \nabla \otimes \rrbracket) (y \oplus xs)) .$$

Let \ominus analogously denote a left-idempotent insertion constructor. A duplicate element at the front of a spine does not influence a catamorphism action. We express this by

$$(\llbracket e \nabla \otimes \rrbracket) \cdot \ominus \cdot (\text{outl } \Delta \ominus) = \otimes \cdot (\text{id} \times (\llbracket e \nabla \otimes \rrbracket)) , \quad \text{CATA-IDEM}$$

or, expressed point-wise,

$$(\llbracket e \nabla \otimes \rrbracket) (x \ominus x \ominus xs) = x \otimes (\llbracket e \nabla \otimes \rrbracket) xs .$$

44 The complexity of **CATA-COMM** and **CATA-IDEM** may make it hard to detect, but both equalities follow a common theme. If we let τ denote the initial algebra for the type \mathbb{T} in insert representation (so that we have $\oplus = \tau \cdot \text{inr}$ and $\ominus = \tau \cdot \text{inr}$ in the case of **CATA-COMM** and **CATA-IDEM**, respectively), the equalities adhere to the pattern

$$(\llbracket f \rrbracket) \cdot L\tau = Rf \cdot H(\llbracket f \rrbracket) ,$$

if we instantiate functions L , R , and H for **CATA-COMM** as $L\tau = (\tau \cdot \text{inr}) \cdot (\text{outl } \Delta ((\tau \cdot \text{inr}) \cdot \text{outr}))$, $Rf = (f \cdot \text{inr}) \cdot ((\text{outl} \cdot \text{outr}) \Delta ((f \cdot \text{inr}) \cdot (\text{id} \times \text{outr})))$, and $H(\llbracket f \rrbracket) = \text{id} \times (\text{id} \times (\llbracket f \rrbracket))$. For equality **CATA-IDEM** we get $L\tau = (\tau \cdot \text{inr}) \cdot (\text{outl } \Delta (\tau \cdot \text{inr}))$, $Rf = f \cdot \text{inr}$, and $H(\llbracket f \rrbracket) = \text{id} \times (\llbracket f \rrbracket)$.

Functions L and R are examples of *transformers* [42, 43], as we will soon see. Transformers are categorical encodings of *terms with variables* of type

H. In the form we are using them here, “terms” L and R play the role of representing *the left- and right-hand side (lhs and rhs)* of a datatype equation, respectively. For our purposes, transformers make the burdensome introduction of non-categorical notions like *signature*, *syntax*, and *variables*—which are otherwise needed to reason about terms—superfluous. Using transformers, we can stay in the categorical realm.

45 Definition. Let $h:A \rightarrow B$ be an F -homomorphism from F -algebra $\alpha:FA \rightarrow A$ to $\beta:FB \rightarrow B$, and let H denote another endofunctor. An arrow $T:\mathbf{Alg}(F) \rightarrow \mathbf{Alg}(H)$ is called a **transformer** of type H if

$$h \cdot \alpha = \beta \cdot Fh \quad \Rightarrow \quad h \cdot T\alpha = T\beta \cdot Hh . \quad \text{TRANSFORMER}$$

Rephrasing **TRANSFORMER** diagrammatically clearly reveals T 's property of preserving the homomorphism h :

$$\begin{array}{ccc} \begin{array}{ccc} FA & \xrightarrow{\alpha} & A \\ \downarrow Fh & & \downarrow h \\ FB & \xrightarrow{\beta} & B \end{array} & \Rightarrow & \begin{array}{ccc} HA & \xrightarrow{T\alpha} & A \\ \downarrow Hh & & \downarrow h \\ HB & \xrightarrow{T\beta} & B \end{array} \end{array}$$

TRANSFORMER

//

46 Verifying the **TRANSFORMER** property for equalities **CATA-COMM** and **CATA-IDEM** is a straightforward but somewhat tedious task. We will establish the validity of **TRANSFORMER** for the left-hand side term $L\tau = (\tau \cdot \text{inr}) \cdot (\text{outl} \Delta ((\tau \cdot \text{inr}) \cdot \text{outr}))$ of type $Hh = \text{id} \times (\text{id} \times h)$ in equality **CATA-COMM** only here. The other terms are shown to be transformers using similar calculational arguments.

So let $\alpha = e \nabla \otimes : FA \rightarrow A$ and $\beta = z \nabla \oplus : FB \rightarrow B$ be two F -algebras and $h:A \rightarrow B$ any F -homomorphism (which does not necessarily need to be a catamorphism as they show up in the equalities) between these as in the above definition:

$$\begin{aligned} & h \cdot L\alpha = L\beta \cdot Hh \\ \equiv & \quad \{ \text{unfold } L \text{ and } H, \text{ product} \} \\ & h \cdot \otimes \cdot (\text{outl} \Delta (\otimes \cdot \text{outr})) \\ = & \quad \oplus \cdot (\text{outl} \Delta (\oplus \cdot \text{outr})) \cdot (\text{id} \times (\text{id} \times h)) \\ \equiv & \quad \{ h \text{ homomorphism} \} \end{aligned}$$

$$\begin{aligned}
& \oplus \cdot (id \times h) \cdot (outl \Delta (\otimes \cdot outr)) \\
= & \oplus \cdot (outl \Delta (\oplus \cdot outr)) \cdot (id \times (id \times h)) \\
\equiv & \{ \text{absorption} \} \\
& \oplus \cdot (outl \Delta (h \cdot \otimes \cdot outr)) \\
= & \oplus \cdot (outl \Delta (\oplus \cdot outr)) \cdot (id \times (id \times h)) \\
\equiv & \{ \text{product} \} \\
& \oplus \cdot (outl \Delta (h \cdot \otimes \cdot outr)) \\
= & \oplus \cdot ((outl \cdot (id \times (id \times h))) \Delta (\oplus \cdot outr \cdot (id \times (id \times h)))) \\
\equiv & \{ \text{product: } outl(f \times g) = f \cdot outl, outr(f \times g) = g \cdot outr \} \\
& \oplus \cdot (outl \Delta (h \cdot \otimes \cdot outr)) = \oplus \cdot ((outl \cdot id) \Delta (\oplus \cdot (id \times h) \cdot outr)) \\
\equiv & \{ h \text{ homomorphism, identity} \} \\
& \oplus \cdot (outl \Delta (h \cdot \otimes \cdot outr)) = \oplus \cdot (outl \Delta (h \cdot \otimes \cdot outr)) \\
\equiv & \{ \text{reflexivity of } = \} \\
& true .
\end{aligned}$$

The transformer notion is the only tool we have to fall back upon to define equations for a datatype. Two transformers sharing the same type encode the lhs and rhs terms of an equality, so it is obvious to simply define a datatype equation as a pair of correctly typed transformers. This is exactly what the following definition states.

47 Definition. A pair of transformers (L, R) defines a **datatype equation** (or **equation** for short) if L and R are of the same type. Equation (L, R) **holds** for algebra α if $L\alpha = R\alpha$. //

48 The equalities **CATA-COMM** and **CATA-IDEM** of Paragraph 43 do not only describe the action of a catamorphism on the non-free **bag** and **set** algebras, they also induce the transformers we need to characterize these algebras themselves. For any algebra α in insert representation, the datatype equation describing the left-commutativity of insertion is given by the transformers of **CATA-COMM**, the left-idempotence equation is derived from the transformers occurring in **CATA-IDEM**:

$$\begin{aligned}
& ((\alpha \cdot inr) \cdot (outl \Delta ((\alpha \cdot inr) \cdot outr)) , \\
& (\alpha \cdot inr) \cdot ((outl \cdot outr) \Delta ((\alpha \cdot inr) \cdot (id \times outr)))) & \text{EQ-COMM} \\
& ((\alpha \cdot inr) \cdot (outl \Delta (\alpha \cdot inr)) , \alpha \cdot inr) . & \text{EQ-IDEM}
\end{aligned}$$

Imposing **EQ-COMM** on the list algebra will give us a **bag** algebra as we were calling for in Paragraph 16: if **EQ-COMM** holds for the algebra $[]^{\text{bag}} \nabla^{\text{bag}}$, then, according to Definition 47, we have

$$\begin{aligned}
& L([]^{\text{bag}} \nabla^{\text{bag}}) = R([]^{\text{bag}} \nabla^{\text{bag}}) \\
\equiv & \quad \{ \text{instantiate } L, R \text{ for } \mathbf{EQ-COMM}, \text{ product} \} \\
& \text{bag} \cdot (\text{outl} \Delta (\text{bag} \cdot \text{outr})) = \text{bag} \cdot ((\text{outl} \cdot \text{outr}) \Delta (\text{bag} \cdot (\text{id} \times \text{outr}))) \\
\equiv & \quad \{ \text{point-wise, apply both sides to } (y, (x, xs)) \} \\
& y^{\text{bag}}(x^{\text{bag}}xs) = x^{\text{bag}}(y^{\text{bag}}xs) .
\end{aligned}$$

The construction of a **set** algebra includes the above step and additionally uses equation **EQ-IDEM** to establish the left-idempotence of constructor set :

$$\begin{aligned}
& L([]^{\text{set}} \nabla^{\text{set}}) = R([]^{\text{set}} \nabla^{\text{set}}) \\
\equiv & \quad \{ \text{instantiate } L, R \text{ for } \mathbf{EQ-IDEM}, \text{ product} \} \\
& \text{set} \cdot (\text{outl} \Delta \text{set}) = \text{set} \\
\equiv & \quad \{ \text{point-wise, apply both sides to } (x, xs) \} \\
& x^{\text{set}}(x^{\text{set}}xs) = x^{\text{set}}xs .
\end{aligned}$$

49 Imposing an equation (L, R) on an \mathbf{F} -algebra $\alpha: \mathbf{F}A \rightarrow A$ induces an equivalence relation Eq on $\mathbf{F}A$ that is a congruence for α [42]. Terms in $\mathbf{F}A$ that belong to the same equivalence class denote the same element of the carrier A —which renders the constructors described by \mathbf{F} non-injective and thus introduces the desired confusion.

More precisely, Eq is the least equivalence relation containing all pairs $((L\alpha)t, (R\alpha)t$ for terms t in $\mathbf{F}A$. Additionally, Eq is a congruence for α so that Eq features a closure property with respect to the constructors encoded in \mathbf{F} . For the insert representation case, which means that α is of the form $e \nabla \otimes$, this boils down to (let x, y, xs, ys be in $\mathbf{F}A$)

$$(x, y) \in Eq \wedge (xs, ys) \in Eq \quad \Rightarrow \quad (x \otimes xs, y \otimes ys) \in Eq .$$

Forcing the equation upon the algebra thus means to work with the *quotient algebra*, denoted $\alpha/(L, R)$. This construction extends to and preserves initial algebras: if $\mathbf{Alg}(\mathbf{F})$ includes an initial object τ , then $\tau/(L, R)$ is initial in the subcategory of $\mathbf{Alg}(\mathbf{F})$ containing all and only those \mathbf{F} -algebras for which (L, R) holds. For a proof of this statement see Fokkinga’s original work on transformers [42, page 122 ff].

50 Let us not proceed without explicitly noting that query processing with catamorphisms over non-free algebras like those for **bag** and **set** introduces a well-definedness condition for any catamorphism application.

To illustrate this caveat, suppose that τ is an initial F -algebra upon which we have forced the equation (L, R) of type H as described in the previous paragraph. Furthermore, let α denote another algebra, implementing the target datatype of a query. With these ingredients, we have

$$\begin{aligned}
 & \text{true} \\
 \equiv & \quad \{ \text{reflexivity of } = \} \\
 & \langle \alpha \rangle = \langle \alpha \rangle \\
 \equiv & \quad \{ \text{equation } (L, R) \text{ holds for } \tau \} \\
 & \langle \alpha \rangle \cdot L\tau = \langle \alpha \rangle \cdot R\tau \\
 \Rightarrow & \quad \{ \text{TRANSFORMER on both sides} \} \\
 & L\alpha \cdot H\langle \alpha \rangle = R\alpha \cdot H\langle \alpha \rangle ,
 \end{aligned}$$

which tells us that (L, R) holds for α as well, at least on the range of $\langle \alpha \rangle$.

In the sequel, which will put the theory into query processing practice, this observation will require us to make sure that a catamorphism applied to a bag or set has as its target an algebra with a left-commutative or left-commutative as well as left-idempotent insertion constructor, respectively. This is in compliance with our initial conjecture that a catamorphism must be indifferent to the confusion in the algebra it is applied to.

A check for the left-commutativity and left-idempotence of a constructor is undecidable in general [14, 20]. However, as the query processing context only requires a rather narrow and fixed spectrum of algebras, we can indicate the relevant properties of the involved constructors beforehand, thus freeing the system from the burden to deduce these. Actually, we find ourselves in an even better position: user-level queries will be mapped into the world of initial algebras and catamorphisms by means of an intermediate language, comprehensions (see Section 2.9), whose programs yield well-defined spine transformer compositions only.

We have now assembled everything that is needed to give the final specifications of the collection formers and their algebras that we will encounter in this text. All that is left is to put the pieces together.

51 Definition. Let $F(E, TE) = 1 + E \times TE$. The **list**, **bag**, and **set collection type formers in insert representation** are defined by the following

initial types induced by F , respectively:

$$\begin{aligned} & (nil^{\text{list}} \nabla cons^{\text{list}} \quad , \text{list}) \\ & ((nil^{\text{list}} \nabla cons^{\text{list}})/\text{EQ-COMM} \quad , \text{bag}) \\ & ((nil^{\text{list}} \nabla cons^{\text{list}})/\text{EQ-COMM-IDEM} \quad , \text{set}) , \end{aligned}$$

where the equation **EQ-COMM-IDEM** is given by

$$(L_C \nabla L_I, R_C \nabla R_I) , \quad \text{EQ-COMM-IDEM}$$

with $(L_C, R_C) = \text{EQ-COMM}$ and $(L_I, R_I) = \text{EQ-IDEM}$. Building the sum of the lhs and rhs of two equations like this retains **TRANSFORMER** (this is easily checked). The resulting equation forms the conjunction of the two original equations which is a direct consequence of the properties of categorical sum:

$$(L_C \nabla L_I) = (R_C \nabla R_I) \iff L_C = R_C \wedge L_I = R_I .$$

We will use the shorthand $(nil^{\text{bag}} \nabla cons^{\text{bag}}) = (nil^{\text{list}} \nabla cons^{\text{list}})/\text{EQ-COMM}$ to render the notation more compact. Analogously for **set**. \parallel

2.8 Monads

52 Now that the gory details of the underlying categorical machinery have all been set, let us again take a bird’s eye view to see how the theory developed so far fits into our overall idea of query processing.

Initial types and catamorphisms are the core concepts for *internal* query representation and manipulation in this text. In Chapter 3 we will thus face the challenge of providing a translation from an *external* query syntax into this formal representation. To be honest, the gap between a user-level SQL or OQL query and its equivalent expressed by means of catamorphisms or, more generally, spine transformers, seems to be quite wide.

Of the complications that will crop up, let us merely highlight two here:

- the presence of *variables* in user-level queries, and
- the arbitrarily deep *nesting* of query clauses.

A query representation based on algebras as we have chosen it in this text is an inherently variable-free formalism. It is not at all obvious how the binding and scoping of variables translates into a catamorphism-based algebraic equivalent, given that catamorphisms are—in the lingo of λ -calculus—*combinators*, *i. e.*, closed expressions that are not dependent on an environment of variables. The issue is indeed relevant: we will meet examples of

OQL queries where binding a formerly free variable completely changes the query semantics. We have to find ways to reflect this phenomenon.

Query nesting poses further problems. A composition of spine transformers is basically a *flat* sequence of operations. Nested queries have to be flattened or *unnested* prior to their translation into catamorphisms and type functors. We are thus better off to choose an internal query representation that facilitates such unnesting and related simplification procedures.

53 These observations suggested the employment of a query representation form that will act as a mediator between the user and the catamorphism-determined world.

It is one of the major claims of this text, that a very general query calculus notation, *monad comprehensions* [117]⁶, can successfully play the role of this intermediate representation [58, 60, 101]. Comprehensions feature variable binding and scoping much like the user query languages. Nesting is naturally expressed by means of comprehensions, and—by far more important for our purposes—comprehension expressions enjoy a *normal form* [38, 58] whose derivation actually implements the unnesting transformations we have just identified as a major problem source.

This notion of a query calculus is general in the sense that it is based upon *monads* [75, 79], algebras that offer just enough structure to interpret the constructs of a calculus in the style of the relational calculus [13, 117]. The initial types in insert representation induce monad instances, as we will shortly see. Putting the monad operations down to the now well-known expressions over initial algebras will ensure that the transition from monad comprehensions to catamorphisms is not too hard to define. We may indeed perceive monad comprehensions as mere syntactic sugar for their defining initial algebra expressions. Monad comprehensions, however, enable us to take the before mentioned hurdles with remarkable ease, which is the ultimate rationale behind their employment in this text.

Last but not least: the rather terse comprehension expressions—whose syntactic appearance closely resembles that of Zermelo-Fraenkel (ZF) expressions [61]—are more easily grasped by the human mind and eye than ever longer sequences of catamorphisms.

In what follows now, we will provide a monad notion based on our initial algebras and define the calculus language on top of it. The calculus will be put to full use in the upcoming chapters.

⁶The title of this thesis is actually a pun on the title of [Wadler's](#) seminal article *Comprehending Monads* [117].

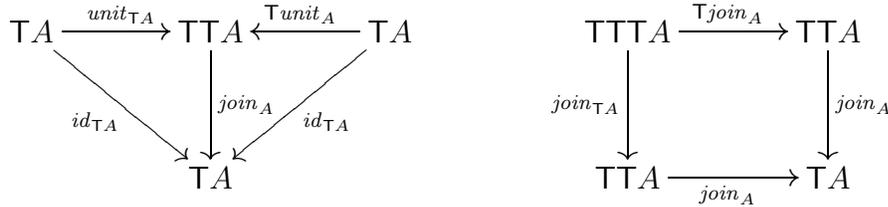
54 Definition. For an endofunctor \mathbb{T} in a category \mathbf{C} (as before, $\mathbf{C} = \mathbf{Set}$ will be the only instantiation of interest for us) we define a **\mathbb{T} -monad** to mean a triple $(\mathbb{T}, \mathit{unit}_A: A \rightarrow \mathbb{T}A, \mathit{join}_A: \mathbb{T}\mathbb{T}A \rightarrow \mathbb{T}A)$ so that

$$\mathit{join}_A \cdot \mathit{unit}_{\mathbb{T}A} = \mathit{id}_{\mathbb{T}A}, \quad \text{MONAD-1}$$

$$\mathit{join}_A \cdot \mathbb{T}\mathit{unit}_A = \mathit{id}_{\mathbb{T}A}, \quad \text{MONAD-2}$$

$$\mathit{join}_{\mathbb{T}A} \cdot \mathit{join}_A = \mathit{join}_A \cdot \mathbb{T}\mathit{join}_A. \quad \text{MONAD-3}$$

The involved types are more easily read off the corresponding diagrams:



If there is the danger of confusing the operations of different monads, we will indicate the monad as a superscript as in $\mathit{join}^{\mathbb{T}}$. Note that the arrow mapping part of functor \mathbb{T} plays the role of the monadic $\mathit{map}^{\mathbb{T}}$ appearing in [Wadler's monad definition \[117\]](#). //

55 In the above definition, subscript A ranges over all objects of \mathbf{C} , so that unit_A and join_A indeed define two families of arrows. For $\mathbf{C} = \mathbf{Set}$, this essentially amounts to saying that unit and join are polymorphic functions with types $\mathit{unit}: \forall A. A \rightarrow \mathbb{T}A$ and $\mathit{join}: \forall A. \mathbb{T}\mathbb{T}A \rightarrow \mathbb{T}A$.

Due to *parametricity* [116], *i. e.*, only by looking at the polymorphic types of unit and join , we can tell that the two monadic operators satisfy the *free theorems* (let $f: A \rightarrow B$ in \mathbf{C} be arbitrary):

$$\begin{aligned} \mathbb{T}f \cdot \mathit{unit}_A &= \mathit{unit}_B \cdot f \\ \mathbb{T}f \cdot \mathit{join}_A &= \mathit{join}_B \cdot \mathbb{T}\mathbb{T}f, \end{aligned}$$

which, to a category theorist anyway, reveals unit and join to be *natural transformations* of types $\mathit{unit}: \mathit{Id} \rightarrow \mathbb{T}$ and $\mathit{join}: \mathbb{T}\mathbb{T} \rightarrow \mathbb{T}$, respectively.

56 Definition. Given two functors of the same type $F, G: \mathbf{C} \rightarrow \mathbf{D}$, a family of arrows $g_A: FA \rightarrow GA$ (for each A in \mathbf{C}) forms a **natural transformation** from F to G (in symbols: $g: F \rightarrow G$), if the *naturality* condition

$$Gf \cdot g_A = g_B \cdot Ff \quad \text{i. e.,} \quad \begin{array}{ccc} FA & \xrightarrow{g_A} & GA \\ Ff \downarrow & & \downarrow Gf \\ FB & \xrightarrow{g_B} & GB \end{array} \quad \text{NATURALITY}$$

is met for each arrow $f:A \rightarrow B$ in \mathbf{C} . Naturality, at least to a certain extent, thus serves as a categorical encoding of the notion of parametric polymorphism. \parallel

57 We will employ monad comprehensions as a *query calculus*, so that they clearly have to support predicates that can act as filters constraining an expression's result to those elements that satisfy the predicate.

The goal to support filters in monad comprehensions (see Section 2.9) will require us to work with slightly extended monads, namely *monads with zero* [117]. To complete the monad introduction, let us define this enriched monad concept here before we go straight on to derive the actual monads that will be useful in defining our query calculus.

58 Definition. If $(\mathbb{T}, \text{unit}, \text{join})$ denotes a \mathbb{T} -monad in \mathbf{C} , we can derive a **\mathbb{T} -monad with zero** $(\mathbb{T}, \text{unit}, \text{join}, \text{zero}_A: 1 \rightarrow \mathbb{T}A)$ from it, provided that

$$\text{join}_A \cdot \text{zero}_{\mathbb{T}A} = \text{zero}_A \quad \text{MONAD-4}$$

$$\text{join}_A \cdot \mathbb{T}\text{zero}_A = \text{zero}_A . \quad \text{MONAD-5}$$

As before, A ranges over the objects of \mathbf{C} . Operator *zero*, just like *unit* and *join*, defines a polymorphic function in **Set** and fulfills law **NATURALITY** (with $f:A \rightarrow B$) which is once more a consequence of parametricity (of *zero*):

$$\mathbb{T}f \cdot \text{zero}_A = \text{zero}_B \cdot f ,$$

i. e., $\text{zero}:\mathbb{K}_1 \dot{\rightarrow} \mathbb{T}$ (observe that the result of *zero* is independent of its argument). \parallel

59 The step from an initial datatype (τ, \mathbb{T}) in insert representation to a \mathbb{T} -monad with zero is actually a small one, as we will now see. Following the construction outlined below we will in particular derive *collection monads*, *i. e.*, **list**, **bag**, and **set** monads according to the collection formers we defined in Definition 51.

60 Lemma. Any initial type $([\]^{\mathbb{T}} \nabla \mathbb{T}, \mathbb{T})$ induced by the bifunctor $F(f, g) = \text{id}_1 + f \times g$ gives rise to a \mathbb{T} -monad $(\mathbb{T}, \text{unit}, \text{join}, \text{zero})$ with zero by defining

$$\begin{aligned} \text{unit} &= \mathbb{T} \cdot (\text{id} \Delta [\]^{\mathbb{T}}) \\ \text{join} &= ([\]^{\mathbb{T}} \nabla \#) \\ \text{zero} &= [\]^{\mathbb{T}} \end{aligned}$$

where function $\# : \mathbb{T}A \times \mathbb{T}A \rightarrow \mathbb{T}A$ is defined to mean

$$xs \# ys = (\llbracket ys \nabla \mathbb{T} \rrbracket) xs .$$

Proof. First, it is easily checked that the so defined monadic operators are correctly typed. Second, observe that *join* is well-defined regardless of \mathbb{T} (respectively \mathbb{T}) being left-commutative or left-idempotent: as $\#$ merely concatenates the spines of its two arguments, $\#$ is commutative or idempotent whenever \mathbb{T} is (in this case we will call the \mathbb{T} -monad itself *commutative* or *idempotent*, respectively). This, in turn, makes **EQ-COMM** or **EQ-IDEM** hold for the algebra $\llbracket \rrbracket^{\mathbb{T}} \nabla \#$, *i. e.*,

$$ys \# (xs \# zs) = xs \# (ys \# zs) \quad \text{and} \quad xs \# (xs \# zs) = xs \# zs ,$$

which ensures the well-definedness of *join* (see Paragraph 50). Observe that $\#$ and *join* commute in the sense that

$$join \cdot \# = \# \cdot (join \times join) ,$$

which essentially is a consequence of the free theorem corresponding to the polymorphic type of $\#$ [116]. We additionally have that $\llbracket \rrbracket^{\mathbb{T}}$ is a left and right unit for $\#$:

$$\begin{aligned} & \llbracket \rrbracket^{\mathbb{T}} \# xs & xs \# \llbracket \rrbracket^{\mathbb{T}} \\ = & \{ \text{unfold } \# \} & = & \{ \text{unfold } \# \} \\ & (\llbracket xs \nabla \mathbb{T} \rrbracket) \llbracket \rrbracket^{\mathbb{T}} & & (\llbracket \rrbracket^{\mathbb{T}} \nabla \mathbb{T} \rrbracket) xs \\ = & \{ \text{CATA-INS-REP} \} & = & \{ \text{CATA-REFLECT} \} \\ & xs & & xs . \end{aligned}$$

The main work we have to do here, namely the verification of the monad laws **MONAD-1** through **MONAD-5**, is carried out calculationally. At several points we will make good use of the catamorphism fusion laws **CATA-FUSION** and **CATA-MAP-FUSION** to shorten the calculations.

MONAD-1:

$$\begin{aligned} & join \cdot unit \\ = & \{ \text{unfold monadic operators} \} \\ & (\llbracket \rrbracket^{\mathbb{T}} \nabla \#) \cdot (\mathbb{T} \cdot (id \nabla \llbracket \rrbracket^{\mathbb{T}})) \\ = & \{ \text{CATA-INS-REP} \} \end{aligned}$$

$$\begin{aligned}
& \# \cdot (id \nabla []^T) \\
= & \{ []^T \text{ is unit for } \# \} \\
& id .
\end{aligned}$$

MONAD-2:

$$\begin{aligned}
& join \cdot Tunit \\
= & \{ \text{unfold, CATA-MAP-FUSION} \} \\
& \langle ([]^T \nabla \#) \cdot F(\text{unit}, id) \rangle \\
= & \{ \text{unfold } F \} \\
& \langle ([]^T \nabla \#) \cdot (id + (\text{unit} \times id)) \rangle \\
= & \{ \text{sum} \} \\
& \langle ([]^T \nabla (\# \cdot (\text{unit} \times id))) \rangle \\
= & \{ \text{unfold } unit \} \\
& \langle ([]^T \nabla (\# \cdot ((\dagger \cdot (id \nabla []^T)) \times id))) \rangle \\
= & \{ \# \text{ is catamorphism on its first argument} \} \\
& \langle ([]^T \nabla (\dagger \cdot (id \times (\# \cdot ([]^T \times id)))) \rangle \\
= & \{ []^T \text{ is unit for } \# \} \\
& \langle ([]^T \nabla (\dagger \cdot (id \times id))) \rangle \\
= & \{ \text{product, CATA-REFLECT} \} \\
& id .
\end{aligned}$$

MONAD-3:

$$\begin{aligned}
& join \cdot join = join \cdot Tjoin \\
\equiv & \{ \text{CATA-MAP-FUSION} \} \\
& join \cdot join = \langle ([]^T \nabla \#) \cdot (id + (join \times id)) \rangle \\
\equiv & \{ \text{sum} \} \\
& join \cdot join = \langle ([]^T \nabla (\# \cdot (join \times id))) \rangle \\
\stackrel{(*)}{\equiv} & \{ \text{CATA-FUSION} \} \\
& true .
\end{aligned}$$

To see the applicability of **CATA-FUSION** at (*), observe that *join* is an *F*-homomorphism from $[]^T \nabla \#$ to $[]^T \nabla (\# \cdot (join \times id))$:

$$join \cdot ([]^T \nabla \#) = ([]^T \nabla (\# \cdot (join \times id))) \cdot F(id, join)$$

$$\begin{aligned}
&\equiv \{ \text{unfold } F \} \\
&\quad \text{join} \cdot ([]^T \nabla \#) = ([]^T \nabla (\# \cdot (\text{join} \times \text{id}))) \cdot (\text{id} + (\text{id} \times \text{join})) \\
&\equiv \{ \text{sum} \} \\
&\quad \text{join} \cdot ([]^T \nabla \#) = []^T \nabla (\# \cdot (\text{join} \times \text{join})) \\
&\equiv \{ \text{sum} \} \\
&\quad (\text{join} \cdot []^T) \nabla (\text{join} \cdot \#) = []^T \nabla (\# \cdot (\text{join} \times \text{join})) \\
&\equiv \{ \text{join is catamorphism, free theorem for } \# \} \\
&\quad \text{true} .
\end{aligned}$$

MONAD-4:

$$\begin{aligned}
&\quad \text{join} \cdot \text{zero} \\
&= \{ \text{unfold} \} \\
&\quad ([]^T \nabla \#) \cdot []^T \\
&= \{ \text{CATA-INS-REP} \} \\
&\quad []^T \\
&= \{ \text{definition of } \text{zero} \} \\
&\quad \text{zero} .
\end{aligned}$$

MONAD-5:

$$\begin{aligned}
&\quad \text{join} \cdot \top \text{zero} \\
&= \{ \text{unfold, CATA-MAP-FUSION} \} \\
&\quad ([]^T \nabla \#) \cdot F(\text{zero}, \text{id}) \\
&= \{ \text{unfold } F, \text{zero} \} \\
&\quad ([]^T \nabla \#) \cdot (\text{id} + ([]^T \times \text{id})) \\
&= \{ \text{sum} \} \\
&\quad ([]^T \nabla (\# \cdot ([]^T \times \text{id}))) \\
&= \{ []^T \text{ is unit for } \# \} \\
&\quad ([]^T, \text{outr}) \\
&= \{ \text{CATA-INS-REP} \} \\
&\quad []^T \\
&= \{ \text{definition of } \text{zero} \} \\
&\quad \text{zero} .
\end{aligned}$$

□

61 Collection monads. Applying the above procedure to the three collection formers of Definition 51 constructs three *collection monads*. In **Set**, the category in which we model the query evaluation process, we specifically get

$$\begin{aligned} (\text{list}, \text{unit}^{\text{list}}, \text{join}^{\text{list}}, \text{zero}^{\text{list}}) &= (\text{list}, \lambda x.[x], \text{flatten}, []) \\ (\text{bag}, \text{unit}^{\text{bag}}, \text{join}^{\text{bag}}, \text{zero}^{\text{bag}}) &= (\text{bag}, \lambda x.\{\{x\}\}, \uplus, \{\{\}\}) \\ (\text{set}, \text{unit}^{\text{set}}, \text{join}^{\text{set}}, \text{zero}^{\text{set}}) &= (\text{set}, \lambda x.\{x\}, \cup, \emptyset), \end{aligned}$$

where square brackets $[\cdot]$ and double curly brackets $\{\{\cdot\}\}$ are used to denote list and bag values, respectively. Function $\text{flatten} : \text{list list } A \rightarrow \text{list } A$ concatenates a list of lists into a flat list as in

$$\text{flatten} [[x_1, x_2], [], [x_3, x_4, x_5], [x_6]] = [x_1, x_2, x_3, x_4, x_5, x_6].$$

The natural transformations $\cup : \text{set set } A \rightarrow \text{set } A$ and $\uplus : \text{bag bag } A \rightarrow \text{bag } A$ flatten a set of sets and a bag of bags in much the same manner. These three monads provide the primary collection abstraction in the discussion ahead.

62 Our use of monads, however, goes further than that. Although the name *insert representation* might suggest so, we can perfectly imagine algebras $e \nabla \otimes$ in insert representation that do *not* model collection type constructors. Instead, we can instantiate e and \otimes in a way that provides us with *aggregation operators* or *quantifiers*. Everything we have said about such algebras until now remains valid, with the difference that a spine over these algebras describes an aggregation or quantification and does not construct a collection value. Consequently, these algebras operate on carrier sets which correspond to *atomic types*. The monads associated with these non-collection algebras then embed aggregation and quantification seamlessly into the soon to be defined monad comprehension calculus (see Section 2.9). Much of the elegance of the approach devised in this text is derived from this uniformity.

63 So let E denote the domain of an atomic type like *Num* or *Bool*; additionally let the identity functor Id instantiate the type functor \mathbb{T} . The endofunctor describing the insert representation constructors $\mathbb{F}_E = \mathbb{K}_1 + \mathbb{K}_E \times \text{Id}$ remains unchanged. We obtain a free non-collection algebra over carrier E with constructors $e : 1 \rightarrow E$ and $\otimes : E \times E \rightarrow E$ (*cf.* Paragraph 19):

$$\begin{aligned} &e \nabla \otimes : \mathbb{F}_E(\mathbb{T}E) \rightarrow \mathbb{T}E \\ = &\{ \text{unfold } \mathbb{F}, \mathbb{T} \} \\ &e \nabla \otimes : 1 + E \times E \rightarrow E. \end{aligned}$$

64 Definition. Let $F_E = K_I + K_E \times \text{Id}$ and $e \nabla \otimes$ as in the previous paragraph. We define **aggregation operators** and **quantifiers in insert representation** by means of the following F_{Num} and F_{Bool} algebras, respectively (ignore the functor annotations until the next paragraph):

$$\begin{array}{lll}
(0 \nabla +) & = & (e \nabla \otimes) / \text{EQ-COMM} \quad (\text{sum}) \\
(1 \nabla *) & = & (e \nabla \otimes) / \text{EQ-COMM} \quad (\text{prod}) \\
(\infty_{Num} \nabla \text{min}_2) & = & (e \nabla \otimes) / \text{EQ-COMM-IDEM} \quad (\text{min}) \\
(-\infty_{Num} \nabla \text{max}_2) & = & (e \nabla \otimes) / \text{EQ-COMM-IDEM} \quad (\text{max}) \\
(\text{true} \nabla \wedge) & = & (e \nabla \otimes) / \text{EQ-COMM-IDEM} \quad (\text{all}) \\
(\text{false} \nabla \vee) & = & (e \nabla \otimes) / \text{EQ-COMM-IDEM} \quad (\text{exists}) .
\end{array}$$

The binary operator min_2 (max_2) returns the smaller (greater) of its two arguments of type Num . Symbol ∞_{Num} shall be perceived as the largest value representable in a concrete implementation of type Num , so that ∞_{Num} acts a unit with respect to min_2 . A dual remark applies to $-\infty_{Num}$ and max_2 . \parallel

65 As the above algebras are not initial, we cannot mimic the construction of Lemma 60 to derive their monads. There is, however, a monad naturally associated with $e \nabla \otimes$ and the functor Id , namely

$$\begin{array}{l}
\text{unit} = \text{id} \\
\text{join} = \text{id} \\
\text{zero} = e .
\end{array}$$

It is easily checked that $(\text{Id}, \text{unit}, \text{join}, \text{zero})$ satisfies laws **MONAD-1** through **MONAD-5** and thus indeed constitutes a monad, namely, adopting a term used by Wadler in [117], an *identity monad* with zero. In the following, we will use the functor aliases listed in Definition 64 to refer to the identity monads associated with the non-collection algebras, *e. g.*, $\text{exists} = (\text{Id}, \text{id}, \text{id}, \text{false})$.

2.9 Monad Comprehensions

66 We are finally ready to keep our promise given in the introduction of this chapter and take a step towards a higher-level query representation, *monad comprehensions* [19, 60, 91, 117, 124]. Due to their syntactic resemblance with the set comprehension notation and thus with the relational calculus, monad comprehension should be particularly accessible to readers used to

these classical query notations. Of the things to come, much may indeed be grasped by applying knowledge acquired in the relational (tuple) calculus domain, although one is trading generality and uniformity for simplicity then.

To start off with a simple example consider the following comprehension to be interpreted in the `set` monad

$$\llbracket x \mid x \leftarrow xs, x \neq 0 \rrbracket^{\text{set}},$$

which is equivalent to the set comprehension $\{x : x \in xs \wedge x \neq 0\}$ (“the set of all x , so that $x \in xs$ and $x \neq 0$ ”). While $x \in xs$ acts as pool from which successive bindings for variable x are generated (exactly one binding for each element in the set xs), predicate $x \neq 0$ restricts the final result set to include only those x that fulfill the specified condition.

More generally, in the \mathbb{T} -monad comprehension

$$\llbracket e \mid q_1, \dots, q_n \rrbracket^{\mathbb{T}},$$

the *qualifiers* q_i are either *generators* $x_i \leftarrow e_i$ or *filters* (expressions of result type `Bool`).

The informal semantics of this expression are as follows: starting with qualifier q_1 , a generator $q_i = x_i \leftarrow e_i$ enriches an initially empty environment of variable bindings by (sequentially) binding x_i to (elements of) its *range* e_i . (If the range of a generator indeed is of an atomic type, the generator is essentially equivalent to a *let*-expression, see Example 71). The binding of x is propagated through the list of qualifiers q_{i+1}, \dots, q_n . Filters, as they are encountered, inhibit further propagation if they evaluate to *false* in the current environment. *Head* e is evaluated in those environments that pass all qualifiers. The results are then injected into the \mathbb{T} -monad using $unit^{\mathbb{T}}$ and finally accumulated via $join^{\mathbb{T}}$.

Let us state the monad comprehension syntax more precisely before we supersede the above informal semantics by a mapping that translates a \mathbb{T} -monad comprehension into a computation over monad \mathbb{T} .

67 Definition. Let the non-terminal symbol *expr* represent the terms of some basic expression language (think of a variant of the λ -calculus). This language is enriched by **monad comprehension syntax** through the following

monad. We define the **T-monad comprehension desugaring scheme** \mathcal{M} (the stylized M serves as a mnemonic for *monad*) to be given by the recursive function below. The occurring identifiers are chosen to fit with those in the abstract comprehension syntax of Definition 67 and represent their respective syntactic category, *e.g.* qs may be replaced by an arbitrary qualifier list.

$$\mathcal{M} \llbracket e \rrbracket^{\top} = \text{unit}^{\top} (\mathcal{M} e) \quad \mathcal{M}\text{-1}$$

$$\mathcal{M} \llbracket e \rrbracket v \leftarrow e' : \top E \rrbracket^{\top} = \top (\lambda v. \mathcal{M} e) (\mathcal{M} e') \quad \mathcal{M}\text{-2}$$

$$\mathcal{M} \llbracket e \rrbracket v \leftarrow e' : \top' E \rrbracket^{\top} = (\top) (\mathcal{M} \llbracket e \rrbracket v \leftarrow e' \rrbracket^{\top'}) \quad \mathcal{M}\text{-3}$$

$$\mathcal{M} \llbracket e \rrbracket e' \rrbracket^{\top} = \text{if } (\mathcal{M} e') \text{ then } \text{unit}^{\top} (\mathcal{M} e) \text{ else } \text{zero}^{\top} \quad \mathcal{M}\text{-4}$$

$$\mathcal{M} \llbracket e \rrbracket qs, qs' \rrbracket^{\top} = \text{join}^{\top} (\mathcal{M} \llbracket \llbracket e \rrbracket qs' \rrbracket^{\top} \rrbracket qs \rrbracket^{\top}) \quad \mathcal{M}\text{-5}$$

$$\mathcal{M} e = e \quad \mathcal{M}\text{-6}$$

∥

70 Notes. (a) Note that, according to branch $\mathcal{M}\text{-5}$, we can arbitrarily divide the comprehension qualifier list into two sublists qs and qs' without affecting the outcome of \mathcal{M} . This is a consequence of qualifier enumeration (denoted by ‘,’) being an associative operation with the empty qualifier list as its unit [117].

(b) The argument to type functor \top in case $\mathcal{M}\text{-2}$ is an anonymous λ -term, *i.e.*, a morphism in category **Set**. Consequently, scheme \mathcal{M} is bound to the assumption that we are operating with **Set** as the underlying category.

(c) Whenever we encounter a generator whose range is of a foreign monad $\top' \neq \top$, branch $\mathcal{M}\text{-3}$ temporarily switches the monad in which the monad comprehension is to be interpreted in to \top' . The results are then coerced to type \top using (\top) . Observe that the cases $\mathcal{M}\text{-2}$ and $\mathcal{M}\text{-3}$ could be merged to read (with \top' not necessarily different to \top)

$$\mathcal{M} \llbracket e \rrbracket v \leftarrow e' : \top' E \rrbracket^{\top} = (\top) (\top' (\lambda v. \mathcal{M} e) (\mathcal{M} e')), \quad \mathcal{M}\text{-2}'$$

which amounts to the same due to law **CATA-REFLECT**.

Coercion is not completely arbitrary though since the well-definedness condition for catamorphisms derived in Paragraph 50 applies: whenever τ' is a left-commutative respective left-idempotent algebra, so has to be τ . This restriction is quite natural, however, as it forbids ambiguous and non-well-defined coercions, like the non-deterministic conversion of **set** value into a list.

(d) Note how branch $\mathcal{M}\text{-4}$ implements the filter q by mapping values failing to fulfill the filter predicate to zero^{\top} so that they do not contribute to the result in the embracing join^{\top} (see **MONAD-4**).

(e) The exchange of qualifiers in \mathcal{M} -5, in interplay with the introduction of λ -abstractions for each generator $v \leftarrow e'$ in \mathcal{M} -2, ensures the visibility of v in each qualifier occurring later in the qualifier list as well as in the head e .

71 Example. A rather broad diversity of computational tasks may be expressed by monad comprehensions and it is the crucial benefit of employing the comprehension notation that this variety is mapped to a completely uniform representation. Quite often we will be able to abstract from the specific monad the query is to be evaluated in and instead reason in the context of *some* monad \mathbb{T} . Such general findings may then be reused during the query translation process in several instantiations and places.

The following examples are, in some sense, teasers that are given here to provide some insight of what the upcoming chapters have in store.

(a) We can now convince ourselves that the initial example of this section, the **set** comprehension $\llbracket x \mid x \leftarrow xs, x \neq 0 \rrbracket^{\text{set}}$, is indeed equivalent to $\{x : x \in xs \wedge x \neq 0\}$. By means of desugaring we get

$$\begin{aligned}
& \mathcal{M} \llbracket x \mid x \leftarrow xs, x \neq 0 \rrbracket^{\text{set}} \\
= & \quad \{ \mathcal{M}\text{-5} \} \\
& \text{join}^{\text{set}} \left(\mathcal{M} \llbracket \llbracket x \mid x \neq 0 \rrbracket^{\text{set}} \mid x \leftarrow xs \rrbracket^{\text{set}} \right) \\
= & \quad \{ \mathcal{M}\text{-2} \} \\
& \text{join}^{\text{set}} \left(\text{set} (\lambda x. \mathcal{M} \llbracket x \mid x \neq 0 \rrbracket^{\text{set}}) xs \right) \\
= & \quad \{ \mathcal{M}\text{-4} \} \\
& \text{join}^{\text{set}} \left(\text{set} (\lambda x. \text{if} (\mathcal{M} (x \neq 0)) \text{ then } (\text{unit}^{\text{set}} (\mathcal{M} x)) \text{ else } \text{zero}^{\text{set}}) xs \right) \\
= & \quad \{ \mathcal{M}\text{-6} \} \\
& \text{join}^{\text{set}} \left(\text{set} (\lambda x. \text{if} (x \neq 0) \text{ then } (\text{unit}^{\text{set}} x) \text{ else } \text{zero}^{\text{set}}) xs \right) \\
= & \quad \{ \text{set monad} \} \\
& \bigcup (\text{set} (\lambda x. \text{if} (x \neq 0) \text{ then } \{x\} \text{ else } \emptyset) xs) .
\end{aligned}$$

The resulting **set**-monad computation implements the set comprehension by mapping each element x of xs to the singleton set $\{x\}$ if $x \neq 0$. Otherwise x is mapped to the empty set \emptyset so that these elements do not contribute to the result during the outer \bigcup .

(b) Monad comprehensions over monad **all** effectively implement universal quantification. Consider the **all**-monad comprehension expression $\llbracket px \mid x \leftarrow xs : \text{set} E \rrbracket^{\text{all}}$ in which p denotes a predicate of type $E \rightarrow \text{Bool}$. Observe that this already constitutes a multi-monad comprehension which was not expressible by earlier related proposals. We calculate as follows:

$$\mathcal{M} \llbracket px \mid x \leftarrow xs : \text{set} E \rrbracket^{\text{all}}$$

$$\begin{aligned}
&= \{ \mathcal{M}\text{-3} \} \\
&\quad (\text{true} \nabla \wedge) (\mathcal{M} \llbracket p x \mid x \leftarrow xs : \text{set} E \rrbracket^{\text{set}}) \\
&= \{ \mathcal{M}\text{-2} \} \\
&\quad (\text{true} \nabla \wedge) ((\text{set} (\lambda x. \mathcal{M} (p x)) (\mathcal{M} xs))) \\
&= \{ \mathcal{M}\text{-6, } \eta\text{-conversion} \} \\
&\quad (\text{true} \nabla \wedge) ((\text{set } p xs)) .
\end{aligned}$$

The resulting expression obviously computes $\forall x \in xs : p x$. Note that this also does the right thing if $xs = \text{zero}^{\text{set}} = \emptyset$ in which case we get *true* as desired.

(c) Monad comprehensions are not bound to denote *collection comprehensions* (as they are understood in [113, 124]). Here is an example of a non-collection comprehension to be interpreted in the **sum** monad:

$$\begin{aligned}
&\mathcal{M} \llbracket (x, y) \mid x \leftarrow 2, y \leftarrow 3 \rrbracket^{\text{sum}} \\
&= \{ \mathcal{M}\text{-5} \} \\
&\quad \text{join}^{\text{sum}} (\mathcal{M} \llbracket \llbracket (x, y) \mid y \leftarrow 3 \rrbracket^{\text{sum}} \mid x \leftarrow 2 \rrbracket^{\text{sum}}) \\
&= \{ \mathcal{M}\text{-2} \} \\
&\quad \text{join}^{\text{sum}} (\text{sum} (\lambda x. (\mathcal{M} \llbracket (x, y) \mid y \leftarrow 3 \rrbracket^{\text{sum}}) 2)) \\
&= \{ \mathcal{M}\text{-2} \} \\
&\quad \text{join}^{\text{sum}} (\text{sum} (\lambda x. (\text{sum} (\lambda y. \mathcal{M} (x, y)) 3)) 2) \\
&= \{ \text{sum monad} \} \\
&\quad \lambda x. ((\lambda y. (x, y)) 3) 2 .
\end{aligned}$$

The final outcome reveals that the **sum** comprehension boils down to a nested *let*-expression:

$$\text{let } x = 2 \text{ in let } y = 3 \text{ in } (x, y) .$$

Note that we could compute the above in any identity monad and that the actual choice of **sum** is arbitrary.

(d) Let us conclude this introduction to monad comprehensions with a more substantial example. The occurrence of a nested **exists**-comprehension inside an outer **sum**-comprehension does not violate the well-definedness condition for multi-monad comprehensions, as the inner comprehension serves as a filter (comprehensions interpreted in monads **exists** and **all** reduce to values of type *Bool*) and not as a generator range. We assume the typings $xs : \text{listNum}$ and $ys : \text{bagNum}$.

$$\mathcal{M} \llbracket x \mid x \leftarrow xs, \llbracket x = y \mid y \leftarrow ys \rrbracket^{\text{exists}} \rrbracket^{\text{sum}}$$

$$\begin{aligned}
&= \{ \mathcal{M}\text{-5} \} \\
&\quad \text{join}^{\text{sum}} (\mathcal{M} \llbracket [x \mid [x = y \mid y \leftarrow ys]]^{\text{exists}} \rrbracket^{\text{sum}} \mid x \leftarrow xs \rrbracket^{\text{sum}}) \\
&= \{ \mathcal{M}\text{-3} \} \\
&\quad \text{join}^{\text{sum}} ((0 \nabla +) (\mathcal{M} \llbracket [x \mid [x = y \mid y \leftarrow ys]]^{\text{exists}} \rrbracket^{\text{sum}} \mid x \leftarrow xs \rrbracket^{\text{list}})) \\
&= \{ \mathcal{M}\text{-2} \} \\
&\quad \text{join}^{\text{sum}} ((0 \nabla +) (\text{list} (\lambda x. \mathcal{M} \llbracket [x \mid [x = y \mid y \leftarrow ys]]^{\text{exists}} \rrbracket^{\text{sum}}) xs)) \\
&= \{ \mathcal{M}\text{-4} \} \\
&\quad \text{join}^{\text{sum}} ((0 \nabla +) (\text{list} (\lambda x. \text{if} (\mathcal{M} \llbracket [x = y \mid y \leftarrow ys]]^{\text{exists}} \rrbracket^{\text{sum}}) \\
&\quad \quad \text{then } (\text{unit}^{\text{sum}} x) \text{ else } \text{zero}^{\text{sum}}) xs)) \\
&= \{ \mathcal{M}\text{-3} \} \\
&\quad \text{join}^{\text{sum}} ((0 \nabla +) (\text{list} (\lambda x. \text{if} ((\text{false} \nabla \vee) (\mathcal{M} \llbracket [x = y \mid y \leftarrow ys]]^{\text{bag}} \rrbracket^{\text{sum}}) \\
&\quad \quad \text{then } (\text{unit}^{\text{sum}} x) \text{ else } \text{zero}^{\text{sum}}) xs)) \\
&= \{ \mathcal{M}\text{-2} \} \\
&\quad \text{join}^{\text{sum}} ((0 \nabla +) (\text{list} (\lambda x. \text{if} ((\text{false} \nabla \vee) (\text{bag} (\lambda y. x = y) ys)) \\
&\quad \quad \text{then } (\text{unit}^{\text{sum}} x) \text{ else } \text{zero}^{\text{sum}}) xs)) .
\end{aligned}$$

The point to note here is that this expressions shows the typical structure we expected a query to exhibit in our representation form. Unfolding the monadic operators discloses the initial monad comprehension to be completely implemented by means of catamorphisms and type functors:

$$((0 \nabla +) (\text{list} (\lambda x. \text{if} ((\text{false} \nabla \vee) (\text{bag} (\lambda y. x = y) ys)) \text{ then } x \text{ else } 0) xs)) .$$

◇

We have reached a point where the foundations for the internal query representation used in this text have all been laid. The forthcoming chapters will put this categorical toolbox to, as we will argue, beneficial use in query translation, transformation, and optimization.

2.10 Union Representation

72 Although the insert representation of algebras has been pervasive in everything we have developed so far, it is not principle to the method. In fact there is at least one alternative class of algebras, the *algebras in union representation*, which has been extensively studied as an alternative foundation for collection (query and programming) languages [13, 14, 20, 110, 111].

A collection value in union representation is built from *singleton collections* which are then merged using an *associative binary union operator*

\oplus . The union representation comes with an additional constant constructor (function), the *empty collection* z . This last constructor is required to be a unit of the binary union so that (z, \oplus) forms a *monoid*. In other words, general binary trees take the place of the right-deep insertion spines we are employing in this text. The elements in the front of these trees are the collection members.

Remarkably few details have to be adapted to turn the material of this chapter into a discussion of the union representation of types. As we have encoded the type constructors by means of the insert representation functor F , it is only this functor that is to be replaced. In principle, everything else may be left untouched! Let us take the union representation view of the theory of datatypes for a minute.

Given a polymorphic type constructor T , its G -algebra in union representation is

$$empty^T \nabla sng^T \nabla union^T : 1 + E + TE \times TE \rightarrow TE ,$$

i. e., $G(E, TE) = 1 + E + TE \times TE$ and $G(f, g) = id_1 + f + g \times g$.

Obvious instances of this generic algebra, for $T = \text{set}$ and $T = \text{list}$ respectively, include $\emptyset \nabla \lambda x. \{x\} \nabla \cup$ and $[] \nabla \lambda x. [x] \nabla \#$. Setting $T = \text{Id}$ yields non-collection instances (with $sng^{\text{Id}} = id$) as in $0 \nabla id \nabla +$ and $true \nabla id \nabla \wedge$.

As before, $\mathbf{Alg}(G)$ has an initial object τ and it this initial algebra that describes the datatype. The catamorphisms in $\mathbf{Alg}(G)$ turn out to be general *tree transformers*: let $z \nabla s \nabla \oplus$ denote any other G -algebra so that (z, \oplus) has the monoid property, then we get (cf. Paragraph 27)

$$\begin{aligned} \langle z \nabla s \nabla \oplus \rangle empty^T &= z \\ \wedge \langle z \nabla s \nabla \oplus \rangle (sng^T x) &= s x \\ \wedge \langle z \nabla s \nabla \oplus \rangle (union^T(xs, ys)) &= (\langle z \nabla s \nabla \oplus \rangle xs) \oplus (\langle z \nabla s \nabla \oplus \rangle ys) . \end{aligned}$$

CATA-UNION-REP

Not surprisingly, the induced type functor T behaves like a *tree map*, since

$$\begin{aligned} &Tf \\ &= \{ \text{TYPE-FUNCTOR} \} \\ &\quad (\tau \cdot G(f, id)) \\ &= \{ \text{unfold } \tau \text{ and union representation functor } G \} \\ &\quad (\langle empty^T \nabla sng^T \nabla union^T \rangle \cdot (id + f + id \times id)) \\ &= \{ \text{sum} \} \\ &\quad (\langle empty^T \nabla sng^T \cdot f \nabla union^T \rangle) , \end{aligned}$$

which tells us that $\mathbb{T}f$ is the specific tree transformer that applies f to the elements of the front of its argument while retaining the structure of the tree.

From the above it is apparent that the \mathbb{G} -catamorphisms $(\cdot)_\mathbb{G}$ generalize the structural recursion operator over collections in union representation *sru* of [13, 110, 111] and the set divide-and-conquer operator *pump* appearing in the functional database language FAD [4] as well as in the work of [Beeri and Kornatzky](#) on algebraic optimization of object query languages [6].

Initial types in union representation induce monads in much the same manner as their insert representation variants. Specifically, for a given type $(empty^\mathbb{T} \nabla sng^\mathbb{T} \nabla union^\mathbb{T}, \mathbb{T})$, obtain the corresponding \mathbb{T} -monad with zero $(\mathbb{T}, unit, join, zero)$ by defining

$$\begin{aligned} unit &= sng^\mathbb{T} \\ join &= (\cdot)_{empty^\mathbb{T} \nabla id \nabla union^\mathbb{T}} \\ zero &= empty^\mathbb{T} \end{aligned}$$

which exactly generates the monads we have already met in Paragraph 61. This implies that an identical notion of monad comprehensions may alternatively be built on top of types in union representation.

What, then, has been our principle rationale for the choice of the insert representation in this text?

The more obvious reason, which we have already mentioned, is related to a core virtue: laziness. The insert representation makes economical use of just two (instead of three) constructors which, in a world of inductive types and structural recursion, has the pertinent effect of rendering definitions more compact and to simplify calculations. A similar observation can be made for datatype equations. Encoding the monoidal structure of $(empty^\mathbb{T}, union^\mathbb{T})$ requires the specification of two equations even for the $\mathbb{T} = \text{list}$ algebra in union representation. Remember that the list algebra is freely generated in insert representation.

Another, and deeper, reason relates the efficiency and expressiveness of programs written in insert and union representation style, respectively: for some fixed type constructor \mathbb{T} , suppose that $(\sigma = empty^\mathbb{T} \nabla sng^\mathbb{T} \nabla union^\mathbb{T}, \mathbb{T})$ denotes the initial type induced by the functor \mathbb{G} . Given this, we can derive an equivalent \mathbb{F} -algebra $\tau = nil^\mathbb{T} \nabla cons^\mathbb{T}$ in insert representation that operates over the same carrier:

$$\begin{aligned} nil^\mathbb{T} &= empty^\mathbb{T} \\ cons^\mathbb{T} &= union^\mathbb{T} \cdot (sng^\mathbb{T} \times id) \end{aligned}$$

This construction preserves the initiality property: τ is initial in $\mathbf{Alg}(\mathbb{F})$ if and only if σ is initial in $\mathbf{Alg}(\mathbb{G})$ [14]. A query in union representation style,

i. e., effectively a \mathbf{G} -catamorphism $(z \nabla s \nabla \oplus)_{\mathbf{G}}$, thus has a simple first-order and efficient translation into insert representation style:

$$(z \nabla (\oplus \cdot (s \times id)))_{\mathbf{F}} .$$

Can we do the converse? We can certainly construct a \mathbf{G} -algebra σ from a given initial \mathbf{F} -algebra $\tau = nil^{\top} \nabla cons^{\top}$:

$$\begin{aligned} empty^{\top} &= nil^{\top} \\ sng^{\top} &= cons^{\top} \cdot (id \triangle nil^{\top}) \\ union^{\top}(xs, ys) &= (ys \nabla cons^{\top}) xs , \end{aligned}$$

and the resulting algebra will be initial in $\mathbf{Alg}(\mathbf{F})$ as well. The bad news is that the above construction is *not* valid for an arbitrary (especially non-initial) source algebra. This implies that we cannot use this scheme to derive a query in union representation style from a given insertion style query as the query's target algebra will not be initial in general [14]. In fact there is no translation from the insertion to the union style that is simple (first-order) and efficient at the same time:

- In [14], [Breazu-Tannen and Subrahmanyam](#) give a simple translation that makes use of the monoidal structure inherent to (id, \cdot) and thus is higher-order. This would require us to extend our system of type constructors by the function type former (\rightarrow) which we are trying to avoid.
- [Suciu and Wong](#) present a first-order translation that is inefficient in that it maps certain polynomial time computable insertion style queries into union style equivalents that require exponential space [111].

The bottom line of this argument is that a query representation based on the insert representation is at least as expressive as its union representation variant. At the same time, an insertion style query is efficiently computable whenever the equivalent union style query is. The converse is not true.

73 Collection Comprehensions, Kleisli Monads, and Ringads. In a series of articles, [Trinder](#) et al. developed a theory of collection types based on an extended monad notion, the *ringad* [113, 123]. This work used the union representation of a collection type constructor \top as its starting point, *i. e.*, the \mathbf{G} -algebra (which is non-free: $empty^{\top}$ is required to be a unit of $union^{\top}$):

$$empty^{\top} \nabla sng^{\top} \nabla union^{\top} : \mathbf{G}(E, \top E) \rightarrow \top E ,$$

with \mathbf{G} denoting the union representation functor of the previous paragraph. The theory then added an iteration abstraction in the form of the binary operator \gg (pronounced *bind*) to this algebra to complete its notion of a *collection type*. Operator \gg is expected to interact with the above algebra as follows (with $f : E \rightarrow \mathbb{T}E'$):

$$\begin{aligned} (\gg f) \cdot \text{empty}^\top &= \text{empty}^\top \\ (\gg f) \cdot \text{sg}^\top &= f \\ (\gg f) \cdot \text{union}^\top &= \text{union}^\top \cdot ((\gg f) \times (\gg f)), \end{aligned}$$

in which $(\gg f)$ denotes a *section* (partial application) of \gg defined to mean $\lambda x.x \gg f$. Note that \gg provides a restricted form of structural recursion that, in contrast to a general \mathbf{G} -catamorphism $(\cdot)_\mathbf{G}$, is always well-defined:

$$(\gg f) = (\text{empty}^\top \nabla f \nabla \text{union}^\top)_\mathbf{G}.$$

Operator \gg has been shown sufficiently expressive to cover a broad range of query primitives.

Trinder observed that the triple $(\mathbb{T}, \text{sg}^\top, \gg)$ forms a *Kleisli triple* [70], an alternative encoding of the monad concept that is, however, equivalent to the one we introduced in Definition 54 which

$$\begin{aligned} \text{sg}^\top &= \text{unit}^\top \\ (\gg f) &= \text{join}^\top \cdot \mathbb{T}f \end{aligned}$$

suggest (for a formal proof of the equivalence see [80]). Based on this *ringad* structure, *i. e.*, the triple $(\mathbb{T}, \text{sg}^\top, \gg)$ in companion with empty^\top and union^\top (these two constructors indeed add to the expressivity as they are not expressible in the monad, see [20]), Trinder postulated monad comprehensions as an effective notation for collection query languages.

Chan and Trinder developed this basic theory of collection types into a so-called *object comprehension* language that incorporated the features needed to communicate with the object type hierarchy of an object-oriented DBMS [24, 25]. However, they never intended to interpret comprehensions in non-collection monads. Query language concepts like quantification and aggregation have instead been added as additional features outside the collection theory domain.

Along similar lines, Buneman, Naqvi, Tannen, and Wong laid the foundation of a theory of collection types by successively enriching Kleisli monads $(\mathbb{T}, \text{sg}^\top, \gg)$ to ringads, which were subsequently enhanced with a notion of equality and conditional expressions [13, 20]. (Note that for the monad extension operator *ext* occurring in this work we have $\text{ext } f \text{ } xs = xs \gg f$.) The

resulting language, if restricted to operate in the `set` monad, was shown to be exactly as expressive as the nested relational algebra of [Schek and Scholl](#) and related proposals for relational non-first normal form algebras [99].

Building on this foundation, [Wong](#) devised the *Kleisli* collection processing framework which incorporated collection monad comprehensions as a sublanguage [124]. This language relies on *monad morphisms*—mappings between two monads \mathbb{T}' and \mathbb{T} that preserve the monadic operators (see Definition 86 and [117])—to allow for generator ranges of type \mathbb{T}' inside a comprehension to be interpreted in collection monad \mathbb{T} . The relevant translation step is (the rest of the translation resembles desugaring scheme \mathcal{M} of Definition 69)

$$\llbracket e \mid v \leftarrow e' : \mathbb{T}' E \rrbracket^{\mathbb{T}} \rightsquigarrow e' \gg^{\mathbb{T}} (\lambda v. \text{sg}^{\mathbb{T}} \cdot e) ,$$

based on a modified monadic *bind* $\gg^{\mathbb{T}}$ defined to mean

$$\begin{aligned} (\gg^{\mathbb{T}} f) \cdot \text{empty}^{\mathbb{T}'} &= \text{empty}^{\mathbb{T}} \\ (\gg^{\mathbb{T}} f) \cdot \text{sg}^{\mathbb{T}'} &= f \\ (\gg^{\mathbb{T}} f) \cdot \text{union}^{\mathbb{T}'} &= \text{union}^{\mathbb{T}} \cdot ((\gg^{\mathbb{T}} f) \times (\gg^{\mathbb{T}} f)) , \end{aligned}$$

so that $(\gg^{\mathbb{T}} \text{sg}^{\mathbb{T}})$ constitutes a monad morphism between the collection monads \mathbb{T}' and \mathbb{T} .

In the form proposed by [Wong](#) however, the monad morphism approach did not enable the interpretation of comprehensions in non-collection monads. The resulting language is thus a hybrid assembled from collection comprehensions and separate structural recursion operators to implement aggregation and quantification. The work of [Buneman, Naqvi, Tannen, and Wong](#) has, nonetheless, been a major source of inspiration for the work we are discussing in this text. It is the group of researchers around [Buneman](#) to which many of the foundational ideas of a categorical treatment of query languages have to be attributed to.

74 Monoid Comprehensions. [Fegaras and Maier](#) fabricated an alternative query calculus, the *monoid comprehension calculus*, based on *monoids* and *homomorphisms* between these [38, 40]. The core of the calculus is, once more, formed by \mathbb{G} -algebras, *i. e.*, types in union representation. The principal query constructs are monoid homomorphisms—or rather homomorphic extensions—between instances of these types. In the lingo of initial algebras, given the initial type $(\tau = \text{empty}^{\mathbb{T}} \nabla \text{sg}^{\mathbb{T}} \nabla \text{union}^{\mathbb{T}}, \mathbb{T})$ and the \mathbb{G} -algebra $\alpha = z \nabla s \nabla \oplus$, the *homomorphic extension* of function f from τ to α is

$$\text{hom}^{\tau \rightarrow \alpha} f = \llbracket z \nabla f \nabla \oplus \rrbracket_{\mathbb{G}} .$$

Like \gg , combinator hom defines a well-behaved structural recursion.⁷ Unlike \gg , hom is additionally capable of coercing types (from τ to α in the case above) while \gg is bound to operate in the \mathbb{T} -monad induced by its initial type only:

$$(\gg f) = \langle \langle empty^{\mathbb{T}} \nabla f \nabla union^{\mathbb{T}} \rangle \rangle_{\mathbb{G}} = hom^{\tau \rightarrow \tau} f .$$

This property of hom enabled [Fegaras and Maier](#) to postulate monoid comprehensions over different union representation types, including non-collection instances like, *e. g.*, $true \nabla id \nabla \wedge$. Let τ' denote the monoid associated with the type constructor \mathbb{T}' , then the relevant monoid comprehension translation rule reads:

$$\llbracket e \parallel v \leftarrow e' : \mathbb{T}' E \rrbracket^{\mathbb{T}} \rightsquigarrow hom^{\tau' \rightarrow \tau} (\lambda v. sng^{\mathbb{T}} \cdot e) e' .$$

If we rewrite the above only slightly, we can spot a (more or less rough) similarity to the catamorphism-based coercion of monadic types (*cf.* [M-2'](#)):

$$\begin{aligned} & hom^{\tau' \rightarrow \tau} (\lambda v. sng^{\mathbb{T}} \cdot e) \\ = & \{ \text{unfold homomorphic extension } hom \} \\ & \langle \langle empty^{\mathbb{T}} \nabla (\lambda v. sng^{\mathbb{T}} \cdot e) \nabla union^{\mathbb{T}} \rangle \rangle \\ = & \{ \text{sum, } v \text{ not free in } sng^{\mathbb{T}} \} \\ & \langle \langle (empty^{\mathbb{T}} \nabla sng^{\mathbb{T}} \nabla union^{\mathbb{T}}) \cdot (id + \lambda v. e + id \times id) \rangle \rangle \\ = & \{ \text{union representation functor } \mathbb{G} \} \\ & \langle \langle (empty^{\mathbb{T}} \nabla sng^{\mathbb{T}} \nabla union^{\mathbb{T}}) \cdot \mathbb{G}(\lambda v. e, id) \rangle \rangle \\ = & \{ \text{CATA-MAP-FUSION} \} \\ & \langle \tau \rangle \cdot \mathbb{T}(\lambda v. e) , \end{aligned}$$

which relates monoid comprehensions and our proposal for multi-monad comprehensions.

2.11 Monad Comprehensions in Functional Programming

75 [Moggi](#) has been the first to establish a connection between the categorical monad notion and the semantics of programming languages [86, 87]. Since these days, [Wadler](#) and others proposed monads as a paradigm after which functional programs and libraries should be structured. Monads have become

⁷Given that \oplus is commutative and idempotent whenever $union^{\mathbb{T}}$ is, respectively.

ubiquitous in functional programming [68, 73, 93, 94, 117, 119, 120, 121]. In spite of the terseness of their interface, monads opened new and exciting ways to construct combinator libraries for a broad variety of tasks, including those that, up to then, were stumbling stones in a purely functional environment: side-effect free I/O, stateful computation, exception handling, and the often notationally cumbersome continuation passing style of programming.

Even better: the intimate connection between monads and comprehensions provided a convenient syntax to express the just mentioned tasks in functional programs *for free*. This had an impact on the design of the functional programming language Haskell 1.4 [91] in which the list monad lost its exceptional role: comprehensions could be defined over an arbitrary monad \mathbb{T} (with the generator ranges bound to be of type \mathbb{T} , too).⁸

Let us close this chapter with the review of two monad instances which could especially affect the optimization of database programming and update languages: the *exception* and *state transformer* monads, respectively.

Recall from Example 36 the definition of the type constructor **maybe** which lifts any type E onto a type with a distinguished constant *nothing*:

$$\mathit{nothing} \nabla \mathit{just} : 1 + E \rightarrow \mathbf{maybe}E .$$

An expression evaluated in the monad induced by **maybe** may be interpreted as a computation that may fail and thus raise an exception (which is indicated by returning *nothing* as the result).

We get hold of the **maybe** monad through

$$\begin{aligned} \mathit{unit}^{\mathbf{maybe}} &= \mathit{just} \\ \mathit{join}^{\mathbf{maybe}} &= (\mathit{nothing} \nabla \mathit{id}) \\ \mathit{zero}^{\mathbf{maybe}} &= \mathit{nothing} . \end{aligned}$$

(Remember that the morphism mapping **maybe** f yields $f x$ for an argument *just* x and returns *nothing* otherwise.)

The meaning of a comprehension evaluated in the **maybe** monad is given by the desugaring scheme \mathcal{M} of Definition 69: a generator $v \leftarrow e : \mathbf{maybe}E$ binds variable v to x if $e = \mathit{just} x$ and propagates the binding. Should e raise an exception, the comprehension evaluates to *nothing* as a whole (the exception is “thrown”); a filter acts like an assertion that, should it fail, also raises an exception. Given this, the comprehension

$$\llbracket x/y \rrbracket x \leftarrow f, y \leftarrow g, y \neq 0 \rrbracket^{\mathbf{maybe}}$$

⁸At the time of writing, this design decision has been reversed on a syntactical level: the now current Haskell 98 definition adopts the equivalent **do**-notation to express monadic computations [65].

computes *just* (f/g) provided that both arguments are well-defined and the denominator is not zero. Otherwise the comprehension indicates failure by evaluating to *nothing*.

Stateful computation and purely functional programming can live together if a representation of the program state (of type S , say) is explicitly passed as an extra argument between function calls. State update is implemented by construction of a fresh copy of the state representation. Passing the state explicitly quickly becomes cumbersome and error-prone, especially if the type of state representation S is subject to change during program development. However, an expression that is evaluated in the **state** monad to be defined in a minute, may act *as if* updatable state were available to the program: the state is passed implicitly and its representation becomes opaque (which is what we desire).

The type constructor we need here is $\mathbf{state}E = S \rightarrow E \times S$, which constructs the type of *state transformers*. A value of type $\mathbf{state}E$ denotes a computation that yields a result of type E (which possibly depends on the state it is evaluated in) and additionally “side-effects” the state by constructing a (modified) copy of it. Both E and the new state are returned. Note that **state** makes use of the function type former (\rightarrow) which we are not covering as a datatype constructor in this text; since we are working in the cartesian-closed category **Set** this is no principle obstacle however (*e. g.*, see [10]).

The ingredients of the **state** monad thus are as follows (**state** gets by without *zero*):

$$\begin{aligned}\mathbf{state}E &= S \rightarrow E \times S \\ \mathbf{state}f &= \lambda s.(f \times id) \cdot s\end{aligned}$$

$$\begin{aligned}unit^{\mathbf{state}}x &= (K_x \Delta id) \\ join^{\mathbf{state}}ss &= apply \cdot ss \quad \text{with } apply(f, x) = f x.\end{aligned}$$

The functoriality of **state** directly follows from the properties of arrow composition. Observe that $unit^{\mathbf{state}}x$ denotes a computation that simply returns x without affecting the state. $join^{\mathbf{state}}ss$ applies the state transformer ss yielding a state transformer which is then applied in a subsequent step (*i. e.*, $join^{\mathbf{state}}$ threads the state through a sequence of transformations).

Based on the above, we can readily implement *mutable references* in a pure functional language. [Ohori](#) used a closely related technique to enrich a referentially transparent database programming language with a notion of *object identity* and *mutable objects* without sacrificing the optimizability

of the language [90]. Assuming that a reference is identified by a unique sequence number, it is sufficient to keep the next sequence number to assign and a list that associates those numbers with the referenced content (of type *String* in this example) in the state, thus $S = Num \times list(Num \times String)$. To unclutter the syntax let us additionally define three primitives that operate on a store of type S : *new* allocates a new sequence number and initializes the content of the newly created reference; *asgn* overwrites the content of a given reference; *deref* accesses the contents of a given reference but does not alter the store:

$$\begin{aligned}
 new & : String \rightarrow stateNum \\
 new\ c & = \lambda(n, l).(n, (n + 1, assoc\ l\ (n, c))) \\
 \\
 asgn & : (Num \times String) \rightarrow state1 \\
 asgn\ (n', c) & = \lambda(n, l).((), (n, assoc\ l\ (n', c))) \\
 \\
 deref & : Num \rightarrow stateString \\
 deref\ n' & = \lambda(n, l).(lookup\ l\ n', (n, l))
 \end{aligned}$$

(the functions *assoc* and *lookup* do the obvious on association lists of type $list(Num \times String)$).

Comprehensions evaluated in the **state** monad exhibit a remarkably “imperative” flavor although we are, of course, still operating in a purely functional environment. The comprehension (applied to the initial state $(0, nil)$ in which no references have been allocated yet)

$$\llbracket z \llbracket x \leftarrow new\ "foo", y \leftarrow new\ "bar", z \leftarrow deref\ y, p \leftarrow asgn\ (x, z) \rrbracket^{state} (0, nil) \rrbracket$$

evaluates to $(\text{"bar"}, (2, [(0, \text{"bar"}), (1, \text{"bar"})]))$ and thus describes the same computation as the imperative program

$$x := \text{"foo"}; \quad y := \text{"bar"}; \quad z := y; \quad x := z.$$

It is our conjecture (which we regrettably cannot investigate more deeply this text) that not only queries but database languages in general may benefit from a monadic interpretation. An appropriately adapted multi-monad comprehension desugaring scheme \mathcal{M} could provide the basis for the optimization theory for database programming languages that speak about collections and stateful computation (*i. e.*, updates) at the same time.

But let us not get carried away. We started out to Comprehend Queries.

Chapter 3

Query Compilation and Normalization

76 A comprehension of queries is not complete if it lacks the apparatus to assign meaning to *user-level* query syntax. The first part of this chapter will connect the internal categorical machinery to the “outside” world of declarative query languages. Thanks to the work we have spent with the introduction of monad comprehensions, we will be able to establish this connection by means of a *syntactic mapping*, \mathcal{Q} , to be defined shortly. Once we have established \mathcal{Q} , the composition $\mathcal{M} \cdot \mathcal{Q}$ will assign a categorical semantics to user queries and thus provides us with a formal, albeit preliminary, query comprehension in terms of catamorphic mappings.

Interpreting a query by way of $\mathcal{M} \cdot \mathcal{Q}$ alone is far from feasible if we hope for an efficient query execution. $\mathcal{M} \cdot \mathcal{Q}$ emits nothing but nested spine transformers from which an underlying query engine can hardly derive anything but naive nested-loop iteration strategies. In certain cases, however, nested-loop iteration may actually be an adequate execution strategy and in these cases an execution plan may be directly “*read off*” the catamorphisms emitted by $\mathcal{M} \cdot \mathcal{Q}$. This is the reason why, from time to time, we invest some initial thought in optimizations which are immediately applicable to catamorphic queries. We will have to say more about this.

Nevertheless, it is our hypothesis that monad comprehensions assign meaning to queries in a form from which more efficient execution strategies are derivable. The *normalization* of the monad comprehension expressions constructed by \mathcal{Q} will provide an important rewriting step towards a query form that facilitates exploitation of the various advanced *join* or *grouping* operators which database query engines have to offer. A number of well-established query unnesting techniques turn out to be normalization steps if we express them in the lingo of monad comprehensions. Normalization will

take up the second part of this chapter.

77 Throughout this text, ODMG’s second version of the object query language OQL will be the source language of choice. In [23, Chapter 4], it is stated that OQL shall be understood as a *functional language* and we will treat it as such: OQL queries are built from basic query formers which, given that the language’s typing rules are obeyed, may be freely composed to form more complex query expressions. None of the query formers initiate side-effects.

It is our ambition to understand almost anything of OQL in terms of the categorical framework we have set up so far. We are, however, falling short of emphasizing the purely object-oriented features of the language: object types will be viewed as abstract data types that allow for equality tests (presumably based on object identity but we are not looking into this). Invocation of side-effect free methods will be modeled as the application of a function which—besides type information—remains opaque to us. Lacking the *object* notion, we will also not reason about inheritance and late binding on method calls. These omissions are no principle obstacles. Fegaras and Maier have sketched a translation from OQL to a monoid comprehension calculus (see Paragraph 74) that has been extended later by Fegaras with a notion of object identity and destructive updates expressed within the monoid calculus itself (*cf.* Section 2.11 and the discussion of the *state* monad) [39, 40].

Given these restrictions, OQL queries may be typed using the orthogonal system of categorical datatype constructors we have discussed in Chapter 2. Put differently, we are querying values of arbitrary but finite nesting depth which—as it was argued by Abiteboul and Kanellakis in [1]—form the core of object-oriented database systems. The adoption of OQL as a source language is, of course, not principle to the method. We are positive that other query languages for complex value databases, *e. g.*, the query fragments of SQL-92 and SQL-3, may be comprehended using the techniques we are developing: OQL constitutes a test case that embraces the SQL dialects in terms of expressiveness of its type system as well as orthogonal applicability of its query formers.

Familiarity with OQL syntax as it was proposed in [23] is assumed from now on. The forthcoming will supply the semantics and additionally—as remarked in the last paragraph—initial catamorphic means for the actual execution of OQL queries. In anticipation of the things to come and to provide reference, Riedel and Scholl as well as Cherniack discuss alternative semantic assignments that are reasonably close to ours [30, 97].

3.1 A Catamorphic Interpretation of OQL

This is the right time to explain a few notational conventions and preliminaries, some of which are consequences of the categorical model of data types we are dealing with. OQL keywords will be typeset in `typewriter` style.

78 Types. (a) The default category is **Set**. As before we will use the names E, E', \dots as well as subscripted variants thereof to indicate variables that range over types, *i. e.*, objects in **Set**. If not stated otherwise, \top and variants denote one of the type functors **list**, **bag**, and **set**. The type functors implement the OQL collection type constructors *list*, *bag*, and *set*, respectively.

(b) We translate OQL's basic types as follows: numerical types are represented by type *Num*; character and string literals are of type *String*; *Bool* implements the type of the two boolean values. Object types map into the abstract type *Obj*.

(c) OQL comes equipped with a *record type constructor with tagged fields* and its corresponding value constructor $\mathbf{struct}(l_1:e_1, \dots, l_n:e_n)$ in which the l_i denote unique field tags and the e_i are field entries (arbitrary OQL expressions). If we have $e_i:E_i$, the thus constructed value has the OQL type $\mathbf{struct}(l_1:E_1, \dots, l_n:E_n)$.

Such record values are represented as right-deep nested products of the form $(e_1, (e_2, (\dots (e_n, ()) \dots)))$ whose types are $E_1 \times (E_2 \times (\dots (E_n \times 1) \dots))$. During this translation we maintain a mapping $\mathit{tag} = \{l_i \mapsto i\}, i = 1 \dots n$, that keeps track of tag positions. We may, however, safely drop the rigid parenthesized notation and write $(e_1, \dots, e_n):E_1 \times \dots \times E_n$ instead, *i. e.*, we may act as if we were operating with n -ary product types and values, the latter also called *tuples* from now on. These simplifications are justified by a simple observation about isomorphisms: in a category with products, objects $E_1 \times (E_2 \times E_3)$ and $(E_1 \times E_2) \times E_3$ are isomorphic by means of the arrows

$$\begin{aligned} (\mathit{id} \times \mathit{outl}) \Delta (\mathit{outr} \cdot \mathit{outr}) &: E_1 \times (E_2 \times E_3) \rightarrow (E_1 \times E_2) \times E_3 \\ (\mathit{outl} \cdot \mathit{outl}) \Delta (\mathit{outr} \times \mathit{id}) &: (E_1 \times E_2) \times E_3 \rightarrow E_1 \times (E_2 \times E_3), \end{aligned}$$

which renders \times associative so that the absence or presence of parentheses does not alter the meaning of a nested product. Likewise, $E \times 1$ and E are isomorphic because we have that $\mathit{outl} \cdot (\mathit{id} \Delta ()) = \mathit{id}_E$ and $(\mathit{id} \Delta ()) \cdot \mathit{outl} = \mathit{id}_{E \times 1}$. This justifies the omission of the $()$ in tuple notation.

(d) All types E built from the basic types *Num*, *String*, *Bool*, *Obj*, and categorical datatype constructors are *equality types* in the ML sense [85], *i. e.*, they admit equality tests by means of the overloaded predicate $=: E \times E \rightarrow \mathit{Bool}$. Equality of products is defined component-wise. Equality on $\top E$ is decided by the initial algebra that induces \top (see Section 2.7).

Types *Num* and *String* are taken to be *ordered*. The type $E_1 \times \dots \times E_n$ is ordered if the E_i are ordered (in which case the ordering of the tuple type is defined lexicographically). If E is ordered, we are admitted to apply the overloaded predicate $<: E \times E \rightarrow Bool$ to check if its argument pair is ordered.

(e) To summarize, the following grammar yields all valid query types ty :

ty	\rightarrow	$Num \mid String \mid Bool \mid Obj$	basic types
		\mid	
		$\text{set } ty \mid \text{bag } ty \mid \text{list } ty$	type functors
		\mid	
		$ty \times ty$	products
		\mid	
		(ty)	

79 Typing OQL fragments. (f) Designed as a syntax-directed mapping, \mathcal{Q} descends the parse tree of its argument query, starting with the root which represents the entire query text. \mathcal{Q} operates in a *bottom-up* manner: subquery fragments are translated before these translations are assembled to emit a translation of the enclosing query. During its walk of the parse tree, \mathcal{Q} consequently encounters *open* OQL fragments, *i. e.*, OQL subqueries containing *free query variables* whose binding sites are located further up the parse tree (the root fragment is expected to be closed).

In the OQL query below, variable y appears free in three subqueries: y (in the **select** clause), **forall** x in $xs: x = y$, and $x = y$. Variable x appears free in $x = y$. The binding sites of y and x are the **from** clause and the **forall** quantifier, respectively:

```

select y
  from ys as y
  where forall x in xs: x = y .

```

To cope with the typing of OQL fragments in order to clearly define query variable scoping, we will interpret open OQL fragments as *functions of their free variables*. To illustrate, the above query will be understood as

```

select ( $\lambda v.v$ ) y
  from ys as y
  where ( $\lambda v.\text{forall } x \text{ in } xs: (\lambda v'.v' = v) x$ ) y ,

```

so that, for example, the **where** clause of an OQL **select-from-where** block is a function of type $E \rightarrow Bool$ (where type E is determined by the query's **from** clause).

(g) Query types will be given and derived using standard notation [22]: the *type judgment* $\Gamma \vdash e: E$ asserts the type of term e to be E , given that

the free variables in e are typed according to the *type environment* Γ (“ e has type E in Γ ”). A type environment is of the form $\emptyset, v_1 : E_1, \dots, v_n : E_n$ (with $n \geq 0$, \emptyset denotes the empty environment).

A *typing rule* derives the type judgment for a term (below the horizontal rule) under the premise of a number (maybe zero) of conditions—typically other type judgments—listed above the rule. The typing rule for OQL’s **forall** quantifier reads (note that v may appear free in p and the predicate is thus typed $\lambda v.p : E \rightarrow Bool$):

$$\frac{\Gamma \vdash e : TE \quad \Gamma, v : E \vdash p : Bool}{\Gamma \vdash \text{forall } v \text{ in } e : p : Bool} \text{ (TY-FORALL)}$$

where $\Gamma, v : E$ denotes Γ after the addition of variable-type assignment $v : E$. To render environment updates more compact, let $\Gamma, \bar{v}_n : \bar{E}_n$ denote the environment $\Gamma, v_1 : E_1, v_2 : E_2, \dots, v_n : E_n$ in which n will be determined by the context. For $n = 0$, this denotes Γ itself.

80 Miscellaneous. (h) The term $e[x/y]$ denotes term e with all free occurrences of variable y in e replaced by term x . Let $e[x/\bar{y}_n]$ be a shorthand for $e[x/y_1][x/y_2] \cdots [x/y_n]$ if the y_i are in context. Likewise, we write $e[\bar{x}_n/\bar{y}_n]$ to abbreviate the replacement $e[x_1/y_1][x_2/y_2] \cdots [x_n/y_n]$.

(i) During the translation, we will identify basic OQL constants with their obvious meaning in the categorical world, *e.g.* we do not distinguish $+$ and $+$ = $\overset{\text{sum}}{!}$ or **true** and $true = []^{\text{all}}$, respectively.

81 OQL translation. We will define \mathcal{Q} by case selection on the various query formers provided by OQL. Much of the simplicity of \mathcal{Q} is due to its *uniform* translation strategy [111]: a query e may be compiled independently from subqueries occurring in it. During the translation of e , the subquery fragments are treated as free variables that may be instantiated later to complete the translation. In the variant described in this text, \mathcal{Q} calls itself recursively to perform these instantiations.

Nevertheless, the strategy may also be used in system contexts that call for *separate compilation* of query fragments. Such situations typically arise if the query language is embedded into a host programming language: parts of the query text are dependent on host language variables or are even computed by host language routines and are not known beforehand. Despite the presence of such dynamically generated query parts, we could translate the statically known query fragments (under the premise that the dynamic parts themselves will be correctly typed) at compile-time of the host program.

This completes the preliminaries, so let us now turn to the actual definition of the OQL translation function \mathcal{Q} .

82 Definition. At several places we provide notes concerning simplifications of the translation. Skipping these should not affect the understanding of \mathcal{Q} .

Query variables. We use $\Gamma, v : E, \Gamma'$ to denote an environment that contains the variable-type-assignment $v : E$.

$$\frac{}{\Gamma, v : E, \Gamma' \vdash v : E} \text{(TY-VAR)}$$

$$\mathcal{Q}v = v \qquad \mathcal{Q}\text{-VAR}$$

Constants of atomic type. Let *Basic* denote one of *Num*, *String*, *Bool*, or *Obj*.

$$\frac{c \in \text{Basic}}{\Gamma \vdash c : \text{Basic}} \text{(TY-BASIC)}$$

$$\mathcal{Q}c = c \qquad \mathcal{Q}\text{-BASIC}$$

Record construction (struct). As already noted in Paragraph 78, we maintain a *tag* = $\{l_i \mapsto i\}$ mapping as a side-effect of the **struct** constructor translation. For mere simplicity we assume that record tags are database-wide uniquely determined.

$$\frac{\Gamma \vdash e_1 : E_1 \quad \cdots \quad \Gamma \vdash e_n : E_n}{\Gamma \vdash \text{struct}(l_1 : e_1, \dots, l_n : e_n) : E_1 \times \cdots \times E_n} \text{(TY-STRUCT)}$$

$$\mathcal{Q}(\text{struct}(l_1 : e_1, \dots, l_n : e_n)) = (\mathcal{Q}e_1, \dots, \mathcal{Q}e_n) \qquad \mathcal{Q}\text{-STRUCT}$$

Record access. For a function f , define its iterated application by $f^0 = id$ and $f^n = f \cdot f^{(n-1)}$.

$$\frac{\Gamma \vdash e : E_1 \times \dots \times E_n \quad \text{tag } l = i \quad (1 \leq i \leq n)}{\Gamma \vdash e.l : E_i} \text{ (TY-DOT)}$$

$$\mathcal{Q}(e.l) = (\text{outl} \cdot \text{outr}^{(i-1)})(\mathcal{Q}e) \quad \mathcal{Q}\text{-DOT}$$

Collection-typed constants (`set`, `bag`, `list`). Let $coll$ be one of `set`, `bag`, or `list` and define $\mathsf{T}^{\text{set}} = \text{set}$, $\mathsf{T}^{\text{bag}} = \text{bag}$, and $\mathsf{T}^{\text{list}} = \text{list}$.

$$\frac{\Gamma \vdash e_1 : E \quad \dots \quad \Gamma \vdash e_n : E}{\Gamma \vdash \text{coll}(e_1, \dots, e_n) : \mathsf{T}^{\text{coll}} E} \text{ (TY-COLL)}$$

$$\begin{aligned} & \mathcal{Q}(\text{coll}(e_1, \dots, e_n)) \\ = & \text{cons}^{\mathsf{T}^{\text{coll}}}(\mathcal{Q}e_1, \text{cons}^{\mathsf{T}^{\text{coll}}}(\dots \text{cons}^{\mathsf{T}^{\text{coll}}}(\mathcal{Q}e_n, \text{nil}^{\mathsf{T}^{\text{coll}}}) \dots)) \end{aligned} \quad \mathcal{Q}\text{-COLL}$$

Arithmetic operators. Let symbol \ominus denote `-` (unary minus) or `abs`, and let \oplus represent one of `+`, `-` (subtraction), `*`, `/`, or `mod`.

$$\frac{\Gamma \vdash e : \text{Num}}{\Gamma \vdash \ominus e : \text{Num}} \text{ (TY-ARITH-1)}$$

$$\frac{\Gamma \vdash e : \text{Num} \quad \Gamma \vdash e' : \text{Num}}{\Gamma \vdash e \oplus e' : \text{Num}} \text{ (TY-ARITH-2)}$$

$$\mathcal{Q}(\ominus e) = \ominus(\mathcal{Q}e) \quad \mathcal{Q}\text{-ARITH-1}$$

$$\mathcal{Q}(e \oplus e') = (\mathcal{Q}e) \oplus (\mathcal{Q}e') \quad \mathcal{Q}\text{-ARITH-2}$$

Equality and comparison operators. Let \otimes stand for one of the infix operators `<`, `<=`, `=`, `!=`, `>=`, or `>`. We then have that $\otimes : E \times E \rightarrow \text{Bool}$ (E must be ordered if \otimes represents one of the relational comparison operators).

$$\frac{\Gamma \vdash e : E \quad \Gamma \vdash e' : E}{\Gamma \vdash e \otimes e' : \text{Bool}} \text{ (TY-COMP)}$$

$$\mathcal{Q}(e \otimes e') = (\mathcal{Q}e) \otimes (\mathcal{Q}e') \quad \mathcal{Q}\text{-COMP}$$

Boolean connectives. We use symbol \otimes to denote either and or or.

$$\frac{\Gamma \vdash e : Bool}{\Gamma \vdash \text{not } e : Bool} \text{ (TY-BOOL-1)}$$

$$\frac{\Gamma \vdash e : Bool \quad \Gamma \vdash e' : Bool}{\Gamma \vdash e \otimes e' : Bool} \text{ (TY-BOOL-2)}$$

$$\mathcal{Q}(\text{not } e) = \neg(\mathcal{Q}e) \quad \mathcal{Q}\text{-BOOL-1}$$

$$\mathcal{Q}(e \otimes e') = (\mathcal{Q}e) \otimes (\mathcal{Q}e') \quad \mathcal{Q}\text{-BOOL-2}$$

Selecting into a set (*select-distinct-from-where*). The types of e_2, \dots, e_n express the possibility to specify *dependent joins* in the **from** clause, *i. e.*, x_1, \dots, x_{i-1} may occur free in e_i (\mathcal{Q} accounts for this by placing the e_i as generator ranges in order of their appearance in the **from** clause). We adopt a widespread convention and assume $\mathcal{Q}p = \text{true}$ should the query expression lack the **where** clause.

$$\frac{\begin{array}{l} \Gamma, \bar{x}_n : \bar{E}_n \vdash e : E \\ \Gamma, \bar{x}_{i-1} : \bar{E}_{i-1} \vdash e_i : \mathbb{T}_i E_i \quad (1 \leq i \leq n) \\ \Gamma, \bar{x}_n : \bar{E}_n \vdash p : Bool \end{array}}{\Gamma \vdash \left(\begin{array}{l} \text{select distinct } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \end{array} \right) : \text{set} E} \text{ (TY-SET-SFW)}$$

$$= \llbracket \mathcal{Q}e \parallel x_1 \leftarrow \mathcal{Q}e_1, \dots, x_n \leftarrow \mathcal{Q}e_n, \mathcal{Q}p \rrbracket^{\text{set}} \quad \mathcal{Q}\text{-SET-SFW}$$

Selecting into a bag (*select-from-where*). Note that, to render the translation well-defined, we require $\mathbb{T}_i \neq \text{set}$ as the **bag** algebra obeys equation **EQ-COMM** only (*cf.* Paragraph 70).

$$\frac{\begin{array}{l} \Gamma, \bar{x}_n : \bar{E}_n \vdash e : E \\ \Gamma, \bar{x}_{i-1} : \bar{E}_{i-1} \vdash e_i : \mathbb{T}_i E_i \quad \mathbb{T}_i \neq \text{set} \quad (1 \leq i \leq n) \\ \Gamma, \bar{x}_n : \bar{E}_n \vdash p : Bool \end{array}}{\Gamma \vdash \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \end{array} \right) : \text{bag} E} \text{ (TY-SFW)}$$

$$\begin{aligned}
& \mathcal{Q} \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \end{array} \right) \\
&= \llbracket \mathcal{Q} e \rrbracket x_1 \leftarrow \mathcal{Q} e_1, \dots, x_n \leftarrow \mathcal{Q} e_n, \mathcal{Q} p \rrbracket^{\text{bag}} \quad \mathcal{Q}\text{-SFW}
\end{aligned}$$

Grouping (select-from-where-group by). OQL's grouping construct partitions the n -fold product of the e_1, \dots, e_n with respect to the grouping criteria g_1, \dots, g_m (accessible in the **select** clause e under names l_1, \dots, l_m): n -tuples that agree on all criteria are collected into a bag (accessible in e by referring to name *partition*). References to these names are replaced with corresponding tuple accesses by function *tuplify*. During the recursive descent of \mathcal{Q} , the thus introduced tuple accesses are subsequently translated into sequences of projections by **\mathcal{Q} -Dot**.

$$\begin{array}{l}
\Gamma, \bar{l}_g : \bar{E}'_g, \text{partition} : \text{bag}(E_1 \times \dots \times E_n) \vdash e : E \\
\Gamma, \bar{x}_{i-1} : \bar{E}_{i-1} \vdash e_i : \top_i E_i \quad \top_i \neq \text{set} \quad (1 \leq i \leq n) \\
\Gamma, \bar{x}_n : \bar{E}_n \vdash p : \text{Bool} \quad \Gamma, \bar{x}_n : \bar{E}_n \vdash g_j : E'_j \quad (1 \leq j \leq m)
\end{array} \quad (\text{TY-GROUP})$$

$$\Gamma \vdash \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{group by } l_1 : g_1, \dots, l_m : g_m \end{array} \right) : \text{bag} E$$

$$\mathcal{Q} \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{group by } l_1 : g_1, \dots, l_m : g_m \end{array} \right)$$

$$= \llbracket \mathcal{Q} (\text{tuplify } e) \rrbracket$$

$$\begin{array}{l}
y \leftarrow \llbracket (\mathcal{Q} g_1, \dots, \mathcal{Q} g_m, \llbracket (x'_1, \dots, x'_n) \rrbracket x'_1 \leftarrow \mathcal{Q} e_1, \dots, x'_n \leftarrow \mathcal{Q} e_n, \\
\mathcal{Q} g_1 = \mathcal{Q} (g_1[\bar{x}'_n/\bar{x}_n]), \dots, \\
\mathcal{Q} g_m = \mathcal{Q} (g_m[\bar{x}'_n/\bar{x}_n]) \rrbracket^{\text{bag}} \rrbracket \\
x_1 \leftarrow \mathcal{Q} e_1, \dots, x_n \leftarrow \mathcal{Q} e_n, \mathcal{Q} p \rrbracket^{\text{bag}} \rrbracket^{\text{bag}}
\end{array} \quad \mathcal{Q}\text{-GROUP}$$

where

$$\text{tuplify } t = t[y.1/l_1] \cdots [y.m/l_m][y.(m+1)/\text{partition}].$$

Constrained Grouping (select-from-where-group by-having). We have to call *tuplify* on both the **select** clause e and p' —a predicate to be evaluated against all partitions—to replace references by name into tuple projections.

$$\begin{array}{c}
\Gamma, \bar{l}_m : \bar{E}'_m, \text{partition} : \text{bag}(E_1 \times \cdots \times E_n) \vdash e : E \\
\Gamma, \bar{x}_{i-1} : \bar{E}_{i-1} \vdash e_i : \top_i E_i \quad \top_i \neq \text{set} \quad (1 \leq i \leq n) \\
\Gamma, \bar{x}_n : \bar{E}_n \vdash p : \text{Bool} \quad \Gamma, \bar{x}_n : \bar{E}_n \vdash g_j : E'_j \quad (1 \leq j \leq m) \\
\Gamma, \bar{l}_m : \bar{E}'_m, \text{partition} : \text{bag}(E_1 \times \cdots \times E_n) \vdash p' : \text{Bool} \\
\hline
\Gamma \vdash \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{group by } l_1 : g_1, \dots, l_m : g_m \\ \text{having } p' \end{array} \right) : \text{bag} E \\
\hline
\mathcal{Q} \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{group by } l_1 : g_1, \dots, l_m : g_m \\ \text{having } p' \end{array} \right) \\
= \llbracket \mathcal{Q}(\text{tuplify } e) \rrbracket \\
\quad y \leftarrow \llbracket (\mathcal{Q} g_1, \dots, \mathcal{Q} g_m, \llbracket (x'_1, \dots, x'_n) \rrbracket \left[\begin{array}{l} x'_1 \leftarrow \mathcal{Q} e_1, \dots, x'_n \leftarrow \mathcal{Q} e_n, \\ \mathcal{Q} g_1 = \mathcal{Q}(g_1[x'_n/\bar{x}_n]), \dots, \\ \mathcal{Q} g_m = \mathcal{Q}(g_m[x'_n/\bar{x}_n]) \end{array} \right] \text{bag}) \rrbracket \\
\quad x_1 \leftarrow \mathcal{Q} e_1, \dots, x_n \leftarrow \mathcal{Q} e_n, \mathcal{Q} p \rrbracket^{\text{bag}}, \\
\mathcal{Q}(\text{tuplify } p') \rrbracket^{\text{bag}} \\
\hline
\mathcal{Q}\text{-HAVING}
\end{array}$$

Sorting (select-from-where-order by). In the query below, let *ord* be a placeholder for **asc** or **desc**—a hint to sort the query result in ascending or descending order, respectively. The type E' of the sort criterion is required to be ordered.

$$\begin{array}{c}
\Gamma, \bar{x}_n : \bar{E}_n \vdash e : E \\
\Gamma, \bar{x}_{i-1} : \bar{E}_{i-1} \vdash e_i : \top_i E_i \quad \top_i \neq \text{set} \quad (1 \leq i \leq n) \\
\Gamma, \bar{x}_n : \bar{E}_n \vdash p : \text{Bool} \quad \Gamma, \bar{x}_n : \bar{E}_n \vdash o : E' \\
\hline
\Gamma \vdash \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{order by } o \text{ ord} \end{array} \right) : \text{list} E \\
\hline
\text{(TY-SORT)}
\end{array}$$

$$\begin{aligned}
& \mathcal{Q} \left(\begin{array}{l} \text{select } e \\ \text{from } e_1 \text{ as } x_1, \dots, e_n \text{ as } x_n \\ \text{where } p \\ \text{order by } o \text{ ord} \end{array} \right) \\
= & \llbracket \mathcal{Q}(e[y/\bar{x}_n]) \rrbracket y \leftarrow (\text{nil}^{\text{list}} \nabla \text{ins}^{\text{list}} \otimes) \llbracket \mathcal{Q} o \rrbracket x_1 \leftarrow \mathcal{Q} e_1, \dots, \\
& \qquad \qquad \qquad x_n \leftarrow \mathcal{Q} e_n, \mathcal{Q} p \rrbracket^{\text{bag}} \rrbracket^{\text{list}} \\
& \qquad \mathcal{Q}\text{-SORT}
\end{aligned}$$

where

$$\otimes = \begin{cases} < : E' \times E' \rightarrow \text{Bool} & \text{if } \text{ord} = \text{asc}, \\ > : E' \times E' \rightarrow \text{Bool} & \text{else.} \end{cases}$$

Notes. Coherent with our quest for a catamorphic model for OQL, $(\text{nil}^{\text{list}} \nabla \text{ins}^{\text{list}} \otimes)$ provides a purely (see below) catamorphic representation of *insertion sort*: F-algebra $\text{nil}^{\text{list}} \nabla (\text{ins}^{\text{list}} \otimes) : \mathbf{F}(E', \text{list} E') \rightarrow \text{list} E'$ is the algebra of sorted lists in insert representation. Parameterized with \otimes , operator ins^{list} inserts an element into an already sorted list, respecting order \otimes . Note that ins^{list} itself may be understood as a catamorphism:

$$\text{ins}^{\text{list}} \otimes (x, xs) = (\llbracket x \rrbracket^{\text{list}} \nabla (\text{swap} \cdot (\otimes \cdot (\text{id} \times \text{hd})))?) xs$$

where

$$\begin{aligned}
\text{hd} \cdot \ddagger &= \text{outl} \\
\text{tl} \cdot \ddagger &= \text{outr} \\
\text{swap} &= \ddagger \nabla (\ddagger \cdot ((\text{hd} \cdot \text{outr}) \Delta (\ddagger \cdot (\text{id} \times \text{tl}))))
\end{aligned}$$

(defined like this, $\text{ins}^{\text{list}} \otimes$ is left-commutative but not left-idempotent).

Catamorphic sorting does not stress the efficient implementation of sorting but rather its algebraic properties; it is these properties that we would like \mathcal{Q} to make explicit (however, see [49] where merge sort is derived from a catamorphic insertion sort by program transformations alone). Applied to algebra $\text{nil}^{\text{list}} \nabla \text{ins}^{\text{list}} \otimes$, Lemma 60 yields a monad of sorted lists and thus its corresponding monad comprehension notion, but we are not following this route here.

Intersection (`intersect`).

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma \vdash e' : \mathbb{T}'E}{\Gamma \vdash e \text{ intersect } e' : \mathbb{T}E} \text{ (TY-INTERSECT)}$$

$$\begin{aligned} & \mathcal{Q}(e \text{ intersect } e') \\ = & \llbracket x \mid x \leftarrow \mathcal{Q}e, \llbracket x = x' \mid x' \leftarrow \mathcal{Q}e' \rrbracket^{\text{exists}} \rrbracket^{\mathbb{T}} \quad \mathcal{Q}\text{-INTERSECT} \end{aligned}$$

Difference (`except`). Note that this translation of `except` does not subtract multiplicities of elements (but rather removes all elements from the result that occur in both e and e') if $\mathbb{T} = \mathbb{T}' = \text{bag}$.

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma \vdash e' : \mathbb{T}'E}{\Gamma \vdash e \text{ except } e' : \mathbb{T}E} \text{ (TY-EXCEPT)}$$

$$\begin{aligned} & \mathcal{Q}(e \text{ except } e') \\ = & \llbracket x \mid x \leftarrow \mathcal{Q}e, \llbracket x \neq x' \mid x' \leftarrow \mathcal{Q}e' \rrbracket^{\text{all}} \rrbracket^{\mathbb{T}} \quad \mathcal{Q}\text{-EXCEPT} \end{aligned}$$

Union and concatenation (`union` and `+`). OQL designates symbol `+` to denote list concatenation if $\mathbb{T} = \mathbb{T}' = \text{list}$. If (τ, \mathbb{T}) and (τ', \mathbb{T}') denote the initial types associated with type functors \mathbb{T} and \mathbb{T}' , define $\mathbb{T}'' = \mathbb{T}$ if τ obeys all equations that τ' obeys, otherwise set $\mathbb{T}'' = \mathbb{T}'$.

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma \vdash e' : \mathbb{T}'E}{\Gamma \vdash e \text{ union } e' : \mathbb{T}''E} \text{ (TY-UNION)}$$

$$\begin{aligned} & \mathcal{Q}(e \text{ union } e') \\ = & \text{join}^{\mathbb{T}''} (\llbracket []^{\mathbb{T}''} \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e)^{\mathbb{T}''} \llbracket []^{\mathbb{T}''} \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e')^{\mathbb{T}''} \llbracket []^{\mathbb{T}''} \rrbracket) \quad \mathcal{Q}\text{-UNION} \end{aligned}$$

Notes. By inlining the definition of $\text{join}^{\mathbb{T}''}$ and with the help of law **ACID-RAIN** (Chapter 6), this initial translation may be simplified to read

$$\llbracket (\llbracket []^{\mathbb{T}''} \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e') \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e) \rrbracket,$$

which, for the case $\mathbb{T}'' = \mathbb{T} = \mathbb{T}'$, further reduces to

$$\begin{aligned} & \llbracket (\llbracket []^{\mathbb{T}''} \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e') \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e) \\ = & \{ \text{CATA-REFLECT} \} \\ & \llbracket \mathcal{Q}e' \nabla \mathbb{T}'' \rrbracket (\mathcal{Q}e) \\ = & \{ \text{definition of } \# \text{ for } \mathbb{T}'' \text{ (Lemma 60)} \} \\ & (\mathcal{Q}e) \# (\mathcal{Q}e'). \end{aligned}$$

Universal quantification (forall-in).

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma, v : E \vdash p : Bool}{\Gamma \vdash \text{forall } v \text{ in } e : p : Bool} \text{ (TY-FORALL)}$$

$$\mathcal{Q}(\text{forall } v \text{ in } e : p) = \llbracket \mathcal{Q}p \parallel v \leftarrow \mathcal{Q}e \rrbracket^{\text{all}} \quad \mathcal{Q}\text{-FORALL}$$

Existential quantification (exists-in).

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma, v : E \vdash p : Bool}{\Gamma \vdash \text{exists } v \text{ in } e : p : Bool} \text{ (TY-EXISTS)}$$

$$\mathcal{Q}(\text{exists } v \text{ in } e : p) = \llbracket \mathcal{Q}p \parallel v \leftarrow \mathcal{Q}e \rrbracket^{\text{exists}} \quad \mathcal{Q}\text{-EXISTS}$$

OQL defines alternative syntactic forms for the quantifiers which we rewrite so that they become subject to $\mathcal{Q}\text{-FORALL}$ or $\mathcal{Q}\text{-EXISTS}$. As before, \odot serves as a placeholder for the operators $<$, $<=$, $=$, $!=$, $>=$, or $>$.

$$\frac{\Gamma \vdash e : E \quad \Gamma \vdash e' : \mathbb{T}E}{\Gamma \vdash e \odot \text{all } e' : Bool} \quad \frac{\Gamma \vdash e : E \quad \Gamma \vdash e' : \mathbb{T}E}{\Gamma \vdash e \odot \text{some } e' : Bool}$$

$$\mathcal{Q}(e \odot \text{all } e') = \mathcal{Q}(\text{forall } v \text{ in } e' : e \odot v)$$

$$\mathcal{Q}(e \odot \text{some } e') = \mathcal{Q}(\text{exists } v \text{ in } e' : e \odot v)$$

$$\frac{\Gamma \vdash e : E \quad \Gamma \vdash e' : \mathbb{T}E}{\Gamma \vdash e \text{ in } e' : Bool}$$

$$\mathcal{Q}(e \text{ in } e') = \mathcal{Q}(\text{exists } v \text{ in } e' : e = v)$$

Test for non-empty collections (exists).

$$\frac{\Gamma \vdash e : \mathbb{T}E}{\Gamma \vdash \text{exists}(e) : \text{Bool}} \text{ (TY-PEEK)}$$

$$\mathcal{Q}(\text{exists}(e)) = \llbracket \text{true} \parallel x \leftarrow \mathcal{Q}e \rrbracket^{\text{exists}} \quad \mathcal{Q}\text{-PEEK}$$

Notes. This translation may suggest to scan collection e completely just to determine whether e contains at least one arbitrary element. There are two points to note here. First, remember that \mathcal{Q} 's purpose is to assign meaning to queries, not execution plans. Second, depending on the execution model implemented by the query engine, $\mathcal{Q}\text{-PEEK}$ actually is as efficient a translation as one may hope for. To see this, consider the calculation

$$\begin{aligned} & (\mathcal{M} \cdot \mathcal{Q})(\text{exists}(e)) \\ = & \quad \{ \mathcal{Q}\text{-PEEK} \} \\ & \mathcal{M} \llbracket \text{true} \parallel x \leftarrow \mathcal{Q}e \rrbracket^{\text{exists}} \\ = & \quad \{ \mathcal{M}\text{-3} \} \\ & (\llbracket \text{false} \nabla \vee \rrbracket) (\mathcal{M} \llbracket \text{true} \parallel x \leftarrow \mathcal{Q}e \rrbracket^{\top}) \\ = & \quad \{ \mathcal{M}\text{-2} \} \\ & (\llbracket \text{false} \nabla \vee \rrbracket \cdot (\top K_{\text{true}})) (\mathcal{Q}e) \\ = & \quad \{ \text{CATA-MAP-FUSION} \} \\ & (\llbracket \text{false} \nabla \vee \rrbracket \cdot (\text{id} + K_{\text{true}} \times \text{id})) (\mathcal{Q}e) \\ = & \quad \{ \text{product, true} \vee x = \text{true} \} \\ & (\llbracket \text{false} \nabla K_{\text{true}} \rrbracket) (\mathcal{Q}e). \end{aligned}$$

A database engine that implements query operators following a *lazy* (or *on-demand*) stream iterator model [18, 54] will demand evaluation of $\mathcal{Q}e$ not beyond the point at which it gets hold of the top-most constructor ($\llbracket \cdot \rrbracket^{\top}$ or \dagger) of the spine of $\mathcal{Q}e$. In terms of functional programming languages, it is sufficient to evaluate $\mathcal{Q}e$ to *weak head normal form* [18, 92]. We will have to say a lot more about streaming in Chapter 6.

Test for singleton collections (unique).

$$\frac{\Gamma \vdash e : \mathbb{T}E}{\Gamma \vdash \text{unique}(e) : \text{Bool}} \text{ (TY-UNIQUE)}$$

$$\mathcal{Q}(\text{unique}(e)) = (= \cdot (K_{\llbracket \cdot \rrbracket^{\top}} \Delta \text{tl})) (\mathcal{Q}e) \quad \mathcal{Q}\text{-UNIQUE}$$

(tl defined as before).

Subset relationship (\leq). The monad comprehension below implements the semantics of the subset relationship predicate \subseteq . The typing allows for the application to arbitrary collection type formers but multiplicity or order of elements is not accounted for.

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \Gamma \vdash e' : \mathbb{T}'E}{\Gamma \vdash e \leq e' : \mathit{Bool}} \text{ (TY-SUBSET)}$$

$$\mathcal{Q}(e \leq e') = \llbracket \llbracket x = x' \mid x' \leftarrow \mathcal{Q} e' \rrbracket^{\text{exists}} \mid x \leftarrow \mathcal{Q} e \rrbracket^{\text{all}} \quad \mathcal{Q}\text{-SUBSET}$$

Summation (sum).

$$\frac{\Gamma \vdash e : \mathbb{T}Num \quad \mathbb{T} \neq \text{set}}{\Gamma \vdash \text{sum}(e) : Num} \text{ (TY-SUM)}$$

$$\mathcal{Q}(\text{sum}(e)) = \llbracket x \mid x \leftarrow \mathcal{Q} e \rrbracket^{\text{sum}} \quad \mathcal{Q}\text{-SUM}$$

Counting (count).

$$\frac{\Gamma \vdash e : \mathbb{T}E \quad \mathbb{T} \neq \text{set}}{\Gamma \vdash \text{count}(e) : Num} \text{ (TY-COUNT)}$$

$$\mathcal{Q}(\text{count}(e)) = \llbracket K_1 x \mid x \leftarrow \mathcal{Q} e \rrbracket^{\text{sum}} \quad \mathcal{Q}\text{-COUNT}$$

Average (avg).

$$\frac{\Gamma \vdash e : \mathbb{T}Num \quad \mathbb{T} \neq \text{set}}{\Gamma \vdash \text{avg}(e) : Num} \text{ (TY-AVG)}$$

$$\begin{aligned} & \mathcal{Q}(\text{avg}(e)) \\ = & \left(/ \cdot ((\lambda v. \mathcal{Q}(\text{sum}(v))) \Delta (\lambda v. \mathcal{Q}(\text{count}(v)))) \right) (\mathcal{Q} e) \quad \mathcal{Q}\text{-AVG} \end{aligned}$$

Notes. After recursive invocation of \mathcal{Q} and subsequent monad comprehension desugaring, the translated query essentially exhibits the typical shape of a *tupling catamorphism* $(\llbracket f \rrbracket) \Delta (\llbracket g \rrbracket)$. Operationally speaking, tupling catamorphisms walk their argument twice but the equational theory for catamorphisms once more provides a hook for optimization, this time in terms of law

BANANA-SPLIT: for a tupling catamorphism $\langle f \rangle_{\mathbb{F}} \Delta \langle g \rangle_{\mathbb{F}}$ whose components are \mathbb{F} -catamorphisms, it is true that

$$\langle f \rangle_{\mathbb{F}} \Delta \langle g \rangle_{\mathbb{F}} = \langle (f \cdot \mathbb{F} \text{outl}) \Delta (g \cdot \mathbb{F} \text{outr}) \rangle_{\mathbb{F}} \quad \text{BANANA-SPLIT}$$

(a proof may be found in several places, including [67] and [42]).

For an **avg** query, **BANANA-SPLIT** yields the catamorphism

$$\langle (K_0 \Delta K_0) \nabla ((+ \cdot \text{id} \times \text{outl}) \Delta (+ \cdot (K_1 \times \text{outr})) \rangle_{\mathbb{F}},$$

which, using a single traversal of $\mathcal{Q}e$, computes a (nominator, denominator) pair from which the average is immediate.

Maximum and minimum (**max** and **min**). Types instantiating E are required to be ordered.

$$\frac{\Gamma \vdash e : \mathbb{T}E}{\Gamma \vdash \text{max}(e) : E} \quad (\text{TY-MAX}) \qquad \frac{\Gamma \vdash e : \mathbb{T}E}{\Gamma \vdash \text{min}(e) : E} \quad (\text{TY-MIN})$$

$$\mathcal{Q}(\text{max}(e)) = \llbracket x \mid x \leftarrow \mathcal{Q}e \rrbracket^{\text{max}} \quad \mathcal{Q}\text{-MAX}$$

$$\mathcal{Q}(\text{min}(e)) = \llbracket x \mid x \leftarrow \mathcal{Q}e \rrbracket^{\text{min}} \quad \mathcal{Q}\text{-MIN}$$

Flattening (**flatten**). If (τ, \mathbb{T}) and (τ', \mathbb{T}') denote the initial types associated with type functors \mathbb{T} and \mathbb{T}' , define $\mathbb{T}'' = \mathbb{T}$ if τ obeys all equations that τ' obeys, otherwise set $\mathbb{T}'' = \mathbb{T}'$.

$$\frac{\Gamma \vdash e : \mathbb{T}\mathbb{T}'E}{\Gamma \vdash \text{flatten}(e) : \mathbb{T}''E} \quad (\text{TY-FLATTEN})$$

$$\mathcal{Q}(\text{flatten}(e)) = \llbracket x' \mid x \leftarrow \mathcal{Q}e, x' \leftarrow x \rrbracket^{\mathbb{T}''} \quad \mathcal{Q}\text{-FLATTEN}$$

Note. An appeal to \mathcal{M} and the functoriality of \mathbb{T}'' reveals OQL's **flatten** operator to simply be the monadic $\text{join}^{\mathbb{T}''}$ if $\mathbb{T}'' = \mathbb{T} = \mathbb{T}'$.

Duplicate elimination (**listtaset** and **distinct**).

$$\frac{\Gamma \vdash e : \text{list}E}{\Gamma \vdash \text{listtaset}(e) : \text{set}E} \quad (\text{TY-DUPELIM-1})$$

$$\frac{\Gamma \vdash e : \text{bag}E}{\Gamma \vdash \text{distinct}(e) : \text{set}E} \quad (\text{TY-DUPELIM-2})$$

$$\mathcal{Q}(\text{listtaset}(e)) = \llbracket x \mid x \leftarrow \mathcal{Q}e \rrbracket^{\text{set}} \quad \mathcal{Q}\text{-DUPELIM-1}$$

$$\mathcal{Q}(\text{distinct}(e)) = \llbracket x \mid x \leftarrow \mathcal{Q}e \rrbracket^{\text{set}} \quad \mathcal{Q}\text{-DUPELIM-2}$$

Extraction (`element` and `first`). OQL designates the keyword `first` for the case $\top = \text{list}$.

$$\frac{\Gamma \vdash e : \top E}{\Gamma \vdash \text{element}(e) : E} \text{ (TY-ELEMENT)}$$

$$\mathcal{Q}(\text{element}(e)) = (\perp \nabla \text{outl})(\mathcal{Q}e) \quad \mathcal{Q}\text{-ELEMENT}$$

Notes. The catamorphism yields \perp should e evaluate to $[\]^\top$. Here, the *bottom* symbol \perp (denoting a completely undefined value) signals a situation in which the query engine should raise an exception according to the specification in [23]. For the above to work we require *outl* to be *non-strict* in the second component of its argument, *i. e.*, $\text{outl}(x, \perp) = x$. \parallel

83 Arrays. As presented in Definition 82, mapping \mathcal{Q} lacks support for OQL’s `array` type former and its associated query primitives, most notably array subscripting via $[\cdot]$. Support for arrays in a setting of categorical datatypes and catamorphic computation raises no principal problems—indeed, a number of efforts in that direction have already been undertaken [5, 16, 40, 77].

An insert representation of arrays is conceivable but leads to a rather awkward handling of arrays as either (a) each insertion affects the array *shape* (dimensionality and extension of each dimension) or (b) the array needs to hold *empty slots* on creation into which array elements are to be inserted later. Arrays are more naturally constructed in a bulk-oriented fashion. In [5], arrays with element type E are constructed by the initial algebra

$$\sigma \nabla \rho : \text{list}E + \text{Num} \times \text{array}E \rightarrow \text{array}E$$

(where `array` denotes the array type functor). An array is constructed from a potentially infinite list of values of type E . Rather than on the array itself, the constructors operate on the shape of the array and thus describe how the array is to be built from the infinite pool of values supplied by the list: an application of the *scalarize* constructor $\sigma[x_0, x_1, \dots]^{\text{list}}$ builds the 0-dimensional (scalar) array consisting of element x_0 only; *redimensioning* ρn raises the dimensionality of the array by 1 (the new dimension has extension n): $((\rho 3) \cdot \sigma)[x_0, x_1, \dots]^{\text{list}}$ constructs a 1-dimensional array containing elements x_0, x_1, x_2 ; $(\rho m) \cdot (\rho n) \cdot \sigma$ denotes the $(n \times m)$ -matrix constructor, and so on.

It is this rather different nature of arrays which led us to exclude their treatment in this text. The details of their coverage would blur much of the conceptual simplicity (*e. g.*, the need for the single insert representation functor $F = K_I + K_E \times \text{Id}$, only) we benefit from in this stage of discussion.

84 The catamorphisms (primarily in the disguise of monad comprehensions) emitted by \mathcal{Q} provide a somewhat unfamiliar view of database queries. Much of this is due to the uniform catamorphic representation of query primitives: the encoding of projections, joins, and selections (*e. g.*, the translation of a **select-from-where** block) is subject to inspection as is the representation of quantification or element extraction.

This is unlike the situation we face with “classical” query formalisms, relational algebra or calculus, say. User-level query language primitives that lie outside the reach—in terms of expressiveness—of the relational operators find their way into the relational query representation as *black box* operators or function calls. During the query rewrite phase these black boxes may be merely moved around as their algebraic properties are not known. Their interaction with the well-understood optimization rules in the relational domain remains unclear. Much may be missed that way.

To prove the catamorphic (monadic) comprehension of queries viable, there are thus two questions on our agenda:

- How is already well-established and provably useful knowledge on query optimization expressed in the catamorphic model? Can we mimic these approaches? We try to find answers to these questions starting with the section on query normalization in the current chapter and the material on query combinators in Chapter 4. Actually, we hope to provide more than just mimicry: in places, the categorical query model can simplify and unify previous ideas.
- Where does the catamorphic comprehension of queries give additional insight? Can the categorical query model answer previously open questions? We explore the now newly opened opportunities for query comprehension throughout this text, especially in Chapters 5 and 6.

Before we go on to work on the items of this agenda, let us spend the next few paragraphs and try to give a taste of what kind of possibilities for optimization show up if queries are tackled using a categorical toolbox.

85 Different placements of OQL’s extraction operator **element** (or **first**, respectively) in a query may obviously influence the quality of a query plan in dramatic ways. Early execution of element extraction can lead to removal of joins or even query unnesting. Join removal, *e. g.*, is possible for the OQL query below (remember the convention that a query expression like fxy denotes a query f containing free variables whose occurrences are replaced by

x, y):

```

element(select f x y
        from xs as x, ys as y) .

```

Computing the cartesian product of xs and ys is wasted work as we are throwing the result away should the product unexpectedly contain more than a single element (in which case the query raises an exception). An `element`-aware query optimizer could emit the equivalent

$$f(\text{element}(xs))(\text{element}(ys))$$

which removes the need for the cartesian product and thus never wastes work. Note that the simplified query raises an exception if and only if the original form would. However, `element` pushdown has a perilous nature:

- The above rewrite does not preserve equivalence if we are computing with sets (replace `select` by `select distinct`): function f might not be one-to-one and duplicate removal could potentially reduce the product to hold exactly one element.
- We must not push `element` extraction beyond a selection as the selection might filter its input so that exactly one element passes (*e. g.*, selection on a key).

This raises at least two questions. How are the safe rewrites recognizable and, in safe cases, how do we obtain the optimized form?

This is where the theory jumps in. Once we have observed that OQL's `element` extraction operators act like *monad morphisms* we can answer both questions.

86 Definition. Given two monads \mathbb{T}, \mathbb{T}' , an arrow $h: \mathbb{T}E \rightarrow \mathbb{T}'E$ is called a **monad morphism** [117] if it interacts with the operations of \mathbb{T} and \mathbb{T}' as follows:

$$\begin{aligned}
 h \cdot \mathbb{T}f &= \mathbb{T}'f \cdot h && \text{MONAD-MORPH-1} \\
 h \cdot \text{unit}^{\mathbb{T}} &= \text{unit}^{\mathbb{T}'} && \text{MONAD-MORPH-2} \\
 h \cdot \text{join}^{\mathbb{T}} &= \text{join}^{\mathbb{T}'} \cdot \mathbb{T}'h \cdot h . && \text{MONAD-MORPH-3}
 \end{aligned}$$

Should \mathbb{T} and \mathbb{T}' be monads with zero, we additionally require

$$h \cdot \text{zero}^{\mathbb{T}} = \text{zero}^{\mathbb{T}'} . \quad \text{MONAD-MORPH-4}$$

//

87 Element extraction acts like a monad morphism from a collection monad \mathbb{T} to the identity monad Id (here we take the freedom to mix OQL keywords and formal notation, having in mind that **Q-ELEMENT** establishes the correspondence $\text{element} = (\perp \nabla \text{outl})$):

$$\begin{aligned} \text{element} \cdot \mathbb{T}f &= f \cdot \text{element} && \text{ELEM-PUSHDOWN-1} \\ \text{element} \cdot \text{unit}^{\mathbb{T}} &= \text{id} && \text{ELEM-PUSHDOWN-2} \\ \text{element} \cdot \text{join}^{\mathbb{T}} &= \text{element} \cdot \text{element} \cdot && \text{ELEM-PUSHDOWN-3} \end{aligned}$$

Two restrictions apply. First, **ELEM-PUSHDOWN-1** does not hold for $\mathbb{T} = \text{set}$ in general but only for constant f (i. e., $f = K_x$ for some x). Second, element is not a morphism between monads with zero since $\text{element} \cdot \text{zero}^{\mathbb{T}} = \perp \neq \text{zero}^{\text{Id}}$.

Observe that element fails to act like a monad morphism in exactly those two situations that we have identified as unsafe pushdown rewritings in Paragraph 85. This suggests to exploit the above three equivalences as rewritings that propagate element through the monad operations and use these with \mathcal{Q} as a basis for safe element pushdown.

For the example query of Paragraph 85 we get

$$\begin{aligned} & (\mathcal{M} \cdot \mathcal{Q}) \left(\text{element} \left(\begin{array}{l} \text{select } f \ x \ y \\ \text{from } x \text{ s as } x, \ y \text{ s as } y \end{array} \right) \right) \\ = & \quad \{ \text{Q-ELEMENT, Q-SFW} \} \\ & \text{element} (\mathcal{M} \llbracket f \ x \ y \llbracket x \leftarrow x \text{ s}, \ y \leftarrow y \text{ s} \rrbracket^{\text{bag}}) \\ = & \quad \{ \mathcal{M-5} \} \\ & (\text{element} \cdot \text{join}^{\text{bag}}) (\mathcal{M} \llbracket \llbracket f \ x \ y \llbracket y \leftarrow y \text{ s} \rrbracket^{\text{bag}} \llbracket x \leftarrow x \text{ s} \rrbracket^{\text{bag}}) \\ = & \quad \{ \mathcal{M-3} \} \\ & (\text{element} \cdot \text{join}^{\text{bag}}) (\text{bag} (\lambda x. \mathcal{M} \llbracket f \ x \ y \llbracket y \leftarrow y \text{ s} \rrbracket^{\text{bag}}) x \text{ s}) \\ = & \quad \{ \mathcal{M-3} \} \\ & (\text{element} \cdot \text{join}^{\text{bag}}) (\text{bag} (\lambda x. (\text{bag} (\lambda y. f \ x \ y) y \text{ s})) x \text{ s}) \\ = & \quad \{ \text{ELEM-PUSHDOWN-3} \} \\ & (\text{element} \cdot \text{element} \cdot \text{bag}) ((\lambda x. (\text{bag} (\lambda y. f \ x \ y) y \text{ s})) x \text{ s}) \\ = & \quad \{ \text{ELEM-PUSHDOWN-1} \} \\ & \text{element} ((\lambda x. (\text{bag} (\lambda y. f \ x \ y) y \text{ s})) (\text{element } x \text{ s})) \\ = & \quad \{ \beta\text{-reduction} \} \\ & (\text{element} \cdot \text{bag}) ((\lambda y. f (\text{element } x \text{ s}) y) y \text{ s}) \\ = & \quad \{ \text{ELEM-PUSHDOWN-1} \} \\ & (\lambda y. f (\text{element } x \text{ s}) y) (\text{element } y \text{ s}) \end{aligned}$$

$$= \{ \beta\text{-reduction} \} \\ f(\text{element } xs)(\text{element } ys)$$

and thus obtain the optimized form we were after. The propagation rules (**ELEM-PUSHDOWN-1** through **ELEM-PUSHDOWN-3**) push element extraction down as far as possible but not beyond filters (which are implemented with the help of monadic zeroes for which no propagation rule applies). We can see this even more easily if we exploit the syntactic sugar we have introduced for monadic computations: comprehension syntax.

In [117], Wadler observed that the action of a monad morphism on a monadic computation may be concisely described by way of comprehension syntax. A monad morphism $h: \mathbb{T}E \rightarrow \mathbb{T}'E$ acts on a \mathbb{T} -monad comprehension as follows:

$$h \llbracket e \parallel qs \rrbracket^{\mathbb{T}} = \llbracket e \parallel h qs \rrbracket^{\mathbb{T}'}$$

MONAD-MORPH-5

where h 's extension on a qualifier list qs is defined by (ε and identifiers v, e, qs, qs' represent the syntactic categories of Definition 67)

$$\begin{array}{ll} h \varepsilon & = \varepsilon & \text{empty} \\ h(v \leftarrow e) & = v \leftarrow h e & \text{generators} \\ h(qs, qs') & = h qs, h qs' & \text{qualifier lists} \\ h e & = e . & \text{filters} \end{array}$$

For $h = \text{element}$, these simple syntactical rewritings exactly express the `element` pushdown optimization and are therefore especially suitable for inclusion into a rule-based optimizer. For our former example we calculate:

$$\begin{aligned} & \text{element} \llbracket f x y \parallel x \leftarrow xs, y \leftarrow ys \rrbracket^{\text{bag}} \\ = & \{ \text{MONAD-MORPH-5} \} \\ & \llbracket f x y \parallel \text{element}(x \leftarrow xs, y \leftarrow ys) \rrbracket^{\text{ld}} \\ = & \{ \text{MONAD-MORPH-5} \} \\ & \llbracket f x y \parallel \text{element}(x \leftarrow xs), \text{element}(y \leftarrow ys) \rrbracket^{\text{ld}} \\ = & \{ \text{MONAD-MORPH-5} \} \\ & \llbracket f x y \parallel x \leftarrow \text{element } xs, y \leftarrow \text{element } ys \rrbracket^{\text{ld}} \\ = & \{ \text{comprehension in ld monad} \} \\ & f(\text{element } xs)(\text{element } ys) \end{aligned}$$

(for the last step remember that a monad comprehension in the `ld`-monad essentially describes a nested *let*-expression, cf. Example 71).

Optimizations like these are indeed relevant. Early execution of element extraction can decrease the nesting depth of a query as the final example shows:

$$\begin{aligned}
& \mathcal{Q} \left(\begin{array}{l} \text{element}(\text{select}(\text{select } f \ x \ y \\ \qquad \qquad \qquad \text{from } ys \ \text{as } y) \\ \qquad \qquad \qquad \text{from } xs \ \text{as } x) \end{array} \right) \\
= & \{ \mathcal{Q} \} \\
& \text{element} \llbracket \llbracket f \ x \ y \mid y \leftarrow ys \rrbracket^{\text{bag}} \mid x \leftarrow xs \rrbracket^{\text{bag}} \\
= & \{ \text{MONAD-MORPH-5} \} \\
& \llbracket \llbracket f \ x \ y \mid y \leftarrow ys \rrbracket^{\text{bag}} \mid \text{element}(x \leftarrow xs) \rrbracket^{\text{ld}} \\
= & \{ \text{MONAD-MORPH-5} \} \\
& \llbracket \llbracket f \ x \ y \mid y \leftarrow ys \rrbracket^{\text{bag}} \mid x \leftarrow \text{element } xs \rrbracket^{\text{ld}} \\
= & \{ \text{comprehension in ld monad} \} \\
& \llbracket f(\text{element } xs) \ y \mid y \leftarrow ys \rrbracket^{\text{bag}} .
\end{aligned}$$

The last monad comprehension has the equivalent OQL form

```

select f(element(xs)) y
  from ys as y ,

```

i. e., a query which simply maps f over collection ys instead of creating a nested bag of bags like the original form did.

88 It is rather straightforward to check that OQL’s duplicate elimination primitives `listtset` and `distinct` act like monad morphisms from the `list` and `bag` monads to the `set` monad, respectively. In other words, the function `listtset = λe.⟦x ∣ x ← e⟧set:listE → setE` fulfills laws **MONAD-MORPH-1** through **MONAD-MORPH-4** and thus law **MONAD-MORPH-5**. An analogous remark applies to `distinct`.

We can exploit this observation to *delay duplicate removal* inside the query engine.

- To delay the removal of duplicates may seem counter-productive as, in general, duplicate elimination reduces operand sizes and thus might save execution costs. Duplicate removal, however, is a rather expensive operation to perform as it requires a sort or hash pass to identify duplicate elements. This, in turn, may disrupt the desirable stream-based processing of queries (to which we will turn in Chapter 6).

- Second, if we take a look at the innards of a query engine implemented on a common platform, actual implementations are almost exclusively list-based. This is partly due to the just mentioned cost inherent to duplicate removal (which manifests itself in the `set` constructor) but additionally enables the engine to reason about interesting (sort) orders [102] of the streams it processes.

The query rewrite that enables delay of duplicate elimination is justified by the equivalence (let $xs:\text{list}E$, we could reason analogously for $xs:\text{bag}E$ and `distinct`)

$$\llbracket e \parallel qs, x \leftarrow xs, qs' \rrbracket^{\text{set}} = \text{listtset} \llbracket e \parallel qs, x \leftarrow xs, qs' \rrbracket^{\text{list}}.$$

The right-hand side evaluates the comprehension in the `list` monad which preserves an potentially interesting sort order of xs until duplicates are finally removed by the outer `listtset` primitive.

We could verify the correctness of this transformation by means of the monad morphism properties of `listtset`. We will not do so here but rather appeal to monad comprehension *unnesting* which will be the focus of the forthcoming section. The correctness of delayed duplicate removal is immediately asserted by comprehension *nesting*:

$$\begin{aligned} & \llbracket e \parallel qs, x \leftarrow xs, qs' \rrbracket^{\text{set}} \\ = & \quad \{ \text{nesting via } \mathcal{M}\text{-NORM-3} \} \\ & \llbracket x' \parallel x' \leftarrow \llbracket e \parallel qs, x \leftarrow xs, qs' \rrbracket^{\text{list}} \rrbracket^{\text{set}} \\ = & \quad \{ \text{definition of } \text{listtset} \} \\ & \text{listtset} \llbracket e \parallel qs, x \leftarrow xs, qs' \rrbracket^{\text{list}}. \end{aligned}$$

Unnesting plays a crucial role during monad comprehension *normalization* to which we will turn now.

3.2 Monad Comprehension Normalization

89 Being a completely syntax-directed translation, \mathcal{Q} rather directly reflects the nested structure of an OQL query in the comprehensions it emits. Deriving anything but nested loop forms from a deeply nested query expression is a considerably hard task and a widely recognized challenge in the query processing literature [35, 48, 69, 105, 107]. In order to make good use of the efficient processing primitives of an underlying query engine we are better off to unravel nested queries.

The monad comprehension calculus provides particularly efficient yet simple handles to attack the problem.

- Different types of query nesting lead to similar, if not identical, nested forms of monad comprehensions. Rather than to maintain and identify a number of special cases of nesting—this route has been taken by numerous approaches, notably Kim’s original and followup work on classifying nested SQL queries [48, 69]—we can concentrate on unnesting the relatively few comprehension forms.
- Much of the actual unnesting work can be achieved by the application of a small set of *normalization rules* to monad comprehensions. Despite its simplicity, comprehension normalization covers a wide range of query unnesting techniques.

The normalization rules will prove to be valuable transformation tools in general (above we have already used one of these rules—more specifically its dual, nesting rule—to justify delayed duplicate removal).

90 The normalization rule set is assembled from a number of sources. Wadler already described the core rule \mathcal{M} -NORM-3 in his original work on endo-monadic comprehensions [117]. We augment this rule with a number of comprehension transformations that have appeared in the work of Wong and Fegaras [38, 124].

Here, we will put these rules to use and confirm their correctness in the multi-monad comprehension framework. An assessment of the soundness of the rules is crucial. The past has shown that relational unnesting is a query transformation already complex enough to let incorrect rewritings—*e. g.*, the infamous count bug of [69]—remain undiscovered for long.

91 Definition. The following set of equivalences is referred to as the **monad comprehension normalization rules** (an equivalence defines a normalization rule if it is read from left to right—interpreted from right to left we obtain the *dual* or *nesting* rule). Let $\top' \neq \text{Id}$.

$$\begin{aligned}
\llbracket e \parallel qs, v \leftarrow \text{zero}^{\top'}, qs' \rrbracket^{\top} &= \text{zero}^{\top} && \mathcal{M}\text{-NORM-1} \\
\llbracket e \parallel qs, v \leftarrow \text{unit}^{\top'} e', qs' \rrbracket^{\top} &= \llbracket e[e'/v] \parallel qs, qs'[e'/v] \rrbracket^{\top} && \mathcal{M}\text{-NORM-2} \\
\llbracket e \parallel qs, v \leftarrow \llbracket e' \parallel qs'' \rrbracket^{\top'}, qs' \rrbracket^{\top} &= \llbracket e[e'/v] \parallel qs, qs'', qs'[e'/v] \rrbracket^{\top} && \mathcal{M}\text{-NORM-3} \\
\llbracket e \parallel qs, \llbracket e' \parallel qs'' \rrbracket^{\text{exists}}, qs' \rrbracket^{\top} &= \llbracket e \parallel qs, qs'', e', qs' \rrbracket^{\top} && \mathcal{M}\text{-NORM-4}
\end{aligned}$$

$$\begin{aligned} & \quad (\top \text{ left-idempotent}) \\ \llbracket \llbracket e \parallel qs' \rrbracket^{\text{ld}} \parallel qs \rrbracket^{\text{ld}} &= \llbracket e \parallel qs, qs' \rrbracket^{\text{ld}} \quad \mathcal{M}\text{-NORM-5} \end{aligned}$$

//

92 Theorem. The set of monad comprehension normalization rules is sound, *i. e.*, the left-hand and right-hand sides of each rule map to equivalent monadic values under \mathcal{M} .

Proof. We will only prove rules $\mathcal{M}\text{-NORM-3}$ and $\mathcal{M}\text{-NORM-5}$ valid here. The structure of the proofs for the other three rules closely follows that of $\mathcal{M}\text{-NORM-3}$ and thus provide few new insights. Along the way of establishing the correctness of $\mathcal{M}\text{-NORM-3}$, however, we will invoke a transformation (validated in Lemma 93) that will be useful in its own right later on.

If we slightly specialize the comprehension qualifier list by setting $qs = \varepsilon$ we can find a rather simple proof for rule $\mathcal{M}\text{-NORM-3}$. Being so far, we establish the proof in a second step by generalizing from $qs = \varepsilon$ to arbitrary qs , assuming the specialized rule to be valid.

Set $qs = \varepsilon$. Let, as usual, τ denote the initial algebra of the initial type (τ, \top) .

$$\begin{aligned} & \mathcal{M} \llbracket e \parallel v \leftarrow \llbracket e' \parallel qs'' \rrbracket^{\top'}, qs' \rrbracket^{\top} \\ = & \quad \{ \mathcal{M}\text{-5} \} \\ & \text{join}^{\top} (\mathcal{M} \llbracket \llbracket e \parallel qs' \rrbracket^{\top} \parallel v \leftarrow \llbracket e' \parallel qs'' \rrbracket^{\top'} \rrbracket^{\top}) \\ = & \quad \{ \mathcal{M}\text{-3} \} \\ & (\text{join}^{\top} \cdot (\tau)) (\mathcal{M} \llbracket \llbracket e \parallel qs' \rrbracket^{\top} \parallel v \leftarrow \llbracket e' \parallel qs'' \rrbracket^{\top'} \rrbracket^{\top'}) \\ = & \quad \{ \mathcal{M}\text{-2} \} \\ & (\text{join}^{\top} \cdot (\tau)) (\top' (\lambda v. \mathcal{M} \llbracket e \parallel qs' \rrbracket^{\top}) (\mathcal{M} \llbracket e' \parallel qs'' \rrbracket^{\top'})) \\ = & \quad \{ \text{FUNCTOR-}\mathcal{M}\text{-FUSION, established below} \} \\ & (\text{join}^{\top} \cdot (\tau)) (\mathcal{M} \llbracket \llbracket e[e'/v] \parallel qs'[e'/v] \rrbracket^{\top} \parallel qs'' \rrbracket^{\top'}) \\ = & \quad \{ \mathcal{M}\text{-3} \} \\ & \text{join}^{\top} (\mathcal{M} \llbracket \llbracket e[e'/v] \parallel qs'[e'/v] \rrbracket^{\top} \parallel qs'' \rrbracket^{\top}) \\ = & \quad \{ \mathcal{M}\text{-5} \} \\ & \mathcal{M} \llbracket e[e'/v] \parallel qs'', qs'[e'/v] \rrbracket^{\top}. \end{aligned}$$

Generalize to arbitrary qualifier list qs :

$$\mathcal{M} \llbracket e \parallel qs, v \leftarrow \llbracket e' \parallel qs'' \rrbracket^{\top'}, qs' \rrbracket^{\top}$$

$$\begin{aligned}
&= \{ \mathcal{M}\text{-5} \} \\
&\quad \text{join}^\top (\mathcal{M} \llbracket [e \parallel v \leftarrow [e' \parallel qs'']]^{\top'} \parallel qs' \rrbracket^\top \parallel qs \rrbracket^\top) \\
&= \{ \mathcal{M}\text{-NORM-3 with } qs = \varepsilon \} \\
&\quad \text{join}^\top (\mathcal{M} \llbracket [e[e'/v] \parallel qs''] \parallel qs'[e'/v] \rrbracket^\top \parallel qs \rrbracket^\top) \\
&= \{ \mathcal{M}\text{-5} \} \\
&\quad \mathcal{M} \llbracket e[e'/v] \parallel qs, qs'', qs'[e'/v] \rrbracket^\top.
\end{aligned}$$

Following the course of the proof for rule $\mathcal{M}\text{-NORM-3}$, establishing rule $\mathcal{M}\text{-NORM-1}$ then is immediate by appealing to the functoriality of \top and law MONAD-4 . A similar remark applies to rule $\mathcal{M}\text{-NORM-2}$ and monad law MONAD-2 . For $\mathcal{M}\text{-NORM-4}$ note that if \top is left-idempotent we have (for arbitrary x)

$$\begin{aligned}
&\text{if } (p \vee q) \text{ then } \text{unit}^\top x \text{ else } \text{zero}^\top \\
&= \{ \text{+ idempotent} \} \\
&\quad (\text{if } p \text{ then } \text{unit}^\top x \text{ else } \text{zero}^\top) \text{+} (\text{if } q \text{ then } \text{unit}^\top x \text{ else } \text{zero}^\top).
\end{aligned}$$

To prove $\mathcal{M}\text{-NORM-5}$ simply observe that in the ld-monad we have $\text{join}^{\text{ld}} = \text{id}$ and thus

$$\begin{aligned}
&\mathcal{M} \llbracket [e \parallel qs']^{\text{ld}} \parallel qs \rrbracket^{\text{ld}} \\
&= \{ \text{join}^{\text{ld}} = \text{id} \} \\
&\quad \text{join}^{\text{ld}} (\mathcal{M} \llbracket [e \parallel qs']^{\text{ld}} \parallel qs \rrbracket^{\text{ld}}) \\
&= \{ \mathcal{M}\text{-5} \} \\
&\quad \mathcal{M} \llbracket e \parallel qs, qs' \rrbracket^{\text{ld}}.
\end{aligned}$$

□

93 Lemma. The application of a type functor $\top \neq \text{ld}$ may be fused with a comprehension evaluated in the associated \top -monad (no variable bound in qs may appear free in f):

$$\top f (\mathcal{M} \llbracket e \parallel qs \rrbracket^\top) = \mathcal{M} \llbracket f e \parallel qs \rrbracket^\top. \quad \text{FUNCTOR-}\mathcal{M}\text{-FUSION}$$

Proof. Induction over qs . The equation holds for $qs = \varepsilon$ due to $\mathcal{M}\text{-1}$ and the functoriality of \top . For non-empty qs , we may safely assume that qs starts with a generator, *i. e.*, $qs = v \leftarrow e', qs'$. Should qs start with a filter p we simply move p to the end of qs . (It is straightforward to see that filters may

take an arbitrary place in the qualifier list as long as they are not moved out of the scope of the free variables they refer to. In our case, p is constant with respect to the variables bound in qs .)

Quite like in the validation of the normalization rules, we will first establish **FUNCTOR- \mathcal{M} -FUSION** for $qs' = \varepsilon$ and later generalize our findings to arbitrary qs' .

Induction base case. Set $qs' = \varepsilon$. In general, $e : \mathbb{T}'E$ with $\mathbb{T} \neq \mathbb{T}'$, E arbitrary:

$$\begin{aligned}
& \mathcal{M} \llbracket f e \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}} \\
= & \quad \{ \mathcal{M}\text{-3} \} \\
& (\tau) (\mathcal{M} \llbracket f e \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}'}) \\
= & \quad \{ \mathcal{M}\text{-2} \} \\
& ((\tau) \cdot \mathbb{T}'(\lambda v. f e)) (\mathcal{M} e') \\
= & \quad \{ \text{CATA-MAP-FUSION} \} \\
& \mathbb{T}(\lambda v. f e) (\mathcal{M} e') \\
= & \quad \{ v \text{ not free in } f, \mathbb{T} \text{ functor} \} \\
& (\mathbb{T}f \cdot \mathbb{T}(\lambda v. e)) (\mathcal{M} e') \\
= & \quad \{ \mathcal{M}\text{-2} \} \\
& \mathbb{T}f (\mathcal{M} \llbracket e \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}}) .
\end{aligned}$$

Assuming that the lemma is true for qs' , perform the induction step. This will complete the proof.

$$\begin{aligned}
& \mathcal{M} \llbracket f e \rrbracket v \leftarrow e', qs' \rrbracket^{\mathbb{T}} \\
= & \quad \{ \mathcal{M}\text{-5} \} \\
& \text{join}^{\mathbb{T}} (\mathcal{M} \llbracket \llbracket f e \rrbracket qs' \rrbracket^{\mathbb{T}} \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}}) \\
= & \quad \{ \text{induction hypothesis} \} \\
& \text{join}^{\mathbb{T}} (\mathcal{M} \llbracket \mathbb{T}f \llbracket e \rrbracket qs' \rrbracket^{\mathbb{T}} \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}}) \\
= & \quad \{ \text{induction base case} \} \\
& (\text{join}^{\mathbb{T}} \cdot \mathbb{T}\mathbb{T}f) (\mathcal{M} \llbracket \llbracket e \rrbracket qs' \rrbracket^{\mathbb{T}} \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}}) \\
= & \quad \{ \text{join}^{\mathbb{T}} \text{ is natural transformation } \mathbb{T}\mathbb{T} \dot{\rightarrow} \mathbb{T} \} \\
& \mathbb{T}f (\text{join}^{\mathbb{T}} (\mathcal{M} \llbracket \llbracket e \rrbracket qs' \rrbracket^{\mathbb{T}} \rrbracket v \leftarrow e' \rrbracket^{\mathbb{T}})) \\
= & \quad \{ \mathcal{M}\text{-5} \} \\
& \mathbb{T}f (\mathcal{M} \llbracket e \rrbracket v \leftarrow e', qs' \rrbracket^{\mathbb{T}}) .
\end{aligned}$$

See Lemma 114 for a related transformation applicable to ld-monad comprehensions. \square

94 Rules \mathcal{M} -NORM-1 through \mathcal{M} -NORM-5 form a confluent and terminating set of rewriting rules which is the main incentive to refer to them as *normalization* rules. Given a monad comprehension, repeated application of the rules derives an equivalent syntactic *normal form*

$$\llbracket e \parallel v_1 \leftarrow e_1, v_2 \leftarrow e_2, \dots, v_n \leftarrow e_n, p \rrbracket^\top$$

for the original comprehension. This is achieved by rules \mathcal{M} -NORM-3 and \mathcal{M} -NORM-4. Normalization gives an unnesting procedure that is *complete* in the sense that an exhaustive application of the rules leads to a comprehension in which all semantically sound unnestings have been performed [40, 41]. If a query is evaluated in idempotent monads only (*cf.* the proviso of \mathcal{M} -NORM-4), then this can go as far as that all e_i are atomic expressions with respect to the monad comprehension syntax given in Definition 67, *i. e.*, database entry points or constants. Nested terms may still occur in the filter p and the comprehension's head e . We will approach these types of nesting when we have introduced *algebraic combinators* in the upcoming Chapter 4.

The comprehension normal form has been used by Fegaras and Maier as a canonical starting point for query optimization [40, 41]. Scanning the comprehension qualifier list from left to right, the algorithm described in [41] introduces *outerjoins* and algebraic grouping operators to cope with nested comprehensions that could not be flattened by normalization.

Unnesting disentangles queries and makes operands of formerly inner queries accessible in the outer enclosing comprehension. This, in turn, provides new possibilities for further rewritings and optimizations that will prove especially useful in the course of Chapters 4 and 5. Unnesting principally facilitates the applicability of other optimizing transformations that may help to, *e. g.*, trade a nested loop for a more efficient join.

Comprehension syntax provides a rather poor variety of syntactical forms, but in the early query translation phase this is more of a virtue than a shortcoming. Monad comprehensions extract and emphasize the *structural* information contained in a query rather than to stress diversity of query clauses. It is this uniformity that facilitates deep query analysis like the completeness result for query unnesting via comprehension normalization which we have mentioned in Paragraph 94. This can lead to new insights and simplifications which is the point the three concluding paragraphs try to make.

95 In [106], Steenhagen, Apers, and Blanken analyzed a class of SQL-like queries which exhibit correlated nesting in the `where`-clause, more specifically

```

select distinct f x
  from xs as x
 where p x z
    with z = (select g x y
              from ys as y
              where q x y) .

```

It is the question whether queries of this class may be rewritten into *flat join queries* of the form

```

select distinct f x
  from xs as x, ys as y
 where q x y
    and p' x v
    with v = g x y .

```

Queries for which such a replacement predicate p' cannot be found have to be processed either (a) using a nested loop strategy, or (b) by grouping, ideally via a *nestjoin* (see Paragraph 108 and [48, 69, 105, 106, 107]). Whether we can derive a flat join query is, naturally, dependent on the nature of the yet unspecified predicate p .

Steenhagen et al. state the following theorem—reproduced here using monadic terminology—which provides a partial answer to the question:

Whenever $p x z$ can be rewritten into $\llbracket p' x v \mid v \leftarrow z \rrbracket^{\text{exists}}$ (for some p') the original query may be evaluated by a flat join.

For space reasons a proof is omitted in [106]. After the original query is compiled via \mathcal{Q} , however, normalization can provide a succinct proof in a few lines (collection monad \mathbb{T} is arbitrary):

$$\begin{aligned}
& \llbracket f x \mid x \leftarrow xs, p x z \rrbracket^{\text{set}} \\
= & \{ \text{unfold } p \} \\
& \llbracket f x \mid x \leftarrow xs, \llbracket p' x v \mid v \leftarrow z \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \mathcal{M}\text{-NORM-4} \} \\
& \llbracket f x \mid x \leftarrow xs, v \leftarrow z, p' x v \rrbracket^{\text{set}} \\
= & \{ \text{unfold } z \} \\
& \llbracket f x \mid x \leftarrow xs, v \leftarrow \llbracket g x y \mid y \leftarrow ys, q x y \rrbracket^{\mathbb{T}}, p' x v \rrbracket^{\text{set}}
\end{aligned}$$

$$= \{ \mathcal{M}\text{-NORM-3} \} \\ \llbracket f x \llbracket x \leftarrow xs, y \leftarrow ys, q x y, p' x (g x y) \rrbracket^{\text{set}} .$$

A flat join between xs and ys can evaluate the resulting comprehension (this will become more apparent in Chapter 4).

But we can say even more and strengthen the statement of the theorem (thus answering an open question that has been put by Steenhagen et al. in [106]):

If p is not rewriteable into an existential quantifier like above, then we can conclude—based on the completeness result for comprehension normalization—that unnesting will in fact be impossible whatsoever.

Note. In [106], the *antijoin* (cf. Chapter 4) is also considered a flat join operator and we will do so in the upcoming chapters, too. This turns $\llbracket \neg p' x v \llbracket v \leftarrow z \rrbracket^{\text{all}} \rrbracket$ into an additionally acceptable form for predicate p .

96 Kim’s seminal work on the unnesting of SQL queries [69] may largely be understood in terms of normalization if queries are interpreted in the monad comprehension calculus. At the same time we can generalize and avoid the treatment of special cases. We additionally gain insight into questions on the validity of these unnesting strategies in the context of a complex object model featuring diverse collection type formers beyond `set`.

In [69], a classification of query nesting was proposed which associates each type of nesting with an equivalence-preserving unnesting transformation (modulo the infamous `count` bug whose treatment requires the introduction of the relational outerjoins or nestjoins [48, 107]) whose primary goal is to rewrite the original nested expression into a flat join query. Monad comprehension normalization readily unnests queries of Kim’s *type J*, i. e., SQL queries of the form

```
select distinct f x
  from xs as x
 where p x in (select g y
              from ys as y
              where q x y) .
```

Note that predicate q refers to the query variable x so that outer and nested query are correlated; uncorrelated (*type N*) queries need no special treatment and are transformed likewise. We will turn to a considerable generalization

of *type D* queries—containing universal quantifiers and thus essentially performing relational division—in Paragraph 116 when the associated notions are in context.

Normalization deduces an equivalent join query for a *type J* expression in two steps:

$$\begin{aligned}
& \llbracket f x \mid x \leftarrow xs, \llbracket p x = v \mid v \leftarrow \llbracket g y \mid y \leftarrow ys, q x y \rrbracket^T \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \mathcal{M}\text{-NORM-3} \} \\
& \llbracket f x \mid x \leftarrow xs, \llbracket p x = g y \mid y \leftarrow ys, q x y \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \mathcal{M}\text{-NORM-4} \} \\
& \llbracket f x \mid x \leftarrow xs, y \leftarrow ys, q x y, p x = g y \rrbracket^{\text{set}} .
\end{aligned}$$

The resulting comprehension is Kim’s *canonical 2-relation query* (i. e., the join)

```

select distinct f x
  from xs as x, ys as y
 where q x y
    and p x = g y .

```

97 Above, we have chosen the collection monads as general possible: this form of *type J full unnesting* is sound only if the outer query is evaluated in the **set** monad; no such restriction is necessary for the inner query.

This restriction on the type formers becomes void, however, if we already stop the unnesting process after the application of rule $\mathcal{M}\text{-NORM-3}$. As we will see in the next chapter, this comes with the extra benefit that the query becomes subject to evaluation by an efficient *semijoin* (\bowtie) while the fully unnested query calls for a general θ -join. Full unnesting may blur the detection of beneficial query execution alternatives.

On these grounds, while we believe in monad comprehension normalization as an important *preprocessing step* that should precede other query optimization phases, we do not subscribe to full comprehension unnesting as it was argued for in [38, 40, 41].

Chapter 4

Combinators

98 Given a simple spine transformer-based query engine as well as the mapping $\mathcal{M} \cdot \mathcal{Q}$ we are already able to rapidly prototype a fully functional complex value database system with OQL support. Based on the work reported in the two previous chapters, a spine transformer database engine has actually been manufactured from scratch and successfully operated in the realm of the *CROQUE* object database project [52, 58]. For efficiency reasons we have discussed in the introduction to Chapter 3, this is, however, clearly *not* the route we want to take. We will nevertheless briefly take up the idea in Chapter 6.

On the other hand, from the viewpoint of the catamorphic query representation, we know that *any* action executed in the query engine may eventually be broken down into walks of a spine. It is only that a classical database query engine implements algorithms—*physical operators* [66]—which allow for a significantly more efficient, possibly parallel, execution of a potentially large series of spine transformations. Typically this efficiency is accomplished by means of more or less involved support data structures, *e. g.*, *indices*. Although these physical operators construct spines they generally do so in completely different ways (random access to spine nodes, block-wise processing, and so on) than in adherence to the primitive catamorphic recursion pattern.

The principle observation that motivates the present chapter is that it is nevertheless possible to *encode the action* of a physical operator using our abstract query representation, *i. e.*, via monad comprehensions and thus catamorphisms. This encoding emphasizes the algebraic properties of the operator rather than its—potentially many—underlying physical realizations. Join \bowtie , for example, is an abstraction that reduces join algorithms like *sort-merge join* or *hash-join* to their algebraic properties (*e. g.*, being commutative).

We refer to these encodings as *combinators* for reasons which will become clear soon. Unlike in the monad comprehension calculus, expressing a query

in terms of combinators provides the query engine with clear hints which algorithms, besides nested loops, may be enlisted to execute the query. The set of combinators encoding the physical algorithms of a query engine forms the *logical algebra* associated with this engine [55].

99 Combinators. We have just said that a *combinator* encodes a series of actions on a spine which, when viewed from the outside, *act like* a possibly complex catamorphic spine transformer. For our aims, this is an operator’s prominent property. Its actual physical realization is going on behind the scenes. This opacity is reflected by the means we use to model operators, namely as *closed expressions*, *i. e.*, *combinators*, in the sense of the λ -calculus [92, 95]. As an example, the relational *selection* operator σ is represented by the closed expression

$$\begin{aligned} \sigma & : (E \rightarrow Bool) \rightarrow \text{set}E \rightarrow \text{set}E \\ \sigma & = \underbrace{\lambda p. \lambda xs. \llbracket x \parallel x \leftarrow xs, px \rrbracket^{\text{set}}}_{\text{closed}} . \end{aligned}$$

As the set of free variables of a combinator expression is empty by definition we do not encounter inter-combinator variable dependencies in a composition of combinators. This facilitates the reordering and transformation of pure combinator expressions, a property that has long been exploited in algebraic query rewriting frameworks. (Although surprisingly many proposed “query algebras” heavily depend on the presence of free variables and thus do not qualify as algebras in the strict sense of the word.)

To reach specific optimization goals, notably the generation of stream-based programs for queries (query execution plans that do not allocate memory to buffer intermediate results), it pays off to break the encapsulation the combinator notion provides. We will not do so until Chapter 6, however.

100 Combinators from monad comprehensions. The path is clear: given a query engine’s logical algebra we encode the actions of its operators as closed monad comprehension expressions or, in rare cases, catamorphisms. The PTIME expressive power of the catamorphic query model makes this possible for query engines implementing classical relational algebra as well as a broad range of advanced query algebras, from various proposals for the non-first normal form (NF²) relational model [1, 12, 35, 99, 100, 105] to algebras for complex value databases [6, 24, 30, 32, 76, 115]. Here, we will focus on a variant of the latter.

The equational theory of monads and the initial algebras the monads were derived from will provide the equational theory for the combinator algebra.

Since a combinator is, in principle, nothing else than a monad comprehension *macro* (built adhering to the closed expression property) its interaction with an enclosing comprehension expression is well-defined.

This paves the way for a step-wise or incremental discipline for the derivation of a pure combinator expression from a given monad comprehension: the combinator definitions serve as *patterns* which may be used to trigger a replacement trading the matched comprehension subexpression for an equivalent combinator. Along the way we are left with *hybrid* query expressions of which parts are either built using comprehension syntax or combinators. If the underlying query engine implements the catamorphism combinator (\cdot) for the initial algebras its data model rests upon—this is actually a light requirement—this step-wise replacement process is guaranteed to generate a pure combinator expression. Falling back onto catamorphisms, however, should clearly be the last resort.

Each replacement step enriches the original query expression by hints on which strategies for query execution may be profitable. This does not yet select the physical operators themselves, but (a) the introduction of a combinator seeds the search for the actual execution plan, and (b) a combinator comes with a limited number of physical implementation alternatives. In this sense we decrease the level of abstraction by placing combinators inside a formerly comprehension-based query.

101 An exploration of a mixture of the calculus and algebra styles of query rewriting over monadic types forms the core of this chapter. Query transformations may be formulated using the style in which their applicability is most easily detected and expressed.

This procedure has much in common with the work of [Nakano](#) as well as [Steenhagen, Apers, Blanken, and de By](#) in the relational and NF² domains, respectively [[88](#), [105](#), [106](#), [107](#)]. There, logical algebra semantics were specified using specific instances of the relational tuple calculus. As already described in the previous paragraph, this opens the scene for a hybrid query rewriting framework. [Steenhagen et al.](#) make significant use of relational calculus to prove correct and justify rewritings, especially when it comes to the processing of complex predicates containing quantifiers. The actual query transformation rules (and thus the actual query rewriting), however, are formulated at the algebra-level.

Here, we employ the monad comprehension calculus also as a tool to calculate on the “semantics level”, *i. e.*, the core query transformation work is performed on the calculus-level. Once again, comprehension normalization greatly simplifies our task. (The explicit treatment of *conjunctive path*

formulas in [105, p. 67 ff.] is, for example, an algebra-level formulation of unnesting rule \mathcal{M} -NORM-4).

It is our aim *not* to restrict our attention to queries that originate from the use of a specific subset of the user-level query language (OQL, say). The `select-from-where` block is the principle query construct, but we believe that studying its efficient translation in isolation from the rest of the query language comes with the danger of missing profitable plans. The uniformity of monad comprehension-based query representation has something to offer here. (Chapter 5 presents extended examples of the treatment of queries other than of the `select-from-where` type.)

4.1 Defining and Detecting Combinators

102 The following definition sets up a combinator algebra for a query engine that can operate over monadic types. It is our assumption that the query engine copes with set, bag, as well as list values and supports similar facilities for the three collection formers. At the combinator level this is reflected by polymorphic typing: the combinators are parametric in their monadic input and result types which implies that they behave similarly for the different type instantiations in which they are going to be used. At the same time this prevents the introduction of a large number of families of combinators in which the members of a family implement a common operator but each member for a specific type only.

103 Definition. Let \mathbb{T}, \mathbb{T}' denote monadic types and let E, E', E'' be placeholders for valid query types (*cf.* Paragraph 78). The list below defines the set of combinators that we will refer to in the sequel. Because we will encounter *curried* (*i. e.*, partial) applications of the combinators, this definition displays them in prefix notation instead of the “classical” infix form. (Fully applied binary operators are also written between their collection-typed arguments, *i. e.*, we consider

$$\bowtie p f xs ys \quad \text{and} \quad xs \underset{p|f}{\bowtie} ys$$

to be equivalent notations.) It goes without saying that a combinator application is sensible only if—given the monadic types of its arguments—its defining monad comprehension is well-defined for these type instantiations.

$$\pi \quad : \quad (E \rightarrow E') \rightarrow \mathbb{T}E \rightarrow \mathbb{T}E'$$

$\pi f xs$	$= \top f xs = \llbracket f x \parallel x \leftarrow xs \rrbracket^\top$	PROJECT
α^\top	$: (E \rightarrow E') \rightarrow \top' E \rightarrow \top E'$	
$\alpha^\top f xs$	$= \llbracket f x \parallel x \leftarrow xs \rrbracket^\top$	AGGREGATE
σ	$: (E \rightarrow Bool) \rightarrow \top E \rightarrow \top E$	
$\sigma p xs$	$= \llbracket x \parallel x \leftarrow xs, px \rrbracket^\top$	SELECT
\times	$: \top E \rightarrow \top' E \rightarrow \top(E \times E')$	
$\times xs ys$	$= \llbracket (x, y) \parallel x \leftarrow xs, y \leftarrow ys \rrbracket^\top$	CROSS
\bowtie	$: (E \rightarrow E' \rightarrow Bool) \rightarrow (E \rightarrow E' \rightarrow E'')$ $\rightarrow \top E \rightarrow \top' E' \rightarrow \top E''$	
$\bowtie p f xs ys$	$= \llbracket f xy \parallel x \leftarrow xs, y \leftarrow ys, pxy \rrbracket^\top$	JOIN
\triangle	$: (E \rightarrow E' \rightarrow Bool) \rightarrow (E \rightarrow \top' E' \rightarrow E'')$ $\rightarrow \top E \rightarrow \top' E' \rightarrow \top E''$	
$\triangle p f xs ys$	$= \llbracket f x \llbracket y \parallel y \leftarrow ys, pxy \rrbracket^{\top'} \parallel x \leftarrow xs \rrbracket^\top$	NESTJOIN
\bowtie	$: (E \rightarrow E' \rightarrow Bool) \rightarrow \top E \rightarrow \top' E' \rightarrow \top E$	
$\bowtie p xs ys$	$= \llbracket x \parallel x \leftarrow xs, \llbracket pxy \parallel y \leftarrow ys \rrbracket^{\text{exists}} \rrbracket^\top$	SEMIJOIN
\bowtie	$: (E \rightarrow E' \rightarrow Bool) \rightarrow \top E \rightarrow \top' E' \rightarrow \top E$	
$\bowtie p xs ys$	$= \llbracket x \parallel x \leftarrow xs, \llbracket \neg pxy \parallel y \leftarrow ys \rrbracket^{\text{all}} \rrbracket^\top$	ANTIJOIN
μ	$: \top \top' E \rightarrow \top E$	
μxss	$= \llbracket x \parallel xs \leftarrow xss, x \leftarrow xs \rrbracket^\top$	UNNEST
ζ	$: (E \times E \rightarrow Bool) \rightarrow \top E \rightarrow \text{list} E$	
$\zeta \odot xs$	$= (\text{nil}^{\text{list}} \nabla \text{ins}^{\text{list}} \odot) xs$	SORT

//

104 Notes. (a) Issues of polymorphism aside, the combinators are standard and should explain themselves. The nestjoin \triangle combines joining with grouping: for each x in xs a group from those y in ys is built that make predicate pxy evaluate to *true*. Note that nestjoin provides an abstraction for nesting in a comprehension's *head*, a case the normalization rules of Definition 91 do not account for. Nestjoin thus plays a central role in nested

query processing. Variants of nestjoin have appeared as, *e.g.*, the *binary grouping* operator of [35] or the *hierarchical join* of [96].

(b) **SEMIJOIN** and **ANTIJOIN** define the *left* variants of the semijoin and antijoin operators, respectively. The right variants \bowtie and \bowtie^c are dually defined. Semijoin $\bowtie p\ x\ y\ s$ emits all x in $x\ s$ for which some (arbitrary) y in $y\ s$ with $p\ x\ y$ can be found. The corresponding antijoin computes the complement.

(c) In a setting of monadic types, combinator α^T has two facets: evaluated in an identity monad, *i. e.*, a monad that has been derived from the algebras of Definition 64, α^T acts like a quantifier or aggregation operator; evaluated in a collection monad the combinator coerces between collection types (*e. g.*, duplicates are eliminated for $T = \text{set}$).

(d) Remember that it has not been our goal to define a *minimal* algebra in the sense that each combinator adds expressive power (indeed, (\cdot) is all what would be needed then). Highly specialized combinators like nestjoin rather flag the presence of an efficient execution algorithm inside the query engine and are thus introduced for optimization purposes only.

105 Partial matches. It is actually only a small step from a combinator definition to its associated detection rule: reading a definition from right to left specifies a rewriting rule that spots the opportunity to *completely* replace the right-hand side by a combinator. Note that such a rule application is valid only if the lhs combinator arguments are closed expressions (this is a requirement inherent to the combinator notion). Rarely, however, will such perfect matches occur.

The detection process is significantly enhanced if we provide (rule-based) knowledge about *partial* matches and replacements. Rules of this type obviously lead to hybrid query expressions as we have discussed them in Paragraph 100. As the thus induced rules exhibit a similar structure let us list only a selection of them below.

Most of the rules rely on the simple syntactic recipe after which monad comprehensions are constructed: as before, let qs, qs', qs'' represent the syntactic category of possibly empty lists of qualifiers (*cf.* Definition 67). The left-hand sides thus specify (qualifier) list patterns that match the specified template of generators and filters. We will assess the rules' correctness in the next paragraph.

Note. The uses of $p\ x$ or $p\ x\ y$ below denote predicates p in which no variables besides x (or x and y , respectively) may occur free. This includes closed predicates. The replacements in the rhs of rules **M-JOIN** and **M-NESTJOIN** implement the deconstruction of the tuple-typed elements found in the join

results. The introduction of pattern-matching capabilities for comprehension generators—which we have not done for simplicity reasons, but see [92, 125]—make such replacements obsolete.

$$\begin{aligned}
\llbracket f e \parallel qs \rrbracket^{\top} &\rightarrow \pi f \llbracket e \parallel qs \rrbracket^{\top} && \mathcal{M}\text{-PROJECT} \\
&&& (f \text{ closed}) \\
\llbracket e \parallel qs : \Gamma' \rrbracket^{\top} &\rightarrow \alpha^{\top} id \llbracket e \parallel qs \rrbracket^{\top} && \mathcal{M}\text{-AGGREGATE} \\
\llbracket e \parallel qs, x \leftarrow xs, &\rightarrow \llbracket e \parallel qs, x \leftarrow \sigma p xs, && \mathcal{M}\text{-SELECT} \\
qs', px, qs'' \rrbracket^{\top} &qs', qs'' \rrbracket^{\top} \\
\llbracket e \parallel qs, x \leftarrow xs, y \leftarrow ys, qs' &\rightarrow \llbracket e[outl v/x][outr v/y] \parallel qs, \\
pxy, qs'' \rrbracket^{\top} &v \leftarrow xs \quad \boxtimes \quad ys, \\
& \quad p|\lambda x.\lambda y.(x,y) \\
&qs'[outl v/x][outr v/y], \\
&qs''[outl v/x][outr v/y] \rrbracket^{\top} && \mathcal{M}\text{-JOIN} \\
\llbracket e \parallel qs, x \leftarrow xs, qs', &\rightarrow \llbracket e[outl v/x] \parallel qs, \\
pxe', qs'' \rrbracket^{\top} &v \leftarrow xs \quad \Delta \quad ys, \\
& \quad q|\lambda x.\lambda y.(x,gy) \\
&qs'[outl v/x], \\
\text{with} & && p[outl v/x][outr v/e'], \\
e' = \llbracket gxy \parallel y \leftarrow ys, qxy \rrbracket^{\top} &qs''[outl v/x] \rrbracket^{\top} && \mathcal{M}\text{-NESTJOIN} \\
&&& \\
\llbracket e \parallel qs, x \leftarrow xs, qs', &\rightarrow \llbracket e \parallel qs, x \leftarrow xs \times ys, && \mathcal{M}\text{-SEMIJOIN} \\
\llbracket pxy \parallel y \leftarrow ys \rrbracket^{\text{exists}}, &qs', qs'' \rrbracket^{\top} && p \\
qs'' \rrbracket^{\top} &&&
\end{aligned}$$

106 Lemma. The introduction of combinators via the rules listed in Paragraph 105 does not alter the meaning of the affected monad comprehensions.

Proof. The proof scheme is so simple that we establish the correctness of rules $\mathcal{M}\text{-PROJECT}$ and $\mathcal{M}\text{-JOIN}$ only here: (a) in the rhs, unfold the combinator definition found in Definition 103, then (b) use monad comprehension normalization to derive the lhs.

$\mathcal{M}\text{-PROJECT}$:

$$\pi f \llbracket e \parallel qs \rrbracket^{\top}$$

$$\begin{aligned}
&= \{ \text{unfold } \pi \text{ via } \mathbf{PROJECT} \} \\
&\quad \llbracket f x \parallel x \leftarrow \llbracket e \parallel qs \rrbracket^\top \rrbracket^\top \\
&= \{ \mathcal{M}\text{-NORM-3} \} \\
&\quad \llbracket (f x)[e/x] \parallel qs \rrbracket^\top \\
&= \{ \text{replacement} \} \\
&\quad \llbracket f e \parallel qs \rrbracket^\top .
\end{aligned}$$

$\mathcal{M}\text{-JOIN}$: we assume $xs : \top' E$ for some E :

$$\begin{aligned}
&\llbracket e[\text{outl } v/x][\text{outr } v/y] \parallel qs, v \leftarrow xs \underset{p|\lambda x.\lambda y.(x,y)}{\bowtie} ys, \\
&\quad qs'[\text{outl } v/x][\text{outr } v/y], \\
&\quad qs''[\text{outl } v/x][\text{outr } v/y] \rrbracket^\top \\
&= \{ \text{unfold } \bowtie \text{ via } \mathbf{JOIN}, \beta\text{-reduction} \} \\
&\llbracket e[\text{outl } v/x][\text{outr } v/y] \parallel qs, v \leftarrow \llbracket (x, y) \parallel x \leftarrow xs, y \leftarrow ys, p x y \rrbracket^{\top'}, \\
&\quad qs'[\text{outl } v/x][\text{outr } v/y], \\
&\quad qs''[\text{outl } v/x][\text{outr } v/y] \rrbracket^\top \\
&= \{ \mathcal{M}\text{-NORM-3} \} \\
&\llbracket e[\text{outl } v/x][\text{outr } v/y][(x, y)/v] \parallel qs, x \leftarrow xs, y \leftarrow ys, p x y, \\
&\quad qs'[\text{outl } v/x][\text{outr } v/y][(x, y)/v], \\
&\quad qs''[\text{outl } v/x][\text{outr } v/y][(x, y)/v] \rrbracket^\top \\
&= \{ \text{replacements, product} \} \\
&\llbracket e \parallel qs, x \leftarrow xs, y \leftarrow ys, p x y, qs', qs'' \rrbracket^\top \\
&= \{ \text{shift filter in qualifier list} \} \\
&\llbracket e \parallel qs, x \leftarrow xs, y \leftarrow ys, qs', p x y, qs'' \rrbracket^\top .
\end{aligned}$$

□

It comes as no surprise that unnesting the rhs of the combinator introduction rules derives their lhs: it is the ultimate point of combinators to capture the structure of a comprehension that has not been fully unnested or, equivalently, rewritten into its nested loops form (*cf.* Paragraph 97).

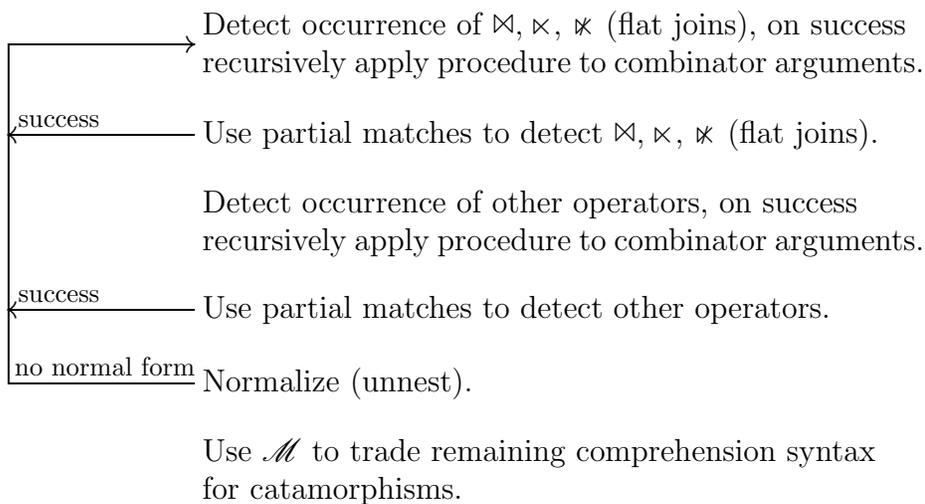
4.2 Combinators From Monad Comprehensions

107 Where is our current position in the greater picture? Prior to the introduction of combinators in this chapter we were dealing with a completely deterministic query translation process, namely $\mathcal{M} \cdot \mathcal{Q}$ followed by

full unnesting. This has changed with our goal to derive combinator queries from monad comprehensions: the set of rewriting rules built from the normalization rules, combinator definitions, and partial combinator patterns (of Paragraph 105) is clearly non-confluent. From the monad comprehension, *e. g.*, $\llbracket f x \parallel x \leftarrow xs, y \leftarrow ys, p y \rrbracket^T$, the rules open a spectrum from purely catamorphic plans (*i. e.*, nested loops), over combinator queries using \bowtie or \bowtie to the most efficient \bowtie -based alternative.

- A query rewriter that gives preference to the application of rules stemming from the combinator definitions will perform a *top-down* derivation—with respect to the syntactic structure of comprehensions—of combinator queries. *Benefit:* early detection of, *e. g.*, a semijoin occurrence (or of a flat join in general), can significantly improve the overall quality of the emitted combinator expression.
- The lhs of partial match rules inspect the innards of a comprehension and, once applied, replace parts of the comprehension (qualifier list), thus working their way *bottom-up*. *Benefit:* the lhs are much more likely to match a given query as they require partial matches only and may rewrite a comprehension such that a combinator definition can subsequently be detected.

This suggests the following sketch of a procedure which can guide a query optimizer in its task to assign priorities to rule groups and to organize its rewriting strategy. Probing the rules of the topmost alternative first, the optimizer only as a last resort falls through to the generation of nested loops via the monad comprehension desugaring scheme \mathcal{M} (Definition 69).



108 Let us pick up the query class of Paragraph 95:

$$\begin{aligned}
 Q \equiv & \text{select } f x \\
 & \text{from } xs \text{ as } x \\
 & \text{where } p x z \\
 & \text{with } z = \left(\begin{array}{l} \text{select } g x y \\ \text{from } ys \text{ as } y \\ \text{where } q x y \end{array} \right) .
 \end{aligned}$$

We already know how to deal with this type of query should predicate p be essentially equivalent to $\llbracket p' x v \rrbracket v \leftarrow z \rrbracket^{\text{exists}}$ for some p' . The theorem sketched in Paragraph 95 indicates the need for grouping in all other cases. To illustrate, let p denote a predicate of general form. A combinator query for Q should ideally employ the nestjoin Δ as this is the combinator which encodes the grouping functionality of the underlying query engine. The calculation goes on as follows (collection monads \mathbb{T}, \mathbb{T}' are arbitrary):

$$\begin{aligned}
 & \mathcal{Q} Q \\
 = & \{ \mathcal{Q}, \text{unfold } z \} \\
 & \llbracket f x \rrbracket x \leftarrow xs, p x \llbracket g x y \rrbracket y \leftarrow ys, q x y \rrbracket^{\mathbb{T}'} \rrbracket^{\mathbb{T}} \\
 = & \{ \mathcal{M}\text{-NESTJOIN} \} \\
 & \llbracket (f x)[\text{outl } v/x] \rrbracket v \leftarrow xs \quad \Delta_{q|\lambda x.\lambda y.(x,y)} \quad ys, (p x)[\text{outl } v/x] \rrbracket^{\mathbb{T}} \\
 = & \{ \text{replacements} \} \\
 & \llbracket (f \cdot \text{outl}) v \rrbracket v \leftarrow xs \quad \Delta_{q|\lambda x.\lambda y.(x,y)} \quad ys, (p \cdot \text{outl}) v \rrbracket^{\mathbb{T}} \\
 = & \{ \mathcal{M}\text{-SELECT} \} \\
 & \llbracket (f \cdot \text{outl}) v \rrbracket v \leftarrow \sigma(p \cdot \text{outl})(xs \quad \Delta_{q|\lambda x.\lambda y.(x,y)} \quad ys) \rrbracket^{\mathbb{T}} \\
 = & \{ \text{PROJECT} \} \\
 & \pi(f \cdot \text{outl})(\sigma(p \cdot \text{outl})(xs \quad \Delta_{q|\lambda x.\lambda y.(x,y)} \quad ys)) .
 \end{aligned}$$

The resulting combinator query is a close variant of the plan that has been proposed in the original work on nestjoin-based query processing [106]. As we do not invoke some sort of *Eureka step*, it is reasonable to assume that the above trail of rewritings will actually be found by the optimizer.

Nestjoin encodes nesting a comprehension's head (see Paragraph 104) so

that its canonical use arises whenever queries of the form

$$\begin{aligned}
 Q' &\equiv \text{select } f x z \\
 &\quad \text{from } xs \text{ as } x \\
 &\quad \text{where } p x \\
 &\quad \text{with } z = \left(\begin{array}{l} \text{select } g x y \\ \text{from } ys \text{ as } y \\ \text{where } q x y \end{array} \right)
 \end{aligned}$$

are to be processed. For $g x y = y$, the nestjoin occurrence in Q' is apparent as soon as \mathcal{Q} has desugared the OQL syntax [107]:

$$\begin{aligned}
 &\mathcal{Q} Q' \\
 = &\quad \{ \mathcal{Q}, \text{unfold } z \} \\
 &\quad \llbracket f x \llbracket y \llbracket y \leftarrow ys, q x y \rrbracket^T \llbracket x \leftarrow xs, p x \rrbracket^T \rrbracket^T \\
 = &\quad \{ \mathcal{M}\text{-PROJECT} \} \\
 &\quad \llbracket f x \llbracket y \llbracket y \leftarrow ys, q x y \rrbracket^T \llbracket x \leftarrow \sigma p xs \rrbracket^T \rrbracket^T \\
 = &\quad \{ \text{NESTJOIN} \} \\
 &\quad (\sigma p xs) \quad \triangle \quad ys . \\
 &\quad \quad \quad q | \lambda x. \lambda y. (x, y)
 \end{aligned}$$

For arbitrary projections g , a *tupling* transformation hand in hand with **FUNCTOR- \mathcal{M} -FUSION** pulls g out of the nested query block, thus paving the way for the use of nestjoin. A generalized nestjoin operator (accepting g as a parameter) that perfectly matches the above query class has been proposed in [106].

Query rewriting in this hybrid framework of combinators and monad comprehensions draws much of its attractiveness on the interleaving of combinator detection and calculus-level rewriting. The uniformity of the monad comprehension calculus—together with its syntactical simplicity—allows to strengthen the rewriting framework by adding only *few* but *generic* (*i. e.*, parametric in the monads) rules. This is the next point we try to make.

109 Example. Both nested existential quantifiers in the following OQL query are not closed (neither p nor q may be evaluated in the scope of one of the quantifiers only):

$$\begin{aligned}
 Q'' &\equiv \text{select distinct } x \\
 &\quad \text{from } xs \text{ as } x \\
 &\quad \text{where exists } y \text{ in } ys: \\
 &\quad \quad (\text{exists } z \text{ in } zs: q x z \text{ and } p y z) .
 \end{aligned}$$

Such variable interdependencies in the source query and its equivalent under mapping \mathcal{Q} may render the introduction of combinators a hard problem. Assuming that the query rewriter is equipped with monad comprehension normalization rules and combinator patterns (we will redo this example shortly), we then calculate as follows:

$$\begin{aligned}
& \mathcal{Q} Q'' \\
= & \{ \mathcal{Q} \} \\
& \llbracket x \mid x \leftarrow xs, \llbracket [q x z \wedge p y z \mid z \leftarrow zs] \text{exists} \mid y \leftarrow ys \rrbracket \text{exists} \rrbracket^{\text{set}} \\
= & \{ \mathcal{M}\text{-NORM-4} \} \\
& \llbracket x \mid x \leftarrow xs, y \leftarrow ys, \llbracket [q x z \wedge p y z \mid z \leftarrow zs] \text{exists} \rrbracket \text{set} \rrbracket \\
= & \{ \text{partial match for } \times \} \\
& \llbracket \text{outl } v \mid v \leftarrow xs \times ys, \\
& \quad \llbracket [q(\text{outl } v) z \wedge p(\text{outl } v) z \mid z \leftarrow zs] \text{exists} \rrbracket \text{set} \rrbracket \\
= & \{ \mathcal{M}\text{-SEMIJOIN} \} \\
& \llbracket \text{outl } v \mid v \leftarrow (xs \times ys) \quad \times \quad zs \rrbracket^{\text{set}} \\
& \quad \lambda v. \lambda z. q(\text{outl } v) z \wedge p(\text{outl } v) z \\
= & \{ \text{PROJECT} \} \\
& \pi \text{ outl } \left((xs \times ys) \quad \times \quad zs \right) \\
& \quad \lambda v. \lambda z. q(\text{outl } v) z \wedge p(\text{outl } v) z
\end{aligned}$$

The introduction of the cross product performs a tupling of xs and ys so that predicates p and q may be evaluated given this stream of tuples and the bindings of z generated from zs . Only then we obtain a constellation a semijoin can implement. Tupling comes at a rather high cost, not least because of the necessary untupling realized by the π combinator. \diamond

110 The rewriting process will clearly benefit from general means to unwind variable interdependencies. Transformations of this type are naturally expressed at the monad comprehension calculus level. The following three lemmata provide such tools. Having these at hand, Paragraph 115 will revisit the above example.

111 Lemma. Let \mathbb{T} denote a commutative monad. Adjacent qualifiers qs', qs'' in a \mathbb{T} -monad comprehension may be exchanged without affecting the comprehension's meaning under \mathcal{M} , given that qs' does not refer to variables bound in qs'' and vice versa [114]:

$$\mathcal{M} \llbracket e \mid qs, qs', qs'', qs''' \rrbracket^{\mathbb{T}} = \mathcal{M} \llbracket e \mid qs, qs'', qs', qs''' \rrbracket^{\mathbb{T}} \quad \text{QUAL-EX}$$

Note that repeated application of **QUAL-EX** gives a more general qualifier exchange theorem.

Proof. We can reuse the proof scheme we have applied to establish the monad comprehension normalization rules in Theorem 92: a proof for the case $qs = \varepsilon$ is found by desugaring the lhs of **QUAL-EX** and an appeal to the commutativity of $\#$ in $join^\top = ([]^\top \nabla \#)$ (cf. Lemma 60). Appealing to \mathcal{M} and the base case then readily establishes the statement for arbitrary qs . \square

A direct consequence of the above and comprehension unnesting rule **M-NORM-5** is the following.

112 Lemma. For **ld**-monads (monads derived from the algebras of Definition 64) we have that

$$\mathcal{M} \llbracket [e \parallel qs']^{\text{ld}} \parallel qs \rrbracket^{\text{ld}} = \mathcal{M} \llbracket [e \parallel qs]^{\text{ld}} \parallel qs' \rrbracket^{\text{ld}} \quad \text{NEST-EX}$$

provided that qs and qs' do not include mutual variable references.

Proof.

$$\begin{aligned} & \mathcal{M} \llbracket [e \parallel qs']^{\text{ld}} \parallel qs \rrbracket^{\text{ld}} \\ = & \{ \text{M-NORM-5} \} \\ & \mathcal{M} \llbracket e \parallel qs, qs' \rrbracket^{\text{ld}} \\ = & \{ \text{QUAL-EX} \} \\ & \mathcal{M} \llbracket e \parallel qs', qs \rrbracket^{\text{ld}} \\ = & \{ \text{M-NORM-5} \} \\ & \mathcal{M} \llbracket [e \parallel qs]^{\text{ld}} \parallel qs' \rrbracket^{\text{ld}} . \end{aligned}$$

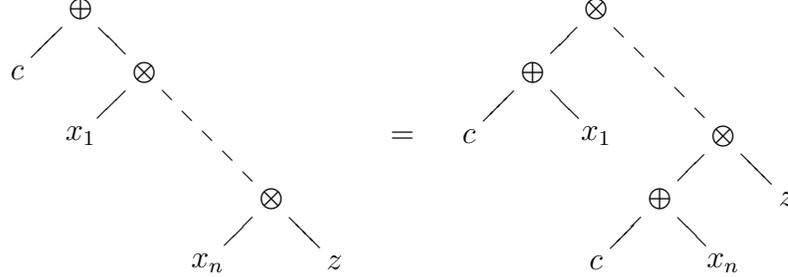
\square

The last of the three “unwinding” tools provides a catamorphic generalization of the *descoping* transformation described in [106] and [15]. Based on the notion of *distributivity* defined below, descoping can pull expressions out of a comprehension’s head.

113 Definition. Let $\alpha = z \nabla \otimes : \mathbf{F}A \rightarrow A$ denote an **F**-algebra, with **F** being the polynomial endofunctor in **Set** describing algebras in insert representation. Function $\oplus : A \times A \rightarrow A$ **distributes** over α if for all $c : A$

$$\oplus \cdot (K_c \Delta (\alpha)) = (z \nabla \otimes \cdot ((\oplus \cdot (K_c \Delta id)) \times id)) , \quad \text{DIST}$$

which, in terms of transformations of a spine $[x_1, \dots, x_n]^{\top}$ (with \top denoting the initial type induced by F), can be depicted as



In this sense, \wedge distributes over *false* $\nabla \vee$, \vee distributes over *true* $\nabla \wedge$, and \max_2 distributes over $\infty \nabla \min_2$, for example. \parallel

114 Lemma. Let $(\tau = z \nabla \otimes, \text{ld})$ denote an initial type in insert representation. Given that \oplus distributes over τ we can pull applications of $\oplus: E \times E \rightarrow E$ out of the scope of a monad comprehension if no variable bound in qs appears free in c (with $e, c: E$):

$$c \oplus \mathcal{M} \llbracket e \parallel qs \rrbracket^{\text{ld}} = \mathcal{M} \llbracket c \oplus e \parallel qs \rrbracket^{\text{ld}}. \quad \text{DESCOPE}$$

Proof. Following the proof scheme of Lemma 93, we assume $qs = v \leftarrow e': \top', qs'$ and first set $qs' = \varepsilon$. If $\top' = \text{ld}$ then **DESCOPE** reads $c \oplus \text{let } v = e' \text{ in } e = \text{let } v = e' \text{ in } c \oplus e$ which is true because v does not appear free in c . For $\top' \neq \text{ld}$ we calculate:

$$\begin{aligned} & c \oplus \mathcal{M} \llbracket e \parallel v \leftarrow e' \rrbracket^{\text{ld}} \\ = & \{ \mathcal{M}\text{-3} \} \\ & c \oplus ((\tau) (\mathcal{M} \llbracket e \parallel v \leftarrow e' \rrbracket^{\top'})) \\ = & \{ \text{DIST} \} \\ & (z \nabla \otimes \cdot (\oplus \cdot (K_c \Delta id) \times id)) (\mathcal{M} \llbracket e \parallel v \leftarrow e' \rrbracket^{\top'}) \\ = & \{ \text{ACID-RAIN} \} \\ & ((\tau) \cdot \top' (\oplus \cdot (K_c \Delta id))) \mathcal{M} \llbracket e \parallel v \leftarrow e' \rrbracket^{\top'} \\ = & \{ \text{FUNCTOR-}\mathcal{M}\text{-FUSION} \} \\ & (\tau) (\mathcal{M} \llbracket c \oplus e \parallel v \leftarrow e' \rrbracket^{\top'}) \\ = & \{ \mathcal{M}\text{-3} \} \\ & \mathcal{M} \llbracket c \oplus e \parallel v \leftarrow e' \rrbracket^{\text{ld}}. \end{aligned}$$

Induction step. Assume that **DIST** holds for $qs = qs'$.

$$\begin{aligned}
& \mathcal{M} \llbracket c \oplus e \parallel v \leftarrow e', qs' \rrbracket^{\text{ld}} \\
= & \{ \mathcal{M}\text{-5} \} \\
& \text{join}^{\text{ld}} (\mathcal{M} \llbracket \llbracket c \oplus e \parallel qs' \rrbracket^{\text{ld}} \parallel v \leftarrow e' \rrbracket^{\text{ld}}) \\
= & \{ \text{induction hypothesis} \} \\
& \mathcal{M} \llbracket c \oplus \llbracket e \parallel qs' \rrbracket^{\text{ld}} \parallel v \leftarrow e' \rrbracket^{\text{ld}} \\
= & \{ \text{induction base case} \} \\
& c \oplus \mathcal{M} \llbracket \llbracket e \parallel qs' \rrbracket^{\text{ld}} \parallel v \leftarrow e' \rrbracket^{\text{ld}} \\
= & \{ \mathcal{M}\text{-NORM-5} \} \\
& c \oplus \mathcal{M} \llbracket e \parallel v \leftarrow e', qs' \rrbracket^{\text{ld}}.
\end{aligned}$$

□

115 Example (continued from Paragraph 109). Now that we are equipped with these calculus-level rewriting rules, let us revisit query Q'' and the problem of correlated quantifiers. An appeal to **NEST-EX** and **DESCOPE** can unwind the nested **exists**-clauses which finally enables the derivation of an efficient semijoin plan for Q'' :

$$\begin{aligned}
& \mathcal{Q} Q'' \\
= & \{ \mathcal{Q} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket \llbracket qxz \wedge pyz \parallel z \leftarrow zs \rrbracket^{\text{exists}} \parallel y \leftarrow ys \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \text{NEST-EX} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket \llbracket qxz \wedge pyz \parallel y \leftarrow ys \rrbracket^{\text{exists}} \parallel z \leftarrow zs \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \text{DESCOPE} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxz \wedge \llbracket pyz \parallel y \leftarrow ys \rrbracket^{\text{exists}} \parallel z \leftarrow zs \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \llbracket p \wedge q \parallel qs \rrbracket^{\text{exists}} = \llbracket p \parallel qs, q \rrbracket^{\text{exists}} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxz \parallel z \leftarrow zs, \llbracket pyz \parallel y \leftarrow ys \rrbracket^{\text{exists}} \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
= & \{ \mathcal{M}\text{-SEMIJOIN} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxz \parallel z \leftarrow zs \quad \times \quad ys \rrbracket^{\text{exists}} \rrbracket^{\text{set}} \\
& \hspace{10em} \lambda z. \lambda y. pyz \\
= & \{ \text{SEMIJOIN} \} \\
& xs \times_q (zs \quad \times \quad ys) . \\
& \hspace{10em} \lambda z. \lambda y. pyz
\end{aligned}$$

Similar mechanical calculations derive the algebraic forms that have been proposed by Bry to efficiently process the *open nested quantified* queries listed in [15, Section 3.2]. ◇

To conclude this chapter on hybrid query rewriting, we briefly review the transformation of another class of queries that has been discussed by Claussen, Kemper, Moerkotte, and Peithner in [34].

116 Example. Queries of the form shown below are considered hard to translate for purely algebraic optimizers mainly because there exists no canonical translation for universal quantification in the `where`-clause:

$$Q''' = \begin{array}{l} \text{select } x \\ \text{from } xs \text{ as } x \\ \text{where forall } y \text{ in } z : qxy \\ \text{with } z = \left(\begin{array}{l} \text{select } y \\ \text{from } ys \text{ as } y \\ \text{where } py \end{array} \right) \end{array}$$

(Claussen et al. distinguish 16 instances in this query class, each instance corresponding to a subset of $\{x, y\}$ of free variables referenced in p and q . Here, we discuss instances in which (a) the universal quantifier, and (b) the quantifier and z are correlated with the outer query block as these were identified as the most challenging in [34].)

Possible algebraic equivalents involve *set difference*, *relational division*, or a combination of grouping and counting to implement the universal quantifier, *e. g.*,

$$xs \setminus (\pi \text{ outl}(xs \bowtie_{-q|\lambda x.\lambda y.(x,y)} (\sigma p ys)))$$

(symbol \setminus denotes set difference).

The derivation of these algebraic forms is tedious, however, and the resulting expressions are judged to be inefficient and too complex to be useful as input to subsequent rewriting phases [34, 88, 105]. The case is even more complex if z is not closed (*e. g.*, if the quantifier predicate py is replaced by pxy). This renders the use of a division combinator impossible.

The universal quantifier poses no particular problem for the monad comprehension calculus. By means of **Q-FORALL**, the universal quantifier is translated into a comprehension to be interpreted in monad `all`. If z is closed, Q''' already emits a perfect match for the rhs of **ANTIJOIN**:

$$\begin{aligned} & \mathcal{Q} Q''' \\ = & \{ \mathcal{Q} \} \\ & \llbracket x \llbracket x \leftarrow xs, \llbracket qxy \llbracket y \leftarrow \llbracket y \llbracket y \leftarrow ys, py \rrbracket^{\text{bag}} \rrbracket^{\text{all}} \rrbracket^{\text{bag}} \\ = & \{ \text{ANTIJOIN} \} \end{aligned}$$

$$\begin{aligned}
& xs \quad \bowtie \quad (\llbracket y \parallel y \leftarrow ys, py \rrbracket^{\text{bag}}) \\
& \lambda x. \lambda y. \neg q \ x \ y \\
= & \quad \{ \text{SELECT} \} \\
& xs \quad \bowtie \quad (\sigma \ p \ ys) . \\
& \lambda x. \lambda y. \neg q \ x \ y
\end{aligned}$$

Should z be correlated with the outer query (replace py by pxy), the derivation of an antijoin-based plan is immediate once p has been moved into the scope of the universal quantifier:

$$\begin{aligned}
& \mathcal{Q} Q''' \\
= & \quad \{ \mathcal{Q} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxy \parallel y \leftarrow \llbracket y \parallel y \leftarrow ys, pxy \rrbracket^{\text{bag}} \rrbracket^{\text{all}} \rrbracket^{\text{bag}} \\
= & \quad \{ \mathcal{M}\text{-NORM-3} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxy \parallel y \leftarrow ys, pxy \rrbracket^{\text{all}} \rrbracket^{\text{bag}} \\
= & \quad \{ \llbracket p \parallel qs, q \rrbracket^{\text{all}} = \llbracket p \vee \neg q \parallel qs \rrbracket^{\text{all}} \} \\
& \llbracket x \parallel x \leftarrow xs, \llbracket qxy \vee \neg pxy \parallel y \leftarrow ys \rrbracket^{\text{all}} \rrbracket^{\text{bag}} \\
= & \quad \{ \text{ANTIJOIN} \} \\
& xs \quad \bowtie \quad ys . \\
& \lambda x. \lambda y. \neg q \ x \ y \wedge p \ x \ y
\end{aligned}$$

The derived combinator queries are identical to those that have been identified as the most efficient forms by Claussen et al. [34]. \diamond

117 Consider a rewriting rule set consisting of (a) defining equations for the algebraic combinators (Definition 103) as well as their associated partial match rules (Paragraph 105), (b) monad comprehension normalization rules (Definition 91), and (c) the calculus-level equivalences of this section.

Thanks to the poor syntactic variety monad comprehensions exhibit, the patterns occurring in this rule set are *simple*. The generic nature (being parametric in the monads) of the rules renders this rule set *compact*. A quite simple heuristic ordering of the rules, sketched in the introduction to this section, has already led to promising rewriting results. We could significantly improve the quality of the rewriting process with only few rules which is mainly due to their *general applicability in many instantiations*. Should the derivation of a pure combinator query fail, this rule set comes with a guaranteed fallback, namely \mathcal{M} , that emits catamorphisms for those query parts that could not be mapped to combinators (these parts are subject to further optimization, see Chapter 6). In this sense, the rule set is *complete*.

We believe that such a comprehension of queries makes a promising foundation upon which to rest a query optimizer for OQL-like declarative languages.

Chapter 5

Comprehending Queries

118 Perhaps the most principle and influential decision in solving a problem is the choice of language in which we represent both the problem and its possible solutions. Choosing the “right” language can turn the concealed or difficult into the obvious or simple. This chapter revisits three problems in the advanced query processing domain. In all three cases, it is our aim to show how a catamorphic and monadic query language can (a) simplify, if not automate, the derivation of proposed solutions to the problem, (b) help to assess the correctness of these solutions, (c) possibly generalize the class of queries described by the problem and thus clarify the applicability of its solution.

The three problems were previously identified and tackled by others. The first two concentrate on the efficient processing of queries featuring grouping and aggregation ([Chatziantoniou and Ross](#), *Groupwise Processing of Relational Queries* [27] and [Cluet and Moerkotte](#), *Efficient Evaluation of Aggregates on Bulk Types* [36]), while the last is concerned with a pure combinator query representation and its normalization ([Cherniack, Malhotra, and Zdonik](#), *Experiences with Query Translation: Object Queries meet DB2*, [31]). We appreciate the work of these authors and the following is certainly not to lecture them. Rather, we like to argue that the techniques developed and applied in this text can contribute to the comprehension of queries.

5.1 Parallelizing Group-By Queries

119 The database backends of decision support systems frequently face SQL queries of the following general type (termed *group queries* in [26, 27]):

$$Q f g \text{ agg } xs \quad \equiv \quad \begin{array}{l} \text{select } f x, \text{ agg}(g x) \\ \text{from } xs \text{ as } x \\ \text{group by } f x . \end{array}$$

Group queries extract a particular dimension or feature—described by function f —from given base data xs and then pair each data point $f x$ in this dimension with aggregated data $\text{agg}(g x)$ associated with that point; agg may be instantiated by one of the SQL aggregate functions, *e. g.*, **sum** or **max**.

There exist a number of execution plans suitable to process the queries in class Q , one of the more efficient but maybe not so obvious variants being a *self- nestjoin* of xs with itself. To find this variant, start with the monad comprehension equivalent for Q (derivable by \mathcal{Q} if its definition is slightly modified to take care of syntactic sugar provided by SQL):

$$Q f g \text{ agg } xs \quad \equiv \quad \llbracket (f x, \llbracket g y \mid y \leftarrow xs, f y = f x \rrbracket^{\text{agg}}) \mid x \leftarrow xs \rrbracket^{\text{set}}$$

(let **agg** denote the **Id**-monad associated with the non-collection algebra implementing the SQL aggregate agg , *e. g.*, **sum** or **max**). Then perform combinator pattern matching as we have proposed in Chapter 4 (we assume $xs : \text{bag } E$ for some E):

$$\begin{aligned} & \llbracket (f x, \llbracket g y \mid y \leftarrow xs, f y = f x \rrbracket^{\text{agg}}) \mid x \leftarrow xs \rrbracket^{\text{set}} \\ = & \quad \{ \text{M-AGGREGATE} \} \\ & \llbracket (f x, \alpha^{\text{agg}} g \llbracket y \mid y \leftarrow xs, f y = f x \rrbracket^{\text{bag}}) \mid x \leftarrow xs \rrbracket^{\text{set}} \\ = & \quad \{ \text{product} \} \\ & \llbracket (f \times \alpha^{\text{agg}} g) (x, \llbracket y \mid y \leftarrow xs, f y = f x \rrbracket^{\text{bag}}) \mid x \leftarrow xs \rrbracket^{\text{set}} \\ = & \quad \{ \text{M-PROJECT} \} \\ & \pi (f \times \alpha^{\text{agg}} g) \llbracket (x, \llbracket y \mid y \leftarrow xs, f y = f x \rrbracket^{\text{bag}}) \mid x \leftarrow xs \rrbracket^{\text{set}} \\ = & \quad \{ \text{NESTJOIN} \} \\ & \pi (f \times \alpha^{\text{agg}} g) (xs \quad \triangle \quad xs) . \\ & \quad \quad \quad \lambda x. \lambda y. f x = f y \mid \lambda x. \lambda y. (x, y) \end{aligned}$$

120 In [26, 27], Chatziantoniou and Ross propose to take a rather different three-step route to execute queries of type Q :

- (a) Separate the data points in dimension f of xs in a preprocessing step, *i. e.*, partition xs with respect to f .
- (b) Evaluate a simplified variant Q' of Q on each partition. In particular, Q' does not need to take care of grouping. Let ps denote a partition of xs , then we have

$$Q' f g \text{ agg } ps \quad \equiv \quad \begin{array}{l} \text{select } f x, \text{ agg}(g x) \\ \text{from } ps \text{ as } x, \end{array}$$

or, equivalently,

$$Q' f g \text{ agg } ps \quad \equiv \quad \llbracket (f x, \llbracket g y \rrbracket y \leftarrow ps \rrbracket^{\text{agg}}) \rrbracket x \leftarrow ps \rrbracket^{\text{set}}.$$

- (c) Finally, join the results obtained in step (b) to form the query response.

This strategy clearly shows its benefit in stage (b): first, since xs has been split into disjoint partitions during the preprocessing step, we may execute Q' on the different partitions in parallel. Second, there is a chance to process Q' in main memory should the partitions of xs fit. Measurements reproduced in [26, 27] show the performance gains in terms of time and I/O costs to compensate for the effort spent in the partitioning and joining stages.

121 In [26, 27], classical relational algebra is the language at which the translation of group queries is targeted. This choice of query representation introduces subtleties.

Relational algebra lacks canonical forms to express the grouping and aggregation operations found in Q . Chatziantoniou and Ross thus propose to understand Q as a *syntactical* query class: the membership of a specific query in this class and thus the applicability of the partitioning strategy is decided by inspection of the SQL parse tree for that query. (Actually, query graphs are extracted from the parse trees and a query graph level criterion decides the group query property.)

Relational algebra falls short to provide idioms that could express the preprocessing, *i. e.*, partitioning, step of the strategy. To remedy the situation, Chatziantoniou and Ross attribute the nodes of the query graphs to indicate which partition is represented by a specific node.

Finally, the core stage (b) of the partitioning strategy has no equivalent at the target language level as well. Classical relational algebra is unable to express the iteration (or parallel application) inherent to this stage.¹

¹Chatziantoniou and Ross implemented this step *on top* of the relational backend and thus outside the relational domain.

Facing this mix of query representations (SQL syntax, query graphs, relational algebra, iteration), it is considerably hard to assess the correctness of [Chatziantoniou and Ross](#)' parallel processing strategy for class Q .

122 Reasoning in the monad comprehension calculus and its associated combinator algebra can significantly simplify the matter. Once expressed using monadic notions, we can construct a correctness proof for the strategy which is basically built from the unfolding of definitions and normalization steps alone (and thus could possibly be conducted automatically, although we consider automated proofs outside the scope of this text). Let us proceed by filling the two gaps (partitioning and iteration) that relational algebra has left open.

First, partitioning the base data collection xs with respect to a function f is expressible in the monad comprehension calculus simply as (let \mathbb{T} denote a collection monad and E' some equality type, *cf.* Paragraph 78)

$$\begin{aligned} \text{partition} & : (E \rightarrow E') \rightarrow \mathbb{T}E \rightarrow \text{set } \mathbb{T}E \\ \text{partition } f \text{ } xs & = \llbracket \llbracket y \rrbracket y \leftarrow xs, f x = f y \rrbracket^{\mathbb{T}} \llbracket x \leftarrow xs \rrbracket^{\text{set}}, \quad \text{PARTITION} \end{aligned}$$

which builds a set of disjunct partitions so that all elements inside one partition agree on feature f .

Second, recall from Section 2.5 that the type or map functors provide the principle iteration abstractions in our categorical setting. At the combinator algebra level, iteration is at our disposal by means of the combinator $\pi f = \mathbb{T}f$ (see **PROJECT**). Map functor $\mathbb{T}f$ also adequately encodes parallel application of f to the elements of its (collection-typed) argument: the applications of f to the spine elements do not interfere and can be evaluated in parallel. See, for example, the work of [Hill](#) in which a complete theory of *data-parallel programming* is developed on top of map functors [62].

With the definition of Q' from Paragraph 120 we can now recast the parallel grouping plan as

$$\mu (\pi (Q' f g \text{agg}) (\text{partition } f \text{ } xs)) .$$

The following rewrite finally establishes the equivalence of this plan and $Q f g \text{agg } xs$ and thus provides a purely calculational proof of the correctness of parallel grouping.

$$\begin{aligned} & \mu (\pi (Q' f g \text{agg}) (\text{partition } f \text{ } xs)) \\ = & \quad \{ \text{UNNEST} \} \\ & \llbracket y \rrbracket ys \leftarrow \pi (Q' f g \text{agg}) (\text{partition } f \text{ } xs), y \leftarrow ys \rrbracket^{\text{set}} \end{aligned}$$

we assume $\otimes = \leq$):

$$Q f \otimes \text{agg } xs \ ys \equiv \begin{array}{l} \text{select } f x (\text{agg } z) \\ \text{from } xs \text{ as } x \\ \\ \text{with } z = \left(\begin{array}{l} \text{select } y \\ \text{from } ys \text{ as } y \\ \text{where } y \otimes x \end{array} \right) . \end{array}$$

To be more precise, assume $xs : \top E$, $ys : \top' E$, $f : E \rightarrow E \rightarrow E'$, $\otimes : E \times E \rightarrow \text{Bool}$ and let agg denote the identity monad realizing agg (see Section 5.1). Q is implemented by the combinator form

$$\pi (\lambda x. f x (\alpha^{\text{agg}} (\sigma (\lambda y. y \otimes x) ys)) xs) .$$

We always have the naive catamorphic interpretation at hand but [Kim](#) already described a superior evaluation strategy for Q in [69]: (a) join collections xs and ys with respect to \otimes , then (b) group the join result, and finally (c) aggregate the groups separately to form the overall query result.

Although an improvement, given that the join in step (a) in general will *not* be an equi-join of xs and ys , we need $O(|xs| \cdot |ys|)$ space to hold the intermediate result. The time complexity of this strategy is clearly the same.

124 [Cluet and Moerkotte](#) realized that an optimizer can reduce both the space and time needed to evaluate queries in Q , provided that the aggregate agg exhibits a *decomposition* property [36]. Importing the work presented in [36] into our system of categorical type constructors is straightforward.

For $ys : \top' E$, the aggregate function $\text{agg} : \top' E \rightarrow E$ is *decomposable*, if we can find functions $\alpha : \top' E \rightarrow E$ and $\oplus : E \rightarrow E \rightarrow E$ so that

$$\text{agg } ys = (\alpha \ ys') \oplus (\alpha \ ys'') \quad \text{with} \quad ys = \text{union}^{\top'} (ys', ys'') .$$

Note that this almost characterizes agg as a homomorphism of algebras in union representation (*cf.* Section 2.10).

The decomposition property of aggregate agg paves the way for an efficient but non-standard evaluation strategy for Q : since agg is decomposable we can *incrementally update* the aggregate value for each group every time we encounter a member of that group [52]. Each group's current aggregate state is kept in a support data structure—the \otimes -table—which forms the core of the algorithm.

The following pseudo-code sketches the implementation of this procedure:

```

new( $t, \otimes$ );
foreach  $x$  in  $xs$ 
    ins( $t, x$ );
sort( $t$ );
foreach  $y$  in  $ys$ 
    update( $t, y$ );
eval( $t$ );

```

$new(t, \otimes)$: Allocate an empty \otimes -table t of type $list(Num \times E \times E)$. For each list member, the first tuple component protocols the multiplicity of element $x:E$ in xs ; x is stored in the second component; the last component reflects the aggregate state of the group x is associated with.

$ins(t, x)$: $t := (1, x, \alpha \text{ empty}^{\top'}) \text{ list } t$.

$sort(t)$: Sort list t on second component with respect to \otimes , remove duplicates but update multiplicity and aggregate components accordingly.

$update(t, y)$: Incremental update: search for tuple (n, x, a) with smallest x such that $y \otimes x$. Replace this tuple in t with $(n, x, a \oplus (\alpha (sng^{\top'} y)))$.

```

eval( $t$ ):  $a' := \alpha \text{ empty}^{\top'}$ ;
foreach  $(n, x, a)$  in  $t$  do
     $a' := a \oplus a'$ ;
    replace tuple by  $(n, x, a')$ ;
od;

```

Note that the actual groups are never constructed. The \otimes -table thus has space requirements of $O(|xs|)$ as there are no more groups than elements in xs . Furthermore, should \otimes impose an order on E , we can exploit binary search to implement the *update* operation on the table. This reduces the time complexity to $O(|xs| \cdot \log |xs| + |ys| \cdot \log |xs|)$.

125 There remains the question of the *effective* applicability of the \otimes -table approach. Among other concluding remarks in [36], Cluet and Moerkotte marked two open issues concerning effectiveness:

- (a) How to enable the optimizer to automatically deduce functions α and \oplus for a given query, specifically for a given aggregate *agg*?

- (b) Clearly, there is a larger class of queries which can benefit from \otimes -table-driven processing. A \otimes -table can support the evaluation of non-equi-joins like $\bowtie(\lambda x.\lambda y.x \otimes y) f xs ys$ and may additionally prove efficient during nestjoin processing (*cf.* query class Q' of Paragraph 108). A syntactical characterization as it is provided by Q seems to be too coarse.

We can provide effective answers to both questions once we have re-expressed Q via our catamorphic (monadic) comprehension of queries.

126 There is an obvious connection between the almost homomorphic nature of a decomposable aggregate function and its catamorphic realization in our world of categorical query semantics. This observation provides an answer to issue (a) raised in the previous paragraph.

Application of mapping $\mathcal{M} \cdot \mathcal{Q}$ makes the catamorphic implementation of Q accessible (let $\tau = e \nabla \otimes : 1 + E \times E \rightarrow E$ denote the algebra in insert representation that is associated with the identity monad \mathbf{agg} , see Definition 64):

$$\begin{aligned}
& (\mathcal{M} \cdot \mathcal{Q})(Q f \otimes \mathbf{agg} xs ys) \\
&= \{ \mathcal{Q} \} \\
& \quad \mathcal{M} \llbracket f x \llbracket y \llbracket y \leftarrow ys, y \otimes x \rrbracket^{\mathbf{agg}} \llbracket x \leftarrow xs \rrbracket^{\top} \\
&= \{ \mathcal{M}\text{-3} \} \\
& \quad \llbracket f x (\llbracket \tau \rrbracket \llbracket y \llbracket y \leftarrow ys, y \otimes x \rrbracket^{\top'}) \llbracket x \leftarrow xs \rrbracket^{\top} ,
\end{aligned}$$

i. e., we have $\mathbf{agg} = (\llbracket \tau \rrbracket) : \top' E \rightarrow E$.

Now note that the \otimes -table algorithm exclusively performs element-wise insertions of elements into the \otimes -table (by means of function *update*) which, in turn, lead to element-wise updates of the aggregate state. This suggests to take an insert representation view of the decomposition property: aggregate function $\mathbf{agg} : \top' E \rightarrow E$ is *decomposable*, if we can find functions $\alpha : \top' E \rightarrow E$ and $\oplus : E \rightarrow E \rightarrow E$ so that

$$\mathbf{agg} ys = y \oplus (\alpha ys') \quad \text{with} \quad ys = y^{\top'} : ys' .$$

Being so far, we can almost immediately read off suitable definitions for α and \oplus :

$$\begin{aligned}
& \mathbf{agg} ys \\
&= \{ ys = y^{\top'} : ys', \mathbf{agg} = (\llbracket \tau \rrbracket) \} \\
& \quad (\llbracket \tau \rrbracket) (y^{\top'} : ys') \\
&= \{ \tau = e \nabla \otimes, \mathbf{CATA}\text{-}\mathbf{INS}\text{-}\mathbf{REP} \}
\end{aligned}$$

$$\begin{aligned}
& y \otimes ((e \nabla \otimes) ys') \\
\Rightarrow & \{ \text{aggregate decomposition} \} \\
& \alpha = (e \nabla \otimes) = (\tau) \\
\wedge \oplus & = \otimes .
\end{aligned}$$

127 To tackle the effective detection of (sub)queries that are subject to \otimes -table-based evaluation, observe that the \otimes -table algorithm is obviously capable of processing queries of the general form

$$\llbracket f x \llbracket y \llbracket y \leftarrow ys, y \otimes x \rrbracket^{\mathsf{T}'} \llbracket x \leftarrow xs \rrbracket^{\mathsf{T}} .$$

The monad comprehension form provides a *structural* characterization of the \otimes -table processible queries. This structural encoding of the query class is more effective and *useful*² than the simple syntactical query pattern Q . The monad comprehension form covers all application scenarios envisioned in [36], including those that were understood as possible future extensions:

(a) (Sub)queries in class Q are detected since $\mathcal{Q}Q$ is an instance of the above monad comprehension form (with $\mathsf{T}' = \mathbf{agg}$, see the previous paragraph).

(b) For the particular queries in which monad T' is instantiated by a collection monad (as opposed to an \mathbf{Id} -monad), we obtain a perfect match with the defining pattern for Δ (**NESTJOIN**), specifically $\Delta \otimes f xs ys$. In fact, this is a rediscovery of a remark in [36] in which the possible applicability of \otimes -tables to the efficient implementation of nestjoin-like operators has been outlined.

(c) Finally, the cases in which T' denotes a collection monad are not special with respect to an automatic derivation of α and \oplus . To see this, simply repeat the calculation of the previous paragraph (with $\tau = \llbracket \rrbracket^{\mathsf{T}'} \nabla \ddagger^{\mathsf{T}'}$):

$$\begin{aligned}
\alpha & = (\llbracket \rrbracket^{\mathsf{T}'} \nabla \ddagger^{\mathsf{T}'}) = id \quad (\mathbf{CATA-REFLECT}) \\
\wedge \oplus & = \ddagger^{\mathsf{T}'}.
\end{aligned}$$

These choices for α and \oplus exactly configure the \otimes -table algorithm for the processing of non-equi-joins and nestjoins as it has been proposed in [36].

²Here, “useful” is used in the sense of Cluet and Moerkotte in [36]: “[Further research should] come up with a useful characterization of the corresponding queries. By useful we mean that an optimization can easily detect the applicability of θ -tables.” (In [36], Cluet and Moerkotte use symbol θ in place of \otimes .)

Note that the above insights are completely based on rather simple pattern matches and unifications with the monad comprehension-based encoding of the query class. It is this encoding which enables a query optimizer to effectively exploit the \otimes -table idea.

5.3 Normalization of KOLA Combinator Queries

In a series of articles, [Cherniack and Zdonik](#) developed the KOLA query algebra as a target algebra for the translation and optimization of OQL queries [28, 29, 30, 31, 32, 33]. Being a successor to the AQUA algebra [108, 109], KOLA follows a design that emphasizes the use of few but generic higher-order query primitives. All KOLA primitives are combinators: KOLA has no notion of variables, scope, and binding, whatsoever. In this sense, KOLA exhibits similarities with the combinator algebra of Chapter 4.

Mainly because the meaning of variables is strongly tied to their scope, [Cherniack and Zdonik](#) argue that programs expressed in variable-based query languages are hard to optimize: the analysis of subqueries is inherently tightly coupled with the inspection of the global variable environment inherited from the enclosing query. There is no true locality of query analysis. This, in turn, complicates the formalization of query transformations and rewriting rules. The applicability of a rule is not only dependent on a structural (syntactical) match with its lhs but, in general, additionally on the current variable environment.

KOLA queries, lacking variables, enable pure *syntactical* rule matching and rewriting. The detection of a query pattern requires local analysis only. In *any* expression of the form (`join` denotes KOLA's θ -join combinator; the predicate is evaluated against pairs built from the elements from collections xs and ys)

$$\text{join}((p \cdot \text{outl}) \wedge q) \text{ } xs \text{ } ys$$

p may be pushed down the expression tree for xs since p is applied after the projection *outl* and thus cannot depend on the second join argument ys (an assurance justified by p 's nature of being a combinator). Here, predicate p 's independence of ys is expressed structurally rather than by (the absence of) variable references. [Cherniack and Zdonik](#) exploited this feature of KOLA to prove a large library of rewriting rules correct with the help of an automated theorem prover [29, 30].

128 Let us emphasize one consequence of the combinator approach. KOLA uses nested pairs to represent tuples and structured objects much like we did

in the preparations to the discussion of mapping \mathcal{Q} in Paragraph 78. A typical KOLA query is swamped with primitives that build and destruct pairs as there are no means to reference a specific tuple or object component besides an explicit sequence of projections on pairs. Similar remarks apply to KOLA rewrite rules. The rule sets that implement a specific global query rewrite—like the KOLA normalization we will examine in detail below—tend to be large: a rewrite is not only concerned with the global transformation goal but also needs to appeal to a variety of auxiliary rules which are primarily designed to, *e. g.*, manipulate nested pairs.

The large number of rules calls for sophisticated rewriting strategies to prevent a KOLA optimizer from getting lost in the vast space of equivalent expressions it can generate from these rule sets. This observation led [Cherniack and Zdonik](#) to the design and implementation of a control language for rule application, KOKO [32, 33].

While there are clear benefits of combinator query languages (*cf.* Paragraph 99), we believe that the combinator paradigm considerably complicates the implementation of certain query translation stages. To exemplify, we will try to comprehend a complex KOLA query rewrite [31] in terms of the monad comprehension calculus. A change in the point of view on query representation can have significant effects on the complexity of the specification and implementation of optimization goals. This is the point we are trying to make before we conclude the current chapter.

129 In [31], [Cherniack et al.](#) describe an OQL frontend which is designed to run on top of the relational DBMS DB2. OQL queries are translated into KOLA, then optimized. A subsequent stage generates DB2 SQL statements from KOLA queries. For the OQL `select-from-where` block

```
select f(x1, ..., xn)
  from e1 as x1, ..., en as xn
 where p(x1, ..., xn),
```

the frontend emits the KOLA form

$$((\text{iterate } p f) \cdot (\text{unnest id } (K f_n)) \cdot \dots \cdot (\text{unnest id } (K f_1)) \cdot \text{singleton}) ()$$

in which the combinators and symbols are defined to mean (here, we recast the combinators as monadic functions to give their precise semantics as well as to prepare the things to come)

$$\text{unnest } f g xs = \llbracket f(x, y) \mid x \leftarrow xs, y \leftarrow g x \rrbracket^{\text{bag}}$$

$$\begin{aligned}
\text{iterate } p f xs &= \llbracket f x \mid x \leftarrow xs, p x \rrbracket^{\text{bag}} \\
\text{singleton} &= \text{unit}^{\text{bag}} \\
\mathbb{K} &= K \\
() &= \text{an arbitrary dummy value} \\
f_i &= \text{KOLA form to represent } e_i.
\end{aligned}$$

This translation scheme is applicable to arbitrary e_i (*e.g.*, path expression, subqueries) and does not restrict the e_i to be constant collections (in the KOLA query above, the f_i may be arbitrarily complex KOLA functions).

The goal of the subsequent optimization stage is to derive a join query from this initial KOLA form.³ Cherniack et al. proceed in two steps to achieve this: (a) *normalize* the initial form, *i. e.*, rewrite it into an equivalent query in which a chain of **unnest** combinators is invoked on the f_i directly, thus rendering the auxiliary construction involving the dummy value $()$ void, and (b) derive a join query from the normalized expression which is more amenable to join detection as the initial form.

The KOLA implementation of this rewrite is too complex to reproduce here as it falls back on 46 rewrite rules and about 100 lines of rule firing control code [31].

130 A monad comprehension semantics for KOLA, in contrast, provides a straightforward, almost mechanical, derivation of a join query from the initial KOLA form. There is nothing more involved than the unfolding of combinator definitions given in the previous paragraph and the application of monad comprehension normalization rules from Definition 91. The simplicity of this process makes external guidance through specific rule firing control code obsolete.

The appeal to comprehension normalization comes at no surprise once we analyze the involved KOLA rewriting rules a bit deeper. A significant ratio of these rules actually implement specific calculus normalization steps. The KOLA rule (with $\text{shr}((x, y), z) = (x, (y, z))$; note that **shr** is one of the two arrows that establish the isomorphism between nested products, *cf.* Paragraph 78)

$$\begin{aligned}
&(\text{unnest}(\text{outr} \cdot \text{shr})(f \cdot \text{outr})) \cdot (\text{unnest } id(\mathbb{K} e)) \cdot \text{singleton} \rightsquigarrow \\
&\mathbb{K}(\text{unnest } id f e),
\end{aligned}$$

³This join query forms the input to the last stage which finally emits code to implement the OQL **select-from-where** block using DB2 SQL.

for example, bundles the application of normalization rules \mathcal{M} -NORM-2 and \mathcal{M} -NORM-3:

$$\begin{aligned}
& ((\text{unnest } (\text{outr} \cdot \text{shr}) (f \cdot \text{outr})) \cdot (\text{unnest } \text{id } (\mathbb{K} e)) \cdot \text{singleton}) e' \\
= & \quad \{ \text{unfold KOLA combinators} \} \\
& \llbracket (\text{outr} \cdot \text{shr}) (x, y) \llbracket x \leftarrow \llbracket (x', y') \llbracket x' \leftarrow \text{unit}^{\text{bag}} e', y' \leftarrow K_e x' \rrbracket^{\text{bag}}, \\
& \quad y \leftarrow (f \cdot \text{outr}) x \rrbracket^{\text{bag}} \\
= & \quad \{ \mathcal{M}\text{-NORM-2} \} \\
& \llbracket (\text{outr} \cdot \text{shr}) (x, y) \llbracket x \leftarrow \llbracket (e', y') \llbracket y' \leftarrow e \rrbracket^{\text{bag}}, y \leftarrow (f \cdot \text{outr}) x \rrbracket^{\text{bag}} \\
= & \quad \{ \mathcal{M}\text{-NORM-3} \} \\
& \llbracket (\text{outr} \cdot \text{shr}) ((e', y'), y) \llbracket y' \leftarrow e, y \leftarrow (f \cdot \text{outr}) (e', y') \rrbracket^{\text{bag}} \\
= & \quad \{ \text{outr, shr} \} \\
& \llbracket (y', y) \llbracket y' \leftarrow e, y \leftarrow f y' \rrbracket^{\text{bag}} \\
= & \quad \{ \text{unnest, K} \} \\
& \mathbb{K} (\text{unnest } \text{id } f e) e'
\end{aligned}$$

(from which the rule follows immediately by the principle of extensionality).

The normalization of the KOLA form goes off as a canonical monad comprehension normalization once we have unfolded the monadic KOLA combinator definitions:

$$\begin{aligned}
& ((\text{iterate } p f) \cdot (\text{unnest } \text{id } (\mathbb{K} f_n)) \cdot \dots \cdot (\text{unnest } \text{id } (\mathbb{K} f_1)) \cdot \\
& \quad \text{singleton}) () \\
= & \quad \{ \text{unfold singleton, unnest} \} \\
& ((\text{iterate } p f) \cdot (\text{unnest } \text{id } (\mathbb{K} f_n)) \cdot \dots \cdot (\text{unnest } \text{id } (\mathbb{K} f_2))) \\
& \quad \llbracket (x_1, y_1) \llbracket x_1 \leftarrow \text{unit}^{\text{bag}} (), y_1 \leftarrow K_{f_1} x_1 \rrbracket^{\text{bag}} \\
= & \quad \{ \mathcal{M}\text{-NORM-2} \} \\
& ((\text{iterate } p f) \cdot (\text{unnest } \text{id } (\mathbb{K} f_n)) \cdot \dots \cdot (\text{unnest } \text{id } (\mathbb{K} f_2))) \\
& \quad \llbracket ((), y_1) \llbracket y_1 \leftarrow f_1 \rrbracket^{\text{bag}} \\
= & \quad \{ \text{unfold unnest} \} \\
& ((\text{iterate } p f) \cdot (\text{unnest } \text{id } (\mathbb{K} f_n)) \cdot \dots \cdot (\text{unnest } \text{id } (\mathbb{K} f_3))) \\
& \quad \llbracket (x_2, y_2) \llbracket x_2 \leftarrow \llbracket ((), y_1) \llbracket y_1 \leftarrow f_1 \rrbracket^{\text{bag}}, y_2 \leftarrow K_{f_2} x_2 \rrbracket^{\text{bag}} \\
= & \quad \{ \mathcal{M}\text{-NORM-3} \} \\
& ((\text{iterate } p f) \cdot (\text{unnest } \text{id } (\mathbb{K} f_n)) \cdot \dots \cdot (\text{unnest } \text{id } (\mathbb{K} f_3))) \\
& \quad \llbracket ((((), y_1), y_2) \llbracket y_1 \leftarrow f_1, y_2 \leftarrow f_2 \rrbracket^{\text{bag}} \\
= & \quad \{ \text{repeat last two steps } n - 3 \text{ times} \} \\
& (\text{iterate } p f) \\
& \quad \llbracket (((\dots ((), y_1) \dots), y_n) \llbracket y_1 \leftarrow f_1, \dots, y_n \leftarrow f_n \rrbracket^{\text{bag}}
\end{aligned}$$

$$= \{ \text{isomorphism of nested products, iterate} \} \\ \llbracket f(y_1, \dots, y_n) \llbracket y_1 \leftarrow f_1, \dots, y_n \leftarrow f_n, p(y_1, \dots, y_n) \rrbracket^{\text{bag}} \rrbracket .$$

The generation of an equivalent KOLA `join` query is now immediate by means of the combinator patterns `JOIN` and `M-JOIN`.

Note, in addition, that the normalized query is exactly what mapping \mathcal{Q} would directly emit for the OQL `select-from-where` block shown in the previous paragraph. We believe that a pure combinator representation renders the early stages of query processing—namely translation and normalization—an unnecessarily hard task. An employment of the monad comprehension calculus clearly shows its elegance here.

Chapter 6

Query Deforestation

131 Recall that a combinator query exhibits a tree-like structure in which the nodes represent abstractions of the spine transformers (physical operators) supplied by the query engine backend (*cf.* the introduction to Chapter 4). It is a demanding task to efficiently schedule the flow of data from the leaves of this processing tree—*i. e.*, the complex values which are actually materialized on persistent storage—to its root which represents the query result.

Query engine models that assign combinators to separate operating system processes or threads connected by inter-process communication (IPC) facilities (to control the flow of execution) as well as disk files (to communicate temporary results) have been found to be inefficient: the operating system scheduler, necessary process context switches, IPC, and I/O of temporary results incur an overhead that dominates the overall query cost by far [53, 54].

Different lines of research [46, 47, 78] thus led to the development of single-process query engines. A combinator query is compiled into a single *iterative* or *recursive procedure* which is then executed inside the monolithic query engine process. Whenever possible, these approaches strive for a *stream-based* or *pipelined* query execution to avoid the I/O of temporary data. The query engine benefits from streaming since data is addressed and brought in from persistent storage only once. Further processing is in-memory and no intermediate writes to the persistent store and subsequent reads of materialized intermediate results occur.

132 The derivation of a stream-based query program from a combinator expression is a problem that has been primarily tackled on the implementation level only. In a sense, the forthcoming material drags this query processing phase out of the hands of the database implementors and instead aims at

making the matter accessible to our categorical tools—and thus comprehensible by an optimizer.

6.1 Streaming

133 Given that the query compiler arranges the query plan as a tree structure of independent combinators, it is not immediate how an efficient streaming procedure may be automatically derived. Combinators consume their input spine as a whole, transform it, and then emit the transformed spine as an intermediate result. We are not given any reference to “reach inside” and collect partial results early in order to immediately pass them to the parent combinator. The possibility to do so would be the key to realize streaming but violates the combinator encapsulation.

These observations led to a spectrum of solutions to this dilemma, ranging from implementation level modifications to the query engine to algebraic transformation of execution plans.

134 In his seminal work on the engineering of query engines for large databases, Graefe proposed an iterator *implementation discipline* for combinators that facilitates streaming [53, 54]. Any combinator implements an interface specifying a function *next* and a constant *empty!*. Application of this iterator discipline does not alter the initial combinator tree but leads the combinators to schedule each other inside the tree structure.

Given a combinator tree, query evaluation is driven by the root combinator *on demand*: repeated calls to *next* request the production of the next stream element (*e.g.*, tuple, complex value, or object) by the combinator’s child node(s). The children recursively forward the *next* call until a leaf combinator can satisfy the request through read access to the data source it is associated with. Each combinator passes stream elements upwards the processing tree, possibly after application of a combinator-specific action to the elements. Combinators iterate this procedure until the special stream element *empty!* signals that the child nodes are exhausted.

Following this discipline, a sketch of the implementation of function *next* for the πf combinator (see **PROJECT**) would thus read (with $empty?x \Leftrightarrow x = empty!$):

```

x := next;
if (not (empty?x)) then
  return (f x);
return empty!;
```

135 The on-demand scheduling of combinators to thread the next stream element (the head of the stream’s tail) through the combinator tree bears close resemblance with the *lazy semantics of evaluation* in functional programming systems [92]. Buneman, Frankel, and Rishiyur actually devised a combinator query language with lazy semantics, FQL, in [18].

A lazy stream xs is represented in its *weak head normal form* [92] $x \parallel s$ in which x denotes the stream’s head element and s represents a *suspension*, *i. e.*, a function that—if ever evaluated—returns the stream $x' \parallel s'$, the tail of xs . A suspension, or *closure*, $\ulcorner f, x \urcorner$ bundles a function f and its argument x . Function f is not applied until the suspension is actually evaluated (*forced*):

$$\text{eval } \ulcorner f, x \urcorner = f x .$$

Iterating over a stream then means to repeatedly evaluate its suspended tail.

In continuation from the previous paragraph, we could implement the π combinator in the suspension model as

$$\begin{aligned} \pi f xs &= \ulcorner f, x \urcorner \parallel \ulcorner \pi f, \text{eval } s \urcorner \\ &\text{with } xs = x \parallel s . \end{aligned}$$

Note that π does not yet apply f to the stream’s head x but rather creates a suspension $\ulcorner f, x \urcorner$. This saves the cost for the application of f to x should the stream element be discarded later on.

Suspensions are, once again, a device that postpones the solution of the streaming problem until query engine implementation. It is the point of this chapter to show that we can benefit from pulling the issue up to the query representation level. At this level, streaming programs are derived by rewriting much like in a transformation-based query optimizer. Let us proceed with a short review of a transformational approach due to Freytag and Goodman [44, 45, 46, 47] before the next section discusses a remarkably simple method to achieve streaming in the catamorphic query model.

136 Guided by observations about the advantages of single-process query engines, Freytag and Goodman set out with the principal goal to transform combinator queries into monolithic *iterative* procedures [44, 45, 46, 47]. The emitted code was designed to be *compiled* and then linked against a relational database backend.

During the derivation of an iterative program form, Freytag and Goodman break the encapsulation inherent to the combinator notion (*cf.* Paragraph 99): given a combinator query Q , combinators in Q are *unfolded*, *i. e.*, replaced by their defining expressions—recursive programs expressed in a subset of LISP—with the principle aim to fuse these with the definitions of

neighboring combinators. To continue our treatment of the π combinator in different streaming models, unfolding $\pi f xs$ would reveal its *control structure* and *action* as the LISP form on the rhs

$$\begin{aligned} \pi f xs &= (if (empty? xs) \\ &\quad empty! \\ &\quad (cons (f (hd xs)) (\pi f (tl xs)))) . \end{aligned}$$

Actual fusion is then performed through simplification of the unfolded LISP forms (all occurring forms are restricted to use a referentially transparent subset of LISP thus facilitating the specification of equivalence-preserving simplification rules). The rules aim to transform the program for Q so that occurrences of the defining expression for Q can be detected. On detection, these subexpressions are replaced by recursive calls to Q itself (*folding*). In a final step, recursion is traded for iteration.

To illustrate, consider the following sketch of the *unfold-simplify-fold* steps for the combinator query $Q f g xs = \pi g (\pi f xs)$:

$$\begin{aligned} &Q f g xs \\ = &\quad \{ \text{unfold defining query for } Q \} \\ &(\pi g (\pi f xs)) \\ = &\quad \{ \text{unfold } \pi g \} \\ &(if (empty? xs) \\ &\quad empty! \\ &\quad (cons (g (hd (\pi f xs))) \\ &\quad\quad (\pi g (tl (\pi f xs))))) \\ = &\quad \{ \text{unfold } \pi f \text{ once} \} \\ &(if (empty? xs) \\ &\quad empty! \\ &\quad (cons (g (hd (if (empty? xs) \\ &\quad\quad empty! \\ &\quad\quad (cons (f (hd xs)) (\pi f (tl xs))))) \\ &\quad\quad (\pi g (tl (\pi f xs))))) \\ = &\quad \{ \text{further unfolding, simplification steps} \} \\ &(if (empty? xs) \\ &\quad empty! \\ &\quad (cons (g (f (hd xs))) \\ &\quad\quad (\underbrace{\pi g (\pi f (tl xs))}_{\text{marked occurrence}}))) \\ = &\quad \{ \text{fold marked occurrence of } Q f g (tl xs) \} \end{aligned}$$

$$\begin{aligned}
 & (if (empty? xs) \\
 & \quad empty! \\
 & \quad (cons (g (f (hd xs))) \\
 & \quad \quad (Q f g (tl xs)))) .
 \end{aligned}$$

Note that the simplified program for Q is equivalent to the combinator query $\pi (g \cdot f) xs$ which, unlike the original query, does not generate any intermediate result.

137 This approach comes with apparent difficulties.

(a) Combinators are encoded by recursive LISP forms that intertwine control structure *and* action. Consequently, simplification rules are not only concerned with the optimization of stream operations like, *e. g.*,

$$(tl (cons e e')) \rightsquigarrow e'$$

but also modify control structure as in

$$(if p e (if p e' e'')) \rightsquigarrow (if p e e'') .$$

The resulting derivations tend to be long, leading to significant transformation effort [44]. We could clearly benefit from a separation of control and action.

(b) As an instance of [Burstall and Darlington](#)'s general *unfold-transform-fold* program transformation strategy [21], the approach suffers from the need for a *fold* step. The transformation algorithm needs a “memory” to spot previously seen expressions in order to be able to successfully “*tie the knot*”. We also have to find an upper bound on the number of unfolding steps to perform. Otherwise, unfolding might never stop for the cases in which we fail to spot expressions to recur. The resulting complexity of the approach has prevented its inclusion into actual query optimizers.

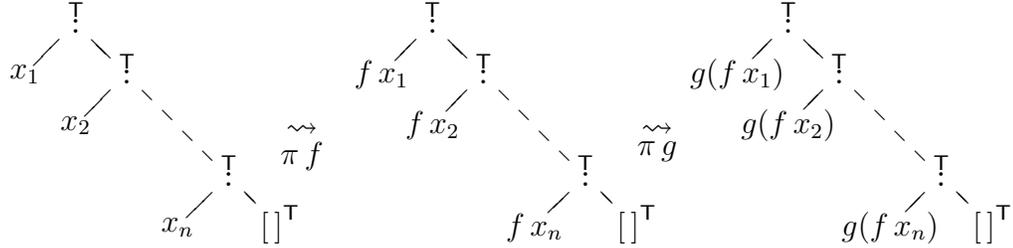
6.2 Cheap Deforestation

138 The catamorphic query discipline provides us with handles to rather elegantly approach these difficulties. First, note that the catamorphism $(e \nabla \otimes)$ separates action (the target algebra $e \nabla \otimes$) from control in that (\cdot) abstracts from the primitive recursion scheme **CATA-INS-REP**. Second, as (\cdot) is the *only* recursive form occurring, we do not need the full power of a general *unfold-fold* transformation strategy.

A general theme of this text recurs once again: a restrictive discipline of expression can lead to deeper insights.

139 The following continues the train of thought we have started in Section 2.6 on spine transformer fusion. To see that fusion is tightly connected with the derivation of streaming programs let us take up the π combinator as an example one last time.

An illustration of the application of query $Q = (\pi g) \cdot (\pi f)$ to $xs = [x_1, \dots, x_n]^\top$ gives the obvious sequence of spine transformations (recall from **PROJECT** that $\pi f = \top f$):



A streaming program for Q would skip the intermediate step but rather directly emit the spine on the right-hand side. Reducing the length of a spine transformer chain reduces the number of necessary spine walks and thus the number of intermediate results to handle. Complete fusion ultimately leads to real stream-based execution.

In the example above, we experience no particular difficulty to derive the streaming program as we only have to appeal to the functoriality hidden in π :

$$\begin{aligned}
 & (\pi g) \cdot (\pi f) \\
 = & \{ \text{PROJECT} \} \\
 & \top g \cdot \top f \\
 = & \{ \text{FUNCTOR} \} \\
 & \top(g \cdot f) \\
 = & \{ \text{PROJECT} \} \\
 & \pi(g \cdot f) .
 \end{aligned}$$

For general catamorphic queries, the following theorem establishes an equivalence which serves as the only but effective fusion tool we have to fall back upon.

140 Acid Rain Theorem [112]. The underlying category is **Set**. In the category of F -algebras $\mathbf{Alg}(F)$, once more let $\tau : FT \rightarrow T$ and $\alpha = FA \rightarrow A$ denote the initial and an arbitrary algebra, respectively. Fix a type D . Then it holds for any function $f : \forall C.(FC \rightarrow C) \rightarrow D \rightarrow C$, that

$$(\alpha)_F \cdot (f \tau) = f \alpha . \qquad \text{ACID-RAIN}$$

Proof. It is remarkable how *parametricity* [116], *i. e.*, the polymorphic type of a function, provides the proof of **ACID-RAIN** for free: the parametricity theorem for the type of f reads (with $\alpha:FA \rightarrow A$, $\beta:FB \rightarrow B$ and $h:B \rightarrow A$ chosen at will)

$$h \cdot \beta = \alpha \cdot Fh \quad \Rightarrow \quad h \cdot (f \beta) = f \alpha .$$

Choosing $\beta = \tau$ and $h = \langle \alpha \rangle_{\mathbb{F}}:T \rightarrow A$ turns the premise of the statement of the parametricity theorem into the characterizing property for catamorphisms **CATA** which, of course, is met by $\langle \alpha \rangle_{\mathbb{F}}$. As a consequence, **ACID-RAIN** holds. \square

The polymorphic type of $f:\forall C.(FC \rightarrow C) \rightarrow D \rightarrow C$ gives insight into the workings of f : as f is indifferent to the actual choice of C , we can be sure that f builds its result of type C solely by application of the constructors for C (which f receives as its first argument).

Law **ACID-RAIN** then validates the intuition that if we manufacture a value using the constructors of τ (this is what $f \tau$ does) and subsequently replace τ 's constructors by those of α (this is the action of $\langle \alpha \rangle_{\mathbb{F}}$), we may rather directly construct the result by using the constructors of α from the start. This way we avoid the computation of the intermediate result of type T —which is what fusion is all about.

141 The functional programming community has devised a family of program transformation techniques, referred to as *deforestation* transformations, which aim to eliminate the production of intermediate (algebraic, and thus tree-shaped) data structures. Wadler characterized a syntactic program class, the *treeless programs*, for which a proof of the guaranteed termination of a *unfold-fold*-based deforestation algorithm could be found [118]. Since then, the criteria which identify the programs that are subject to deforestation have been relaxed [82], but the inherent problems of controlling the *unfold-fold* steps remain.

On the contrary, Takano and Meijer's **ACID-RAIN** theorem exploits a rather fixed producer ($f \tau$) and consumer ($\langle \alpha \rangle_{\mathbb{F}}$) scenario to establish a *cheap* one-step transformation that fuses adjacent consumers and producers. **ACID-RAIN** constitutes a categorical generalization of earlier work on the cheap (or shortcut) deforestation of *listful* programs (in which *foldr* takes the role of the list consumer, see Paragraph 28) [50, 51] and the deforestation of programs over mutually recursive sum-of-product types [37, 74].

142 There are three major driving forces behind this chapter.

(a) In continuation from Paragraph 98, pick up the intriguing idea of a purely spine transformer-based query engine (which means to essentially ignore the material on combinator algebras of Chapter 4). Query deforestation, applied after the syntactic mapping $\mathcal{M} \cdot \mathcal{Q}$ can thus lead to an (at least moderately) efficient, purely transformation-based, implementation of a complex value database engine [59, 60].

(b) We have seen that the combinator pattern matching process of Chapter 4 can fail in that monad comprehension-based parts of a query survive combinator introduction. Comprehension desugarer \mathcal{M} then acts as a fallback which emits catamorphisms, *i. e.*, nested loops, for these query parts. Subsequent deforestation can further optimize these subqueries through reduction of the loop nesting depth. In this sense, deforestation may be perceived as an *a posteriori execution plan optimization*, to be performed after algebraic query rewriting.

(c) Unlike the query engine level techniques (iterator models, suspended evaluation of a stream's tail) which realize streaming *dynamically* at query runtime, deforestation derives streaming plans *statically*, *i. e.*, at query compile time. A deforestation step via law **ACID-RAIN** cancels a spine transformer pair so that formerly separate query parts become adjacent. This, in turn, can reveal previously hidden opportunities for further optimizations. In this sense, contrary to point (b), we can regard deforestation also as a combinator level and thus algebraic optimization tool.

To enable deforestation as a combinator level optimization, as suggested in (c), means to allow the optimizer to desugar (selected) combinator definitions (Definition 103) via mapping \mathcal{M} in order to unbox their catamorphic implementations. Deforestation will especially pay off for those combinators the query engine is likely to implement through simple looping anyway, *e. g.*, π and α . Let us close this chapter with two related examples.

143 Example. Under the proviso $xs : \top Num$ with $\top \neq \text{set}$, naive evaluation of the following combinator query obviously wastes work:

$$((\alpha^{\text{sum}} K_1) \cdot (\pi g) \cdot (\pi f)) xs$$

(note that \mathcal{Q} yields the aggregate $\alpha^{\text{sum}} K_1 : \top Num \rightarrow Num$ as a translation of OQL's `count` clause). Since the query effectively counts the nodes of xs ' spine only, there is no point in mapping g or f over xs at all. The π combinators could be discarded.

Note that Graefe's iterator discipline yields a streaming program for the above query but the resulting plan will nevertheless apply f and g to all elements that pass the π combinators.

At query runtime, the suspension model will create the closure $\overline{g}, \overline{f}, \overline{x}$ but will not force it provided that combinator K_x has a non-strict implementation.

Deforestation, however, realizes this optimization statically:

$$\begin{aligned}
& (\alpha^{\text{sum}} K_1) \cdot (\pi g) \cdot (\pi f) \\
= & \{ \text{AGGREGATE}, \text{PROJECT}, \mathcal{M} \} \\
& \langle 0 \nabla (+ \cdot (K_1 \times id)) \rangle \cdot \top g \cdot \top f \\
= & \{ \text{TYPE-FUNCTOR}, \text{ACID-RAIN} \} \\
& \langle 0 \nabla (+ \cdot (K_1 \times id)) \rangle \cdot \langle ([]^\top \nabla !) \cdot F(g \cdot f, id) \rangle \\
= & \{ \text{ACID-RAIN} \} \\
& \langle (0 \nabla (+ \cdot (K_1 \times id))) \cdot F(g \cdot f, id) \rangle \\
= & \{ \text{sum} \} \\
& \langle 0 \nabla (+ \cdot ((K_1 \cdot (g \cdot f)) \times id)) \rangle \\
= & \{ K_x \cdot f = K_x \} \\
& \langle 0 \nabla (+ \cdot (K_1 \times id)) \rangle \\
= & \{ \text{AGGREGATE} \} \\
& \alpha^{\text{sum}} (K_1) .
\end{aligned}$$

The resulting plan does not rely on lazy evaluation semantics at runtime and may be executed by a strict query engine (which is the common case, as the typical imperative implementation languages have strict semantics). \diamond

144 Example. Recall that we have adopted a catamorphic representation of the sort operator ζ (see **SORT** and **2-SORT**). In interaction with deforestation this sets the scene to reveal and remove superfluous sorting stages in a query plan by means of query rewriting. Traditionally, optimizations of this kind are detected through the inspection of the *physical sort order(s)* of a query expression. As sort orders have no algebraic denotation, their inspection requires extra code outside the core query rewriting system [53, 55].

The catamorphic encoding, however, provides us with enough clue to comprehend ζ 's action at the algebraic level. Let **agg** denote the **ld**-monad associated with the algebra $e \nabla \otimes$ over carrier E (*i. e.*, \otimes is left-commutative, *cf.* Definition 64). With $\ominus : E \times E \rightarrow \text{Bool}$, we then have that

$$\begin{aligned}
& (\alpha^{\text{agg}} f) \cdot (\zeta \ominus) \\
= & \{ \text{AGGREGATE}, \text{SORT} \} \\
& \langle e \nabla (\otimes \cdot (f \times id)) \rangle \cdot \langle nil^{\text{list}} \nabla ins^{\text{list}} \ominus \rangle
\end{aligned}$$

$$\begin{aligned}
&= \{ \text{ACID-RAIN} \} \\
&\quad (e \nabla (\otimes \cdot (f \times id))) \\
&= \{ \text{AGGREGATE} \} \\
&\quad \alpha^{\text{agg}} f
\end{aligned}$$

(a similar calculation asserts the idempotence of ς : $(\varsigma \otimes) \cdot (\varsigma \otimes) = \varsigma \otimes$). \diamond

Chapter 7

Finalization

145 Datatypes. The consistent comprehension of queries as mappings between datatypes, *i. e.*, their algebras of value constructors, has been our guide rail throughout the entire text.

Once this was all set, almost everything else was a consequence of the datatype centric approach: if a query f is a mapping between types A and B , how can we go about and represent f inside a computer or, more specifically, a query engine? Under the proviso that A is finite we could exploit a lookup table to encode f . But what if A is infinite (as are the domains we define query languages over)? Then, the first thing we need is a finite (recursive) description F of the elements in A . This is just what the endofunctors encoding the insert representation of objects in category **Set** did for us. The type of f will thus actually be $f:FA \rightarrow B$. Second, for the encoding of f to be finite, too, it has to track the description F of A , which is nothing else than the hand-waving way of stating that f is a homomorphism (or, for initial A , a catamorphism). This describes the basic understanding of queries in this text already fairly well.

146 Catamorphisms. Much of the conciseness and elegance of this treatment of queries was drawn on the uniform encoding F , the insert representation, of bulk values as well as aggregates and quantifiers. This made level to represent a diversity of query constructs, from parallel application via type functors to sorting (ς), using a single program form, namely the catamorphism combinator (\cdot) .

The sorting combinator ς provides a canonical example of our intention behind this approach: at the time of its introduction as a catamorphic insertion sort this may have appeared intriguing or interesting at best, but a weird thing to do in practice (sorting in the context of very large databases calls for external sorting techniques that operate on secondary storage). This,

however, was not the point at that time. Rather, we could benefit from this encoding as it characterized the act of sorting on the algebraic level. Sorting, being just another spine transformer, thus became accessible to algebraic program transformation which is a big win for a query optimizer since it facilitates purely algebraic—*i. e.*, calculational—query rewriting (maintenance of an extra sortedness attribute is superfluous in such a setup).

Last but not least, catamorphisms are a deeply understood program form in the theory of *Constructive Algorithmics*. With the three fusion laws, notably law **ACID-RAIN**, the work of this community provided input that proved to be most useful in many places.

147 Monads and the Monad Comprehension Calculus. In some sense, this text used monads in the role that sets play in the relational calculus. We consider it a feature not a bug of the monad notion that it comes with just enough internal structure that is needed to interpret a query calculus. The resulting monad comprehension calculus is poor with respect to the variety of syntactic forms it offers but this ultimately led to a discipline in query compilation that extracted the core structure inherent to a query. No obfuscation caused by syntactic sugar (of which approaches that rewrite user level syntax, *e. g.*, OQL or SQL, suffer) remained.

Being completely parametric in the monad an expression of the calculus is evaluated in, the number of different query forms we encountered was significantly reduced: the \mathbb{T} -monad comprehension $\llbracket f\ x \parallel x \leftarrow xs \rrbracket^{\mathbb{T}}$ can describe parallel application of f , duplicate elimination, aggregation, or a quantifier ranging over xs , dependent on the actual choice of monad \mathbb{T} . This uniformity enabled us to spot useful and sometimes unexpected dualities between query constructs, *e. g.*, the close connection of the class of flat join queries and the queries evaluated in the **exists** monad (*e. g.*, existential quantification) in Section 3.2.

The terseness of the calculus additionally had a positive impact on the size of the rule sets necessary to express complex query rewrites. Rewriting rules could be established by appealing to the abstract monad notion in general and then used in many instantiations.

We have found this comprehension of queries based on catamorphisms and monads to cover, simplify, and generalize many of the proposed views of database queries. Even better, however, it enabled us to offer a terse and thus elegant account of issues in query optimization that were cumbersome or impossible to express in a way that is effectively accessible for a query optimizer.

148 We have no doubt that one can and should further travel down this road.

(a) The catamorphic spine transformers have turned out to be expressive enough to grasp modern declarative user query languages. We could, however, push the limits of expressiveness of this approach through substitution of catamorphisms by *paramorphisms* [83]. Paramorphisms provide a significant generalization of the catamorphism notion: *any* function definable on an initial datatype has a paramorphic equivalent. Let $\tau: FT \rightarrow T$ denote the initial algebra for functor F . Any function $f: T \rightarrow A$ can be computed by the paramorphism $outl \cdot (\alpha \Delta (\tau \cdot F outl))$ for a suitably chosen F -algebra $\alpha: F(A \times T) \rightarrow A$. As paramorphisms are particular tupling catamorphisms they enjoy calculational properties similar to those of catamorphisms, *e. g.*, a variant of law **CATA-FUSION** can be established for paramorphisms, too. Paramorphisms thus provide full expressiveness *and* effective control over program form. This could lay a foundation of the comprehension of full-fledged database programming languages and their optimization.

(b) The monads we have encountered in this text have been induced by the type system and the constructs of the user query language. Of course, this is not principle to the method. We envision a query engine in which further aspects, especially implementation level issues, of query processing are encoded in the monad framework and thus become subject to analysis and transformation just like queries. This could spawn research heading into several directions. As already pointed out in the closing paragraph of Chapter 2, the referentially transparent treatment of stateful computation in the **state** monad can lead to an optimizable model of database updates and transactions. Certain tree-shaped datatypes are open to a monadic interpretation as well. It is conceivable to develop a deeper understanding of lookups in tree-shaped index structures (think of *B-trees*) and how these lookups could be efficiently interleaved with query execution. Monadic tree datatypes also make level for a comprehension of queries over semi-structured data in the spirit of [17].

(c) Different levels of sophistication are equally expressible in this framework because the query translation process itself is completely compositional. A query compiler can be rapidly prototyped and look as simple as the composition $\mathcal{M} \cdot \mathcal{Q}$. Throughout this text, however, we have repeatedly observed our model to assimilate other advanced approaches to query optimization with relative ease. Its abstract foundation, catamorphisms and monads, make this understanding of queries a lowest common denominator of these approaches, in a sense. This provides an ideal playground in which the interaction of a diversity of optimization techniques may be studied, much like we did in Chapter 5. There have been other optimization techniques knocking on the

door in the course of Chapter 5 and we are positive that we can further improve our comprehension of queries if we let them in.

References

- [1] Serge Abiteboul and Paris C. Kanellakis. Query Languages for Complex Object Databases. *ACM SIGACT News*, 21(3):9–18, 1990.
- [2] Andrea Asperti and Giuseppe Longo. *Categories, Types, and Structures*. Foundations of Computing Series. The MIT Press, 1991.
- [3] Roland Backhouse. An Exploration of the Bird-Meertens Formalism. Technical Report CS8810, University of Groningen, Department of Mathematics and Computer Science, 1988.
- [4] François Bancilhon, Ted Briggs, Setrag Khoshafian, and Patrick Valduriez. FAD, a Powerful and Simple Database Language. In *Proc. of the 13th Int'l Conference on Very Large Data Bases (VLDB)*, pages 97–105, Brighton, England, September 1987.
- [5] Colin R. Banger and David B. Skillicorn. Flat Arrays as a Categorical Data Type. Department of Computing and Information Science, Queen's University, Kingston, Canada, November 1992.
- [6] Catriel Beeri and Yoram Kornatzky. Algebraic Optimization of Object-Oriented Languages. In Paris C. Kanellakis and Serge Abiteboul, editors, *Proc. of the 3rd Int'l Conference on Database Theory (ICDT)*, number 470 in Lecture Notes in Computer Science (LNCS), pages 72–88, Paris, France, December 1990. Springer Verlag.
- [7] Richard S. Bird. An Introduction to the Theory of Lists. In Manfred Broy, editor, *Logic of Programming and Calculi of Design*, volume 36 of *NATO ASI Series*, pages 5–42. Springer Verlag, April 1987.
- [8] Richard S. Bird. Algebraic Identities for Program Calculation. *The Computer Journal*, 32(2):122–126, 1989.
- [9] Richard S. Bird. Functional Algorithm Design. In *Mathematics of Program Construction (MPC)*, number 947 in Lecture Notes in Computer Science (LNCS), pages 2–17. Springer Verlag, 1995.

- [10] Richard S. Bird and Oege de Moor. *Algebra of Programming*. Prentice Hall International Series in Computer Science. Prentice Hall Europe, 1997.
- [11] Richard S. Bird and Lambert Meertens. Two Exercises Found in a Book on Algorithmics. In Lambert Meertens, editor, *Proc. of the IFIP TC2/WG 2.1 Working Conference on Program Specification and Transformation*, pages 451–457, Bad Tölz, Germany, April 1986. Elsevier Science Publishers B.V. (North Holland).
- [12] Peter A. Boncz, Fred Kwakkel, and Martin L. Kersten. High Performance Support for OO Traversals in Monet. In *Proc. of the 14th British National Conference on Databases (BNCOD)*, pages 152–169, Edinburgh, UK, 1996.
- [13] Val Breazu-Tannen, Peter Buneman, and Limsoon Wong. Naturally Embedded Query Languages. In *Proc. of the Int'l Conference on Database Theory (ICDT)*, pages 140–154, Berlin, Germany, October 1992.
- [14] Val Breazu-Tannen and Ramesh Subrahmanyam. Logical and Computational Aspects of Programming with Sets/Bags/Lists. In *Proc. of the 18th Int'l Colloquium on Automata, Languages and Programming*, pages 60–75, Madrid, Spain, July 1991.
- [15] François Bry. Towards an Efficient Evaluation of General Queries: Quantifier and Disjunction Processing Revisited. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 193–204, Portland, Oregon, USA, June 1989.
- [16] Peter Buneman. The Fast Fourier Transform as a Database Query. Technical Report MS-CIS-93-37/L&C 60, University of Pennsylvania, CIS Department, March 1993.
- [17] Peter Buneman, Susan Davidson, Gerd Hillebrand, and Dan Suciu. A Query Language and Optimization Techniques for Unstructured Data. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 505–516, Montreal, Canada, June 1996.
- [18] Peter Buneman, Robert E. Frankel, and Nikhil Rishiyur. An Implementation Technique for Database Query Languages. *ACM Transactions on Database Systems*, 7(2):164–186, June 1982.

- [19] Peter Buneman, Leonid Libkin, Dan Suciu, Val Tannen, and Limsoon Wong. Comprehension Syntax. *ACM SIGMOD Record*, 23:87–96, March 1994.
- [20] Peter Buneman, Shamim Naqvi, Val Tannen, and Limsoon Wong. Principles of Programming with Complex Objects and Collection Types. *Theoretical Computer Science*, 149(1):3–48, 1995.
- [21] Rod M. Burstall and John Darlington. A Transformation System for Developing Recursive Programs. *Journal of the ACM*, 24(1):44–67, January 1977.
- [22] Luca Cardelli. Type Systems. In *Handbook of Computer Science and Engineering*, chapter 103. CRC Press, 1997.
- [23] Rick G. Cattell and Douglas K. Barry, editors. *The Object Database Standard: ODMG 2.0*. Morgan Kaufmann Publishers, San Francisco, California, 1997. Release 2.0.
- [24] Daniel Chan. *Object-oriented Language Design and Processing*. PhD thesis, Computer Science Department, University of Glasgow, 1994.
- [25] Daniel Chan and Philip Trinder. Object Comprehensions: A Query Notation for Object-Oriented Databases. In *Proc. of the 12th British National Conference on Databases (BNCOD)*, pages 55–72, Guildford, UK, 1994.
- [26] Damianos Chatziantoniou. *Optimization of Complex Aggregate Queries in Relational Databases*. PhD thesis, Columbia University, June 1997.
- [27] Damianos Chatziantoniou and Kenneth A. Ross. Groupwise Processing of Relational Queries. In *Proc. of the 23rd Int'l Conference on Very Large Data Bases (VLDB)*, pages 476–485, Athens, Greece, August 1997.
- [28] Mitch Cherniack. Form(ers) Over Function(s): The KOLA Reference Manual. Technical report, Department of Computer Science, Brown University, 1996.
- [29] Mitch Cherniack. Translating Queries into Combinators. Technical report, Department of Computer Science, Brown University Database Group, 1996.
- [30] Mitch Cherniack. *Building Query Optimizers with Combinators*. PhD thesis, Brown University, May 1999.

- [31] Mitch Cherniack, Ashok Malhotra, and Stanley B. Zdonik. Experiences with Query Translation: Object Queries meet DB2. Technical Report A149, Department of Computer Science, Brown University Database Group, 1998.
- [32] Mitch Cherniack and Stanley B. Zdonik. Rule Languages and Internal Algebras for Rule-Based Optimizers. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 401–412, June 1996.
- [33] Mitch Cherniack and Stanley B. Zdonik. Changing the Rules: Transformation for Rule-Based Optimizers. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, Seattle, Washington, June 1998.
- [34] Jens Claussen, Alfons Kemper, Guido Moerkotte, and Klaus Peithner. Optimizing Queries with Universal Quantification in Object-Oriented Databases. In *Proc. of the 23rd Int'l Conference on Very Large Databases (VLDB)*, pages 286–295, Athens, Greece, August 1997.
- [35] Sophie Cluet and Guido Moerkotte. Nested Queries in Object Bases. In *Proc. of the 4th Int'l Workshop on Database Programming Languages (DBPL)*, September 1993.
- [36] Sophie Cluet and Guido Moerkotte. Efficient Evaluation of Aggregates on Bulk Types. In *Proc. of the 5th Int'l Workshop on Database Programming Languages (DBPL)*, Gubbio, Italy, September 1995.
- [37] Leonidas Fegaras. Efficient Optimization of Iterative Queries. In *Proc. of the 4th Int'l Workshop on Database Programming Languages (DBPL)*, pages 200–225, New York City, USA, August 1993.
- [38] Leonidas Fegaras. A Uniform Calculus for Collection Types. Technical Report 94-030, Oregon Graduate Institute of Science & Technology, 1994.
- [39] Leonidas Fegaras. Optimizing Queries with Object Updates. *Journal of Intelligent Information Systems*, 12(2/3):219–242, 1999. Special Issue on Functional Approach to Intelligent Information Systems.
- [40] Leonidas Fegaras and David Maier. Towards an Effective Calculus for Object Query Languages. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, May 1995.

- [41] Leonidas Fegaras and David Maier. Optimizing Object Queries Using an Effective Calculus. Unpublished work, 1998.
- [42] Maarten M. Fokkinga. *Law and Order in Algorithmics*. PhD thesis, University of Twente, Enschede, 1992.
- [43] Maarten M. Fokkinga. Datatype Laws without Signatures. *Mathematical Structures in Computer Science*, 6:1–32, 1996.
- [44] Johann Christoph Freytag. *Translating Relational Queries into Iterative Programs*. PhD thesis, IBM Almaden Research Center, San Jose, USA, April 1987.
- [45] Johann Christoph Freytag and Nathan Goodman. Rule-Based Translation of Relational Queries into Iterative Programs. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 206–214, Washington, D.C., USA, May 1986.
- [46] Johann Christoph Freytag and Nathan Goodman. Translating Aggregate Queries into Iterative Programs. In *Proc. of the 12th Int'l Conference on Very Large Data Bases (VLDB)*, pages 138–146, Kyoto, Japan, August 1986.
- [47] Johann Christoph Freytag and Nathan Goodman. On the Translation of Relational Queries into Iterative Programs. *ACM Transactions on Database Systems*, 14(1):1–27, March 1989.
- [48] Richard A. Ganski and Harry K. T. Wong. Optimization of Nested SQL Queries Revisited. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 23–33, San Francisco, USA, 1987.
- [49] Jeremy Gibbons. The Third Homomorphism Theorem. *Journal of Functional Programming*, 6(4):657–665, 1996.
- [50] Andrew J. Gill. *Cheap Deforestation for Non-strict Functional Languages*. PhD thesis, Department of Computing Science, University of Glasgow, January 1996.
- [51] Andrew J. Gill, John Launchbury, and Simon L. Peyton Jones. A Short Cut to Deforestation. In *Proc. of the ACM Conference on Functional Programming and Computer Architecture (FPCA)*, pages 223–232, Copenhagen, Denmark, April 1993.

- [52] Dieter Gluche, Torsten Grust, Christof Mainberger, and Marc H. Scholl. Incremental Updates for Materialized OQL Views. In François Bry, Raghu Ramakrishnan, and Kotagiri Ramamohanarao, editors, *Proc. of the 5th Int'l Conference on Deductive and Object-Oriented Databases (DOOD'97)*, number 1341 in Lecture Notes in Computer Science (LNCS), pages 52–66, Montreux, Switzerland, December 1997. Springer Verlag.
- [53] Goetz Graefe. Encapsulation of Parallelism in the Volcano Query Processing System. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 102–111, Atlantic City, New Jersey, USA, 1990.
- [54] Goetz Graefe. Query Evaluation Techniques for Large Databases. *ACM Computing Surveys*, 25(2):73–170, 1993.
- [55] Goetz Graefe and William J. McKenna. The VOLCANO Optimizer Generator: Extensibility and Efficient Search. In *Proc. IEEE Conference on Data Engineering*, page 209 ff., April 1993.
- [56] Malcolm Grant. Homomorphisms and Promotability. In Jan L. A. van de Snepscheut, editor, *Proc. of the Int'l Conference on Mathematics of Program Construction*, number 375 in Lecture Notes in Computer Science (LNCS), pages 335–347, Groningen, The Netherlands, June 1989. Springer Verlag.
- [57] Malcolm Grant. Data Structures and Program Transformation. *Science of Computer Programming*, 14:255–279, 1990.
- [58] Torsten Grust, Joachim Kröger, Dieter Gluche, Andreas Heuer, and Marc H. Scholl. Query Evaluation in CROQUE—Calculus and Algebra Coincide. In Carol Small, Paul Douglas, Roger Johnson, Peter King, and Nigel Martin, editors, *Proc. of the 15th British National Conference on Databases (BNCOD15)*, number 1271 in Lecture Notes in Computer Science (LNCS), pages 84–100, London, Birkbeck College, July 1997. Springer Verlag.
- [59] Torsten Grust and Marc H. Scholl. Query Deforestation. Technical Report 68/1998, Faculty of Mathematics and Computer Science, Database Research Group, University of Konstanz, June 1998.
- [60] Torsten Grust and Marc H. Scholl. How to Comprehend Queries Functionally. *Journal of Intelligent Information Systems*, 12(2/3):191–218,

- March 1999. Special Issue on Functional Approach to Intelligent Information Systems.
- [61] Seymour Hayden and John F. Kennison. *Zermelo-Fraenkel Set Theory*. Merrill, Columbus, Ohio, 1968.
- [62] Jonathan M.D. Hill. *Data-Parallel Lazy Functional Programming*. PhD thesis, University of London, Queen Mary and Westfield College, September 1994.
- [63] Jonathan M.D. Hill and Keith Clarke. An Introduction to Category Theory Monads and their Relationship to Functional Programming. Technical Report QMW-DCS-681, Department of Computer Science, Queen Mary & Westfield College, August 1994.
- [64] C.A.R. Hoare. Notes on an Approach to Category Theory for Computer Scientists. In Manfred Broy, editor, *Constructive Methods in Computing Science*, volume F 55 of *NATO ASI Series*, pages 245–305. Springer Verlag, 1989.
- [65] John Hughes and Simon L. Peyton Jones (editors). Haskell 98: A Non-strict, Purely Functional Language. <http://haskell.org/definition/>, February 1999.
- [66] Matthias Jarke and Jürgen Koch. Query Optimization in Database Systems. *ACM Computing Surveys*, 16(2):111–152, June 1984.
- [67] Johan Jeuring. *Theories for Algorithm Calculation*. PhD thesis, University of Utrecht, 1993.
- [68] Mark P. Jones and Luc Duponcheel. Composing Monads. Technical report, Yale University, Department of Computer Science, New Haven, Connecticut, USA, December 1993.
- [69] Won Kim. On Optimizing an SQL-like Nested Query. *ACM Transactions on Database Systems*, 7(3):443–469, September 1982.
- [70] Heinrich Kleisli. Every Standard Construction is Induced by a Pair of Adjoint Functors. *Proc. of the American Mathematical Society (AMS)*, 16:544–546, 1965.
- [71] Joachim Lambek. A Fixpoint Theorem for Complete Categories. *Mathematische Zeitschrift*, 103:151–161, 1968.

- [72] Joachim Lambek. Least Fixpoints of Endofunctors of Cartesian Closed Categories. *Mathematical Structures in Computer Science*, 3(2):229–257, 1993.
- [73] John Launchbury and Simon L. Peyton Jones. Lazy Functional State Threads. In *Proc. of the ACM Conference on Programming Languages Design and Implementation (PLDI)*, Orlando, Florida, June 1994. ACM Press.
- [74] John Launchbury and Tim Sheard. Warm Fusion: Deriving Build-Catas from Recursive Definitions. In *Proc. of the ACM Conference on Functional Programming and Computer Architecture (FPCA)*, La Jolla, California, June 1995.
- [75] Kazem Lellahi and Val Tannen. A Calculus for Collections and Aggregates. In *Proc. of the 7th Conference on Category Theory and Computer Science (CTCS)*, number 1290 in Lecture Notes in Computer Science (LNCS), pages 261–280, S. Margherita Ligure, Italy, September 1997. Springer Verlag.
- [76] Theodore W. Leung, Gail Mitchell, Bharathi Subramanian, Stanley B. Zdonik, and other. The AQUA Data Model and Algebra. In *Proc. of the 4th Int'l Workshop on Database Programming Languages (DBPL)*, September 1993.
- [77] Leonid Libkin, Rona Machlin, and Limsoon Wong. A Query Language for Multidimensional Arrays: Design, Implementation, and Optimization Techniques. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 228–239, June 1996.
- [78] Daniel F. Liewen and David J. DeWitt. A Transformation Based Approach to Optimizing Loops in Database Programming Languages. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 91–100, San Diego, California, June 1992.
- [79] Saunders Mac Lane. *Categories for the Working Mathematician*, volume 5 of *Graduate Texts in Mathematics*. Springer Verlag, 1971.
- [80] Ernest G. Manes. *Algebraic Theories*, volume 26 of *Graduate Texts in Mathematics*. Springer Verlag, 1976.
- [81] Ernest G. Manes and Michael A. Arbib. *Algebraic Approaches to Program Semantics*. Texts and Monographs in Computer Science. Springer Verlag, 1986.

- [82] Simon David Marlow. *Deforestation for Higher-Order Functional Programs*. PhD thesis, Department of Computing Science, University of Glasgow, September 1995.
- [83] Lambert Meertens. Paramorphisms. *Formal Aspects of Computing*, 4: 413–424, 1992.
- [84] Erik Meijer, Marten M. Fokkinga, and Ross Paterson. Functional Programming with Bananas, Lenses, Envelopes and Barbed Wire. In *Proc. of the ACM Conference on Functional Programming and Computer Architecture (FPCA)*, number 523 in Lecture Notes in Computer Science (LNCS), pages 124–144, Cambridge, USA, 1991. Springer Verlag.
- [85] Robert Milner, Mads Tofte, Robert Harper, and David Mac Queen. *The Definition of Standard ML – Revised*. The MIT Press, 1997.
- [86] Eugenio Moggi. Computational Lambda-Calculus and Monads. In *Proc. of the IEEE Symposium on Logic in Computer Science*, pages 14–23, Pacific Grove, California, USA, June 1989.
- [87] Eugenio Moggi. Notions of Computations and Monads. *Information and Computation*, 93(1):55–92, 1991.
- [88] Ryohei Nakano. Translation with Optimization from Relational Calculus to Relational Algebra Having Aggregate Functions. *ACM Transactions on Database Systems*, 15(4):518–557, December 1990.
- [89] Piergiorgio Odifreddi. *Classical Recursion Theory*, volume 125 of *Studies in Logic and the Foundation of Mathematics*. Elsevier Science Publishers B.V. (North Holland), 1989.
- [90] Atsushi Ohori. Representing Object Identity in a Pure Functional Language. In Paris C. Kanellakis and Serge Abiteboul, editors, *Proc. of the 3rd Int’l Conference on Database Theory (ICDT)*, number 470 in Lecture Notes in Computer Science (LNCS), pages 41–55, Paris, France, December 1990. Springer Verlag.
- [91] John Peterson and Kevin Hammond (editors). Report on the Programming Language Haskell. <http://haskell.org/definition/haskell-report-1.4.ps.gz>, April 1997.
- [92] Simon L. Peyton Jones. *The Implementation of Functional Programming Languages*. Prentice Hall International Series in Computer Science. Prentice Hall International (UK) Ltd., 1987.

- [93] Simon L. Peyton Jones and John Launchbury. State in Haskell. *Lisp and Symbolic Computation*, 8(4):293–341, December 1995.
- [94] Simon L. Peyton Jones and Philip Wadler. Imperative Functional Programming. In *20th ACM Symposium on Principles of Programming Languages (POPL)*, pages 71–84, Charlotte, North Carolina, January 1993.
- [95] György E. Révész. *Lambda-Calculus, Combinators, and Functional Programming*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press, 1988.
- [96] Christian Rich, Arnon Rosenthal, and Marc H. Scholl. Reducing Duplicate Work in Relational Join(s): A Unified Approach. In *Proc. of the Int'l Conference on Information Systems and Management of Data (CISMOD)*, pages 87–102, Delhi, India, October 1993.
- [97] Holger Riedel and Marc H. Scholl. A Formalization of ODMG Queries. In *Proc. of the 7th Int'l Conference on Data Semantics (DS-7)*, Leysin, Switzerland, October 1997.
- [98] David E. Rydeheard and Rod. M. Burstall. *Computational Category Theory*. Prentice Hall International Series in Computer Science. Prentice Hall International (UK) Ltd., 1988.
- [99] Hans-Jörg Schek and Marc H. Scholl. The Relational Model with Relation-Valued Attributes. *Information Systems*, 11(2):137–147, 1986.
- [100] Marc H. Scholl. Extensions to the Relational Data Model. In Peri Loucopoulos and Roberto Zicari, editors, *Advances in Conceptual Modelling and CASE: An Integrated View of Information Systems Development*, pages 163–182. John Wiley & Sons, New York, 1992.
- [101] Marc H. Scholl and Torsten Grust. Hybrid Strategies for Query Translation and Optimisation. Technical Report 72/1998, Faculty of Mathematics and Computer Science, Database Research Group, University of Konstanz, October 1998. Also appeared as ESPRIT *Pastel* Research Report RT2R1.
- [102] Patricia G. Selinger, Morton M. Astrahan, Donald M. Chamberlin, Raymond A. Lorie, and Thomas G. Price. Access Path Selection in a Relational Database Management System. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 23–34, Boston, Massachusetts, USA, 1979.

- [103] Mike Spivey. A Categorical Approach to the Theory of Lists. In Jan L. A. van de Snepscheut, editor, *Proc. of the Int'l Conference on Mathematics of Program Construction*, number 375 in Lecture Notes in Computer Science (LNCS), pages 399–408, Groningen, The Netherlands, June 1989. Springer Verlag.
- [104] Mike Spivey. Category Theory for Functional Programming. Technical Report PRG-TR-7-93, Computing Laboratory, Oxford University, 1993.
- [105] Hennie J. Steenhagen. *Optimization of Object Query Languages*. PhD thesis, Department of Computer Science, University of Twente, 1995.
- [106] Hennie J. Steenhagen, Peter M.G. Apers, and Henk M. Blanken. Optimization of Nested Queries in a Complex Object Model. In *Proc. of the 4th Int'l Conference on Extending Database Technology (EDBT)*, pages 337–350, Cambridge, UK, March 1994.
- [107] Hennie J. Steenhagen, Peter M.G. Apers, Henk M. Blanken, and Rolf A. de By. From Nested-Loop to Join Queries in OODB. In *Proc. of the 20th Int'l Conference on Very Large Databases (VLDB)*, pages 618–629, Santiago, Chile, September 1994.
- [108] Bharathi Subramanian, Stanley B. Zdonik, Theodore W. Leung, and Scott L. Vandenberg. Ordered Types in the AQUA Data Model. In *Proc. of the 4th Int'l Workshop on Database Programming Languages (DBPL)*, September 1993.
- [109] Bharati Subramanian, Theodore W. Leung, Scott L. Vandenberg, and Stanley B. Zdonik. The AQUA Approach to Querying Lists and Trees in Object-Oriented Databases. In *Proc. of the 11th Intl. Conference on Data Engineering*, 1995.
- [110] Dan Suciu and Val Breazu-Tannen. A Query Language for NC. In *Proc. of the 13th ACM Symposium on Principles of Databases (PODS)*, pages 167–178, Minneapolis, Minnesota, May 1994.
- [111] Dan Suciu and Limsoon Wong. On Two Forms of Structural Recursion. In Georg Gottlob and Moshe Y. Vardi, editors, *Proc. of the 5th Int'l Conference on Database Theory (ICDT)*, number 893 in Lecture Notes in Computer Science (LNCS), pages 111–124, Prague, Czech Republic, January 1995. Springer Verlag.

- [112] Akihiko Takano and Erik Meijer. Shortcut Deforestation in Calculational Form. In *Proc. of the ACM Conference on Functional Programming and Computer Architecture (FPCA)*, pages 306–313, La Jolla, USA, June 1995. ACM Press.
- [113] Philip Trinder. Comprehensions, a Query Notation for DBPLs. In *Proc. of the 3rd Int'l Workshop on Database Programming Languages (DBPL)*, pages 55–68, Nafplion, Greece, 1991.
- [114] Philip Trinder and Philip Wadler. Improving List Comprehension Database Queries. In *Proc. of TENCEN'89*, pages 186–192, Bombay, India, November 1989.
- [115] Bennet Vance. Towards an Object-Oriented Query Algebra. Technical Report 91–008, Oregon Graduate Institute of Science & Technology, January 1992.
- [116] Philip Wadler. Theorems for Free! In *Proc. of the 4th Int'l Conference on Functional Programming and Computer Architecture (FPCA)*, London, England, September 1989.
- [117] Philip Wadler. Comprehending Monads. In *Conference on Lisp and Functional Programming*, pages 61–78, June 1990.
- [118] Philip Wadler. Deforestation: Transforming Programs to Eliminate Trees. *Theoretical Computer Science*, 73:231–248, 1990.
- [119] Philip Wadler. The Essence of Functional Programming. In *19th ACM Symposium on Principles of Programming Languages (POPL)*, Albuquerque, New Mexico, January 1992.
- [120] Philip Wadler. Monads for Functional Programming. In Johan Jeuring and Erik Meijer, editors, *Advanced Functional Programming*, number 925 in Lecture Notes in Computer Science (LNCS). Springer Verlag, 1995.
- [121] Philip Wadler. How to Declare an Imperative. *ACM Computing Surveys*, 29(3):240–263, 1997.
- [122] Philip Wadler and Stephen Blott. How to make Ad-hoc Polymorphism less Ad hoc. In *16th ACM Symposium on Principles of Programming Languages (POPL)*, Austin, Texas, January 1989.

- [123] David A. Watt and Philip Trinder. Towards a Theory of Bulk Types. Technical Report FIDE/91/26, Computer Science Department, University of Glasgow, 1991.
- [124] Limsoon Wong. *Querying Nested Collections*. PhD thesis, University of Pennsylvania, Philadelphia, August 1994.
- [125] Limsoon Wong. *The Kleisli/CPL Extensible Query Optimizer—Programmer Guide*. Institute of Systems Science, Singapore, April 1996.

Index

- abstract datatype, 17
- ACID-RAIN, 132
- aggregate (α), **99**, 134, 135
- aggregation operator, **46**, 114, 118
 - decomposable, 118, 120
- Alg** (category of algebras and homomorphisms), 15
- algebra, 6
 - carrier of an, 14
 - categorical, 14
 - collection, 37
 - combinator, 98
 - free, 32
 - initial, 16, 17
 - logical (query), 96
 - non-collection, 46
 - polymorphic, 24
- all, 46
- all (OQL), 75
- and (OQL), 70
- antijoin (\bowtie), 92, **99**
- AQUA, 122
- array, 79
 - shape of an, 79
- array, 79
- array (OQL), 79
- arrow, *see* morphism
- avg (OQL), 77
- bag, 38
- bag (OQL), 69
- banana brackets ($\lfloor \cdot \rfloor$), 17
- BANANA-SPLIT, 78
- bifunctor, 12
- bind (\gg), 56
- bintree, 26
- black box, 80
- Bool*, 22
- bottom (\perp), 79
- calculation, 7
- calculus
 - monad comprehension, 39
 - monoid comprehension, 57
 - query, 6
 - relational (tuple), 47, 97
- cancellation, 10
- carrier, 14, 24
- Cat** (category of functors and morphisms), 9
- CATA, 17
- CATA-COMM, 33
- CATA-FUSION, 29
- CATA-IDEM, 33
- CATA-INS-REP, 20
- CATA-MAP-FUSION, 29
- CATA-REFLECT, 21
- CATA-UNION-REP, 53
- catamorphism, 3, 6, **17**, 20, 131
 - $\lfloor \cdot \rfloor$, 17
 - tupling, 77, 139
- category, 8
 - Alg**, 15
 - Cat**, 9
 - default, 8
 - product, 12

- Set**, 9
 - theory, 6, 7
- category theory, 3
- closed expression (in the λ -calculus), 96
- closure, 8, *see* suspension
- codomain, 9
- coercion, 49
- collection monad, 45
- combinator, 38, **96**, 128
- commutative diagram, 8, 18
- composition, 8, 9
 - functional, 28
- comprehension
 - endo-monadic, 48
 - monad, 6, 39, **46**
 - desugaring (\mathcal{M}), 49
 - filters in a, 41, 48
 - generators in a, 48
 - normal form of a, 39
 - normalization, **90**, 125
 - normalization rules, 86
 - qualifiers in a, 48
 - syntax, 47, 83
 - unnesting, 85, 86
 - monoid, 57
 - multi-monad, 7, **48**, 49
 - object, 56
- confusion, 16, 17, 32
- cons* (:), 1, 13
- Constructive Algorithmics, 3, 7
- constructor, 1, 5, 16, 133
 - collection type, 13
 - polymorphic, 24
 - , 13
 - datatype, 9
 - [], 13
- control structure, 1, 130
- coproduct, *see* sum
- count** (OQL), 77
- cross (X), 99
- database engine, *see* query engine
- datatype, 5
 - abstract, 17
 - equation for a, 32, **35**
 - former, 16
 - polymorphic, 24
- decomposition (of an aggregate), 118
- deforestation, 31, 133, 135
 - cheap, 32, 133
- DESCOPE, 108
- descoping, 107
- destructor, 18
- diagram, 8
 - commutative, 18
 - pasting, 18
- DIST, 107
- distinct** (OQL), 78
 - delay of, 84
- domain, 9
- duplicate removal, 84
 - delay of, 84
- element, 14
- element** (OQL), 79
 - pushdown of, 80
- element type, 13
- empty set (\emptyset), 9
- empty tuple (), 65, 124
- endofunctor, 9
 - fixpoint of, 19
- EQ-COMM, 35
- EQ-COMM-IDEM, 38
- EQ-IDEM, 35
- equality type, 65
- equation, 6, 32, 35
- equational reasoning, 7
- Eureka step, 2, 31, 104
- except** (OQL), 74
- exists, 46
- exists** (OQL), 76

- exists-in (OQL), 75
- ext*, 56
- extensionality, 7
- F-algebra, 14
 - polymorphic, 24
- F-catamorphism, 17
- F-homomorphism, 15
- filter, *see* comprehension
- filter*, 22
- first (OQL), 79
 - pushdown of, 80
- flatten*, 45
- flatten (OQL), 78
- fold step, 130, 131
- foldr*, 21
- forall-in (OQL), 75
- FQL, 129
- free theorem, 40, 133
- functional programming
 - lazy evaluation in, 129
 - monads in, 59
- FUNCTOR, 9
- functor, 9
 - all, 46
 - array, 79
 - bag, 38
 - bintree, 26
 - composition, **9**, 12
 - constant (**K**), 12
 - exists, 46
 - identity (**Id**), **9**, 12
 - list, 26, **38**
 - max, 46
 - maybe, 27, 59
 - min, 46
 - polynomial, 12
 - prod, 46
 - set, 38
 - state, 60
 - sum, 46
 - type, 25, 116
- FUNCTOR- \mathcal{M} -FUSION, 88
- fusion, 28, 130, 132
- group query, 114
- Haskell, 59
- hd*, 73
- HOM, 15
- hom*, 57
- homomorphism, 2, 15, 118
 - monoid, 57
- Id, 9, 12
- id*, 8
- in (OQL), 75
- incremental update, 118
- infix conditional (?), 23
- initial algebra, 16
- initial type, 26
- initiality, 9, 17
- injection, 10
 - left (*inl*), 10
 - right (*inr*), 10
- inl*, 10
- inr*, 10
- insert representation, **12**, 120
 - aggregation in, 46
 - collection type in, 37
 - quantifier in, 46
- insertion sort, 73
- intersect (OQL), 74
- isomorphism, 9
- iterator, 115, 128
- join*, 40
- join (\bowtie), 95, **99**
- junc (∇), 10
- junk, 16, 17
- K*, 22, 124
- K**, 12

- Kleisli triple, 56
- KOLA, 122
- Lambek's Lemma, **18**, 32
- lazy evaluation, 76, 129, 135
- least fixpoint semantics, 19
- left-commutative, 14, 33
- left-idempotent, 14, 33
- length*, 22, 30
- let* (local binding), 51
- list, 10, 26, 38
- list (functor), 9
- list (OQL), 69
- listtset (OQL), 78
 - delay of, 84
- \mathcal{M} , 49
- \mathcal{M} -NORM-1, 87
- \mathcal{M} -NORM-2, 87
- \mathcal{M} -NORM-3, 87
- \mathcal{M} -NORM-4, 87
- \mathcal{M} -NORM-5, 87
- map*, 9, 25, 30
- map functor, *see* type functor
- max, 46
- max (OQL), 78
- maybe, 27, 59
- min, 46
- min (OQL), 78
- ML, 65
- monad, 40
 - all, 46
 - bag, 45
 - \gg , 56
 - collection, 45
 - commutative, 42
 - $\#$, 42
 - exception, 59
 - exists, 46
 - idempotent, 42
 - identity (**ld**), 46
 - join*, 40
 - Kleisli, 56
 - list, 45
 - map*, 40
 - max, 46
 - maybe, 59
 - min, 46
 - prod, 46
 - set, 45
 - state, 60
 - state transformer, 59
 - sum, 46
 - unit*, 40
 - zero*, 41
- monad comprehension, *see* comprehension
- monad morphism, 57, **81**, 83
- MONAD-1, 40
- MONAD-2, 40
- MONAD-3, 40
- MONAD-4, 41
- MONAD-5, 41
- MONAD-MORPH-1, 81
- MONAD-MORPH-2, 81
- MONAD-MORPH-3, 81
- MONAD-MORPH-4, 81
- MONAD-MORPH-5, 83
- monoid, 53, 57
- morphism, 8
 - composition, 8
 - identity (*id*), 8
 - type of a, 8
- mutator, 5
- natural transformation (\dashv), 40
- NATURALITY, 40
- naturality, 40
- NEST-EX, 107
- nested loops, 102, 103, 134
- nestjoin (\triangle), 91, **99**, 104, 114, 121
- nil* ($[]$), 13

- normal form, 39
 - weak head, 76, 129
- normalization, 63
 - KOLA, 123, 124
 - monad comprehension, 86, **90**, 125
 - rules for, 86
- Num*, 22
- Obj*, 65
- object, 8
 - element of an, 14
 - initial (*1*), 9
 - product (\times), 10
 - source (*src*), 8
 - sum (+), 10
 - target (*tgt*), 8
- object identity, 60, 64
- OQL, 64, 68
 - arrays in, 79
 - variables in, 66
- or (OQL), 70
- orthogonality, 28
- outl*, 10
- outr*, 10
- parallel application, 115, 116
- parametricity, 40, 133
- paramorphism, 139
- parser, 18
- PARTITION, 116
- partition*, 116
- partitioning, 115, 116
- pattern
 - combinator, 97
 - partial, 100, 103
- pipelining, *see* streaming
- point-free, 7
- point-wise, 7
- polymorphism, 6, 24, 40, 98, 133
 - free theorem and, 40
 - naturality and parametric, 41
 - parametricity and, 3, 40
- polynomial functor, 12
- primitive recursive function, 21
- prod**, 46
- product, 10
 - cartesian, 11
 - nested, 65, 122, 124
 - \times , 10
 - Δ , 10
- project (π), **99**, 116, 128–130, 132
- projection, 10, 123
 - left (*outl*), 10
 - right (*outr*), 10
- \mathcal{Q} , 63, **68**
- QUAL-EX, 106
- qualifier, *see* comprehension
- quantifier, 46
 - existential, 51
 - universal, 50, 110
- query
 - nesting in a, 38
 - unnesting of a, 39
- query engine
 - monolithic, 127
 - spine transformer-based, 95, 134
- query former, 64
- quotient algebra, 36
- record, 65
 - tag, 65
- recursion
 - primitive, 21
 - structural, 2, 18, 21
 - sri*, 21
 - sru*, 54
- referential transparency, 5
- relational algebra, 21
 - nested, 21

- relational calculus, 47, 97
- rewrite, 7
- ringad, 55
- select (σ), 96, **99**
- select-distinct-from-where (OQL), 70
- select-from-where (OQL), 70
- select-from-where-group by (OQL), 71
- select-from-where-group by-having (OQL), 72
- select-from-where-order by (OQL), 72
- semijoin (\bowtie), 93, **99**, 106, 109
- set, 38
- Set** (category of sets and total functions), 9
- set (OQL), 69
- shortcut deforestation, *see* deforestation
- signature, 6, 34
- some (OQL), 75
- sort (ς), **99**, 135
- source, *see* object
- spine, 13, 20
 - transformation of a, **20**, 25, 127, 132
 - walk of a, 18, 28, 30, 132
- spine transformer, **20**, 127, 132
- split (Δ), 10
- Squiggol, 7, 17
- src*, 8
- state, 60
- streaming, 28, 76, 127, **128**
- String*, 61
- struct (OQL), 68
- structural recursion, 18, 54
- sum, 10
 - ∇ , 10
 - $+$, 10
- sum, 46
- sum (OQL), 77
- suspension ($\ulcorner f, x \urcorner$), 129
- T-monad, 40
- tag (of a disjoint sum), 11
- tag (of a record field), 65
- target, *see* object
- tgt*, 8
- \otimes -table, 118
- tl*, 73
- TRANSFORMER, 34
- transformer, 33
- tree transformation, 20
- tupling, 77, 105, 106
- type
 - equality, 65
 - initial, 26
 - polymorphic, 133
 - record, 65
 - unit*, 14
- type constructor, 12
- type environment (Γ), 67
- type functor, **25**, 116
- type judgment, 66
- TYPE-FUNCTOR, 26
- typing rule, 67
- union (OQL), 74
- union representation, **52**, 118
- unique, 76
- unit* type, 14
- unit* (monadic injection), 40
- unnest (μ), 99
- unnesting, 39, 63, 90
 - monad comprehension, 85, 86, **90**, 102
 - SQL query, 92
- variable, 6, 33, 38

free, 66, 96, 110
interdependency, 106
scope of a, 122
variable environment, 122

weak head normal form, 76, 129

zero, 41
ZF-expression, *see* comprehension