

# Multi-Model Query Compilation Using Bi-Algebra Comprehensions

Ryan Wisnesky  
Conexus AI

September 2, 2020

## Abstract

In this paper we extend the monad comprehension calculus, a well-studied user-facing query formalism sufficient to encode (nested) relational algebra and admitting an easily-optimized iterator-based execution model, to arbitrary user-defined abstract data types. Following Erwig, our formalism is parameterized by a bi-cartesian category  $L$  representing an ambient type theory and we represent an implementation of abstract data type  $F$  as an  $F$ -bi-algebra on  $L$ . Using the internal language of the free bi-cartesian category, we define  $L$ -independent implementations of bi-algebras for sets, graphs, infinite streams, finite maps, state, and people (representing a business object), and demonstrate how to translate bi-algebras from one choice of  $L$  to another; translation from one bi-algebra to another within the same  $L$  is accomplished with natural transformations in the usual way.

## 1 Introduction

In this paper we are concerned with three recent trends in database systems and theory:

- the development, beginning in the 1990s [19, 11], of query formalisms, compilers, and runtimes based on the interplay between structural recursion and comprehension syntax, as exemplified by systems such as Microsoft’s Dyrad-LINQ [20] and Data Parallel Haskell [5]; and
- the development, beginning in the 2000s, of data integration formalisms based on first-order logic, used to mediate between semantically heterogeneous database schemas, as exemplified in systems such as IBM’s Clio [12] and Microsoft’s Rondo [13]; and
- the development, beginning in the 2010s, of “big data” query formalisms, compilers, and runtimes based on graph walks over various distributed data stores, as exemplified by systems such as Apache Tinkerpop [14] and foreshadowed by mixed SQL/Hadoop runtimes such as Hadapt [3].

The above trends are related to each other as follows:

- although query compilers based on structural recursion and comprehension syntax are often used in relational data integration systems, including Clio and Rondo, using them with graph databases runs into non-trivial problems [8] with purely functional graph representations; and
- although data integration formalisms based on first-order logic can mediate between schemas in the same data model, using them to mediate between data models themselves runs into non-trivial problems with representing logic using logic [1]; and
- although graph-based “big data” systems such as Apache TinkerPop organically developed techniques similar to structural recursion, comprehension syntax and logic-based schema mappings, because of the problems described above, many of these techniques are purely informal, or formalized operationally [14].

## 1.1 Contributions

In this paper we address the challenges above by extending the monad comprehension calculus, a well-studied user-facing query formalism sufficient to encode (nested) relational algebra and admitting an easily-optimized iterator-based execution model, to arbitrary user-defined abstract data types. Following Erwig [7], our formalism is parameterized by a bi-cartesian category  $L$ , representing an ambient type theory, and we represent an implementation of abstract data type  $F$  as an  $F$ -bi-algebra on  $L$ . After reviewing the internal language of bi-cartesian categories, we use that syntax to define  $L$ -independent implementations of bi-algebras for sets, graphs, infinite streams, finite maps, state, and people (representing a business object), and demonstrate how to translate bi-algebras from one choice of  $L$  to another; translating from one bi-algebra to another within the same  $L$  is accomplished with natural transformations in the usual way. Our formalism has many uses:

- to generalize compilers that are based on structural recursion and comprehensions, such as LINQ [20] and Data Parallel Haskell [5], to handle graphs; and
- to generalize data integration formalisms based on first-order logic to account for morphisms between underlying data models, such as between graphs and relations; and
- to provide a rigorous mathematical foundation / intermediate form for existing graph processing systems, such as Apache Tinkerpop [14].

To exemplify these uses, we prove theorems relating our formalism to that of algebraic databases [16] and algebraic property graphs [17], demonstrate how our formalism can interpret the Apache Tinkerpop Gremlin family of languages [14], and conclude with a discussion of current work on distributed runtimes for our formalism .

## 1.2 Outline

This paper is structured as follows:

- In section 2, we describe a standard formalism (which we call CAML) for query processing with structural recursion and co-recursion, similar to Grust’s monad comprehension calculus [11, 19, 18].
- In section 3, we describe how to extend CAML to allow for arbitrary user-defined abstract data types, following Erwig’s bi-algebra approach [9, 7], and we describe how to migrate bi-algebras from one programming language to another by way of CAML.
- In section 4, we implement a set of bi-algebras and morphisms for common abstract data types such as sets and graphs, following Erwig’s inductive graph approach [8], and demonstrate how to interpret a fragment of Apache Tinkerpop Gremlin [14] in our formalism.
- In section 5, we relate our formalism to those of algebraic databases [16] and property graphs [17].

We conclude in section 6 with some discussion on distributed runtimes for our formalism. Knowledge of type theory and category theory at the level of [2] and knowledge of functional query languages at the level of [11] are required to understand this paper.

# 2 Categorical Abstract Machine Language (CAML)

In this section we describe a minimal, standard formalism (which we call CAML) for query processing with structural recursion on finite lists and structural co-recursion on infinite streams, similar to Wong’s nested relational calculus [19], Grust’s monad comprehension calculus [11], Curien’s Categorical Abstract Machine Language [6], and many others.

## 2.1 Syntax

A *type*  $t$  is defined as a term in the context-free grammar:

$$t ::= \text{Int} \mid 0 \mid 1 \mid t + t \mid t \times t \mid t \Rightarrow t \mid \text{List } t \mid \text{Stream } t$$

A  $(t, t'$ -typed) *expression*  $e : t \rightarrow t'$  is defined as a term in the context-free grammar (where we omit the  $|$  symbols between productions on different lines):

$$\begin{aligned}
e : t \rightarrow t' & ::= \text{id}_t : t \rightarrow t \mid (e' : t' \rightarrow t'') \circ (e : t \rightarrow t') : t \rightarrow t'' \\
\text{const}_c \ (c \in \{0, 1, -1, \dots\}) & : 1 \rightarrow \text{Int} \mid \text{add}, \text{mul} : \text{Int} \times \text{Int} \rightarrow \text{Int} \mid \text{neg} : \text{Int} \rightarrow \text{Int} \\
\text{eq}_t & : t \times t \rightarrow 1 + 1 \mid \delta_{t, t', t''} : t \times (t' + t'') \rightarrow (t \times t') + (t \times t'') \\
\text{tt}_t : t \rightarrow 1 & \mid \text{fst}_{t, t'} : t \times t' \rightarrow t \mid \text{snd}_{t, t'} : t \times t' \rightarrow t' \mid \langle e : t \rightarrow t', e' : t \rightarrow t'' \rangle : t \rightarrow t' \times t'' \\
\text{ff}_t : 0 \rightarrow t & \mid \text{inl}_{t, t'} : t \rightarrow t + t' \mid \text{inr}_{t, t'} : t' \rightarrow t + t' \mid \langle e : t' \rightarrow t \mid e' : t'' \rightarrow t \rangle : t' + t'' \rightarrow t \\
& \Lambda(e : t \times t' \rightarrow t'') : t \rightarrow t' \Rightarrow t'' \mid \text{eval}_{t, t'} : (t' \Rightarrow t) \times t' \rightarrow t \\
\text{in}_t : 1 + (t \times \text{List } t) & \rightarrow \text{List } t \mid \text{fold}(e : \tau \times (1 + (t \times t')) \rightarrow t') : \tau \times \text{List } t \rightarrow t' \\
\text{out}_t : \text{Stream } t & \rightarrow t \times \text{Stream } t \mid \text{unfold}(e : t \rightarrow t' \times t) : t \rightarrow \text{Stream } t'
\end{aligned}$$

We abbreviate:

$$\begin{aligned}
\text{Bool} & := 1 + 1 \quad \text{true}_t := \text{inl} \circ \text{tt}_t \quad \text{false}_t := \text{inr} \circ \text{tt}_t \quad f \otimes g := \langle f \circ \text{fst}, g \circ \text{snd} \rangle \quad f \oplus g := \langle \text{inl} \circ f \mid \text{inr} \circ g \rangle \\
\text{fold}'(f) & := \text{fold}(f \circ \text{snd}) \circ (\text{tt}, \text{id}) \quad \text{unfold}'(f) := \text{unfold}(\langle \text{fst} \circ f, \langle \text{fst}, \text{snd} \circ f \rangle \rangle)
\end{aligned}$$

As an example, to sum the length of a list we may write:  $\text{fold}'(\text{const}_0 \mid \text{add}) : \text{List } \text{Int} \rightarrow \text{Int}$ .

So that our formalism can be used without function types, we use parameterized folds, which additionally require an input parameter of type  $\tau$  compared to the form usually seen in the literature; parameterized folds can be defined in terms of the other expressions, but only in a higher-order setting (requiring  $\text{eval}$  and  $\Lambda$ ). Similarly, the distributivity expression  $\delta$  can be defined in terms of the other expressions, but only in a higher-order setting, because not all bi-cartesian categories are distributive, but all closed ones are [10].

The fragment of CAML that does not contain function types or stream types can express exactly the primitive recursive functions on lists; dually, the fragment of CAML that does not contain function types or list types can express exactly the primitive co-recursive function on streams; in the presence of function types, CAML can express non-primitive-recursive functions on lists and streams, such as the Ackermann's function.

## 2.2 Axiomatic Semantics

The axiomatic semantics of CAML is defined as the smallest equivalence relation  $\approx$  that is also a congruence (i.e., for which  $f \approx f'$  implies  $g \circ f \approx g \circ f'$ , etc):

$$\begin{aligned}
\text{id} \circ f & \approx f \quad f \circ \text{id} \approx f \quad f \circ (g \circ h) \approx (f \circ g) \circ h \quad \text{eq} \circ \langle e, f \rangle \approx \text{true} \text{ iff } e \approx f \\
f : t \rightarrow 1 & \approx \text{tt}_t \quad \text{fst} \circ \langle f, g \rangle \approx f \quad \text{snd} \circ \langle f, g \rangle \approx g \quad \langle \text{fst} \circ f, \text{snd} \circ f \rangle \approx f \\
f : 0 \rightarrow t & \approx \text{ff}_t \quad \langle f \mid g \rangle \circ \text{inl} \approx f \quad \langle f \mid g \rangle \circ \text{inr} \approx g \quad \langle f \circ \text{inl} \mid f \circ \text{inr} \rangle \approx f \\
\text{eval} \circ (\Lambda(f) \otimes \text{id}) & \approx f \quad \Lambda(\text{eval} \circ (g \otimes \text{id})) \approx g \quad \langle \text{inl} \circ \text{snd} \mid \text{inr} \circ \text{snd} \rangle \circ \delta \approx \text{snd} \quad \langle \text{fst} \mid \text{fst} \rangle \circ \delta \approx \text{fst} \\
\delta \circ (\text{id} \otimes \text{inl}) & \approx \text{inl} \quad \delta \circ (\text{id} \otimes \text{inr}) \approx \text{inr} \quad \delta \circ (f \otimes \langle \text{inr} \circ g \mid \text{inl} \circ h \rangle) \approx \langle \text{inr} \circ (f \otimes g) \mid \text{inl} \circ (f \otimes h) \rangle \circ \delta \\
\text{fold}(f) \circ (\text{id} \otimes \text{in}_t) & \approx f \circ \langle \text{fst}, (\text{tt}_1 \oplus (\text{id}_t \otimes \text{fold}(f))) \circ \alpha \rangle \quad g \approx \text{fold}(f) \text{ if } g \circ (\text{id} \otimes \text{in}_t) \approx f \circ \langle \text{fst}, (\text{tt}_1 \oplus (\text{id}_t \otimes g)) \circ \alpha \rangle \\
g & \approx \text{unfold}(f) \text{ if } g \circ \text{unfold}(f) \approx ((\text{id}_t \otimes g)) \circ f \quad \text{out}_t \circ \text{unfold}(f) \approx (\text{id}_t \otimes \text{unfold}(f)) \circ f
\end{aligned}$$

with isomorphism:  $\alpha : \Gamma \times (1 + (A \times B)) \rightarrow 1 + (A \times (\Gamma \times B))$  and equations stating that  $\text{Int}$  forms a commutative ring in the obvious way, our choice of base types, primitives, and their axioms being ad-hoc in this paper.

The *syntactic category* CAML is defined as the category formed by taking types as objects and equivalence classes of expressions modulo  $\approx$  as morphisms. These axioms, minus lists, streams, arithmetic, and equality, were shown to be decidable in [15].

**Remark on Type Isomorphism.** Product and sum types have the pleasing property of admitting a complete (and decidable) axiomatization as the least congruence  $\sim$  for which  $(0, +)$  and  $(1, \times)$  are commutative monoids,  $\times$  distributes over  $+$ , and  $t \times 0 \sim 0$ . This axiomatization allows mediating between isomorphic product and sum types automatically, a common operation in query processing. Surprisingly, although the type constructor  $\Rightarrow$  satisfies Tarski's "high school algebra axioms" such as  $a \Rightarrow b \times c \sim (a \Rightarrow b) \times (a \Rightarrow c)$ , type isomorphism with  $+$  and  $\times$  and  $\Rightarrow$  together is not recursively axiomatizable, and its decidability is unknown [10].

## 2.3 Denotational Semantics

The set-theoretic semantics  $\llbracket t \rrbracket$  types  $t$  is standard and can be recursively defined as a functor from CAML to the category of sets:

- $\text{Int}$  denotes the set of integers,  $\{0, 1, -1, 2, -2, \dots\}$ .
- $0$  denotes the empty set  $\{\}$ , also written  $\emptyset$ , and  $1$  denotes the singleton set  $\{\emptyset\}$ .
- $t \times t'$  denotes the Cartesian product of the sets  $\llbracket t \rrbracket$  and  $\llbracket t' \rrbracket$ , where e.g,

$$\{1, 2\} \times \{2, 3\} = \{(1, 2), (1, 3), (2, 2), (2, 3)\}.$$

- $t + t'$  denotes the disjoint union of the sets  $\llbracket t \rrbracket$  and  $\llbracket t' \rrbracket$ , where e.g,

$$\{1, 2\} + \{2, 3\} = \{(\emptyset, 1), (\emptyset, 2), (\{\emptyset\}, 2), (\{\emptyset\}, 3)\},$$

where we have chosen  $\emptyset$  and  $\{\emptyset\}$  as left and right “labels”.

- $t \Rightarrow t'$  denotes the set of functions from  $\llbracket t \rrbracket$  to  $\llbracket t' \rrbracket$ .
- $\text{List } t$  denotes the set of finite lists of  $\llbracket t \rrbracket$ .
- $\text{Stream } t$  denotes the set of infinite lists (so-called *streams*) with elements in  $\llbracket t \rrbracket$ .

The set-theoretic semantics  $\llbracket e \rrbracket$  of an expression  $e : t \rightarrow t'$  is a function  $\llbracket t \rrbracket \rightarrow \llbracket t' \rrbracket$  is also standard:

- $\text{id}_t : t \rightarrow t$  denotes the identity function on  $\llbracket t \rrbracket$ , and  $f \circ g$  denotes the function composition of  $\llbracket f \rrbracket$  with  $\llbracket g \rrbracket$ ; i.e.,  $\llbracket f \circ g \rrbracket(x) := \llbracket f \rrbracket(\llbracket g \rrbracket(x))$ .
- $\text{const}_c : 1 \rightarrow \text{Int}$  denotes the constant function returning integer  $c$ , and  $\text{add}, \text{mul} : \text{Int} \times \text{Int} \rightarrow \text{Int}$  denote integer addition and multiplication, respectively, and  $\text{neg} : \text{Int} \rightarrow \text{Int}$  denotes negation (subtraction from zero).
- $\text{eq}_t : t \times t \rightarrow \text{Bool}$  denotes the characteristic function of  $\{(x, x) \mid x \in \llbracket t \rrbracket\}$ ; i.e., the  $\llbracket t \rrbracket$ -equality function.
- $\text{tt}_t : t \rightarrow 1$  denotes the function sending any element of  $\llbracket t \rrbracket$  to  $\emptyset$ , and  $\text{ff}_t : 0 \rightarrow t$  denotes the only function from  $\emptyset$  to  $\llbracket t \rrbracket$ .
- $\text{fst}_{t,t'}$  and  $\text{snd}_{t,t'}$  denote the first and second projections from  $\llbracket t \times t' \rrbracket$ , respectively, and  $\langle f : t \rightarrow t', g : t \rightarrow t'' \rangle$  denotes the function taking each  $x \in \llbracket t \rrbracket$  to the pair  $(\llbracket f \rrbracket(x), \llbracket g \rrbracket(x))$ .
- $\text{inl}_{t,t'}$  and  $\text{inr}_{t,t'}$  denote the first and second injections to  $\llbracket t + t' \rrbracket$ , respectively, and  $\langle f : t' \rightarrow t \mid g : t'' \rightarrow t \rangle$  denotes the case-analysis function taking each  $(\emptyset, x) \in \llbracket t' + t'' \rrbracket$  to  $\llbracket f \rrbracket(x)$  and each  $(\{\emptyset\}, y) \in \llbracket t' + t'' \rrbracket$  to  $\llbracket g \rrbracket(y)$ .
- $\delta_{t,t',t''} : t \times (t' + t'') \rightarrow (t \times t') + (t \times t'')$  denotes the function taking  $(x, (l, y))$  to  $(l, (x, y))$ .
- $\Lambda(e : t \times t' \rightarrow t'') : t \rightarrow t' \Rightarrow t''$  denotes the “currying” function taking each  $x \in \llbracket t \rrbracket$  to the function taking each  $y \in \llbracket t' \rrbracket$  to  $\llbracket e \rrbracket(x, y)$ , and  $\text{eval}_{t,t'}$  denotes the function taking  $(f : \llbracket t \rrbracket \rightarrow \llbracket t \rrbracket, g : \llbracket t' \rrbracket)$  to  $f(g)$ .
- $\text{in}_t$  denotes the initial algebra for lists of type  $t$  in insert presentation, and  $\text{fold}(e : \tau \times (1 + t \times t') \rightarrow t') : \tau \times \text{List } t \rightarrow t'$  denotes structural iteration with a parameter of type  $\tau$ .
- $\text{out}_{\tau,t}$  denotes the final co-algebra for streams in insert presentation, and  $\text{unfold}(e : t' \rightarrow t \times t') : t' \rightarrow \text{Stream } t$  denotes structural co-iteration.

**Theorem 1.** If  $f, g : t \rightarrow t'$  are such that  $f \approx g$ , then  $\llbracket f \rrbracket = \llbracket g \rrbracket$  as set-theoretic functions  $\llbracket t \rrbracket \rightarrow \llbracket t' \rrbracket$ .

**Proof.** In Coq, by induction on  $f \cong g$ .  $\square$

## 2.4 Lambda calculus Form

A more user-friendly syntax for CAML than the variable-free form above is that of a  $\lambda$ -calculus [6] that allows typed variables, which we write as  $v:t$  for variable  $v$  and type  $t$ . A *context* is a list of distinct variable-type pairs:

$$\Gamma ::= - \mid \Gamma, v:t \ (v \notin \Gamma)$$

A *typing judgment*, also called an *expression-in-context*,  $\Gamma \vdash e : t$  is a relation on contexts  $\Gamma$ , terms  $e$ , and types  $t$  and is inductively defined using the inference rules:

$$\begin{array}{c} \frac{\Gamma \vdash v : t \quad v \neq v'}{\Gamma, v:t \vdash v : t} \quad \frac{}{\Gamma, v:t \vdash v : t} \quad \frac{c \in \{\dots, 0, \dots\}}{\Gamma \vdash \mathbf{K}_c : \mathbf{Int}} \quad \frac{\Gamma \vdash e : t \quad \Gamma \vdash f : t}{\Gamma \vdash \mathbf{eq} \ e \ f : \mathbf{Bool}} \quad \frac{}{\Gamma \vdash \mathbf{tt} : \mathbf{1}} \\ \\ \frac{\Gamma \vdash e : t \times t'}{\Gamma \vdash \mathbf{fst} \ e : t} \quad \frac{\Gamma \vdash e : t \times t'}{\Gamma \vdash \mathbf{snd} \ e : t'} \quad \frac{\Gamma \vdash e : t \quad \Gamma \vdash e' : t'}{\Gamma \vdash \langle e, e' \rangle : t \times t'} \quad \frac{\Gamma \vdash e : \mathbf{Int} \quad \Gamma \vdash e' : \mathbf{Int}}{\Gamma \vdash e + / * e' : \mathbf{Int}} \\ \\ \frac{\Gamma \vdash e : \mathbf{Int}}{\Gamma \vdash - \ e : \mathbf{Int}} \quad \frac{\Gamma \vdash e : \mathbf{0}}{\Gamma \vdash \mathbf{ff}_t \ e : t} \quad \frac{\Gamma \vdash e : t}{\Gamma \vdash \mathbf{inl}_t \ e : t' + t} \quad \frac{\Gamma \vdash e : t'}{\Gamma \vdash \mathbf{inr}_t \ e : t + t'} \\ \\ \frac{\Gamma, v : t \vdash e'' : t'' \quad \Gamma, v : t' \vdash e' : t'' \quad \Gamma \vdash e : t + t' \quad v \notin \Gamma}{\Gamma \vdash \mathbf{case}_v \ e \ \mathbf{of} \ e' \mid e'' : t''} \quad \frac{\Gamma, v:t \vdash e : t' \quad v \notin \Gamma}{\Gamma \vdash \lambda v:t.e : t \Rightarrow t'} \\ \\ \frac{\Gamma \vdash e' : t' \quad \Gamma \vdash e : t' \Rightarrow t}{\Gamma \vdash e e' : t} \quad \frac{\Gamma \vdash e : \mathbf{1} + t \times \mathbf{List} \ t}{\Gamma \vdash \mathbf{in} \ e : \mathbf{List} \ t} \quad \frac{\Gamma \vdash e : \mathbf{Stream} \ t}{\Gamma \vdash \mathbf{out} \ e : t \times \mathbf{Stream} \ t} \\ \\ \frac{\Gamma, v : \mathbf{1} + t \times t' \vdash e : t' \quad v \notin \Gamma \quad \Gamma \vdash f : \mathbf{List} \ t}{\Gamma \vdash \mathbf{fold} \ (\lambda v.e) \ f : t'} \quad \frac{\Gamma, v : t \vdash e : t \times t' \quad v \notin \Gamma \quad \Gamma \vdash f : t}{\Gamma \vdash \mathbf{unfold} \ (\lambda v.e) \ f : \mathbf{Stream} \ t'} \end{array}$$

We have altered the typeface of keywords above to ensure that they are distinct from the combinator form of CAML. The translation from the above syntax to the variable-free syntax is standard: contexts are translated into product types:

$$[-] := \mathbf{1} \quad [\Gamma, v : t] := [\Gamma] \times t$$

and a typing judgment  $\Gamma \vdash e : t$  is translated into an expression  $[e] : [\Gamma] \rightarrow [t]$  in combinator form by translating each variable into a projection:

$$\begin{aligned} [\Gamma, v' : t' \vdash v : t] &:= [\Gamma \vdash v : t] \circ \mathbf{fst}_{[\Gamma], [t]} & [\Gamma, v : t \vdash v : t] &:= \mathbf{snd}_{[\Gamma], [t]} & [\Gamma \vdash \mathbf{K}_c : \mathbf{Int}] &:= \mathbf{const}_c \circ \mathbf{tt}_{[\Gamma]} \\ [\Gamma \vdash \mathbf{eq} \ e \ f : \mathbf{Bool}] &:= \mathbf{eq} \circ \langle [e], [f] \rangle & [\Gamma \vdash e + / * e' : \mathbf{Int}] &:= \mathbf{plus/mul} \circ \langle [e], [f] \rangle & [\Gamma \vdash - \ e : \mathbf{Int}] &:= \mathbf{neg} \circ [e] \\ [\Gamma \vdash \mathbf{fst} \ e : t] &:= \mathbf{fst}_{[t], [t']} \circ [\Gamma \vdash e : t \times t'] & [\Gamma \vdash \mathbf{snd} \ e : t'] &:= \mathbf{snd}_{[t], [t']} \circ [\Gamma \vdash e : t \times t'] \\ [\Gamma \vdash \langle e, e' \rangle : t \times t'] &:= \langle [\Gamma \vdash e : t], [\Gamma \vdash e' : t'] \rangle & [\Gamma \vdash \mathbf{tt} : \mathbf{1}] &:= \mathbf{tt}_{[\Gamma]} & [\Gamma \vdash \mathbf{ff} \ e : t] &:= \mathbf{ff}_{[t]} \circ [\Gamma \vdash e : \mathbf{0}] \\ [\Gamma \vdash \mathbf{case}_v \ e \ \mathbf{of} \ e' \mid e'' : t''] &:= \langle [\Gamma, v:t \vdash e' : t''] \mid [\Gamma, v:t' \vdash e'' : t''] \rangle \circ \delta \circ \langle \mathbf{id}_{[\Gamma]}, [\Gamma \vdash e : t + t'] \rangle \\ [\Gamma \vdash \mathbf{inl}_{t'} \ e : t + t'] &:= \mathbf{inl}_{[t], [t']} \circ [\Gamma \vdash e : t] & [\Gamma \vdash \mathbf{inr}_{t'} \ e : t'] &:= \mathbf{inr}_{[t], [t']} \circ [\Gamma \vdash e : t'] \\ [\Gamma \vdash \lambda x : t.e : t \Rightarrow t'] &:= \Lambda[\Gamma, x : t \vdash e : t'] & [\Gamma \vdash e' e : t'] &:= \mathbf{eval} \circ \langle [\Gamma \vdash e' : t \Rightarrow t'], [\Gamma \vdash e : t] \rangle \\ [\Gamma \vdash \mathbf{in} \ e : \mathbf{List} \ t] &:= \mathbf{in}_{[t]} \circ [\Gamma \vdash e : \mathbf{1} + (t \times \mathbf{List} \ t)] & [\Gamma \vdash \mathbf{out} \ e : t \times \mathbf{List} \ t] &:= \mathbf{out}_{[t]} \circ [\Gamma \vdash e : \mathbf{Stream} \ t] \\ [\Gamma \vdash \mathbf{fold} \ (\lambda v.e) \ f : t'] &:= \mathbf{fold}_{\Gamma, t, t'} [\Gamma, v : \mathbf{1} + (t \times t') \vdash e : t'] \circ \langle \mathbf{id}_{[\Gamma]}, [\Gamma \vdash f : \mathbf{List} \ t] \rangle \\ [\Gamma \vdash \mathbf{unfold} \ (\lambda v.e) \ f : \mathbf{Stream} \ t] &:= \mathbf{unfold}'_{\Gamma, t, t'} [\Gamma, v : t \times t' \vdash e : t \times t'] \circ \langle \mathbf{id}_{[\Gamma]}, [\Gamma \vdash f : \mathbf{Stream} \ t] \rangle \end{aligned}$$

## 2.5 The Set Monad

Structural recursion (and co-recursion) has much to recommend it as in intermediate form for query compilation. In particular, folds (and unfolds) are easily optimized [11]. However, we run into complications when generalizing fold to *non-free* collection types. For example, although we may represent sets as “lists up to permutation”, to do so we require that whenever we fold a function  $f$  through a list  $l$  representing a set, for every permutation  $l'$  of  $l$ , that  $f(l) = f(l')$ , a potentially undecidable property. For this reason, a restricted form of fold, known as *comprehension*, is often proposed as a query language for non-free inductive datatypes, which then compiles into structural recursion in a way known to respect the defining equations of the non-free type, such as invariance under permutation.

To exemplify such non-free types, we add types and expressions for the *set monad* to CAML. The types and expressions of CAML are extended with:

$$\begin{aligned}
 t & ::= \dots \mid \mathbf{Set} \ t \\
 e & ::= \dots \mid \eta_t : t \rightarrow \mathbf{Set} \ t \mid \emptyset_t : \mathbf{1} \rightarrow \mathbf{Set} \ t \mid \Upsilon_t : \mathbf{Set} \ t \times \mathbf{Set} \ t \rightarrow \mathbf{Set} \ t \mid \mu_t : \mathbf{Set} \ (\mathbf{Set} \ t) \rightarrow \mathbf{Set} \ t \\
 & \quad \mathbf{map}(e : t \rightarrow t') : \mathbf{Set} \ t \rightarrow \mathbf{Set} \ t' \mid \rho_{t,t'} : t \times \mathbf{Set} \ t' \rightarrow \mathbf{Set} \ (t \times t')
 \end{aligned}$$

and the axiomatic semantics with:

$$\begin{aligned}
 \mathbf{map}(\mathbf{snd}) \circ \rho & \approx \mathbf{snd} \quad \rho \circ (\mathbf{id} \otimes \eta) \approx \eta \mathbf{map}(f) \circ \emptyset' \approx \emptyset' \quad \rho \circ \langle e, \emptyset' \rangle \approx \emptyset' \quad \mu \circ \emptyset' \approx \emptyset' \\
 \rho \circ (\mathbf{id} \otimes \mu) & \approx \mu \circ \mathbf{map}(\rho) \circ \rho \quad \mathbf{map}(f \otimes g) \circ \rho \approx \rho \circ (f \otimes \mathbf{map}(g)) \quad \mathbf{map}(\theta) \circ \rho \approx \rho \circ (\mathbf{id} \otimes \rho) \circ \theta \\
 \Upsilon \circ \langle e, f \rangle & \approx \Upsilon \circ \langle f, e \rangle \quad \Upsilon \circ \langle e, e \rangle \approx e \quad \Upsilon \circ \langle e, \emptyset' \rangle \approx e \quad \Upsilon \circ \langle e, \Upsilon \circ \langle f, g \rangle \rangle \approx \Upsilon \circ \langle \Upsilon \circ \langle e, f \rangle, g \rangle
 \end{aligned}$$

where  $\theta := \langle \mathbf{fst} \circ \mathbf{fst}, \langle \mathbf{snd} \circ \mathbf{fst}, \mathbf{snd} \rangle \rangle$  and  $\emptyset' := \emptyset \circ \mathbf{tt}$ . The set-theoretic semantics is extended with:

- $\llbracket \emptyset_t \rrbracket$  denotes the function out of  $\llbracket \mathbf{1} \rrbracket$  returning the empty set,  $\llbracket \eta_t \rrbracket$  denotes the function taking each  $x \in \llbracket t \rrbracket$  to the singleton set  $\{x\}$ , and  $\llbracket \Upsilon_t \rrbracket$  denotes the function taking each pair  $(x, y)$  of sets of  $\llbracket t \rrbracket$  to their union  $x \cup y$ . These three operations are called the *union presentation of sets*.
- $\llbracket \mathbf{map}(f) \rrbracket$  denotes the function taking each set  $\{s_1, \dots, s_n\}$  to the set  $\{\llbracket f \rrbracket(s_1), \dots, \llbracket f \rrbracket(s_n)\}$ , and  $\llbracket \mu_t \rrbracket$  denotes the function taking a set of sets of  $\llbracket t \rrbracket$ , say  $X$ , to the “n-ary-union”  $\bigcup_{x \in X} x$ , returning a set of  $\llbracket t \rrbracket$ . Collectively these two operations, and  $\emptyset$ , form the *set monad with zero*.
- $\llbracket \rho_{t,t'} \rrbracket$  denotes the function taking each pair  $(x, \{s_1, \dots, s_n\}) \in \llbracket t \rrbracket \times \llbracket \mathbf{Set} \ t' \rrbracket$  to the set of pairs  $\{(x, s_1), \dots, (x, s_n)\}$ . Such an operation is often called a *strength*, and a monad with a strength, such as  $\mathbf{Set}$ , is called a *strong monad*.

The set monad can be implemented using lists considered up to permutation and duplication as:

$$\begin{aligned}
 \llbracket \emptyset \rrbracket & := \mathbf{in}' \circ \mathbf{inl} \circ \mathbf{tt} \quad \llbracket \eta \rrbracket := \mathbf{in}' \circ \mathbf{inr} \circ \langle \mathbf{id}, \llbracket \emptyset \rrbracket \circ \mathbf{tt} \rangle \quad \llbracket \Upsilon \rrbracket := \mathbf{fold}(\langle \mathbf{fst} \mid \mathbf{in}' \circ \mathbf{inr} \circ \mathbf{snd} \rangle \circ \delta) \circ \mathbf{swap} \\
 \llbracket \mu_t \rrbracket & := \mathbf{fold}'(\langle \llbracket \emptyset \rrbracket \mid \llbracket \Upsilon \rrbracket \rangle) \quad \llbracket \mathbf{map}_t \rrbracket := \mathbf{fold}'(\langle \llbracket \emptyset \rrbracket \mid \mathbf{in}' \circ \mathbf{inr} \circ \langle f \circ \mathbf{fst}, \mathbf{snd} \rangle \rangle) \quad \llbracket \rho_t \rrbracket := \mathbf{fold}(\langle \llbracket \emptyset \rrbracket \circ \mathbf{tt} \mid \mathbf{in}' \circ \mathbf{inr} \circ \omega \rangle \circ \delta)
 \end{aligned}$$

where  $\omega := \langle \langle \mathbf{fst}, \mathbf{fst} \circ \mathbf{snd} \rangle, \mathbf{snd} \circ \mathbf{snd} \rangle$  is an isomorphism and so is  $\mathbf{swap} := \langle \mathbf{snd}, \mathbf{fst} \rangle$ . Although not necessary, we use  $\mathbf{swap}$  in  $\llbracket \Upsilon \rrbracket$  so that it preserves input list order, which makes our proofs easier.

**Theorem 2.** Let  $\triangleright$  denote the function taking lists to sets and  $\mathbf{map}$  denote the map function on lists (i.e.,  $\mathbf{map}(f)(x_1, \dots, x_n) = (f(x_1), \dots, f(x_n))$ ). Then:

$$\begin{aligned}
 \triangleright \llbracket \llbracket \emptyset \rrbracket \rrbracket (x) & = \llbracket \llbracket \emptyset \rrbracket \rrbracket (x) \quad \triangleright \llbracket \llbracket \eta \rrbracket \rrbracket (x) = \llbracket \llbracket \eta \rrbracket \rrbracket (x) \quad \triangleright \llbracket \llbracket \Upsilon \rrbracket \rrbracket (x, y) = \llbracket \llbracket \Upsilon \rrbracket \rrbracket (\triangleright x, \triangleright y) \\
 \triangleright \llbracket \llbracket \mu \rrbracket \rrbracket (x) & = \llbracket \llbracket \mu \rrbracket \rrbracket (\triangleright \mathbf{map}(\triangleright, x)) \quad \triangleright \llbracket \llbracket \mathbf{map}(f) \rrbracket \rrbracket (x) = \llbracket \llbracket \mathbf{map}(f) \rrbracket \rrbracket (\triangleright x) \quad \triangleright \llbracket \llbracket \rho \rrbracket \rrbracket (x, y) = \llbracket \llbracket \rho \rrbracket \rrbracket (x, \triangleright y)
 \end{aligned}$$

**Proof.** In Coq. □

## 2.6 Comprehensions as Queries

Although we may use the set monad as described above directly, a comprehension is defined as an expression in the  $\lambda$ -calculus form of CAML that de-sugars into the set monad. Comprehension syntax has many concrete presentations, as well as a convenient normal form, and in this paper we adopt the “do/bind notation” of Haskell [18], and extend the  $\lambda$ -calculus version of CAML with:

$$\frac{}{\Gamma \vdash \emptyset_t : \text{Set } t} \quad \frac{\Gamma \vdash e : t}{\Gamma \vdash \text{sng } e : \text{Set } t} \quad \frac{\Gamma \vdash e : \text{Set } t \quad \Gamma \vdash e' : \text{Set } t}{\Gamma \vdash e \vee e' : \text{Set } t}$$

$$\frac{\Gamma \vdash e : \text{Set } t \quad \Gamma, v : t \vdash f : \text{Set } t' \quad v \notin \Gamma}{\Gamma \vdash \text{do } v:t \leftarrow e; f : \text{Set } t'}$$

and corresponding translations:

$$\begin{aligned} [\Gamma \vdash \emptyset_t : \text{Set } t] &:= \emptyset_{[t]} \circ \text{ott}_{[\Gamma]} \quad [\Gamma \vdash \text{sng } e : \text{Set } t] := \eta \circ [\Gamma \vdash e : t] \quad [\Gamma \vdash e \vee e' : \text{Set } t] := \gamma \circ \langle [\Gamma \vdash e], [\Gamma \vdash e'] \rangle \\ [\Gamma \vdash \text{do } v:t \leftarrow e; f : \text{Set } t'] &:= \mu \circ \text{map}([v:t, \Gamma \vdash f]) \circ \rho \circ \langle \text{id}, [\Gamma \vdash e] \rangle \end{aligned}$$

Set-theoretically,  $\text{do } v \leftarrow e; f$  denotes “flat mapping”  $f$  through  $e$ ; i.e.,  $\cup_{v \in e} f(v)$ .

**Theorem 3.** Let  $\text{subst}(v \mapsto e, f)$  indicate the capture-avoiding substitution of  $e$  for variable  $v$  in  $f$  and  $\text{weaken}(v, f)$  the weakening of  $f$  to include and ignore fresh variable  $v$ . Then:

$$\begin{aligned} [[\Gamma \vdash \text{do } v \leftarrow e; \text{sng } v]] &= [[\Gamma \vdash e]] \quad [[\Gamma \vdash \text{do } v \leftarrow \text{sng } e; f]] = [[\Gamma \vdash \text{subst}(v \mapsto e, f)]] \\ [[\Gamma \vdash \text{do } y \leftarrow (\text{do } x \leftarrow m; f) ; g]] &= [[\Gamma \vdash \text{do } x \leftarrow m ; (\text{do } y \leftarrow f; \text{weaken}(x, g))]] \\ [[\Gamma \vdash \text{do } v \leftarrow e; \emptyset]] &= [[\Gamma \vdash \emptyset]] = [[\Gamma \vdash \text{do } v \leftarrow \emptyset; f]] \end{aligned}$$

**Proof.** In Coq. □

Additional equations are provable, such as commutativity and idempotency of do/bind, but we single out the equations above because are true in any monad with zero (lists, bags, etc), such as each type in the Boom hierarchy [4], and they can be oriented to form a *strongly normalizing* re-write system [11]. So far, then, our CAML compiler consists of comprehension normalization (representing a simple query optimizer) and then translation into recursion (representing a simple query compiler).

## 3 Extending CAML with Bi-Algebra ADTs

## 4 A Standard Library of Bi-Algebra ADTs

## 5 Application: Generic Compilation of CQL

## 6 Conclusion

## References

- [1] Suad Alagić and Philip A. Bernstein. A model theory for generic schema management. In Giorgio Ghelli and Gösta Grahne, editors, *Database Programming Languages*. Springer Berlin Heidelberg, 2002.
- [2] Steve Awodey. *Category Theory*. Oxford University Press, Inc., 2nd edition, 2010.
- [3] Kamil Bajda-Pawlikowski, Daniel J. Abadi, Avi Silberschatz, and Erik Paulson. Efficient processing of data warehousing queries in a split execution environment. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data*, SIGMOD '11, page 1165–1176, New York, NY, USA, 2011. Association for Computing Machinery.
- [4] Alexander Bunkenburg. *The Boom Hierarchy*, pages 1–8. Springer London, London, 1994.

- [5] Manuel M. T. Chakravarty, Roman Leshchinskiy, Simon Peyton Jones, Gabriele Keller, and Simon Marlow. Data parallel haskell: A status report. In *Proceedings of the 2007 Workshop on Declarative Aspects of Multicore Programming*, DAMP '07, page 10–18, New York, NY, USA, 2007. Association for Computing Machinery.
- [6] G. Cousineau, P.-L. Curien, and M. Mauny. The categorical abstract machine. *Science of Computer Programming*, 8(2):173 – 202, 1987.
- [7] Martin Erwig. Categorical programming with abstract data types. In *Proceedings of the 7th International Conference on Algebraic Methodology and Software Technology*, AMAST '98, page 406–421, Berlin, Heidelberg, 1999. Springer-Verlag.
- [8] Martin Erwig. Inductive graphs and functional graph algorithms. *J. Funct. Program.*, 11(5):467–492, September 2001.
- [9] Martin Erwig. Comprehending adts. [https://www.researchgate.net/publication/228520206\\_Comprehending\\_ADTs](https://www.researchgate.net/publication/228520206_Comprehending_ADTs), 01 2011.
- [10] Marcelo Fiore, Roberto Di Cosmo, and Vincent Balat. Remarks on isomorphisms in typed lambda calculi with empty and sum types. *Annals of Pure and Applied Logic*, 141(1), 2006.
- [11] Torsten Grust and Marc H. Scholl. How to comprehend queries functionally. *Journal of Intelligent Information Systems*, 12(2):191–218, Mar 1999.
- [12] Laura M. Haas, Mauricio A. Hernández, Howard Ho, Lucian Popa, and Mary Roth. Clio grows up: From research prototype to industrial tool. In *Proceedings of the 2005 ACM SIGMOD International Conference on Management of Data*, SIGMOD '05, page 805–810, New York, NY, USA, 2005. Association for Computing Machinery.
- [13] Sergey Melnik, Erhard Rahm, and Philip A. Bernstein. Developing metadata-intensive applications with rondo. *Journal of Web Semantics*, 1(1):47 – 74, 2003.
- [14] Marko A. Rodriguez. The gremlin graph traversal machine and language (invited talk). *Proceedings of the 15th Symposium on Database Programming Languages - DBPL 2015*, 2015.
- [15] Gabriel Scherer. Deciding equivalence with sums and the empty type. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages*, POPL 2017, page 374–386, New York, NY, USA, 2017. ACM.
- [16] Patrick Schultz, David I. Spivak, and Ryan Wisnesky. Algebraic model management: A survey. In Phillip James and Markus Roggenbach, editors, *Recent Trends in Algebraic Development Techniques*. Springer International Publishing, 2017.
- [17] Joshua Shinavier and Ryan Wisnesky. Algebraic property graphs, 2019.
- [18] Philip Wadler. Comprehending monads. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming*, LFP '90, page 61–78, New York, NY, USA, 1990. ACM.
- [19] Limsoon Wong. *Querying Nested Collections*. PhD thesis, University of Pennsylvania, Philadelphia, PA, USA, 1994. AAI9503855.
- [20] Yuan Yu, Michael Isard, Dennis Fetterly, Mihai Budiu, Úlfar Erlingsson, Pradeep Kumar Gunda, and Jon Currey. Dryadlinq: A system for general-purpose distributed data-parallel computing using a high-level language. In *OSDI'08: Eighth Symposium on Operating System Design and Implementation*, December 2008.

## Contents

<b>1 Introduction</b>	<b>1</b>
1.1 Contributions . . . . .	2
1.2 Outline . . . . .	2



<b>2</b>	<b>Categorical Abstract Machine Language (CAML)</b>	<b>2</b>
2.1	Syntax . . . . .	2
2.2	Axiomatic Semantics . . . . .	3
2.3	Denotational Semantics . . . . .	4
2.4	Lambda calculus Form . . . . .	5
2.5	The Set Monad . . . . .	6
2.6	Comprehensions as Queries . . . . .	7
<b>3</b>	<b>Extending CAML with Bi-Algebra ADTs</b>	<b>7</b>
<b>4</b>	<b>A Standard Library of Bi-Algebra ADTs</b>	<b>7</b>
<b>5</b>	<b>Application: Generic Compilation of CQL</b>	<b>7</b>
<b>6</b>	<b>Conclusion</b>	<b>7</b>