

# Towards Representation-Independent Logic and Data

CONSTANTINE THEOCHARIS, University of St Andrews, UK

Abstract TODO

CCS Concepts: • **Theory of computation** → **Categorical semantics**; **Abstraction**; *Operational semantics*; • **Software and its engineering** → *Translator writing systems and compiler generators*.

Additional Key Words and Phrases: Representation, Algebraic Data Types, Categorical Semantics, Compilation, Low-Level Data Structures

## ACM Reference Format:

Constantine Theocharis. 2024. Towards Representation-Independent Logic and Data. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 8 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

There is a fundamental gap between programs written for human convenience and programs written for machine efficiency. It can be argued that a large segment of programming language design aims to bridge this gap, by providing abstraction techniques that allow programmers to write programs in a high-level language, and then compile them to a low-level language that can be executed efficiently.

However, as the source programs become increasingly abstract and high-level, and the compilation techniques become increasingly sophisticated in order to keep up with the larger gap, it is very difficult to predict the final low-level code output. As such, there can be an observed performance ceiling for programs written in high-level languages, which can only be broken by "getting one's hands dirty" and dropping down to low-level code, where human intuition and insight is more effective at selecting the correct algorithms and data structures for a given task.

This is unfortunate, however, because making this transition leads to a loss of abstraction, and thus a loss of many things including ease of readability, ease of modification, or even correctness guarantees. It appears as if there is a trade-off at play. But what if we could have the best of both worlds? What if we can write programs in a high-level, expressive, abstract language with all the fancy features under the sun, and at the same time be able to *specify* how these programs are to be represented in a lower-level sense, but without requiring us to repeat ourselves in terms of business logic and data?

### 1.1 The high-level idea

This setup of having two sites of "programming" is not a new concept. In fact, it could be argued that this is exactly how the process of standard compilation works. We have a high-level language, which is compiled into a low-level language, which is then executed on a machine.

However, the process of compilation is usually not considered to be part of the high-level program itself. Rather, it is the job of the compiler to figure out exactly what is the optimal representation of

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference acronym 'XX, June 03–05, 2018, Woodstock, NY*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/XXXXXXX.XXXXXXX>

a high-level program in the low-level language. This is especially the case when the discrepancy between the levels becomes great, for example in functional languages such as Haskell or ML. In these languages, the programmer has very little control over the eventual representation of their program in the low-level language, unless they use specialised primitives which are meant to mirror the low-level capabilities of the target system.

This is what we are aiming to change with this formalism. We envision a process of compilation in which the programmer has a lot more control over the representation of their program in the low-level language, while still keeping a clear separation between that and the core logic of the high-level program itself.

## 1.2 Translation of logic and data

In the present categorical presentation of this process, we want each program to construct its own customised compilation functor  $C$ . This functor will be constructed by the programmer, just like the rest of the high-level logic of the program.

One might remark that this sounds like a lot of work for the programmer. However, we will see that the compilation functor can be constructed automatically based on a set of core "data type representations" that are provided by the programmer. These data type representations will be the only thing that the programmer will have to write, and they will be reusable across many different programs.

Furthermore, there will be a trade-off at play. The more custom parts of the compilation functor that are written, the more control they have over the eventual representation, but the more coupling there is between the high-level program and the low-level representation. Crucially however, this coupling is never present in the high-level program itself, but rather sits alongside it.

## 2 PRELIMINARIES

### 3 TECHNIQUE

### 4 CATEGORICAL MODEL

Let  $\mathbb{S}$  and  $\mathbb{L}$  be two 2-categories. Eventually we will require additional structure to be present in these, but for now we want to consider them as context categories for two different type theories. We will consider 2-morphisms to be rewriting rules, and 1-morphisms to be terms, in the style of [CITE].

We want to view  $\mathbb{S}$  as the "high-level", nice category that we want to write our programs in. On the other hand,  $\mathbb{L}$  is the "low-level", ugly category that we actually execute our programs in. It is ugly from the point of view of the programmer, but nice from the point of view of a machine, because it closely follows the way the machine actually executes the program.

Since we have a high-level and a low-level category, we want to have a way to translate between them. We will do this by defining two functors:

- A functor  $C : \mathbb{S} \longrightarrow \mathbb{L}$ , called the *compilation* functor. It takes a program in  $\mathbb{S}$  and compiles it to a program in  $\mathbb{L}$ .
- A functor  $Q : \mathbb{L} \longrightarrow \mathbb{S}$ , called the *quoting* functor. it allows us to opaquely operate on program fragments from  $\mathbb{L}$  within  $\mathbb{S}$ .

Importantly, both of these functors are *lax*, meaning that they only respect functoriality up to 2-morphisms. In other words, compiling a high-level program  $C(f \circ g)$  is not necessarily the same as compiling  $C(f)$  and then  $C(g)$  separately, but there is a 2-morphism  $C(f) \circ C(g) \Rightarrow C(f \circ g)$  implying that one evaluates to the other.

The quoting functor  $Q$  is left-adjoint to the compilation functor  $C$ , in the sense that compiling a quoted program yields the same program. Commonly this adjunction is strict.

#### 4.1 Features of $\mathbb{S}$

Since  $\mathbb{S}$  is the site of our high-level programs, we want it to have a lot of the things that we expect from a programming language. In particular, we want it to be cartesian closed, and we want it to admit fixpoints of endofunctors. Furthermore, we expect that there exists a terminal object  $1 \in \text{Ob}(\mathbb{S})$ .

#### 4.2 Features of $\mathbb{L}$

We require much less convenience in the structure of  $\mathbb{L}$ , instead relying on the functor  $C$  to translate human-friendly programs into machine-friendly ones.

Commonly  $\mathbb{L}$  will not be closed, but for this paper we will assume it is. It might sometimes be a Kleisli category with respect to some monad (probably a monad holding the state of the stack/heap/environment/etc).

We will commonly want to restrict  $\mathbb{L}$  to only be monoidal in some less-than-cartesian way, for example to model a linear logic for the tracking and preservation of resources.

#### 4.3 The category $\text{Repr}_{\mathbb{L}}^{\mathbb{S}}$

The data about how to translate high-level programs into low-level programs will be stored in a category  $\text{Repr}_{\mathbb{L}}^{\mathbb{S}}$ .

Consider the following diagram:

$$\begin{array}{ccc}
 & \text{Int} & \\
 & \curvearrowright & \\
 \mathbb{S} & \xleftarrow{\mu\text{Src}} & \text{Repr}_{\mathbb{L}}^{\mathbb{S}} \xrightarrow{\text{Tar}} \mathbb{L} \\
 & \xrightarrow{\sigma} & \downarrow \text{Src} \\
 & \xrightarrow{\mu} & [\mathbb{S}, \mathbb{S}]
 \end{array} \tag{1}$$

The category  $\text{Repr}_{\mathbb{L}}^{\mathbb{S}}$  is the category of representations of  $\mathbb{S}$  in  $\mathbb{L}$ . There are a few key functors here:

- $\text{Src} : \text{Repr}_{\mathbb{L}}^{\mathbb{S}} \longrightarrow [\mathbb{S}, \mathbb{S}]$  is the source functor. It takes a representation of  $\mathbb{S}$  in  $\mathbb{L}$  and returns the original high-level inductive type as an endofunctor in  $[\mathbb{S}, \mathbb{S}]$ . Can be considered one of the bundle projections.
- $\text{Tar} : \text{Repr}_{\mathbb{L}}^{\mathbb{S}} \longrightarrow \mathbb{L}$  is the target functor. It takes a representation of  $\mathbb{S}$  in  $\mathbb{L}$  and returns the low-level representation of the given high-level inductive type. Can be considered the other bundle projection.
- $\mu : [\mathbb{S}, \mathbb{S}] \longrightarrow \mathbb{S}$  is the functor that takes an endofunctor in  $[\mathbb{S}, \mathbb{S}]$  and returns the least fixpoint of that endofunctor in  $\mathbb{S}$ .
- $\sigma : \mathbb{S} \longrightarrow \text{Repr}_{\mathbb{L}}^{\mathbb{S}}$  is the functor that takes an object in  $\mathbb{S}$  and returns the representation of that object in  $\mathbb{L}$ . It is a section of  $\mu\text{Src}$ .

Each object  $R \in \text{Repr}_{\mathbb{L}}^{\mathbb{S}}$  contains the following data:

- An object  $\mu\text{Src}_R \in \mathbb{S}$  which is the source type.
- An endofunctor  $\text{Src}_R : \mathbb{S} \longrightarrow \mathbb{S}$  whose fixpoint is  $\mu\text{Src}_R$ .
- An object  $\text{Tar}_R \in \mathbb{L}$  which is the low-level representation of  $\mu\text{Src}_R$ .
- An object  $\text{Int}_R \in \mathbb{S}$  which represents an intermediate descriptive object that captures information about  $\mu\text{Src}_R$  during a translation process. In simple cases the intermediate object will be something like  $\text{Int}_R = \mu\text{Src}_R + Q\text{Tar}_R$ .
- A morphism

$$\text{IntAlg}_R \in \mathbb{S}(\text{Src}_R(\text{Int}_R), \text{Int}_R)$$

which calculates the structure of the intermediate descriptive object by a fold over the syntactical term structure of the high-level program. For example, this can condense sequences of constructors of the high-level inductive type into a single constructor of the intermediate descriptive object, depending on what the algebra dictates.

- A morphism

$$\text{IntCoAlg}_R \in \prod_{L \in \text{Repr}_{\mathbb{L}}^{\mathbb{S}}} \mathbb{S}(\text{Int}_L^{\text{Src}_R(\text{Int}_R)}, \text{Int}_L^{\text{Int}_R})$$

which calculates the structure of a function on the intermediate object, given a function on the source object. It is almost a coalgebra, if the bound of the product was over all of  $\mathbb{S}$  rather than just the image of  $\text{Int}$ . Syntactically, it is used to transform pattern matches on the high-level inductive type into pattern matches on the intermediate descriptive object.

- A morphism

$$\text{Comp}_R \in \mathbb{S}(\text{Int}_R, Q\text{Tar}_R)$$

that "compiles" the intermediate descriptive object into the low-level representation (quoted).

- A morphism

$$\text{Decomp}_R \in \mathbb{S}(Q\text{Tar}_R, \text{Int}_R)$$

that "decompiles" the low-level representation into the intermediate descriptive object.

Each morphism  $f \in \text{Repr}_{\mathbb{L}}^{\mathbb{S}}(R, R')$  contains the following data:

- TODO: Basically a morphism of each of the above data

From all this data, we should be able to construct the following functions (morphisms in  $\text{Set}$ ):

- A function

$$\text{int}_{\sigma, \Gamma, T} : \mathbb{S}(\Gamma, T) \rightarrow \mathbb{S}(Q\text{Tar}_{\sigma\Gamma}, \text{Int}_{\sigma T})$$

which applies all the  $\text{IntAlg}$  and  $\text{IntCoAlg}$  morphisms by folding and unfolding over the syntactical term structure of the high-level program.

- A function

$$\text{comp}_{\sigma, \Gamma, T} : \mathbb{S}(Q\text{Tar}_{\sigma\Gamma}, \text{Int}_{\sigma T}) \rightarrow \mathbb{S}(Q\text{Tar}_{\sigma\Gamma}, Q\text{Tar}_{\sigma T})$$

which applies the  $\text{Comp}$  morphism to the intermediate descriptive object.

- A function

$$\text{unQ}_{\sigma, \Gamma, T} : \mathbb{S}(Q\text{Tar}_{\sigma\Gamma}, Q\text{Tar}_{\sigma T}) \rightarrow \mathbb{L}(\text{Tar}_{\sigma\Gamma}, \text{Tar}_{\sigma T})$$

which essentially unquotes the result, yielding a term in the low-level language.

Finally, for a given section  $\sigma$  of the bundle of representations, we should be able to define the functor  $C$  as follows:

$$\begin{aligned} C_{\sigma} &: \mathbb{S} \longrightarrow \mathbb{L} \\ C_{\sigma}(T) &:= \text{Tar}_{\sigma T} \\ C_{\sigma}(f) &:= \text{unQ}_{\sigma, \Gamma, T} \circ \text{comp}_{\sigma, \Gamma, T} \circ \text{int}_{\sigma, \Gamma, T}(f) \end{aligned}$$

NOTICE: We have not at all defined what the internal language of  $\mathbb{S}$  or  $\mathbb{L}$  looks like. An advantage of this formalism is that the compilation process in terms of algebras and coalgebras can be defined independently of the actual syntax and semantics of the languages (modulo some requirements such as cartesian closedness for  $\mathbb{S}$ ).

#### 4.4 Properties of the compilation functor

The compilation functor should be coherent with respect to operational reduction in both languages.

We can draw the following lax-commutative square:

$$\begin{array}{ccc}
 \mathbb{S} & \xrightarrow{\text{eval}} & \mathbb{S} \\
 \downarrow c & \nearrow & \downarrow c \\
 \mathbb{L} & \xrightarrow{\text{run}} & \mathbb{L}
 \end{array} \tag{2}$$

### 5 EXAMPLES

#### 5.1 Manual boxing and sequences in the simply-typed lambda calculus

We will define and work in the internal languages of  $\mathbb{S}$  and  $\mathbb{L}$ . By internal language we mean that the objects of each category are the types, and the morphisms are the terms within a given context.

For example, a morphism in  $\mathbb{S}$

$$f \in \mathbb{S}(\Gamma, A)$$

will correspond to a term inside a context in the internal language of  $\mathbb{S}$ ,

$$\Gamma \vdash f : A$$

#### 5.2 Definition of $\mathbb{S}$

$x, y$	(variables)
$l, m$	(labels)
$A, B ::= A \rightarrow B \mid A \times B \mid A + B \mid \mu x.A \mid 1 \mid (l, A) \mid QC$	(types)
$t, u, v ::= \lambda x.t \mid x \mid t \ u \mid (t, u) \mid p_1(t) \mid p_2(t) \mid \text{inl}(t) \mid \text{inr}(t) \mid \text{case}(t, x.u, y.v) \mid q(t) \mid \text{let}(t, u, x.v) \mid \text{letrec}(t, y.u, x.v) \mid (l, t) \mid \star$	(terms)

Why this language?

- We want to be able to express inductive data types  $\mu x.A$ , for lists, numbers, trees, etc.
- We also want to be able to handle quoted terms of the lower language:  $QC/q(t)$ , so that we can define translation functions in the language.
- We want to be able to explicitly label certain types and their inhabitants  $(l, A)/(l, t)$ , even though they might be functionally identical to some others. This is so that we can consider them as different objects in the category, and thus have the compilation functor produce different results for each one of them.
- We want to be able to express the usual constructs of a functional language: functions, pairs, sums, recursion, etc.

Still to do: typing rules, operational semantics.

Still to do: need an extra calculus of representations so that we can define the extra data from  $\text{Repr}_{\mathbb{L}}^{\mathbb{S}}$  i.e. algebras and coalgebras.

### 5.3 Definition of $\mathbb{L}$

$x, y$	(variables)
$n$	(finite natural numbers, word size)
$s ::= \emptyset \mid [s, t]$	(sequences)
$S ::= \emptyset \mid [S, C]$	(type sequences)
$C, D ::= C \multimap D \mid \Sigma x_n. C \mid W \mid S[n] \mid C^n \mid \Box C \mid \mathbf{I}$	(types)
$t, u ::= \lambda x. t \mid x \mid t u \mid \langle n, t \rangle \mid \text{letpair}(t, u, x. y. v) \mid n \mid \text{box}(t)$ $\mid \text{letbox}(t, u, x. v) \mid s \mid \text{letseq}(t, u, x. v) \mid \text{num}(t) \mid \star$	(terms)

Why this language?

- Once again we want a language based on the lambda calculus—for now we will not go too low-level.
- We want this language to be able to represent boxed types, because we want to have lower-level control over the memory representation of inductive data types (usually in the form  $(l, \mu x. A)$ ).
- We allow a restricted form of dependent types in the form of  $\Sigma x_n. C$  types parameterised over some word size  $W$ . This is because sequence types  $C^n$  are indexed by  $W$ , and if we have some sequence that is of some runtime size, we want to be able to represent the fact that some stored size is the size of that sequence.
- We want to be able to represent sequences  $C^n$  for the reason above. We can also have sequence types which can be used to model disjoint unions (same as unions in  $C$ ). The only difference here is we can represent tagged unions by using the  $\Sigma$  type:  $\Sigma x_2. [A, B][x]$  is kind of like  $A + B$  in the high-level language.
- Due to the presence of boxed types, we want the language to be linear. In other words we do not allow weakening or contraction of the context (unclear if we want to allow contraction or not—we do not have explicit destructors so maybe we really want an affine system).
- We don't care about labels here since that is only a concern during compilation.
- We have different let expressions for the different linear types, to be able to extract their inner data all at once or not at all.
- The category  $\mathbb{L}$  is still a monoidal closed category (specifically a symmetric monoidal closed category), so we still allow lambdas to capture variables (i.e. closures). Compiling these to boxed closures is a different matter.

Still to do: typing rules, operational semantics.

### 5.4 Lists to arrays and natural numbers to big unsigned integers

$$\begin{aligned} \text{List}(A) &:= \mu X. (A \times X + 1) && \text{(lists)} \\ \text{Array}(A) &:= \Sigma n. \Box A^n && \text{(arrays)} \end{aligned}$$

$$\begin{aligned} \text{Nat} &:= \mu X. (1 + X) && \text{(numbers)} \\ \text{BigUInt} &:= \Sigma n. \Box W^n && \text{(big unsigned integers)} \end{aligned}$$

Fully write out the representations of these types in  $\mathbb{S}$  and  $\mathbb{L}$ .  
Find less trivial examples.

### 5.5 $\mathbb{S}$ and $\mathbb{L}$ in one using dependent types

We should be able to embed the data of  $\text{Repr}_{\mathbb{L}}^{\mathbb{S}}$  inside  $\mathbb{S}$  if the latter has sufficient quantification structure. In other words, its internal language is some kind of dependent type theory.

## 6 RELATED WORK

TODO: BibLatex

Bit Stealing Made Legal: <https://dl.acm.org/doi/pdf/10.1145/3607858>

Type fusion: <https://www.cs.ox.ac.uk/ralf.hinze/publications/AMAST10.pdf>

Unrolling lists: <https://dl.acm.org/doi/10.1145/182409.182453>

An automatic object inlining optimization and its evaluation: <https://dl.acm.org/doi/10.1145/349299.349344>

Selection of representations for data structures: <https://dl.acm.org/doi/pdf/10.1145/872736.806944>

Linear/non-Linear Types For Embedded Domain-Specific Languages: <https://core.ac.uk/download/pdf/214213829.pdf>

Staged compilation with two-level type theory: <http://arxiv.org/abs/2209.09729>

## 7 CONCLUSION

### 7.1 Ideas about future work

- Algebra-coalgebra pairs as a way to interpret inductive data types and their recursive control-flow in low-level categories. (this work)
- Coherence conditions on an algebra-coalgebra pair to ensure that a chosen representation is faithful. (this work ?)
- Custom shortcuts for derived transformations based on the algebra-coalgebra pairs, for fine-tuning the representation of compound operations.
- Algebra-coalgebra pair generators for equivalence classes of isomorphic data types, to automatically generate representations of commonly seen structures.
- Solving for the representation that optimises some metric of a chosen set of operations (e.g. space complexity, time complexity, constant factors etc.) through various static and dynamic techniques
- Restriction of the context category of the source language in terms of its monoidal structure, to prevent certain low-level operations from being expressible at all.
- Relaxing well-foundedness, to model more complicated control-flow structures such as coroutines, continuations, and so on.

## ACKNOWLEDGMENTS

Acknowledgments

## 8 [TEMP] NOTES

## 9 DOCUMENT PLAN

We probably need to have the following sections:

- Introduction
- Preliminaries – category theory, conventions, basic structure of the problem
- Categorical model of (compilation?) Maybe it should be called something else
- An example with the simply-typed lambda calculus and a language with explicit boxing
- Embedding the representations into the high-level language – dependent types
- Related work
- Conclusion + future

- Would be nice to have an artifact: Agda formalisation?

Received 28 February 2024