

# General Stage Untangling with Applications to Algorithm Derivation

Nicolas Feltman et al.

December 11, 2013

## 1 Introduction

Often when programming high-performance systems, a programmer will wish to specialize a bivariate function  $f : \alpha \times \beta \rightarrow \gamma$  to one of its inputs,  $x : \alpha$ , while leaving the other,  $y : \beta$ , for later. This need may arise because  $x$  varies at a higher frequency than  $y$ , and we would wish to avoid repeated calculation. It may also be the case that  $x$  is available at a time when computational resources are less expensive than when  $y$  becomes available. It's even possible that knowledge of  $x$  allows further code optimizations that would have otherwise been unwise to speculate.

Consider two examples: the discrete Fourier transform (DFT) and a raytracer. In the DFT case we say that  $x : \mathbb{N}$  is the dimension of the transformation,  $y : \mathbb{R}^x$  is the data to be transformed, and  $f$  simply computes the DFT. Much of the structure of modern DFT algorithms depends on  $x$  (for both mathematical decomposition and cache performance reasons) but not on  $y$ , so a great deal of work could be done to specialize the algorithm to a particular number of dimensions. In the raytracer case, we say that  $x : \mathbf{Geom}$  is a specification of the scene geometry (perhaps a list of triangles),  $y : \mathbf{Ray}$  is a ray cast from some point in the scene, and  $f$  computes the first intersection of the ray and the geometry. In modern raytracers, we may cast millions of rays without varying the scene geometry, so it could be advantageous to precompute as much of the algorithm as possible based on only scene geometry.

We say that the problems above exhibit a *staged* structure, in that they can be broken down into parts that run at different points in time (*i.e.* stages). There are a variety of techniques to take advantage of this staged structure. The general theme is to do some work with the immediate input ( $x$ ), producing an intermediate result, and then finish the rest of the computation once the delayed input ( $y$ ) is available. We consider three techniques: program generation, partial evaluation, and pass separation. Each technique has broader applicability than the previous, but is less aggressive of an optimization.

### 1.1 Program Generation

A common solution is to write code that writes code. Essentially, second-stage computations are represented by second-stage code, and second-stage code is a value that can be passed around and manipulated by first-stage code. To use the terminology from above: this first stage code is a function that takes in  $x$  and outputs second stage function, which itself takes in  $y$  and outputs the final answer. Note that this is essentially a manual process, in that the programmer must know the function  $f$  which is being specialized. Using the equational notation of [7], where “ $\llbracket h \rrbracket z$ ” means the “the code  $h$  interpreted as a function and applied to input  $z$ ,” we say that a programmer defines a program generator  $g$  that satisfies

$$\llbracket \llbracket g \rrbracket x \rrbracket y = \llbracket f \rrbracket (x, y)$$

for all  $x$  and  $y$ . When well supported by the language ([3], [some metaml citation]) this technique is often called *metaprogramming*. Of course, even without first class support, one can represent the second stage code using strings or other data types. Defined so broadly, program generation includes all compilers, where  $f$  is a hypothetical interpreter,  $g$  is the actual compiler,  $x$  is a program in the source language, and  $y$  is the input to  $x$ . Program generation has also been used to implement the DFT example in the form of FFTW [4].

The benefit of the program generation technique is full control, even up to the ability to use domain-specific optimizations. This comes at the cost of essentially being a manual operation.

## 1.2 Partial Evaluation

The idea of partial evaluation is to *automatically* specialize the code of  $f$  to the supplied input, in a manner similar to standard evaluation. In particular this means defining a program  $pe$ , called a partial evaluator, to which we can pass the implementation of a function  $f$  along with  $f$ 's immediate input  $x$ , and get back a version of  $f$  that is *specialized* to  $x$ . This specialized version, also called a *residual*, is often denoted as  $f_x$ . Equationally, we have

$$\begin{aligned}\llbracket pe \rrbracket (f, x) &= f_x \\ \llbracket f_x \rrbracket y &= \llbracket f \rrbracket (x, y)\end{aligned}$$

for all  $f$ ,  $x$ , and  $y$ . As seminally observed by [5], partially evaluating an interpreter on a source program is equivalent to compiling that program to whatever language the interpreter was written in! This has the benefit of being an automatic process, although the effective compiler is only as good as the partial evaluator. Much of the subsequent work in partial evaluators has been with the purpose of chasing this goal.

Partial evaluation is a more general technique than program generation, in that the partial evaluator is defined once and works for all  $f$ s that need to be specialized. Of course, the partial evaluator has no domain knowledge of the program it is splitting, so it cannot be nearly as aggressive as a program generator might be. The two techniques are similar in that they both require the first input  $x$  to do their specializing action.

## 1.3 Pass Separation

The final technique, pass separation, is conceptually the simplest. The idea is to define a program called  $ps$ , for *pass separator*. It works by cleaving the function  $f : \alpha \times \beta \rightarrow \gamma$  into two functions  $f_1 : \alpha \rightarrow \tau$  and  $f_2 : \beta \times \tau \rightarrow \gamma$  for some type  $\tau$ , where  $f_1$  builds a data structure (of type  $\tau$ ) from  $x$ , and  $f_2$  consumes that data structure as well as  $y$  to produce the standard output. In analogy to partial evaluation, we call  $f_2$  the *residual*. Equationally, this is

$$\begin{aligned}\llbracket ps \rrbracket f &= (f_1, f_2) \\ \llbracket f_2 \rrbracket (\llbracket f_1 \rrbracket x, y) &= \llbracket f \rrbracket (x, y)\end{aligned}$$

for all  $f$ ,  $x$ , and  $y$ . This technique, performed manually, is a common exercise for every programmer. It can be considered the general form of common compiler optimizations such as loop hoisting, an example considered later. Like partial evaluation, pass separation is an automatic operation, but unlike both partial evaluation and program generation, pass separation does not require the value of the first argument,  $x$ , to properly specialize  $f$ . This makes pass separation the most widely applicable of the staging techniques, at the cost of being able to perform aggressive optimizations in the residual that might depend on the first argument.

In this work we present a pass separation algorithm for a fragment of ML. Specifically, our formulation can separate terms containing sums and recursion, which have not appeared previously in the stage separation literature. We anticipate that the algorithm will also be able to also separate this language extended with first-class functions, although this is left as

future work. Additionally, our presentation is more explicitly type-motivated than those that have come before.

## 2 Binding-Time Analysis and Stage Untangling

The pass separation problem can naturally be decomposed into two subproblems:

- Decide which parts of the definition of  $f$  belong to which stage and produce an annotated version of  $f$ .
- Use the annotated version of  $f$  to produce  $f_1$  and  $f_2$ .

An analogous form of this decomposition exists in the partial evaluation literature, wherein the first subproblem is known as *binding-time analysis*. The second subproblem, in which the difference between pass separation and partial evaluation is more manifest, does not seem to yet have a name. We call it *stage untangling*. Unlike previous work in pass separation, we focus entirely on the stage untangling problem by assuming that  $f$  is given to us in an annotated form.

Previous research has observed that the pass separation problem is inherently ambiguous in that the equations defined above do not fully specify the definition of  $pe$ . That is, there can exist multiple ways to partition  $f$  into  $f_1$  and  $f_2$  that satisfy  $\llbracket f_2 \rrbracket (\llbracket f_1 \rrbracket x, y) = \llbracket f \rrbracket (x, y)$ . For a trivial example, we can set  $f_1$  to be the identity and  $f_2$  to be  $f$ . Fortunately, this ambiguity is contained entirely within the binding-time analysis portion of pass separation; stage untangling is entirely determined.

## 3 $\lambda^{12}$

(Introduce the language, talk about its type system and maybe its evaluation.)

## 4 Examples

Some examples to get us comfortable with the language.

### 4.1 loop hoisting

## 5 Stage Separation Algorithm

Talk about the goal of the transformation. Introduce the general form with open variables and a context, which is *not* the simplistic  $f_1, f_2$  thing. Stress the first image/second image/boundary type stuff.

### 5.1 Type translation

### 5.2 Term translation

State the main theorems we want to be true. Refer to the

### 5.3 Functions

Talk about why the current set up for functions alleviates the boundary-type guessing problem.

## 5.4 Issue with sums

Sums result in dependant types, which is cool.

## 6 Implementation

Not much to discuss. It exists. Talk about how crufty the naive transformation is.

## 7 Related Work in Stage Separation

Our starting language,  $\lambda^{12}$ , is essentially a two-stage version of  $\lambda^\circ$  from [1], but beefed up with products, sums, and fixed-points. That work motivates  $\lambda^\circ$  as an image of linear temporal logic under the Curry-Howard correspondence, and it makes the argument that  $\lambda^\circ$  is a good model for binding time analysis. Although  $\lambda^\circ$  was designed with partial evaluation in mind, it turns out to be a good model for pass separation as well. In contrast, [2] presents  $\lambda^\square$ , which is based off of branching temporal logic. As that work shows,  $\lambda^\square$  is a good system for program generation, although we find the closed-code requirement to be more restricting than we need.

Pass separation was introduced in [8]. They focused on motivating pass separation as a technique for compiler generation, and hinted at wider applicability. Their approach was not mechanized or automated in any sense, and their examples were separated by hand. We believe that our work represents the fulfillment of their prediction that “the [pass separation] approach will elude full automation for some time.”

[6] also uses a pass separation technique for generating compiler/evaluator pairs from interpreters specified as term-rewrite systems. This heavy restriction on the form of the input prevents generalization of their style of pass separation. [I don’t know enough about term-rewrite systems yet to be able to comment further.] Their method, while amenable to mechanization, was not actually implemented.

[9] represents the first attempt at fully automatic pass separation for code approaching general-purpose. In particular, they implemented pass separation for a C-like shading language including basic arithmetic and if statements. The main goal of their work, like ours, was to minimize recomputation. That said, their starting language was sufficiently restricted for recursive boundary types to be inexpressible. Given that memory is at a premium in shading languages, so much of that work was also focused on minimizing the memory-footprint of the boundary type.

## References

- [1] Rowan Davies. A temporal-logic approach to binding-time analysis. In *LICS*, pages 184–195, 1996.
- [2] Rowan Davies and Frank Pfenning. A modal analysis of staged computation. *J. ACM*, 48(3):555–604, 2001.
- [3] Zachary DeVito, James Hegarty, Alex Aiken, Pat Hanrahan, and Jan Vitek. Terra: A multi-stage language for high-performance computing. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI ’13*, pages 105–116, New York, NY, USA, 2013. ACM.
- [4] Matteo Frigo and Steven G. Johnson. The design and implementation of FFTW3. *Proceedings of the IEEE*, 93(2):216–231, 2005. Special issue on “Program Generation, Optimization, and Platform Adaptation”.

- [5] Y. Futamura. Partial evaluation of computation process – an approach to a compiler-compiler. *Systems, Computers, Controls*, 2(5):45–50, 1971.
- [6] John Hannan. Operational semantics-directed compilers and machine architectures. *ACM Trans. Program. Lang. Syst.*, 16(4):1215–1247, 1994.
- [7] Neil D. Jones. An introduction to partial evaluation. *ACM Comput. Surv.*, 28(3):480–503, September 1996.
- [8] Ulrik Jørring and William L. Scherlis. Compilers and staging transformations. In *Proceedings of the 13th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, POPL '86, pages 86–96, New York, NY, USA, 1986. ACM.
- [9] Todd B. Knoblock and Erik Ruf. Data specialization. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation*, PLDI '96, pages 215–225, New York, NY, USA, 1996. ACM.