

Stage-Splitting a Modal Language

Name1

Affiliation1

Email1

Name2 Name3

Affiliation2/3

Email2/3

Abstract

This is the text of the abstract.

Categories and Subject Descriptors CR-number [subcategory]: third-level

General Terms term1, term2

Keywords keyword1, keyword2

1. Introduction

We often find ourselves in situations where the input to a computation arrives in two parts, at different points in time.

The natural response is to try and do some part of the computation early, and the rest of the computation once the second input arrives.

Since the computation runs in two parts, the total cost can be split as well, into $m + n$.

Additionally, we can run the second stage multiple times. Cost is now $m + bn$.

Jorring et al. ([5]) identify three classes of staging techniques: meta-programming, partial evaluation, and stage-splitting. The first two of these have received significantly more attention than the third. Countless meta-programming systems exist (...twelve thousand citations...[3]), and their background theory and type-systems are well understood ([2]). Partial evaluation, too, is well-understood. Partial evaluation systems exist... This paper explores both theory and applications for stage-splitting.

2. Stage-Splitting Definition and Comparison to Partial Evaluation

First, we review the definition of partial evaluation. Informally, a partial evaluator takes the code for a function f , as well as the first-stage input x to that function, and produces the code for a version of that function *specialized* to the first input, often called f_x . This f_x function can then be evaluated with the second-stage input to produce the same final answer that f would have. The goal of the process is that f_x should be cheaper to evaluate than f , although this can't be guaranteed for all inputs. We now state this theorem more formally: a partial evaluator is some function p such that,

$$\forall f, x. \exists f_x. [p(f, x) = f_x \text{ and } \forall y. \llbracket f \rrbracket(x, y) = \llbracket f_x \rrbracket(y)]$$

where (borrowing notation from [4]), $\llbracket f \rrbracket$ means the mathematical function corresponding to the code given by f .

Informally, we define stage-splitting to be the process of taking some function f into two other functions, f_1 and f_2 , where f_1 computes a partial result from the first-stage input, and f_2 uses that partial result and the second-stage input to compute a final result which is the same as if we had just run the original f on both inputs. Again, more formally, a stage-splitter is some s such that,

$$\forall f. \exists f_1, f_2. [s(f) = (f_1, f_2) \text{ and } \forall x, y. \llbracket f \rrbracket(x, y) = \llbracket f_2 \rrbracket(\llbracket f_1 \rrbracket(x), y)]$$

We first discuss a few similarities between partial-evaluation and stage-splitting. First off, both techniques have the same form of input, namely a bivariate function where the first input comes at stage one, and the second input comes at stage two.

Again in both cases, the governing equations are too weak to fully determine the definitions of p and s . Indeed, both admit completely trivial definitions. Consider the stage-splitter which always returns the identity for f_1 and f for f_2 , or analogously the partial evaluator which always returns an f_x that just closes over the input x and internally calls f once y is available. The ambiguity of these equations (modulo standard program equivalence of the outputs) can be resolved by adding annotations to f to clearly specify the parts of the computation that are first stage and the parts that are second stage. Later, we show that the same annotations suffice for both partial evaluation and stage-splitting.

The differences between stage-splitting and partial evaluation are likewise evident from these governing equations. For instance in partial evaluation, the existential f_x depends on x , which means that the partial evaluator cannot be run until x is known. Moreover, if one wishes to specialize f for multiple x 's, then the partial evaluator must be run several times. Depending on the use case and cost of partial evaluation, this may be prohibitively expensive. Alternatively, a stage-splitter need only be run once, and this can be done entirely before any x is known.

2.1 Partial Evaluator from Stage-Splitter

We can recover a valid partial evaluator from a stage-splitter by stage-splitting the input function f into f_1 and f_2 , computing $\llbracket f_1 \rrbracket(x)$ to obtain \bar{x} , and then returning an f_x such that $\llbracket f_x \rrbracket(y) = \llbracket f_2 \rrbracket(\bar{x}, y)$. Note that this does not mean that stage-splitting is a strict generalization of partial evaluation. In practice, partial evaluators easily perform optimizations (such as branch elimination, discussed later) which are beyond the scope of stage-splitting, and would require further technology than has been developed here. It is best to think of stage-splitting as simply the first half of partial evaluation, where the back half is an optimizer. [Might be able to come up with a futamura projection-like statement here, which would be really really really cool.]

2.2 Stage-Splitter from Partial Evaluator

Likewise, we can easily recover a stage-splitter from a partial evaluator. If p is a valid partial evaluator, then we can define a stage

$$\begin{aligned}
w &::= 1 \mid 2 \\
\tau &::= \text{unit} \mid \text{int} \mid \text{bool} \mid \tau \times \tau \mid \tau + \tau \\
&\quad \mid \tau \rightarrow \tau \mid \bigcirc \tau \mid \alpha \mid \mu \alpha. \tau \\
e &::= () \mid i \mid b \mid x \mid \lambda x : \tau. e \mid e e \\
&\quad \mid (e, e) \mid \text{pi1 } e \mid \text{pi2 } e \mid \text{inl } e \mid \text{inr } e \\
&\quad \mid \text{case } e \text{ of } x.e \mid x.e \mid \text{if } e \text{ then } e \text{ else } e \\
&\quad \mid \text{let } x = e \text{ in } e \mid \text{next } e \mid \text{prev } e \mid \text{hold } e \\
\Gamma &::= \bullet \mid \Gamma, x : \tau @ w
\end{aligned}$$

Figure 1. λ^{12} Syntax

$$\begin{aligned}
&\frac{}{\text{unit type @ } w} \text{unit} & \frac{}{\text{int type @ } w} \text{int} \\
&\frac{}{\text{bool type @ } w} \text{bool} & \frac{A \text{ type @ } w \quad B \text{ type @ } w}{A \times B \text{ type @ } w} \times \\
&\frac{A \text{ type @ } w \quad B \text{ type @ } w}{A + B \text{ type @ } w} + \\
&\frac{A \text{ type @ } w \quad B \text{ type @ } w}{A \rightarrow B \text{ type @ } w} \rightarrow & \frac{A \text{ type @ } 2}{\bigcirc A \text{ type @ } 1} \bigcirc
\end{aligned}$$

Figure 2. λ^{12} Valid Types

$$\begin{aligned}
&\frac{}{() \text{ val @ } w} \text{unit} & \frac{}{i \text{ val @ } w} \text{int} & \frac{}{b \text{ val @ } w} \text{bool} \\
&\frac{v_1 \text{ val @ } w \quad v_2 \text{ val @ } w}{(v_1, v_2) \text{ val @ } w} \times & \frac{v \text{ val @ } w}{\text{inl } v \text{ val @ } w} +_1 \\
&\frac{v \text{ val @ } w}{\text{inr } v \text{ val @ } w} +_2 & \frac{}{\lambda x : A. e \text{ val @ } w} \rightarrow \\
&\frac{}{\Xi \vdash \hat{y} \text{ val @ } 1} \bigcirc
\end{aligned}$$

Figure 3. λ^{12} Valid Values

splitter s such that $s(f) = (f_1, f_2)$, where

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

This implicitly requires that the languages in which f_1 and f_2 are expressed are strong enough to write a partial evaluator, but that is the case in this paper. A stage-splitter defined this way leaves much to be desired. Firstly, partial evaluation of f may be too expensive for the context in which f_1 needs to run. Additionally, the intermediate data structure created this way may be much larger than necessary, as it would contain all of the residual code.

3. A Two-Stage Modal Language

Every expression in λ^{12} can be understood as having a *stage*, either 1 or 2, in addition to a type. Just as a type describes *which* values an expression reduces to, the stage tells us *when* it will be fully reduced. The reason for a two-stage design (as opposed to [1], which allows an infinite number of stages) will be discussed later.

The only types that exist at stage 2 are products, sums, functions, and base types. In addition to this usual set, stage 1 also has one more type: $\bigcirc A$, where A is a stage-2 type. This can be understood as *promising* a value of type A at stage 2. We'll see shortly that the promise given by $\bigcirc A$ can only be treated opaquely at stage one; it cannot be redeemed for an actual A until the second stage.

Although there is a strict one-way dependence between the two stages at the type level, their terms are defined mutually recursively: **next** embeds stage-2 expressions into stage 1, whereas **prev** embeds stage-1 expressions into stage 2. These **next** and **prev** constructs are the only expressions that interact with the stage of a term, and we surround their arguments with $\{\}$ braces in the concrete syntax to clearly indicate the stage boundaries with in a program. The problem of how to turn an unstaged program into a staged one by inserting these constructs is called *binding time analysis*. We do not consider that problem here, and instead assume that it has been solved either automatically or by a programmer.

The grammar and type system for our language, λ^{12} , are given in Figures 1 and 4. We annotate the typing judgment and context variables with stages; only \bigcirc and its introductory and eliminatory forms **next** and **prev** affect the stage of a term or type.

Specifically, **next**, given an argument with type A at stage 2, forms a $\bigcirc A$ at stage 1. That is, it forms the promise of a future A out of a construction for an A at the next timestep. This promise can only be redeemed by the **prev** construct. Since **prev** operates at stage 2, this ensures no violation of causality.

The **hold** construct is a mechanism to wrap stage 1 integers so that they can be used at stage 2. Although **hold** can be implemented from our other features, we provide it as a primitive to simplify our examples.

3.1 Staged Evaluation

Fast exponentiation is a method for calculating b^p in $\log p$ time:

$$fexp(b, p) = \begin{cases} 1 & p = 0 \\ fexp(b, p/2)^2 & p \text{ even} \\ b \cdot fexp(b, p-1) & p \text{ odd} \end{cases}$$

Ignoring for the moment how to define recursive functions in λ^{12} , we can implement staged fast exponentiation as follows:¹

```

let fexp (b : $int, e : int) : $int =
  if e == 0 then
    next{1}
  else if (e mod 2) == 0 then
    next{let x = prev{fexp(b,p/2)} in x*x}
  else
    next{prev{b} * prev{fexp(b,p-1)}}

```

Although the code for **fexp** looks very much like the unstaged mathematical definition above, it is in fact a staged program: the **ifs** and exponent decomposition are all at stage 1, since they occur within **prev** blocks. Essentially, the type system is sufficient to prove the observation that all of the decomposition of the exponent can be done at the stage 1 without requiring us to disrupt the elegant functional structure of the code.

This example also illustrates how stage 1 and stage 2 code can be arbitrarily nested. Ordinary evaluation always eliminates outermost redexes first, but this conflicts with the precept that we evaluate all stage 1 code before stage 2 code, in the case that stage 1 code is inside stage 2 code.

A suitably staged dynamic semantics for λ^{12} must therefore evaluate *all* the stage 1 subexpressions of a term before its stage 2 subexpressions.

¹ We render \bigcirc as \$ in code examples.

$$\begin{array}{c}
\frac{}{\Gamma \vdash () : \text{unit} @ w} \text{unit} \quad \frac{}{\Gamma \vdash i : \text{int} @ w} \text{int} \quad \frac{}{\Gamma \vdash b : \text{bool} @ w} \text{bool} \quad \frac{x : A @ w \in \Gamma}{\Gamma \vdash x : A @ w} \text{hyp} \\
\\
\frac{A \text{ type} @ w \quad \Gamma, x : A @ w \vdash e : B @ w}{\Gamma \vdash \lambda x : A. e : A \rightarrow B @ w} \rightarrow I \quad \frac{\Gamma \vdash e_1 : A \rightarrow B @ w \quad \Gamma \vdash e_2 : A @ w}{\Gamma \vdash e_1 e_2 : B @ w} \rightarrow E \\
\\
\frac{\Gamma \vdash e_1 : A @ w \quad \Gamma, x : A @ w \vdash e_2 : B @ w}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : B @ w} \text{let} \quad \frac{\Gamma \vdash e_1 : A @ w \quad \Gamma \vdash e_2 : B @ w}{\Gamma \vdash (e_1, e_2) : A \times B @ w} \times I \quad \frac{\Gamma \vdash e : A \times B @ w}{\Gamma \vdash \text{pi1 } e : A @ w} \times E_1 \\
\\
\frac{\Gamma \vdash e : A \times B @ w}{\Gamma \vdash \text{pi2 } e : B @ w} \times E_2 \quad \frac{\Gamma \vdash e : A @ w}{\Gamma \vdash \text{inl } e : A + B @ w} +I_1 \quad \frac{\Gamma \vdash e : B @ w}{\Gamma \vdash \text{inr } e : A + B @ w} +I_2 \\
\\
\frac{\Gamma \vdash e_1 : A + B @ w \quad \Gamma, x_2 : A @ w \vdash e_2 : C @ w \quad \Gamma, x_3 : B @ w \vdash e_3 : C @ w}{\Gamma \vdash \text{case } e_1 \text{ of } x_2.e_2 \mid x_3.e_3 : C @ w} +E \\
\\
\frac{\Gamma \vdash e_1 : \text{bool} @ w \quad \Gamma \vdash e_2 : A @ w \quad \Gamma \vdash e_3 : A @ w}{\Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : A @ w} \text{if} \quad \frac{\Gamma \vdash e : A @ 2}{\Gamma \vdash \text{next } e : \bigcirc A @ 1} \bigcirc I \quad \frac{\Gamma \vdash e : \bigcirc A @ 1}{\Gamma \vdash \text{prev } e : A @ 2} \bigcirc E \\
\\
\frac{\Gamma \vdash e : \text{int} @ 1}{\Gamma \vdash \text{hold } e : \bigcirc \text{int} @ 2} \text{hold}
\end{array}$$

Figure 4. λ^{12} Static Semantics

Consider how we could achieve this in the following stage 2 term:

```

prev{
let x = (next {4+5}, 7+8) in
next{
  prev{#1 x} * prev{#1 x} * prev{hold (#2 x)}
}
}

```

One option is to duplicate the contents of the first `next` expression, to produce

`(4+5) * (4+5) * 15`

Here we are essentially treating `next{e}` as a stage 1 value, so long as all the `prev` subterms of `e` are reduced. This is the approach taken by [1].

We take a different approach in our dynamic semantics, which is to let-bind the contents of the first `next` expression to some variable (here, `x1`) and duplicate that variable,

```

let x1 = 4+5 in
x1 * x1 * 15

```

We think that this interpretation makes more sense from the point of view of cost, since the contents of a `next` always appear once in the residual, regardless of how many times the result is used. This is *nearly* equivalent to the first expression, because they have different termination behavior when a non-terminating stage-2 expression is declared but never used.

How do our semantics achieve this behavior? In the original term,

```

prev {
let x = (next {4+5}, 7+8) in
next{
  prev{#1 x} * prev{#1 x} * prev{hold (#2 x)}
}
}

```

we must start by evaluating the binding for `x`. We can't evaluate inside the `next` statement, but must still substitute some result

into the body of the `let`. We replace the contents of `next` with a variable and make an explicit substitution, which then floats up to the top of the containing `prev`:

```

prev{
[4+5/x1]
let x = (next{x1}, 15) in
next{
  prev{#1 x} * prev{#1 x} * prev{hold (#2 x)}
}
}

```

We're now free to perform the stage-1 substitution for `x`.

```

prev{
[4+5/x1]
next{
  prev{#1 (next{x1}, 15)}
  * prev{#1 (next{x1}, 15)}
  * prev{hold (#2 (next{x1}, 15))}
}
}

```

To evaluate the remaining `next`, we must first partially evaluate the body by finding all of the contained stage-1 terms and reducing them:

```

prev{
[4+5/x1]
next{ x1 * x1 * 15 }
}

```

The final step in evaluating the outer `prev` is reifying the contained substitutions into `let` statements, yielding

```

let x1 = 4+5 in
x1 * x1 * 15

```

Once again, the point here is that we did not duplicate the `4 + 5` in our residual. Instead we lifted it out into a substitution, leaving a hatted variable [**TODO**: add hats above] in its place, and we reified the substitution into a `let`-binding at the enclosing `prev`.

$$\begin{array}{c}
\frac{\cdot}{\Xi; \Gamma \vdash \lambda x:A.e \Downarrow_1 [\cdot, \lambda x.e]} \rightarrow I \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, \lambda x.e'] \quad \Xi; \Gamma \vdash e_2 \Downarrow_1 [\xi_2, v_2] \quad \Xi, \text{dom}(\xi_1), \text{dom}(\xi_2); \Gamma \vdash [v_2/x]e' \Downarrow_1 [\xi', v']}{\Xi; \Gamma \vdash e_1 e_2 \Downarrow_1 [\xi_1 \circ \xi_2 \circ \xi', v']} \rightarrow E \Downarrow_1 \\
\frac{\cdot}{\Xi; \Gamma \vdash () \Downarrow_1 [\cdot, ()]} \text{unit} \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, v_1] \quad \Xi; \Gamma \vdash e_2 \Downarrow_1 [\xi_2, v_2]}{\Xi; \Gamma \vdash (e_1, e_2) \Downarrow_1 [\xi_1 \circ \xi_2, (v_1, v_2)]} \times I \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, (v_1, v_2)]}{\Xi; \Gamma \vdash \text{pi1 } e \Downarrow_1 [\xi, v_1]} \times E_1 \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, (v_1, v_2)]}{\Xi; \Gamma \vdash \text{pi2 } e \Downarrow_1 [\xi, v_2]} \times E_2 \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, v]}{\Xi; \Gamma \vdash \text{inl } e \Downarrow_1 [\xi, \text{inl } v]} + I_1 \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, v]}{\Xi; \Gamma \vdash \text{inr } e \Downarrow_1 [\xi, \text{inr } v]} + I_2 \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, \text{inl } v] \quad \Xi, \text{dom}(\xi_1); \Gamma \vdash [v/x_2]e_2 \Downarrow_1 [\xi_2, v_2]}{\Xi; \Gamma \vdash \text{case } e_1 \text{ of } x_2.e_2 \mid x_3.e_3 \Downarrow_1 [\xi_1 \circ \xi_2, v_2]} + E_1 \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, \text{inr } v] \quad \Xi, \text{dom}(\xi_1); \Gamma \vdash [v/x_3]e_3 \Downarrow_1 [\xi_3, v_3]}{\Xi; \Gamma \vdash \text{case } e_1 \text{ of } x_2.e_2 \mid x_3.e_3 \Downarrow_1 [\xi_1 \circ \xi_3, v_3]} + E_2 \Downarrow_1 \quad \frac{\cdot}{\Xi; \Gamma \vdash i \Downarrow_1 [\cdot, i]} \text{int} \Downarrow_1 \\
\frac{\cdot}{\Xi; \Gamma \vdash b \Downarrow_1 [\cdot, b]} \text{bool} \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, v_1] \quad \Xi, \text{dom}(\xi_1); \Gamma \vdash [v_1/x]e_2 \Downarrow_1 [\xi_2, v_2]}{\Xi; \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow_1 [\xi_1 \circ \xi_2, v_2]} \text{let} \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, \text{true}] \quad \Xi; \Gamma \vdash e_2 \Downarrow_1 [\xi_2, v]}{\Xi; \Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \Downarrow_1 [\xi_1 \circ \xi_2, v]} \text{if}_T \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e_1 \Downarrow_1 [\xi_1, \text{false}] \quad \Xi; \Gamma \vdash e_3 \Downarrow_1 [\xi_3, v]}{\Xi; \Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \Downarrow_1 [\xi_1 \circ \xi_3, v]} \text{if}_F \Downarrow_1
\end{array}$$

Figure 5. λ^{12} Dynamic Semantics: Core

$$\begin{array}{c}
\frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{next } e \Downarrow_1 [z \mapsto \text{let } y = q \text{ in } z, \hat{y}]} \circ I \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, \hat{y}] \quad [\xi, y] \xRightarrow{R} q}{\Xi; \Gamma \vdash \text{prev } e \Downarrow_2 q} \circ E \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, i]}{\Xi; \Gamma \vdash \text{hold } e \Downarrow_1 [z \mapsto \xi(\text{let } y = i \text{ in } z), \hat{y}]} \text{hold} \Downarrow_1 \quad \frac{\cdot}{\Xi; \Gamma \vdash \hat{y} \Downarrow_1 [\cdot, \hat{y}]} \text{hattedvar} \Downarrow_1 \\
\frac{[\xi, q_2] \xRightarrow{R} q'}{[(y \mapsto q_1) \circ \xi, q_2] \xRightarrow{R} \text{let } y = q_1 \text{ in } q'} \quad \frac{\cdot}{[\cdot, q] \xRightarrow{R} q}
\end{array}$$

Figure 6. λ^{12} Dynamic Semantics: next and prev

$$\begin{array}{c}
\frac{\Xi; \Gamma, x \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \lambda x:A.e \Downarrow_2 \lambda x:A.q} \rightarrow I \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e_1 \Downarrow_2 q_1 \quad \Xi; \Gamma \vdash e_2 \Downarrow_2 q_2}{\Xi; \Gamma \vdash e_1 e_2 \Downarrow_2 q_1 q_2} \rightarrow E \Downarrow_2 \quad \frac{\cdot}{\Xi; \Gamma \vdash () \Downarrow_2 ()} \text{unit} \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_2 q_1 \quad \Xi; \Gamma \vdash e_2 \Downarrow_2 q_2}{\Xi; \Gamma \vdash (e_1, e_2) \Downarrow_2 (q_1, q_2)} \times I \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{pi1 } e \Downarrow_2 \text{pi1 } q} \times E_1 \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{pi2 } e \Downarrow_2 \text{pi2 } q} \times E_2 \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{inl } e \Downarrow_2 \text{inl } q} + I_1 \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{inr } e \Downarrow_2 \text{inr } q} + I_2 \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_2 q_1 \quad \Xi; \Gamma, x_2 \vdash e_2 \Downarrow_2 q_2 \quad \Xi; \Gamma, x_3 \vdash e_3 \Downarrow_2 q_3}{\Xi; \Gamma \vdash \text{case } e_1 \text{ of } x_2.e_2 \mid x_3.e_3 \Downarrow_2 \text{case } q_1 \text{ of } x_2.q_2 \mid x_3.q_3} + E_1 \Downarrow_2 \quad \frac{\cdot}{\Xi; \Gamma \vdash i \Downarrow_2 i} \text{int} \Downarrow_2 \quad \frac{\cdot}{\Xi; \Gamma \vdash b \Downarrow_2 b} \text{bool} \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e_1 \Downarrow_2 q_1 \quad \Xi; \Gamma, x \vdash e_2 \Downarrow_2 q_2}{\Xi; \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow_2 \text{let } x = q_1 \text{ in } q_2} \text{let} \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e_1 \Downarrow_2 q_1 \quad \Xi; \Gamma \vdash e_2 \Downarrow_2 q_2 \quad \Xi; \Gamma \vdash e_3 \Downarrow_2 q_3}{\Xi; \Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \Downarrow_2 \text{if } q_1 \text{ then } q_2 \text{ else } q_3} \text{if}_T \Downarrow_2
\end{array}$$

Figure 7. λ^{12} Dynamic Semantics: Speculation

3.2 Dynamics

The algorithm described above creates three different kinds of expressions which cannot be evaluated further at a particular stage:

- Partial values (pvals) are stage 1 terms whose stage 1 subexpressions have all been fully evaluated. In the example above, `next {x1 * x1 * 15}`.
- Residuals (reses) are stage 2 terms whose stage 1 subexpressions have all been fully evaluated. In the example above, `4+5`.
- Values (vals) are stage 2 terms which are fully evaluated; these are the results of a computation after both stages have been completed. In the example above, the term evaluates to 1215.

The \Downarrow_1 judgment takes an open stage 1 term to a pair of an environment ξ mapping newly-created variables to stage 2 residuals, and a partial value v which may refer to those variables. This judgment necessarily traverse stage 2 subexpressions in order to find their stage 1 subexpressions; we call this process *speculation*. Speculation is implemented by the \Downarrow_2 judgment, which takes a stage 2 term to a residual.

TODO: finish explaining how the nested evaluation actually works. explain reification, hatted variables, and contexts.

A function f which accepts inputs at both stages 1 and 2 can be given a type of the form $A \rightarrow \bigcirc(B \rightarrow C)$.² Given its stage 1 argument $a : A$, we can evaluate the partially-applied function: $\cdot; \cdot \vdash f a \Downarrow_1 [\xi, v]$. The result is an environment ξ and a partial value v of type $\bigcirc(B \rightarrow C)$. Next, we reify this environment into a sequence of `let`-bindings enclosing v , via $[\xi, v] \xRightarrow{R} f_a$. A canonical forms theorem on v ensures that the resulting f_a has type $B \rightarrow C$ in λ^2 . Finally, given a stage 2 argument $b : B$, we can stage-2 evaluate the ultimate result of the function, $f_a b \Downarrow_2 c$.

TODO: make sure this is true (haven't written the canonical forms theorems yet)

That this sequence of evaluations is in fact staged follows from our characterizations of partial values, residuals, and values, that \Downarrow_1 outputs a partial value, and that \xRightarrow{R} outputs an expression in λ^2 .

Remark 3.1. For any $e : A @ 1$ containing no `next` subexpressions, \Downarrow_1 will always compute an empty environment, and a partial value identical to the result of call-by-value evaluation of e .

TODO: If we change the highlighted rules to something slightly different, then we get precisely the semantics from [1]. This essentially bypasses the environment bookkeeping in \Downarrow_1 , by inlining residuals instead of hoisting them in `let`-bindings.

3.3 Metatheory

TODO: terminology for the contexts?

Definition 3.2. Contexts Ξ and Γ are well-formed (Ξ wf, Γ wf) if they contain only stage-2 variables.

Definition 3.3. An environment ξ is well-formed (ξ wf) if either:

1. $\xi = \cdot$; or
2. $\xi = \xi', x \mapsto e$ where $\Xi, \text{dom}(\xi'); \Gamma \vdash e : B @ 2$ and $\Xi, \text{dom}(\xi'); \Gamma \vdash e \text{ 1-val} @ 2$.

Theorem 3.4. If $\Xi; \Gamma \vdash e : A @ 1$ then Ξ wf, Γ wf, and A type $@ 1$.

Theorem 3.5. If $\Xi; \Gamma \vdash e \Downarrow_1 [\xi, v]$ then:

1. $\Gamma \vdash e : A @ 1$ for some 1-type A ;

2. ξ wf;
3. $\Xi, \text{dom}(\xi); \cdot \vdash v : A @ 1$; and
4. $\Xi, \text{dom}(\xi); \cdot \vdash v \text{ 1-val} @ 1$.

Theorem 3.6. If $\Xi; \Gamma \vdash e \Downarrow_2 q$ then:

1. $\Gamma \vdash e : A @ 2$ for some 2-type A ;
2. $\Xi; \Gamma \vdash q : A @ 2$; and
3. $\Xi; \Gamma \vdash q \text{ 2-val} @ 2$.

² We can rewrite `fexp` in this form, or simply apply the following higher-order function which makes the adjustment:

TODO: write the function here

4. Splitting Algorithm

[Present the splitting judgement. Give statements of type and value correctness for splitting. Give all of the splitting rules. Talk through a few of them.]

5. Examples for Staged Pipelines

Give the gist of one-to-one pipeline example (like client/server). Then talk about a one-to-many pipeline. Then talk about a many-to-one pipeline like spark. It clear how to target something like this for known base types on the boundary, and for product types. But sums on the boundary are hard! We leave many-to-one as future-work.

6. Examples of Algorithm Derivation

Fast exponent example.

```
let exp (b : $int, e : int) : $int =
  if e == 0 then
    next{1}
  else if (e mod 2) == 0 then
    next{let x = prev{exp(b,e/2)} in x*x}
  else
    next{prev{b} * prev{exp (b,e-1)}}
```

splits into

```
let exp (b, e) =
  (((), roll (
    if e == 0 then
      inL ()
    else
      inR (
        if (e mod 2) == 0 then
          inL (#2 (exp (b,e/2)))
        else
          inR (#2 (exp (b,e-1)))
        )
      ))
  )
```

and

```
let exp ((b, e), p) =
  case unroll p of
    () => 1
  | d =>
    case d of
      r => let x = exp ((b,()),r) in x*x
    | r => b * exp ((b,()),r)
```

Quickselect example.

```
let qs (l :  $\mu\alpha.()$  + int *  $\alpha$ , i: $int) =
  case unroll l of
    () => next {0}
  | (h,t) =>
    let (left,right,n) = partition h t in
    next{
      let n = prev{hold n} in
      case compare prev{i} n of
        () (*LT*) => prev {qs left i}
      | () (*EQ*) => prev {hold h}
      | () (*GT*) =>
        prev {qs right next{prev{i}-n-1}}
    }
}
```

Things to try: an interpreter which, partially evaluated, does cps or something.

For each of these examples, talk about what partial evaluation would do and why that might be bad.

[Meta-ML eases off on this restriction but does not (I think?) eliminate it.]

[What's going on with names and necessity?]

[Our work bears a lot of similarity to ML5, which also uses a modal type system. The difference is that we target stages systems (each stage talks to the next), whereas they target distributed ones (all stages talk to all others). The type systems reflect this directly in the world accessibility relation. There might be some analogue of stage-splitting in the ML5 work, but I have not yet isolated it (might be buried in CPS conversion).]

7. Related Work

Our stage-splitting algorithm was first suggested in [5] under the name *pass separation*. They essentially proposed that a function f could be split into two others, f_1 and f_2 , such that $f(x,y) = f_2(f_1(x),y)$. They did not distinguish binding time analysis from stage splitting, and so pass separation inherits the former's ambiguity. The main goal of [5] was to motivate pass separation and other staging transformations as a powerful way to think about compilation and optimization. Accordingly, their approach was entirely informal, with no implementation realized. Moreover, they predicted that "the [pass separation] approach will elude full automation for some time."

Implementations of the stage-splitting algorithm have appeared in the literature exclusively (and coincidentally) in the context of graphics pipelines. The first of these ([6]) uses a binding time analysis to separate those parts of graphics shaders that are input-invariant from those which are not, and then uses a stage splitting algorithm to factor that into two shaders, thereby minimizing recomputation. Their shaders are written in a C-like language with basic arithmetic and if statements. Although their analysis does not give an explicit account of the type-level behavior of the splitting algorithm, it effectively can synthesize product and sum boundary type.

Like the previous example, the Spark language ([7]) uses staging to minimize recomputation in real-time rendering applications. But instead of using a binding-time analysis, Spark allows the programmer to manually target stages of the graphics pipeline. Since the modern graphics pipeline is inherently a many-to-one system, this is difficult to reconcile with sum types on the boundary. Fortunately, Spark has a set of syntactic restrictions that prevent sum boundary types. Spark does not clearly identify this conflict, but the authors did note that first-stage if statements were difficult to provide meaning to [need quote].

[RTSL and SH]

[Discuss Yong's recent paper here. It does some pretty sophisticated binding time analysis, with a somewhat straightforward splitting after that. They have the same many-to-one use case as Spark, but syntactic restrictions prevent sum types on the boundary, sort of. If we wanted to faithfully represent their system in ours, we would need some mechanisms for abstraction over stage, which we do not have.]

Davies ([1]) explored the connection between linear temporal logic and its corresponding type system [circle](which we adapted into [circle sub 2]), and showed the equivalence between [circle] and existing systems for binding time analysis. That work provided β and η rules for the next and prev operators, but did not consider a full dynamic semantics for the whole language. Whereas [name of our type system] is appropriate for stage-splitting and partial evaluation, [2] provides a similar system, [square], that is appropriate for meta-programming. The main difference is that terms inside a [prev] operator do not see any stage-2 bindings declared outside of

$$\begin{array}{c}
\frac{}{\Gamma \vdash () : \text{unit} \xrightarrow{1} [((), ()), \text{--}()] } \text{unit} \xrightarrow{1} \quad \frac{}{\Gamma \vdash () : \text{unit} \xrightarrow{2} [((), \text{--}())] } \text{unit} \xrightarrow{2} \quad \frac{}{\Gamma \vdash i : \text{int} \xrightarrow{1} [(i, ()), \text{--}()] } \text{int} \xrightarrow{1} \\
\frac{}{\Gamma \vdash i : \text{int} \xrightarrow{2} [((), \text{--}i)] } \text{int} \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : A \xrightarrow{1} [c_1, l_1.r_1] \quad \Gamma \vdash e_2 : B \xrightarrow{1} [c_2, l_2.r_2]}{\Gamma \vdash (e_1, e_2) : A \times B \xrightarrow{1} \left[\begin{array}{l} (\text{let } (y_1, z_1) = c_1 \text{ in let } (y_2, z_2) = c_2 \text{ in } ((y_1, y_2), (z_1, z_2))), \\ (l_1, l_2).(r_1, r_2) \end{array} \right]} \times I \xrightarrow{1} \\
\frac{\Gamma \vdash e_1 : A \xrightarrow{2} [p_1, l_1.r_1] \quad \Gamma \vdash e_2 : B \xrightarrow{2} [p_2, l_2.r_2]}{\Gamma \vdash (e_1, e_2) : A \times B \xrightarrow{2} [(p_1, p_2), (l_1, l_2).(r_1, r_2)]} \times I \xrightarrow{2} \\
\frac{\Gamma \vdash e : A \times B \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{pi1 } e : A \xrightarrow{1} [\text{let } (y, z) = c \text{ in } (\text{pi1 } y, z), l.\text{pi1 } r]} \times E_1 \xrightarrow{1} \\
\frac{\Gamma \vdash e : A \times B \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{pi2 } e : B \xrightarrow{1} [\text{let } (y, z) = c \text{ in } (\text{pi2 } y, z), l.\text{pi2 } r]} \times E_2 \xrightarrow{1} \quad \frac{\Gamma \vdash e : A \times B \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{pi1 } e : A \xrightarrow{2} [p, l.\text{pi1 } r]} \times E_1 \xrightarrow{2} \\
\frac{\Gamma \vdash e : A \times B \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{pi2 } e : B \xrightarrow{2} [p, l.\text{pi2 } r]} \times E_2 \xrightarrow{2} \quad \frac{x : A @ 1 \in \Gamma}{\Gamma \vdash v : A \xrightarrow{1} [(x, ()), \text{--}()] } \text{hyp} \xrightarrow{1} \quad \frac{x : A @ 2 \in \Gamma}{\Gamma \vdash v : A \xrightarrow{2} [((), \text{--}x)] } \text{hyp} \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : A \xrightarrow{1} [c_1, l_1.r_1] \quad \Gamma, x : A @ 1 \vdash e_2 : B \xrightarrow{1} [c_2, l_2.r_2]}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : B \xrightarrow{1} \left[\begin{array}{l} \left(\begin{array}{l} \text{let } (x, z_1) = c_1 \text{ in} \\ \text{let } (y, z_2) = c_2 \text{ in} \\ (y, (z_1, z_2)) \end{array} \right), \\ (l_1, l_2).\text{let } x = r_1 \text{ in } r_2 \end{array} \right]} \text{let} \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e_1 : A \xrightarrow{2} [p_1, l_1.r_1] \quad \Gamma, x : A @ 2 \vdash e_2 : B \xrightarrow{2} [p_2, l_2.r_2]}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : B \xrightarrow{2} [(p_1, p_2), (l_1, l_2).\text{let } x = r_1 \text{ in } r_2]} \text{let} \xrightarrow{2} \quad \frac{\Gamma \vdash e : A \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{next } e : \bigcirc A \xrightarrow{1} [((), p), l.r]} \bigcirc I \xrightarrow{1} \\
\frac{\Gamma \vdash e : \bigcirc A \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{prev } e : A \xrightarrow{2} [\text{pi2 } c, l.r]} \bigcirc E \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : \text{bool} \xrightarrow{1} [c_1, l_1.r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{1} [c_2, l_2.r_2] \quad \Gamma \vdash e_3 : A \xrightarrow{1} [c_3, l_3.r_3]}{\Gamma \vdash \left(\begin{array}{l} \text{if } e_1 \\ \text{then } e_2 \\ \text{else } e_3 \end{array} \right) : A \xrightarrow{1} \left[\begin{array}{l} \left(\begin{array}{l} \text{let } (y_1, z_1) = c_1 \text{ in} \\ \text{if } y_1 \\ \text{then let } (y_2, z_2) = c_2 \text{ in } (y_2, (z_1, \text{inl } z_2)) \\ \text{else let } (y_3, z_3) = c_2 \text{ in } (y_3, (z_1, \text{inl } z_3)) \end{array} \right), \\ (l_1, l_b).(r_1; \text{case } l_b \text{ of } l_2.r_2 \mid l_3.r_3) \end{array} \right]} \text{if} \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e_1 : \text{bool} \xrightarrow{2} [p_1, l_1.r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{2} [p_2, l_2.r_2] \quad \Gamma \vdash e_3 : A \xrightarrow{2} [p_3, l_3.r_3]}{\Gamma \vdash \left(\begin{array}{l} \text{if } e_1 \\ \text{then } e_2 \\ \text{else } e_3 \end{array} \right) : A \xrightarrow{2} \left[(p_1, p_2, p_3), (l_1, l_2, l_3). \left(\begin{array}{l} \text{if } r_1 \\ \text{then } r_2 \\ \text{else } r_3 \end{array} \right) \right]} \text{if} \xrightarrow{2}
\end{array}$$

Figure 8. Term Splitting

$$\begin{array}{c}
\frac{\Gamma, x : A @ \mathbb{1} \vdash e : B \xrightarrow{1} [c, l.r]}{\Gamma \vdash \lambda x : A. e : A \rightarrow B \xrightarrow{1} [\lambda x. c, \dots (\lambda(x, l). r)]} \rightarrow I \xrightarrow{1} \quad \frac{\Gamma, x : A @ \mathbb{2} \vdash e : B \xrightarrow{2} [p, l.r]}{\Gamma \vdash \lambda x : A. e : A \rightarrow B \xrightarrow{2} [p, l. (\lambda x. r)]} \rightarrow I \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : A \rightarrow B \xrightarrow{1} [c_1, l_1. r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{1} [c_2, l_2. r_2]}{\Gamma \vdash e_1 e_2 : B \xrightarrow{1} [\dots]} \rightarrow E \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e_1 : A \rightarrow B \xrightarrow{2} [p_1, l_1. r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{2} [p_2, l_2. r_2]}{\Gamma \vdash e_1 e_2 : B \xrightarrow{2} [(p_1, p_2), (l_1, l_2). r_1 r_2]} \rightarrow E \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e : A \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{inl } e : A + B \xrightarrow{1} [\text{let } (y, z) = c \text{ in } (\text{inl } y, z), l.r]} + I_1 \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e : A \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{inr } e : A + B \xrightarrow{1} [\text{let } (y, z) = c \text{ in } (\text{inr } y, z), l.r]} + I_2 \xrightarrow{1} \quad \frac{\Gamma \vdash e : A \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{inl } e : A + B \xrightarrow{2} [l. \text{inl } r]} + I_1 \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e : A \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{inr } e : A + B \xrightarrow{2} [l. \text{inr } r]} + I_1 \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : A + B \xrightarrow{1} [c_1, l_1. r_1] \quad \Gamma, x_2 : A @ \mathbb{1} \vdash e_2 : C \xrightarrow{1} [c_2, l_2. r_2] \quad \Gamma, x_3 : B @ \mathbb{1} \vdash e_3 : C \xrightarrow{1} [c_3, l_3. r_3]}{\Gamma \vdash \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ x_2. e_2 \\ | x_3. e_3 \end{array} \right) : C \xrightarrow{1} \left[\begin{array}{c} \left(\begin{array}{c} \text{let } (y_1, z_1) = c_1 \text{ in} \\ \text{case } y_1 \text{ of} \\ x_2. \text{let } (y_2, z_2) = c_2 \text{ in } (y_2, (z_1, \text{inl}(z_2))) \\ | x_3. \text{let } (y_3, z_3) = c_3 \text{ in } (y_3, (z_1, \text{inr}(z_3))) \end{array} \right), \\ (l_1, l_b). \left(\begin{array}{c} \text{let } z = r_1 \text{ in} \\ \text{case } l_b \text{ of} \\ l_2. \text{let } x_2 = z \text{ in } r_2 \\ | l_3. \text{let } x_3 = z \text{ in } r_3 \end{array} \right) \end{array} \right]} + E \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e_1 : A + B \xrightarrow{2} [p_1, l_1. r_1] \quad \Gamma, x_2 : A @ \mathbb{2} \vdash e_2 : C \xrightarrow{2} [p_2, l_2. r_2] \quad \Gamma, x_3 : B @ \mathbb{2} \vdash e_3 : C \xrightarrow{2} [p_3, l_3. r_3]}{\Gamma \vdash \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ x_2. e_2 \\ | x_3. e_3 \end{array} \right) : C \xrightarrow{2} \left[(p_1, p_2, p_3), (l_1, l_2, l_3). \left(\begin{array}{c} \text{case } r_1 \text{ of} \\ x_2. r_2 \\ | x_3. r_3 \end{array} \right) \right]} + E \xrightarrow{2}
\end{array}$$

Figure 9. Sum and Function Splitting

it. They note that where [circle] corresponds to the non-branching temporal logic, [square] corresponds to the branching version.

References

- [1] R. Davies. A temporal-logic approach to binding-time analysis. In *LICS*, pages 184–195, 1996.
- [2] R. Davies and F. Pfenning. A modal analysis of staged computation. *J. ACM*, 48(3):555–604, 2001.
- [3] Z. DeVito, J. Hegarty, A. Aiken, P. Hanrahan, and J. 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. ISBN 978-1-4503-2014-6. . URL <http://doi.acm.org/10.1145/2491956.2462166>.
- [4] N. D. Jones. An introduction to partial evaluation. *ACM Comput. Surv.*, 28(3):480–503, Sept. 1996. ISSN 0360-0300. . URL <http://doi.acm.org/10.1145/243439.243447>.
- [5] U. Jørring and W. 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. . URL <http://doi.acm.org/10.1145/512644.512652>.
- [6] T. B. Knoblock and E. 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. ISBN 0-89791-795-2. . URL <http://doi.acm.org/10.1145/231379.231428>.