

Stage-Splitting a Modal Language

Nicolas Feltman Carlo Angiuli Umut Acar Kayvon Fatahalian

Carnegie Mellon University
{nfeltman,cangiuli,umut,kayvonf}@cs.cmu.edu

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. ([6]) 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...[4]), and their background theory and type-systems are well understood ([3]). 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 [5]), $\llbracket f \rrbracket$ 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

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

Figure 1. λ^{12} Syntax

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

Figure 2. λ^{12} Valid Types

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. The λ^{12} Language

In this section, we describe λ^{12} , a simple two-stage language used as the basis for defining and analyzing stage splitting. Every expression in λ^{12} has a statically known *stage*, either $\mathbb{1}$ or $\mathbb{2}$, which defines when it is fully reduced. Stage- $\mathbb{1}$ and stage- $\mathbb{2}$ types include products, sums, functions, and base types. In addition to these types, stage- $\mathbb{1}$ expressions may also be of type $\bigcirc A$, representing an enclosed stage- $\mathbb{2}$ expression of type A . (The type system of λ^{12} is adopted directly from that of the modal language λ^\bigcirc [2], but is restricted to two stages for simplicity.) $\bigcirc A$ terms must be treated opaquely by stage- $\mathbb{1}$ code as they cannot be reduced (to values with type A) until stage $\mathbb{2}$. Like in λ^\bigcirc there is no reciprocal modality.

Although there is a strict one-way dependence between the two stages at the type level, stage- $\mathbb{1}$ and stage- $\mathbb{2}$ terms are defined mutually recursively: the `next` construct embeds stage- $\mathbb{2}$ expressions into stage $\mathbb{1}$, whereas `prev` embeds stage- $\mathbb{1}$ expressions into stage $\mathbb{2}$. The `next` and `prev` constructs are the only expressions that interact with the stage of a term, and we surround their arguments with braces in λ^{12} syntax to clearly indicate stage boundaries within a program.

The grammar and type system for λ^{12} is given in Figures 1 and 3. (Typing judgments and context variables are annotated with stages.) Only \bigcirc and its introductory and eliminatory forms `next` and `prev` affect the stage of a term or type. Specifically, given an argument with type A at stage $\mathbb{2}$, `next` forms a $\bigcirc A$ at stage $\mathbb{1}$. Stage $\mathbb{2}$ expressions can obtain the original stage $\mathbb{2}$ argument via the `prev` construct. Since `prev` operates at stage $\mathbb{2}$, this ensures no violation of causality [1]. The `hold` construct serves to wrap stage $\mathbb{1}$ integers for use in stage $\mathbb{2}$. Although it is possible to implement `hold` from other λ^{12} features, we provide it as a core primitive to simplify our examples.

3.1 Staged Evaluation

As a simple example of expressing staged programs in λ^{12} , consider the following fast exponentiation algorithm, which calculates b^p in $\log p$ time:

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

The stage- $\mathbb{1}$ term below defines the function `fexp` that implements the fast exponentiation algorithm:

```

let fexp (b :  $\bigcirc$ int, e : int) :  $\bigcirc$ 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 `if` predicates and exponent decomposition are all stage- $\mathbb{1}$ terms, since they occur within `prev` blocks. We note that our work does not consider the problem of *binding time analysis*—transforming an unstaged program into a staged λ^{12} program by inserting appropriate `next`, `prev`, and `hold` constructs. We assume that binding time analysis is performed manually by a programmer (as done above) or via automatic analysis of unstaged code to arrive at a valid λ^{12} program used as input for subsequent splitting.

A key attribute of the `fexp` code example is the nesting of stage- $\mathbb{1}$ and stage- $\mathbb{2}$ expressions. Ordinary term evaluation eliminates outermost redexes first, however in the case that stage $\mathbb{1}$ expressions are contained inside stage $\mathbb{2}$ ones (such as in the example above), this strategy conflicts with the precept of staged execution: that all stage $\mathbb{1}$ code be evaluated before the evaluation of stage $\mathbb{2}$ code. Accordingly, a suitably staged dynamic semantics for λ^{12} must evaluate *all* the stage $\mathbb{1}$ subexpressions of a term before its stage $\mathbb{2}$ subexpressions.

3.2 Non-Duplicating First-Stage Evaluation

To understand how this behavior can be achieved, consider evaluation of the stage- $\mathbb{2}$ term below, which contains nested `next` and `prev` blocks:

```

prev{
  let x = (next {1+2}, 3+4) in
  next{ prev{#1 x} * prev{#1 x} * hold{#2 x} }
}

```

Our goal is to evaluate away all of the stage- $\mathbb{1}$ portions of this code, but there is more than one reasonable way to do this. One option, taken by [2], is to treat `next{1+2}` as a value for the purposes of stage- $\mathbb{1}$ evaluation, consequently duplicating it during substitution for x . This produces:

$$\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}_{A+B} e : A + B @ w} +I_1 \quad \frac{\Gamma \vdash e : B @ w}{\Gamma \vdash \text{inr}_{A+B} 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 : [\mu\alpha.\tau/\alpha]\tau @ w}{\Gamma \vdash \text{roll}_{\mu\alpha.\tau} e : \mu\alpha.\tau @ w} \mu I \\
\\
\frac{\Gamma \vdash e : \mu\alpha.\tau @ w}{\Gamma \vdash \text{unroll } e : [\mu\alpha.\tau/\alpha]\tau @ w} \mu E \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 \quad \frac{\Gamma \vdash e : \text{int} @ 1}{\Gamma \vdash \text{hold } e : \text{int} @ 2} \text{hold}
\end{array}$$

Figure 3. λ^{12} Static Semantics

(1+2) * (1+2) * 7

We take a different approach. Instead of duplicating the contents of the first `next` expression, we bind them to some variable (here, `y`) and duplicate only a reference to that variable. This produces:

`let y = 1+2 in y * y * 7`

Achieving this behavior mechanically requires us to resolve a contradiction: we must substitute for `x` at stage 1, but we cannot evaluate inside the `next` block bound to it. Our solution is to replace the `next` with a new variable and create an explicit substitution (shown with a \mapsto) binding that variable to the `next`'s contents. This substitution then floats up to the top of the containing `prev` block:

```
prev {
  [y ↦ 1+2]
  let x = (ŷ, 7) in
  next{prev{#1 x} * prev{#1 x} * hold{#2 x}}
}
```

Note that the variable left behind is rendered with a hat. By treating this *hatted variable* as a value, we're now free to perform the stage-1 substitution for `x`.

```
prev {
  [y ↦ 1+2]
  next{
    prev{#1 (ŷ, 7)} * prev{#1 (ŷ, 7)} *
    hold{#2 (ŷ, 7)}
  }
}
```

To evaluate the remaining `next`, we must first partially evaluate the body by finding all of the contained stage-1 terms and reducing them. As a rule, these will reduce to hatted variables, and `prev` removes the hat, leaving the variable in place:

```
prev {
  [y ↦ 1+2]
  next{ y * y * 7 }
}
```

Once again, we lift the contents of the `next` into a substitution:

```
prev {
  [y ↦ 1+2]
  [z ↦ y*y*7]
  ŷ
}
```

Finally, when evaluating the outer `prev`, we must *reify* the contained substitutions into let statements, yielding

```
let y = 1+2 in
let z = y * y * 7 in z
```

Thus we have evaluated all of the stage 1 expressions of this program without duplicating the contents of `next` blocks.

3.3 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 that have been fully evaluated, but which may contain hatted variables. In the example above, $(\hat{y}, 7)$.
- Residuals (reses) are stage 2 terms whose stage 1 subexpressions have all been fully evaluated. In the example above, $(1+2)$ and $(\text{let } z = y*y*7 \text{ in } z)$.
- 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 63.

The \Downarrow_1 judgment takes an open stage 1 term to a *future environment* ξ and a partial value v . The future environment is a mapping from newly-created variables (which then may appear—hatted—in v) to residuals. It is the formal manifestation of the floating substitutions in our example. For all of the normal features of λ^{12} (Figure 4), first stage evaluation has the same behavior and effect on values as standard evaluation, and the final future environment is gotten by merging the future environments of the subterms.

When this judgment encounters a `next` block (Figure 5), it searches into the block's stage 2 content to find any contained

$$\begin{array}{c}
\frac{}{\Xi; \Gamma \vdash () \Downarrow_1 [\cdot, ()]} \text{unit} \Downarrow_1 \quad \frac{}{\Xi; \Gamma \vdash i \Downarrow_1 [\cdot, i]} \text{int} \Downarrow_1 \quad \frac{}{\Xi; \Gamma \vdash b \Downarrow_1 [\cdot, b]} \text{bool} \Downarrow_1 \quad \frac{}{\Xi; \Gamma \vdash \lambda x:A.e \Downarrow_1 [\cdot, \lambda x.e]} \rightarrow \text{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 \text{E} \Downarrow_1 \\
\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 \text{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 \text{E}_1 \Downarrow_1 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, (v_1, v_2)]}{\Xi; \Gamma \vdash \text{pi2 } e \Downarrow_1 [\xi, v_2]} \times \text{E}_2 \Downarrow_1 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_1 [\xi, v]}{\Xi; \Gamma \vdash \text{inl } e \Downarrow_1 [\xi, \text{inl } v]} + \text{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]} + \text{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]} + \text{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]} + \text{E}_2 \Downarrow_1 \\
\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 \quad \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 \\
\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 4. λ^{12} Dynamic Semantics: Core

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

Figure 5. λ^{12} Dynamic Semantics: next and prev

$$\begin{array}{c}
\frac{}{\Xi; \Gamma \vdash () \Downarrow_2 ()} \text{unit} \Downarrow_2 \quad \frac{}{\Xi; \Gamma \vdash i \Downarrow_2 i} \text{int} \Downarrow_2 \quad \frac{}{\Xi; \Gamma \vdash b \Downarrow_2 b} \text{bool} \Downarrow_2 \quad \frac{x \in \Gamma}{\Xi; \Gamma \vdash x \Downarrow_2 x} \text{hyp} \Downarrow_2 \\
\frac{\Xi; \Gamma, x \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \lambda x:A.e \Downarrow_2 \lambda x:A.q} \rightarrow \text{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 \text{E} \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)} \times \text{I} \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{pi1 } e \Downarrow_2 \text{pi1 } q} \times \text{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 \text{E}_2 \Downarrow_2 \quad \frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{inl } e \Downarrow_2 \text{inl } q} + \text{I}_1 \Downarrow_2 \\
\frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{inr } e \Downarrow_2 \text{inr } q} + \text{I}_2 \Downarrow_2 \quad \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} + \text{E}_1 \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 6. λ^{12} Dynamic Semantics: Speculation

stage 1 subexpressions and evaluate them in place. We call this search process *speculation*. It is implemented by the \Downarrow_2 judgment, which takes a stage 2 term to a residual. Once the contents of the `next` block are speculated into a residual (q), we then return a fresh hatted variable (\hat{y}), along with a future environment which maps that variable to our residual ($y \mapsto q$).

At all of the normal features (Figure 6), speculation does nothing but recursively speculate into every subexpression. Once speculation finds a `prev` block, we resume stage 1 evaluation of the contents, which produces a future context and (by canonical forms) a hatted variable (\hat{y}). The context is then reified (using the \Rightarrow^R judgment) into a series of `let` bindings, with y —stripped of its hat—at the bottom.

Because the inputs to evaluation and speculation are sometimes open on variables, we adorn our judgments with two contexts. The first of these (Ξ) keeps track of those hatted variables that may appear in the input. Since hatted variables can only be introduced to terms via substitution (*i.e.* they do not appear in the original program), this context is augmented only at stage 1 substitution sites. The second context (Γ) keeps track of stage 2 variables, which are declared in original program.

If we change the `next` and `prev` rules to

$$\frac{\Xi; \Gamma \vdash e \Downarrow_2 q}{\Xi; \Gamma \vdash \text{next } e \Downarrow_1 [\cdot, \text{next } q]} \quad \frac{\Xi; \Gamma \vdash e \Downarrow_1 [\cdot, \text{next } v]}{\Xi; \Gamma \vdash \text{prev } e \Downarrow_2 v}$$

and treat `next` q as a partial value, where q is a residual, then we get precisely the semantics from [2]. This essentially bypasses the environment bookkeeping in \Downarrow_1 , by inlining residuals instead of hoisting them in `let`-bindings. From this it's clear that the semantics of [2] and ours always produce the same value when they both terminate. But moreover, because a reified residual will always be evaluated whether or not its result is consumed, our semantics terminates strictly less often than that of [2].

3.4 A Partial Evaluation System

A function f which accepts inputs at both stages 1 and 2 can be given a type of the form $A \rightarrow \bigcirc(B \rightarrow C)$ ¹. 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 $[\xi, v] \Rightarrow^R f_a$. A canonical forms theorem on v ensures that the resulting f_a has type $B \rightarrow C$ in the residual language. 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 \Rightarrow^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 .

3.5 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 ξ' wf, $\Xi, \text{dom}(\xi'); \Gamma \vdash e : B @ 2$, and $\Xi, \text{dom}(\xi'); \Gamma \vdash e$ res.

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

Theorem 3.5. If $\Xi; \Gamma \vdash e \Downarrow_1 [\xi, v]$ and $\Gamma \vdash e : A @ 1$ then

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

Theorem 3.6. If $\Xi; \Gamma \vdash e \Downarrow_2 q$ and $\Gamma \vdash e : A @ 2$ then

1. $\Xi; \Gamma \vdash q : A @ 2$; and
2. $\Xi; \Gamma \vdash q$ val.

¹ We can rewrite `fexp` in this form, or simply apply the following higher-order function which makes the adjustment:

TODO: code written, but listings don't work in footnotes! commented out in source until we get working

4. Splitting Algorithm

The basic idea is that we want to take a term which contains interleaved stage 1 and stage 2 code and untangle it into two separate terms: one with all of the stage 1 code and one with all the stage 2 code. To facilitate communication between the stages, there is a data structure passed between the stages, which is an output of the stage 1 term and input to the stage 2 term.

As with the dynamics, the form of the splitting statement depends on the external stage of the term. We consider splitting for stage 2 terms first, as it has a simpler form. For example, take $\text{hold}\{1+2\} < 5$. We can split this into two separate terms, $1+2$ and $\text{fn } 1 \Rightarrow 1 < 5$, which are called the *precomputation* and *residual*, respectively. For this splitting to be correct, the residual applied to the precomputation should yield the same result as the original term.

The form of splitting for stage 1 terms is more complicated. To see why, consider the term $(\text{next}\{\text{hold}\{1+2\} < 5\}, 3*4)$. As before, the comparison operation is part of stage 2, and the addition operation is a stage 1 precomputation, the result of which will eventually become input to the residual. The multiplication operation, however, is neither part of the residual nor a precomputation to support it. Instead it is the *immediate result* of stage 1, because it is available to the stage 1 context around our original term. For example:

```
let x = (next{hold{1+2} < 5}, 3*4) in
next{ if prev{#1 x} then 7 else hold {#2 x} }
```

To support this, stage 1 splitting produces two outputs: the *combined result* and the *residual*, where the combined result reduces to a tuple containing the immediate result and the precomputation.

4.1 Formal Setup

We now discuss the formal definition of splitting. The splitting algorithm comprises two judgments, $\overset{1}{\rightsquigarrow}$ and $\overset{2}{\rightsquigarrow}$. The $\overset{1}{\rightsquigarrow}$ judgment sends a stage 1 term (e) to a combined term (c) and residual ($l.r$), while $\overset{2}{\rightsquigarrow}$ sends a stage 2 term (e) to a precomputation (p) and a residual ($l.r$). For concision, we actually represent residuals as a term open on a single variable, rather than as functions. For both judgments, we use a context (Γ) to keep track of the open variables of e .

The rules² for splitting `next`, `prev`, and `hold` are given in Figure 7. The rule for `next` simply tuples up the precomputation of its subexpression with a trivial immediate result, while the rule for `prev` projects the combined result of its subexpression to just the precomputation. The `hold` rule treats the entire combined result of its subexpression as a precomputation, and projects out the integer result in the residual³.

The rules of stage 2 splitting are given in Figure 8. Every rule works by bundling the precomputations of the constituent parts, and then unbundling them with a pattern. The rules of stage 1 splitting are given in Figure 7.

4.2 Speculation

Notice in the stage 2 rules for `ifs`, `cases`, and `functions`, the precomputation is lifted out from within branches. This is the manifestation of the speculation behavior from the semantics.

4.3 Stage 1 Divergence

Consider the case of stage 1 `if` and `case` expressions. In both cases we have enough information at stage 1 to know what branch to

²The rules are written using patterns, including the open variable of the residual.

³The residual of an integer expression is usually trivial, but we have to include it here for termination purposes.

take, so there's no need to speculate. Instead, we evaluate the stage 1 portion of only the active branch, and then inject precomputation into a sum type. Then in the residual, we case on that sum and resume the stage 2 portion of the correct branch.

4.4 Stage 1 Functions

The trick when splitting stage 1 functions is that the contents themselves may be multi-stage. We handle this by splitting them into two functions: one in the immediate result which handles all the stage 1 content of the original, and one in the residual which handles all of the stage 2 content. Note that stage 1 λ -expressions themselves have only a trivial precomputation.

Thus the boundary type appears as output of the first stage and input of the second stage. This makes it hard to type the output, because the boundary types are not represented in the original function type.

4.5 Boundary Type Worst Case

TODO: add the example which is the worst case for figuring out boundary types

4.6 Value Splitting, Etc

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 :  $\bigcirc$ int, e : int) :  $\bigcirc$ 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)
```

$$\begin{array}{c}
\frac{\cdot}{\Gamma \vdash () : \text{unit}} \xrightarrow{\mathbb{1}} [(), (), -()] \quad \text{unit} \xrightarrow{\mathbb{1}} \quad \frac{\cdot}{\Gamma \vdash i : \text{int}} \xrightarrow{\mathbb{1}} [(i, ()), -()] \quad \text{int} \xrightarrow{\mathbb{1}} \quad \frac{x : A @ \mathbb{1} \in \Gamma}{\Gamma \vdash v : A} \xrightarrow{\mathbb{1}} [(x, ()), -()] \quad \text{hyp} \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e_1 : A \xrightarrow{\mathbb{1}} [c_1, l_1.r_1] \quad \Gamma \vdash e_2 : B \xrightarrow{\mathbb{1}} [c_2, l_2.r_2]}{\Gamma \vdash (e_1, e_2) : A \times B \xrightarrow{\mathbb{1}} \left[\left(\begin{array}{c} \text{let } ((y_1, z_1), (y_2, z_2)) = (c_1, c_2) \text{ in} \\ ((y_1, y_2), (z_1, z_2)) \end{array} \right), (l_1, l_2).(r_1, r_2) \right]} \times I \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e : A \times B \xrightarrow{\mathbb{1}} [c, l.r]}{\Gamma \vdash \text{pi1 } e : A \xrightarrow{\mathbb{1}} \left[\begin{array}{c} \text{let } (y, z) = c \text{ in} \\ (\text{pi1 } y, z) \end{array} \right]} \times E_1 \xrightarrow{\mathbb{1}} \quad \frac{\Gamma \vdash e : A \times B \xrightarrow{\mathbb{1}} [c, l.r]}{\Gamma \vdash \text{pi2 } e : B \xrightarrow{\mathbb{1}} \left[\begin{array}{c} \text{let } (y, z) = c \text{ in} \\ (\text{pi2 } y, z) \end{array} \right]} \times E_2 \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e_1 : A \xrightarrow{\mathbb{1}} [c_1, l_1.r_1] \quad \Gamma, x : A @ \mathbb{1} \vdash e_2 : B \xrightarrow{\mathbb{1}} [c_2, l_2.r_2]}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : B \xrightarrow{\mathbb{1}} \left[\left(\begin{array}{c} \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 \right]} \text{let} \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma, x : A @ \mathbb{1} \vdash e : B \xrightarrow{\mathbb{1}} [c, l.r]}{\Gamma \vdash \lambda x : A. e : A \rightarrow B \xrightarrow{\mathbb{1}} [\lambda x. c, -(\lambda(x, l).r)]} \rightarrow I \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e_1 : A \rightarrow B \xrightarrow{\mathbb{1}} [c_1, l_1.r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{\mathbb{1}} [c_2, l_2.r_2]}{\Gamma \vdash e_1 e_2 : B \xrightarrow{\mathbb{1}} \left[\begin{array}{c} \text{let } ((y_1, z_1), (y_2, z_2)) = (c_1, c_2) \text{ in} \\ \text{let } (y_3, z_3) = y_1 y_2 \text{ in} \\ (y_3, (z_1, z_2, z_3)) \end{array} \right], (l_1, l_2, l_3).(r_1 (r_2, l_3))} \rightarrow E \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e : A \xrightarrow{\mathbb{1}} [c, l.r]}{\Gamma \vdash \text{inl } e : A + B \xrightarrow{\mathbb{1}} [\text{let } (y, z) = c \text{ in } (\text{inl } y, z), l.r]} +I_1 \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e : A \xrightarrow{\mathbb{1}} [c, l.r]}{\Gamma \vdash \text{inr } e : A + B \xrightarrow{\mathbb{1}} [\text{let } (y, z) = c \text{ in } (\text{inr } y, z), l.r]} +I_2 \xrightarrow{\mathbb{1}} \\
\\
\frac{\Gamma \vdash e_1 : A + B \xrightarrow{\mathbb{1}} [c_1, l_1.r_1] \quad \Gamma, x_2 : A @ \mathbb{1} \vdash e_2 : C \xrightarrow{\mathbb{1}} [c_2, l_2.r_2] \quad \Gamma, x_3 : B @ \mathbb{1} \vdash e_3 : C \xrightarrow{\mathbb{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{\mathbb{1}} \left[\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), \right.} +E \xrightarrow{\mathbb{1}} \\
\\
\left. \left. \begin{array}{c} (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] \right. \\
\\
\frac{\Gamma \vdash e_1 : \text{bool} \xrightarrow{\mathbb{1}} [c_1, l_1.r_1] \quad \Gamma \vdash e_2 : A \xrightarrow{\mathbb{1}} [c_2, l_2.r_2] \quad \Gamma \vdash e_3 : A \xrightarrow{\mathbb{1}} [c_3, l_3.r_3]}{\Gamma \vdash \left(\begin{array}{c} \text{if } e_1 \\ \text{then } e_2 \\ \text{else } e_3 \end{array} \right) : A \xrightarrow{\mathbb{1}} \left[\left(\begin{array}{c} \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_3 \text{ in } (y_3, (z_1, \text{inl } z_3)) \end{array} \right), \right.} \text{if} \xrightarrow{\mathbb{1}} \\
\\
\left. \left. \begin{array}{c} (l_1, l_b).(r_1; \text{case } l_b \text{ of } l_2.r_2 \mid l_3.r_3) \end{array} \right] \right.
\end{array}$$

Figure 7. Stage $\mathbb{1}$ Splitting

$$\begin{array}{c}
\frac{}{\Gamma \vdash i : \text{int} \xrightarrow{2} [(), _i]} \text{int} \xrightarrow{2} \quad \frac{}{\Gamma \vdash () : \text{unit} \xrightarrow{2} [(), _{}]} \text{unit} \xrightarrow{2} \quad \frac{x : A @ 2 \in \Gamma}{\Gamma \vdash v : A \xrightarrow{2} [(), _x]} \text{hyp} \xrightarrow{2} \\
\\
\frac{\Gamma, x : A @ 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 \text{I} \xrightarrow{2} \quad \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 \text{E} \xrightarrow{2} \\
\\
\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 \text{I} \xrightarrow{2} \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 \text{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 \text{E}_2 \xrightarrow{2} \quad \frac{\Gamma \vdash e : A \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{inl } e : A + B \xrightarrow{2} [l.\text{inl } r]} + \text{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]} + \text{I}_1 \xrightarrow{2} \\
\\
\frac{\Gamma \vdash e_1 : A + B \xrightarrow{2} [p_1, l_1.r_1] \quad \Gamma, x_2 : A @ 2 \vdash e_2 : C \xrightarrow{2} [p_2, l_2.r_2] \quad \Gamma, x_3 : B @ 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]} + \text{E} \xrightarrow{2} \\
\\
\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} \\
\\
\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}{c} \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}{c} \text{if } r_1 \\ \text{then } r_2 \\ \text{else } r_3 \end{array} \right) \right]} \text{if} \xrightarrow{2}
\end{array}$$

Figure 8. Stage 2 Splitting

$$\begin{array}{c}
\frac{\Gamma \vdash e : \bigcirc A \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{prev } e : A \xrightarrow{2} [\text{pi2 } c, l.r]} \bigcirc \text{E} \xrightarrow{2} \quad \frac{\Gamma \vdash e : A \xrightarrow{2} [p, l.r]}{\Gamma \vdash \text{next } e : \bigcirc A \xrightarrow{1} [((), p), l.r]} \bigcirc \text{I} \xrightarrow{1} \\
\\
\frac{\Gamma \vdash e : \text{int} \xrightarrow{1} [c, l.r]}{\Gamma \vdash \text{hold } e : \text{int} \xrightarrow{2} [c, (a, l).(r; a)]} \text{hold} \xrightarrow{2}
\end{array}$$

Figure 9. next and prev Splitting

$$\begin{array}{c}
\frac{}{() \xrightarrow{1} [(), ()]} \text{unit} \xrightarrow{1} \quad \frac{}{i \xrightarrow{1} [i, ()]} \text{int} \xrightarrow{1} \quad \frac{}{b \xrightarrow{1} [b, ()]} \text{bool} \xrightarrow{1} \quad \frac{v_1 \xrightarrow{1} [u_1, w_1] \quad v_2 \xrightarrow{1} [u_2, w_2]}{(v_1, v_2) \xrightarrow{1} [(u_1, u_2), (w_1, w_2)]} \times \xrightarrow{1} \\
\\
\frac{v \xrightarrow{2} v'}{\hat{y} \xrightarrow{1} [(), y]} \bigcirc \xrightarrow{1} \quad \frac{x : A \vdash e : B \xrightarrow{1} [c, l.r]}{\lambda x : A. e \xrightarrow{1} [\lambda x.c, \lambda(x, l).r]} \rightarrow \xrightarrow{1}
\end{array}$$

Figure 10. Value Splitting

Quickselect example.

```
let qs (l :  $\mu\alpha.()$  + int *  $\alpha$ , i:  $\bigcirc$ int) =
  case unroll l of
  () => next {0}
  | (h,t) =>
    let (left,right,n) = partition h t in
    next{
      let n = hold{n} in
      case compare prev{i} n of
      () (*LT*) => prev {qs left i}
      | () (*EQ*) => 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 [6] 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 [6] 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 ([7]) 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 ([?]) 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 ([2]) 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, [3] 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 it. They note that where [circle] corresponds to the non-branching temporal logic, [square] corresponds to the branching version.

References

- [1] A. Cave, F. Ferreira, P. Panangaden, and B. Pientka. Fair reactive programming. *SIGPLAN Not.*, 49(1):361–372, Jan. 2014. ISSN 0362-1340. . URL <http://doi.acm.org/10.1145/2578855.2535881>.
- [2] R. Davies. A temporal-logic approach to binding-time analysis. In *LICS*, pages 184–195, 1996.
- [3] R. Davies and F. Pfenning. A modal analysis of staged computation. *J. ACM*, 48(3):555–604, 2001.
- [4] 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>.
- [5] 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>.
- [6] 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>.
- [7] 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>.