





Exact Recursive Probabilistic Programming*

DAVID CHIANG and COLIN MCDONALD, University of Notre Dame, USA CHUNG-CHIEH SHAN, Indiana University, USA

Recursive calls over recursive data are useful for generating probability distributions, and probabilistic programming allows computations over these distributions to be expressed in a modular and intuitive way. Exact inference is also useful, but unfortunately, existing probabilistic programming languages do not perform exact inference on recursive calls over recursive data, forcing programmers to code many applications manually. We introduce a probabilistic language in which a wide variety of recursion can be expressed naturally, and inference carried out exactly. For instance, probabilistic pushdown automata and their generalizations are easy to express, and polynomial-time parsing algorithms for them are derived automatically. We eliminate recursive data types using program transformations related to defunctionalization and refunctionalization. These transformations are assured correct by a linear type system, and a successful choice of transformations, if there is one, is guaranteed to be found by a greedy algorithm.

CCS Concepts: • Software and its engineering \rightarrow Formal language definitions; Recursion; • Mathematics of computing \rightarrow Probability and statistics.

Additional Key Words and Phrases: probabilistic programming, recursive types, linear types

ACM Reference Format:

David Chiang, Colin McDonald, and Chung-chieh Shan. 2023. Exact Recursive Probabilistic Programming. *Proc. ACM Program. Lang.* 7, OOPSLA1, Article 98 (April 2023), 31 pages. https://doi.org/10.1145/3586050

1 INTRODUCTION

Probability models in natural language processing (NLP), as well as computational biology and other fields, often involve recursive computations on recursive data structures. Sequences of words (or other basic units) have been modeled, for example, using higher-order Markov chains [Chen and Goodman 1999], hidden Markov models [Rabiner 1989], finite automata and transducers [Mohri 1997], and chain-structured conditional random fields (CRFs) [Lafferty et al. 2001]. Neural models that use CRFs [Huang et al. 2015] and connectionist temporal classification [Graves et al. 2006], which are just finite transducers, are at or near the state of the art in a number of NLP tasks. Trees, used for representing syntactic structures and sometimes other linguistic structures, have been modeled using probabilistic context-free grammars (PCFGs) [Booth and Thompson 1973], tree automata and transducers, or tree-structured CRFs, again performing at or near the state of the art in constituency parsing [Stern et al. 2017].

In order to support applications such as optimization of parameters by gradient-based methods, it is important that algorithms for these models can efficiently produce exact results. Typically

*This material is based upon work supported by the National Science Foundation under Award Nos. CCF-2019266 and CCF-2019291.

Authors' addresses: David Chiang, dchiang@nd.edu; Colin McDonald, cmcdona8@nd.edu, Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN, USA; Chung-chieh Shan, ccshan@indiana.edu, Department of Computer Science, Indiana University, Bloomington, IN, USA.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2023 Copyright held by the owner/author(s).

2475-1421/2023/4-ART98

https://doi.org/10.1145/3586050

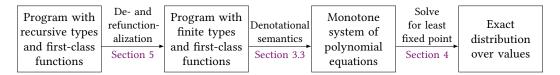


Fig. 1. How a PERPL program is compiled and evaluated. By "recursive type", we mean types such as String that are defined recursively and inhabited by an infinite number of values. By "finite type", we mean types such as Bool that are inhabited by a finite number of values.

they use dynamic programming, but more generally, they can be thought of as solving systems of linear or polynomial equations. Traditionally, these algorithms are coded by hand, which is tedious and error-prone. Worse, small changes to the model can often require nontrivial changes to the algorithm. For this reason, various toolkits have been developed, for example, for finite-state machines [Riley et al. 2009] and tree transducers [May and Knight 2006]. These provide reliable implementations of algorithms, but at the expense of flexibility: small changes to the model can often take it outside the class of model for which the toolkit was designed.

Probabilistic programming languages (PPLs) promise to avoid both manual coding and inflexible toolkits. The promise begins with an intuitive way to express a distribution, namely by describing a generative process: a run of a probabilistic program appears to generate a random outcome by making choices and scoring branches, but those probabilistic side effects are just notation for a distribution whose associated quantities—such as total weight or expected value—need to be computed. A distribution expressed in a PPL can be changed easily by editing the program, without needing to manually redesign inference algorithms.

In PPLs, recursive models such as those in NLP are naturally represented as recursive calls and data. For example, to find the probability of a string under a PCFG, it would be natural to write code that recursively generates random trees and sums the probabilities of trees yielding that string. Indeed, PCFGs served as the central example in the first PPL paper [Koller et al. 1997]. However, the example ultimately did not work out, and efficient inference on these representations remains a longstanding challenge. The general inference methods offered by PPLs today, often based on sampling, are not suitable in many applications because they take too long, don't compute gradients, or exhibit too much variance. Some PPLs do support exact inference—beginning with the "stochastic Lisp" of Koller et al. [1997], and more recently with languages like IBAL [Pfeffer 2005], Fun [Borgström et al. 2011], FSPN [Stuhlmüller and Goodman 2012], Bernoulliprob [Claret et al. 2013], PSI [Gehr et al. 2016, 2020], Hakaru [Walia et al. 2019], Dice [Holtzen et al. 2020], SPPL [Saad et al. 2021], and ProbZelus [Atkinson et al. 2022]—but they impose severe limitations on recursion.

In this paper, we introduce PERPL (Probabilistic Exact Recursive Programming Language). The denotation of a PERPL program is given by a system of polynomial equations, which can be solved efficiently using numerical methods. Unlike many general-purpose PPLs, PERPL performs exact inference, in the sense that even if it solves the equations iteratively (by Newton's method), any desired level of accuracy is guaranteed by a known number of iterations [Stewart et al. 2015]. Moreover, it can differentiate the equations and solve for the derivatives, permitting gradient-based optimization. And unlike other PPLs that do support exact inference, PERPL can express first-class, higher-order functions, unbounded recursive calls, and some recursive data structures. For example, a PCFG parser written in PERPL appears to generate (infinitely many) trees and sum the probabilities of those (exponentially many) trees that yield a given string; yet it compiles to a cubic-sized system of equations whose solution is equivalent to the CYK algorithm. We also show programs using stacks or stacks-of-stacks for which PERPL automatically derives polynomial-time algorithms.

We demonstrate that these asymptotic speedups result in orders-of-magnitude differences in a benchmark.

The key to our approach is a linear type system [Girard 1987; Walker 2005], which enables first-class functions to be used efficiently and recursive types to be eliminated correctly using defunctionalization [Danvy and Nielsen 2001; Reynolds 1972] and refunctionalization [Danvy and Millikin 2009], even in the presence of probability effects. Figure 1 depicts our pipeline of steps that turns a naturally expressed probabilistic program into an exact distribution over return values.

Our contributions are

- (1) a nonstandard denotational semantics for a PPL, in which a λ -abstraction from type τ_1 to type τ_2 denotes not a distribution over a set of functions, which has cardinality $|\tau_2|^{|\tau_1|}$, but a distribution over a set of pairs, which has cardinality $|\tau_1| \cdot |\tau_2|$ (Section 3.3);
- (2) using defunctionalization and refunctionalization to eliminate recursive types (Section 5);
- (3) showing that a linear type system ensures that both of the above are correct, in spite of probability effects (Theorems 3.14 and A.5);
- (4) compiling programs with unbounded recursion into systems of equations and solving them; in particular, loops can be evaluated exactly and directly, not iteratively (Section 4.2);
- (5) compiling natural probabilistic models into polynomial-time parsers for context-free grammars (Section 5.4), pushdown automata (Section 6.1), and their generalization to tree-adjoining grammars and embedded pushdown automata (Section 6.3);
- (6) translating programs with affinely used functions to programs with only linearly used functions (Appendix B.2);
- (7) an open-source implementation based on a compositional, semantics-preserving translation from PERPL to *factor graph grammars* [Chiang and Riley 2020] (Appendix C.2).

2 MOTIVATION

In this section, we present a sequence of examples with increasing demands on expressivity. First, a single coin flip (Example 2.1); second, unbounded loops and recursive calls (Examples 2.2 and 2.3); third, recursive data types (Example 2.4) and their elimination (Examples 2.5 and 2.6). The last of these will also demonstrate the need for (linearly used) first-class functions.

Example 2.1 (Probability). PERPL, like other PPLs, has a mechanism for nondeterministically sampling from probability distributions. To sample from a Bernoulli distribution, we can write:

```
define flip : Bool =
  amb (factor p in true) (factor q in false)
flip
```

where p and q are (metavariables for literal) nonnegative real weights. The distribution flip returns **true** with weight p and **false** with weight q. To get a Bernoulli distribution, we would require p + q = 1, but in general, weights do not have to sum to one. Then the final expression flip calls the distribution, so the value of the whole program is **true** with weight p and **false** with weight q.

Example 2.2 (Loops). To simulate a fair coin using an unfair coin, flip the unfair coin twice. If the two flips disagree, take the first flip; otherwise, repeat [von Neumann 1951]. We can express this unbounded process as:

```
define fair : Bool =
  let x = flip in let y = flip in if x = y then fair else x
fair
```

PERPL turns this program into the linear equations $t = pq + (p^2 + q^2)t$ and $f = qp + (p^2 + q^2)f$, then solves for t and f. If p + q = 1, then the answer is $t = f = \frac{1}{2}$.

Example 2.3 (Recursive calls). The PCFG [Booth and Thompson 1973]

$$S \xrightarrow{p} SS$$
 $S \xrightarrow{q} a$ (1)

defines a distribution over trees, under which the total probability of all *finite* trees is the least nonnegative solution of

$$z = pz^2 + q. (2)$$

If p + q = 1, then $z = \min(1, \frac{1-p}{p})$. This probability is computed by the following PERPL program.

define gen : Unit = if flip then let () = gen in let () = gen in () $--S \rightarrow SS$ else () $--S \rightarrow a$ gen

Although this program always returns (), the weight associated with () is the desired probability z. Notably, this program makes recursive calls of unbounded depth. The denotational semantics of this program will turn out to be the equation (2), which can be solved using standard methods. (In this case, the quadratic formula gives a closed-form solution; in general, we might need to resort to iterative methods like Newton's method, which is guaranteed to converge to the correct answer.)

Example 2.4 (Recursive data). In CFG parsing, we want to know the total weight of all derivations of a CFG that yield a given string. The following program defines a type String for strings over the alphabet $\{a\}$, a function gen that generates strings from the CFG (1), and a function equal that tests whether the generated string is the given one. The program returns a distribution over Booleans: the weight of **true** is the total weight of all CFG derivations that yield the string aaa, and the weight of **false** is the total weight of all other (finite) CFG derivations.

```
data Nonterminal = S
data Terminal = A
data String = Nil | Cons Terminal String
define gen (lhs: Nonterminal) (acc: String) : String =
   case lhs of S → if flip then gen S (gen S acc) else Cons A acc
define equal (xs: String) (ys: String) : Bool =
   case xs of
    Nil → case ys of Nil → true | Cons y _ → false
    Cons x xs → case ys of Nil → false | Cons y ys → x = y and equal xs ys
equal (gen S Nil) (Cons A (Cons A (Cons A Nil)))
```

Processing pipelines like this are common, and writing them in a PPL enables intuitive expression and modular reuse. In particular, Bayesian inference is easy to express. For instance, to predict the next word conditioned on a given prefix, we can leave the gen function as is, change equal to

where **fail** denotes the zero distribution, and change the last line of code to next_word (gen S Nil) (Cons A (Cons A (Cons A Nil)))

Because the String type is recursive and infinite, it is not trivial for PERPL to convert a program like Example 2.4 to a finite system of equations. The rest of this section shows how PERPL manages to compile this program to efficient exact inference. One way to proceed is to apply a whole-program transformation that replaces the recursive type by a nonrecursive data structure that represents all the places in the code where it is constructed. This is closely related to *defunctionalization* [Danvy and Nielsen 2001; Reynolds 1972].

Example 2.5 (Defunctionalization). In Example 2.4, there were two uses of String: for the strings generated by gen, and for the input string *aaa* to be parsed. These can be automatically distinguished. Below, we rename them to GenString and InputString, respectively. Defunctionalization changes the latter into a nonrecursive type.

We introduced a new type Position with four values corresponding to the four places in the original program where an InputString was constructed, all in the last line. A Position can be thought of as a potential InputString, or a tail of the input string. It is converted into an actual InputString by the function shift.

Due to this transformation, neither Position nor InputString is recursive, although GenString still is. If we read these programs as call-by-value random generators, then the transformation has interleaved the producer and consumer of the input string and deforested it away.

Our example still uses the recursive type GenString. We can get rid of it using *refunctionalization* [Danvy and Millikin 2009; Danvy and Nielsen 2001], a whole-program transformation that replaces the recursive type by the computation that consumes it.

Example 2.6 (Refunctionalization). Example 2.5 only consumes a GenString in one place: the case xs expression in equal. Refunctionalization introduces functions gen_nil and gen_cons whose bodies are the branches of this case expression, and replaces the constructors of GenString with these functions.

```
define gen_nil : Position \multimap Bool = \lambdays: Position. case shift ys of InputNil \rightarrow true | InputCons y \_ \rightarrow false
```

```
define gen_cons (x: Terminal) (xs: Position \multimap Bool) : Position \multimap Bool = \lambdays: Position. case shift ys of InputNil \to false

InputCons y ys \to x = y and xs ys

define gen (lhs: Nonterminal) (acc: Position \multimap Bool) : Position \multimap Bool = case lhs of S \to if flip then gen S (gen S acc) else gen_cons A acc

(gen S gen_nil) 0
```

All recursive types are gone now. However, refunctionalization has changed values formerly of type GenString to type Position \multimap Bool, necessitating the use of first-class functions.

The denotational semantics of this program will be a system of equations, and solving these equations is equivalent to the CYK algorithm for CFG parsing. The parser can be generalized to an arbitrary CFG, whether in Chomsky normal form, by extending Nonterminal, Terminal, and gen.

Not all programs of interest are structured as producer-consumer pipelines. In Section 6, we will give examples of programs that use stacks, whose production (by pushes) and consumption (by pops) follow no fixed order.

3 A PROBABILISTIC PROGRAMMING LANGUAGE

We now define PERPL more formally. It has two main distinctive features:

- probability effects, which allow a program to nondeterministically pursue multiple branches, each possibly with a different weight or probability, and
- a linear type system, in which values of certain types must be consumed exactly once.

3.1 Syntax and Typing

The syntax rules are shown in Figure 2.

- 3.1.1 Probability. Evaluating an expression may involve making a probabilistic choice. The expression $\operatorname{amb}\ e_1\ e_2$ chooses between e_1 and e_2 , so, for example, (amb true false, amb true false) splits the current branch of computation into four branches, in which the pair evaluates to (true, true), (true, false), (false, true), and (false, false), each with weight 1. The expression factor w in e, where w is a nonnegative number, multiplies the weight of the current branch by w and evaluates e. The expression fail has weight 0, so it terminates the current branch.
- 3.1.2 Linearity. We use a linear type system for two reasons, both illustrated in Section 2. First, defunctionalization, by changing some computations into data structures, delays them. For example, the calls to Cons in Example 2.4 are delayed to inside shift in Example 2.5. The fact that ys is consumed only once ensures that any probabilistic effects in the input string would not be duplicated by the reordering.

Second, refunctionalization wraps computations inside λ -abstractions (as seen in Example 2.6), and it would be prohibitively expensive to let a λ -abstraction from type τ_1 to type τ_2 denote a distribution over $|\tau_2|^{|\tau_1|}$ functions. But if it is used only once (as xs and acc are in Example 2.6), it can denote a distribution over input–output pairs. Since there are only $|\tau_1| \cdot |\tau_2|$ such pairs, this is far more manageable.

For both of these reasons (detailed below in Sections 5.3 and 3.3, respectively), we decree that functions must be used *linearly* [Girard 1987; Walker 2005], in that, once introduced, they must be used exactly once. For example, the following program is not well-typed, because it calls f twice.

```
let f = \lambda x. amb x (not x) in (f true, f true)
```

```
p := e \mid \text{define } x = e ; p
Programs
                           e := x \mid \lambda x : \tau . e \mid e \mid e \mid amb \mid e \mid fail \mid factor \mid w \mid in \mid e \mid (e, ..., e) \mid \langle e, ..., e \rangle
Expressions
                                  | \mathbf{let} (x, \dots, x) = e \mathbf{in} \ e \ | \ e.i \ | \ c.e \ | \ \mathsf{case} \ e.\mathsf{of} \ c.x \to e \ | \cdots | \ c.x \to e
                           \tau ::= \tau \multimap \tau \mid \mathbf{Unit} \mid \tau \otimes \cdots \otimes \tau \mid \tau \& \cdots \& \tau \mid c \tau \oplus \cdots \oplus c \tau
Types
                           c \in C C is an infinite set of constructors
Constructors
Weights
                          w \in \mathbb{Q}_{>0}
                                                                 let x_1 = e_1 in e' \equiv (\lambda x_1.e') e_1
Syntactic sugar
                             Bool \equiv True Unit \oplus False Unit
                                                                                          true ≡ True ()
                                                                                                                             false \equiv False()
               if e then e_1' else e_2' \equiv \mathsf{case}\ e of \mathsf{True}\ u \to \mathsf{let}\ () = u in e_1' | \mathsf{False}\ u \to \mathsf{let}\ () = u in e_2'
                         e_1 and e_2 \equiv \text{if } e_1 then e_2 else false e_2 \equiv \text{if } e_2 not e_3 \equiv \text{if } e_3 not e_4 \equiv \text{if } e_4 then false else true
```

Fig. 2. Syntax. We use a taller vertical bar (|) to separate different right-hand sides of a BNF production, and a shorter vertical bar (|) to separate the branches of a **case** expression.

3.1.3 Algebraic Data Types. Because of linearity, there are two kinds of tuples [Abramsky 1993]. In a multiplicative tuple (e_1,\ldots,e_n) , of type $\tau_1\otimes\cdots\otimes\tau_n$, all components are computed regardless of demand, and they are all consumed together. In the common case where n=0, we write () for the empty tuple and **Unit** for its type. In an additive tuple $\langle e_1,\ldots,e_n\rangle$, of type $\tau_1\otimes\cdots\otimes\tau_n$, just one component is computed depending on which one the context demands, and only the demanded one is consumed. Additive tuples also allow n=0, but we don't need to notate that case in this paper.

Disjoint union types $\tau ::= c \ \tau \oplus \cdots \oplus c \ \tau$ and terms $e ::= c \ e$ are tagged with *constructors* c drawn from an infinite set C. In our example programs, **data** type declarations can be understood as introducing aliases for unions of multiplicative tuples. In Example 2.5 for instance, Position is an alias for $0 \ \text{Unit} \oplus 1 \ \text{Unit} \oplus 2 \ \text{Unit} \oplus 3 \ \text{Unit}$, inhabited by terms like 2 (). In turn, InputString is an alias for InputNil Unit \oplus InputCons (Terminal \otimes Position), inhabited by terms like InputCons (A (), 2 ()). Here the constructors are 0, 1, 2, 3, InputNil, InputCons, and A.

3.1.4 Global Definitions. The **define** keyword introduces a global definition, which can recursively use any global variable. Because global variables are nonlinear, they allow us to express recursion, as Examples 2.2 to 2.6 demonstrated. Global variables are 'call-by-name' in the sense that they are bound to computations rather than values: each use of a global variable evaluates to a fresh copy of its definition's right-hand side. So the following program is well-typed, because each use of f creates a new λ -expression that is used once:

```
define f = \lambda x. amb x (not x); (f true, f true)
```

But the following is not well-typed, because g is used twice.

```
define f = \lambda x. amb x (not x); let g = f in (g true, g true)
```

And in the following program, b samples from **amb true false** each time it is evaluated (in other words, its nondeterminism is "hot"), so the program has four branches, not two:

```
define b = amb true false; (b, b)
```

3.1.5 Typing. The typing rules are shown in Figure 3. Typing judgements for expressions e are of the form Γ ; $\Delta \vdash e : \tau$, where Γ and Δ are typing contexts for nonlinear and linear bindings, respectively. Nonlinear (or *intuitionistic*) bindings can be used any number of times, whereas linear bindings must each be used exactly once [Barber 1996]. In these typing rules, global definitions always

Fig. 3. Typing rules. We write · for an empty typing context.

populate the nonlinear context, whereas local bindings only enter the linear context. It is useful in practical applications (including the examples in Section 2) to relax the linearity requirement to allow linear bindings to be used affinely (zero times or once), or to allow local nonlinear bindings under certain circumstances. Please see Appendix B for more details.

An **amb** expression can be used in any linear context, and both branches are type-checked in that linear context. Since **fail** is like **amb** but with 0 branches, it can also be used in any linear context.

The difference between multiplicative and additive tuples is reflected in their distinct typing rules. Whereas the components of an additive tuple $\langle e_1, \ldots, e_n \rangle$ are all type-checked in the same linear context Δ , the components of a multiplicative tuple (e_1, \ldots, e_n) are type-checked by partitioning the linear context $\Delta_1, \ldots, \Delta_n$. Thus, the empty tuple () requires the empty linear context.

For a program to type check under the global context Γ , the **define**s of the program must provide global variables with types exactly as promised in Γ . To enforce this consistency, the judgement for a program p takes the form $\Gamma \vdash p : \Gamma'; \tau$ to track the assumed context Γ and the provided context Γ' separately, and the judgement for a complete program $p : \tau$ requires Γ and Γ' to be identical.

Next, we give a standard operational semantics and a slightly nonstandard denotational semantics for PERPL, then prove that they agree thanks to the linear use of functions and additive tuples.

3.2 Operational Semantics

Our operational semantics is defined by a one-step reduction judgement $\gamma \vdash e \Longrightarrow E'$, shown in Figure 4. The expression being reduced is e. The reduction result E' is a distribution over expressions, written as a set of weight-expression pairs; this is used to reduce **amb**, **fail**, and **factor**. The global environment $\gamma = (x_1 = e_1, \dots)$ maps global names to their definitions; this is used to reduce global names in (4a). Our evaluation contexts allow evaluation under λ and on both sides of application, so the order of probabilistic choices is left unspecified; this is benign, as denotations are preserved regardless of the order (Theorem 3.14).

$$\gamma \vdash x \Longrightarrow \{(1,e)\} \text{ if } (x=e) \in \gamma$$
 (4a)

$$\gamma \vdash \qquad (\lambda x_1.e') \ v_1 \Longrightarrow \{(1, e'\{x_1 := v_1\})\} \tag{4b}$$

$$\gamma \vdash$$
 amb $e_1 e_2 \Longrightarrow \{(1, e_1), (1, e_2)\}$ (4c)

$$\gamma \vdash \qquad \qquad \mathsf{fail} \Longrightarrow \emptyset$$
 (4d)

$$\gamma \vdash \qquad \qquad \text{factor } w \text{ in } e \Longrightarrow \{(w, e)\}$$
 (4e)

$$\gamma \vdash$$
 let $(x_1, ..., x_n) = (v_1, ..., v_n)$ in $e' \Longrightarrow \{(1, e'\{x_1 := v_1, ..., x_n := v_n\})\}$ (4f)

$$\gamma \vdash \langle e_1, \dots, e_n \rangle . i \Longrightarrow \{(1, e_i)\}$$
(4g)

$$\gamma \vdash \mathsf{case}\ c_i \ v \ \mathsf{of}\ c_1 \ x_1 \to e_1' \ | \cdots \ | \ c_n \ x_n \to e_n' \Longrightarrow \big\{ \big(1, e_i' \{ x_i \coloneqq v \} \big) \big\}$$
 (4h)

$$\frac{\gamma \vdash e \Longrightarrow E'}{\gamma \vdash C[e] \Longrightarrow \{(w, C[e']) \mid (w, e') \in E'\}} \tag{4i}$$

Syntactic values

$$v := \lambda x.e \mid (v, \ldots, v) \mid c v \mid \langle e, \ldots, e \rangle$$

Evaluation contexts $C := [\cdot] | \lambda x : \tau.C | Ce | eC |$ factor w in C

$$| (e, \dots, e, C, e, \dots, e) | \mathbf{let} (x, \dots, x) = C \mathbf{in} e | \mathbf{let} (x, \dots, x) = e \mathbf{in} C$$

$$| C.i | c C | \mathbf{case} C \mathbf{of} c_1 x \to e | \dots | c_n x \to e$$

Fig. 4. Operational semantics.

If *E* is a distribution over expressions, then we can reduce *E* by reducing any element of *E*:

Definition 3.1. If $E = \{(w, e)\} + E_0$ and $\gamma \vdash e \Longrightarrow E'$, then $\gamma \vdash E \Longrightarrow w \cdot E' + E_0$ (where \cdot and + denote scaling and addition of distributions). We also write \Longrightarrow^* for a sequence of zero or more reductions.

For example, starting with the expression **factor** w_1 **in factor** w_2 **in** e, we can apply Definition 3.1 to (4e) twice to get

$$\gamma \ \vdash \ \big\{ \big(1, \mathsf{factor} \ w_1 \ \mathsf{in} \ \mathsf{factor} \ w_2 \ \mathsf{in} \ e \big) \big\} \ \Longrightarrow \ \big\{ \big(w_1, \mathsf{factor} \ w_2 \ \mathsf{in} \ e \big) \big\} \ \Longrightarrow \ \big\{ \big(w_1 w_2, e \big) \big\}. \ (3)$$

We conclude this section with a proof sketch of type soundness.

Lemma 3.2 (Linear substitution preserves typing). Suppose Γ ; Δ_0 , x_1 : $\tau_1 \vdash e' : \tau'$ and Γ ; $\Delta_1 \vdash e_1 : \tau_1$. Then Γ ; Δ_0 , $\Delta_1 \vdash e' \{x_1 := e_1\} : \tau'$.

PROOF. By induction on the typing derivation of e'.

Definition 3.3. We say that a global environment $\gamma = (x_1 = e_1, ...)$ is *well-typed* for a context $\Gamma = (x_1 : \tau_1, ...)$ iff $\Gamma; \cdot \vdash e_i : \tau_i$ for all i.

Proposition 3.4 (Reduction preserves typing). Let γ be well-typed for Γ . If Γ ; $\Delta \vdash e : \tau$ and $\gamma \vdash e \Longrightarrow E'$, then Γ ; $\Delta \vdash e' : \tau$ for all $(w, e') \in E'$.

PROOF. By induction on the derivation of $\gamma \vdash e \Longrightarrow E'$. Case (4a) uses the well-typedness of γ . Cases (4b), (4f), (4h) use Lemma 3.2. Case (4i) uses the induction hypothesis.

PROPOSITION 3.5 (PROGRESS). Let γ be well-typed for Γ . If Γ ; $\cdot \vdash e : \tau$ then either e is a value or $\gamma \vdash e \Longrightarrow E'$ for some E'.

PROOF. By induction on the typing derivation of e. The benign nondeterminism mentioned above is not used.

3.3 Denotational Semantics

In this section, we define the denotation of a complete program $p : \tau$, which is a distribution over the set denoted by the type τ . First we define a *distribution* over a set X to be a mapping $X \to [0, \infty]$. This makes sense in measure theory whenever X is countable; for us, X is always finite. Note that weights can be irrational or infinite (unlike in the syntax and operational semantics). If χ_1, χ_2 are distributions over X, we define the (complete) partial order $\chi_1 \le \chi_2$ iff for all $x \in X$, $\chi_1(x) \le \chi_2(x)$.

We handle recursive calls in our denotational semantics in a standard way. Without the global environment (our sole source of recursive calls), we would just let each expression e denote a mapping from environments to distributions over semantic values, where an environment maps each free variable of e to its semantic value. With the global environment, we have to first define this denotation relative to an interpretation of each global name as a distribution. Thus, the denotation of a program's global definitions maps an interpretation of each global name to another interpretation of each global name. Finally, we take the least fixed point of this monotonic map [Kozen 1981].

Definition 3.6. A type τ denotes a finite set $[\![\tau]\!]$ of semantic values, defined by induction on τ :

$$\begin{bmatrix} \tau_1 \otimes \cdots \otimes \tau_n \end{bmatrix} = \llbracket \tau_1 \rrbracket \times \cdots \times \llbracket \tau_n \rrbracket
 \llbracket c_1 \tau_1 \oplus \cdots \oplus c_n \tau_n \rrbracket = \{ c_i v \mid 1 \le i \le n, v \in \llbracket \tau_i \rrbracket \}
 \llbracket \tau_1 \multimap \tau_2 \rrbracket = \{ v_1 \mapsto v_2 \mid v_1 \in \llbracket \tau_1 \rrbracket, v_2 \in \llbracket \tau_2 \rrbracket \}
 \llbracket \tau_1 \otimes \cdots \otimes \tau_n \rrbracket = \{ i : v \mid 1 \le i \le n, v \in \llbracket \tau_i \rrbracket \}$$

The first two cases are unsurprising: the tuple type $\tau_1 \otimes \cdots \otimes \tau_n$ denotes a set of tuples, and the union type $c_1 \tau_1 \oplus \cdots \oplus c_n \tau_n$ denotes a disjoint union, whose elements $c_i v$ are just pairs of c_i and v.

The last two cases are where our denotational semantics is nonstandard. One might expect $\llbracket \tau_1 \multimap \tau_2 \rrbracket$ to be the set of functions $\llbracket \tau_2 \rrbracket^{\llbracket \tau_1 \rrbracket}$. But the denotation of an expression of type $\tau_1 \multimap \tau_2$ will involve a distribution over $\llbracket \tau_1 \multimap \tau_2 \rrbracket$, and that would have $|\llbracket \tau_2 \rrbracket|^{|\llbracket \tau_1 \rrbracket|}$ weights and lead to programs that require exponential time and space. Linearity permits a more efficient way. Intuitively, when a function is created, it nondeterministically guesses what (one) argument value it will receive. Later, applying the function (once) is just a matter of unifying the actual argument with this guess. The semantic value $v_1 \mapsto v_2$ is just a suggestively notated pair of v_1 and v_2 . Hence, the cardinality $\llbracket \tau_1 \multimap \tau_2 \rrbracket$ is only $\llbracket \tau_1 \rrbracket | \cdot \llbracket \tau_2 \rrbracket$, enabling PERPL to handle many programs in polynomial time.

Similarly, when an additive tuple is created, it nondeterministically guesses which (one) component will be demanded. The semantic value i:v is just a pair of i and v, and the cardinality $\llbracket \tau_1 \& \cdots \& \tau_n \rrbracket \mid$ is not $\llbracket \tau_1 \rrbracket \mid \cdots \mid \llbracket \tau_n \rrbracket \mid$ but only $\llbracket \tau_1 \rrbracket \mid + \cdots + \llbracket \tau_n \rrbracket \mid$.

Example 3.7. Since **Unit** is the type of 0-tuples, the set [[Unit]] is $\{()\}$. Since **Bool** is syntactic sugar for True **Unit** \oplus False **Unit**, the set [[Bool]] is $\{True(), False()\}$.

Recall from Section 3.1.3 that Position = \emptyset Unit \oplus 1 Unit \oplus 2 Unit \oplus 3 Unit, so [Position]] = $\{\emptyset (), 1 (), 2 (), 3 ()\}$. Then [Position \multimap Position] has $4 \cdot 4 = 16$ members, not $4^4 = 256$:

```
 [\![ \text{Position} \multimap \text{Position} ]\!] = \big\{ \big(\emptyset \, (), \emptyset \, () \big), \big(\emptyset \, (), 1 \, () \big), \big(\emptyset \, (), 2 \, () \big), \big(\emptyset \, (), 3 \, () \big), \\ \big(1 \, (), \emptyset \, () \big), \big(1 \, (), 1 \, () \big), \big(1 \, (), 2 \, () \big), \big(1 \, (), 3 \, () \big), \\ \big(2 \, (), \emptyset \, () \big), \big(2 \, (), 1 \, () \big), \big(2 \, (), 2 \, () \big), \big(2 \, (), 3 \, () \big), \\ \big(3 \, (), \emptyset \, () \big), \big(3 \, (), 1 \, () \big), \big(3 \, (), 2 \, () \big), \big(3 \, (), 3 \, () \big) \big\}.
```

Furthermore, $[\![$ (Position \multimap Position) \multimap Position $]\!]$ has $4\cdot 4\cdot 4=64$ members, not $4^{4^4}\approx 10^{154}$.

$$[\![x]\!]_{\eta}(\delta, v) = \begin{cases} \eta(x)(v) & x \in \operatorname{dom} \eta \\ \mathbb{I}[v = \delta(x)] & x \in \operatorname{dom} \delta \end{cases}$$
 (5a)

$$[\![\lambda x_1.e']\!](\delta_0, v_1 \mapsto v') = [\![e']\!](\delta_0 \cup \{(x_1, v_1)\}, v')$$
(5b)

$$[\![e_0 \ e_1]\!](\delta_0 \cup \delta_1, v') = \sum_{v_1 \in [\![\tau_1]\!]} [\![e_0]\!](\delta_0, v_1 \mapsto v') \cdot [\![e_1]\!](\delta_1, v_1)$$
(5c)

$$[\![amb \ e_1 \ e_2]\!](\delta, v) = [\![e_1]\!](\delta, v) + [\![e_2]\!](\delta, v)$$
(5d)

$$\llbracket \mathbf{fail} \rrbracket (\delta, v) = 0 \tag{5e}$$

$$\llbracket \mathbf{factor} \ w \ \mathbf{in} \ e \rrbracket (\delta, v) = w \cdot \llbracket e \rrbracket (\delta, v) \tag{5f}$$

$$[\![(e_1,\ldots,e_n)]\!](\delta_1\cup\cdots\cup\delta_n,(v_1,\ldots,v_n)) = [\![e_1]\!](\delta_1,v_1)\cdots[\![e_n]\!](\delta_n,v_n)$$
(5g)

$$[\![\mathtt{let} \ (x_1, \dots, x_n) = e \ \mathtt{in} \ e']\!] (\delta \cup \delta', v') = \sum_{v_1, \dots, v_n} [\![e]\!] (\delta, (v_1, \dots, v_n)) \cdot [\![e']\!] (\delta' \cup \{(x_1, v_1), \dots, (x_n, v_n)\}, v') \quad (5h)$$

$$[\![\langle e_1, \dots, e_n \rangle]\!](\delta, i : v) = [\![e_i]\!](\delta, v)$$

$$(5i)$$

$$\llbracket e.i \rrbracket (\delta, v) = \llbracket e \rrbracket (\delta, i : v) \tag{5j}$$

$$\llbracket c \, e \rrbracket (\delta, c' \, v) = \begin{cases} \llbracket e \rrbracket (\delta, v) & c' = c \\ 0 & c' \neq c \end{cases}$$
 (5k)

Fig. 5. Denotational semantics. In $[\![e]\!](\delta, v)$, the domain of δ is always the set of free variables in e, so in (5c), the left-hand side uniquely determines δ_0 and δ_1 on the right-hand side, and similarly in (5g), (5h), and (5l).

A global typing context Γ denotes the set of all mappings from variables in Γ to distributions over semantic values. That is, $\llbracket \Gamma \rrbracket$ contains all possible η such that for each $x : \tau \in \Gamma$, we have a distribution $\eta(x) : \llbracket \tau \rrbracket \to [0, \infty]$. Thus, $\llbracket \Gamma \rrbracket$ is generally uncountable.

A local typing context Δ denotes the set of all mappings from variables in Δ to semantic values. That is, $[\![\Delta]\!]$ contains all possible δ such that if $x : \tau \in \Delta$, then $\delta(x) \in [\![\tau]\!]$. Thus, $[\![\Delta]\!]$ is a Cartesian product of finite sets and itself finite.

Example 3.8. In Example 2.3, the global typing context is $\Gamma = (flip : Bool, gen : Unit)$. Since Bool has two values and Unit has one, each $\eta \in \llbracket \Gamma \rrbracket$ stores 3 weights:

$$\eta(\mathsf{flip})\big(\mathsf{True}\,()\big) \in [0,\infty] \qquad \eta(\mathsf{flip})\big(\mathsf{False}\,()\big) \in [0,\infty] \qquad \eta(\mathsf{gen})\big(()\big) \in [0,\infty]$$

There are no local variables in Example 2.3, so the local typing context Δ is everywhere empty, and the unique $\delta \in [\![\Delta]\!]$ is also empty.

A typing judgement Γ ; $\Delta \vdash e : \tau$ denotes a mapping that takes an $\eta \in \llbracket \Gamma \rrbracket$ to a distribution over $\llbracket \Delta \rrbracket \times \llbracket \tau \rrbracket$. We write this distribution as $\llbracket \Gamma; \Delta \vdash e : \tau \rrbracket_{\eta}$, but as $\llbracket e \rrbracket_{\eta}$ or even $\llbracket e \rrbracket$ for short, and we define it compositionally on the typing derivation, using the equations in Figure 5. The Iverson bracket $\llbracket \Gamma \rrbracket$ is 1 if its contents are true, 0 otherwise.

Example 3.9. Consider the expression amb (factor p in true) (factor q in false) from Example 2.1. We build up the denotation of this expression as follows.

$$[()](\emptyset, ()) \stackrel{(5g)}{=} 1$$

$$[\![\![amb \ (factor \ p \ in \ true) \ (factor \ q \ in \ false)]\!] (\emptyset, c \ ()) \stackrel{(5d)}{=} \begin{cases} p & c = \text{True} \\ q & c = \text{False}. \end{cases}$$

Example 3.10. Moving on to Example 2.3, assume that η is given.

[if flip **then let** () =
$$gen$$

in let () = gen in () else ()]
$$_{\eta}(\emptyset, ()) \stackrel{\text{(5l)}}{=} \eta(\text{flip})(\text{True}()) \cdot \eta(\text{gen})(()) \cdot \eta(\text{gen})(()) \cdot 1$$

+ $\eta(\text{flip})(\text{False}()) \cdot 1$. (7)

Given a set of global definitions $\gamma = (x_1 = e_1 : \tau_1, \dots, x_n = e_n : \tau_n)$, we define its denotation $[\![\gamma]\!]$ to be the global environment η that is the least solution to the equations

$$\eta(x_i)(v_i) = \llbracket e_i \rrbracket_{\eta}(\emptyset, v_i) \tag{8}$$

for each $i=1,\ldots,n$ and each v_i in $\llbracket \tau_i \rrbracket$. Finally, the program **define** $x_1=e_1$; ...; **define** $x_n=e_n$; e_n denotes the distribution $\llbracket e \rrbracket_{\llbracket y \rrbracket}(\emptyset, v)$ over semantic values v.

Example 3.11. For Examples 2.1 and 2.3, the equations are

$$\eta(\text{flip})(c()) \stackrel{(8)}{=} [\![\text{amb (factor } p \text{ in true}) \text{ (factor } q \text{ in false})]\!] (\emptyset, c())$$
 (9)

$$\begin{cases}
p & c = \text{True} \\
q & c = \text{False}
\end{cases} (10)$$

$$\eta(\text{gen})(()) \stackrel{(8)}{=} [\text{if flip then let } () = \text{gen in let } () = \text{gen in } () \text{ else } ()](\emptyset, ())$$
(11)

$$\eta(\text{gen})\big(()\big) \stackrel{\text{(8)}}{=} \left[\!\!\left[\text{if flip then let } () = \text{gen in let } () = \text{gen in } () \text{ else } () \right]\!\!\right] \left(\emptyset, ()\right) \\
\stackrel{\text{(7)}}{=} \eta(\text{flip})\big(\text{True } ()\big) \cdot \eta(\text{gen})\big(()\big) \cdot \eta(\text{gen})\big(()\big) \cdot 1 + \eta(\text{flip})\big(\text{False } ()\big) \cdot 1. \\
(12)$$

In all, there are three equations in three unknowns, $\eta(flip)(True())$, $\eta(flip)(False())$, and $\eta(\text{gen})(())$. The whole program in Example 2.3, of type **Unit**, denotes a single weight, which we write as z for short:

$$z = [\![gen]\!](\emptyset, ()) \stackrel{\text{(5a)}}{=} \eta(gen)(()). \tag{13}$$

Putting (10), (12), and (13) together, we get Equation (2) as promised:

$$z = pz^2 + q. (14)$$

3.4 Soundness

We justify our denotational semantics by showing that it is consistent with our standard operational semantics, thanks to the linear type system.

Definition 3.12. We say that the denotation $[\Gamma; \Delta \vdash e : \tau]_{\eta}$ is deterministic if, for each $\delta \in [\![\Delta]\!]$, it assigns weight 1 to just one $v \in [\tau]$ and weight 0 to all other $v \in [\tau]$.

Lemma 3.13 (Linear substitution preserves denotation).

Suppose Γ ; $\Delta_0, x_1 : \tau_1 \vdash e' : \tau'$ and Γ ; $\Delta_1 \vdash e_1 : \tau_1$ (as in Lemma 3.2). Then

$$[\![e'\{x_1 := e_1\}]\!]_{\eta}(\delta_0 \cup \delta_1, v') = \sum_{v_1 \in [\![\tau_1]\!]} [\![e']\!]_{\eta}(\delta_0 \cup \{(x_1, v_1)\}, v') \cdot [\![e_1]\!]_{\eta}(\delta_1, v_1)$$
(15)

for all $\eta \in \llbracket \Gamma \rrbracket$, $\delta_0 \in \llbracket \Delta_0 \rrbracket$, $\delta_1 \in \llbracket \Delta_1 \rrbracket$, and $v' \in \llbracket \tau' \rrbracket$.

Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.

PROOF. By induction on the typing derivation of e'. See Appendix A.1 for more details.

THEOREM 3.14 (REDUCTION PRESERVES DENOTATION). If $\Gamma: \Delta \vdash e : \tau$ and $\gamma \vdash e \Longrightarrow E'$ in the operational semantics, then in the denotational semantics, for all $\delta \in [\![\Delta]\!]$ and $v \in [\![\tau]\!]$,

$$\llbracket e \rrbracket_{\llbracket \gamma \rrbracket}(\delta,v) = \sum_{(w,e') \in E'} w \cdot \llbracket e' \rrbracket_{\llbracket \gamma \rrbracket}(\delta,v).$$

PROOF. By induction on the reduction derivation of e. We show the case $\gamma \vdash (\lambda x_1.e') e_1 \Longrightarrow \{(1, e'\{x_1 := e_1\})\}$ where e_1 is a syntactic value:

4 INFERENCE

The denotation of a program under the above semantics is the least fixed point of a system of equations, which has a particular form that has been well studied in other contexts.

Definition 4.1. A monotone system of polynomial equations, or MSPE [Esparza et al. 2008; Etessami and Yannakakis 2009], is a system of equations

$$z_1 = P_1(z_1, \dots, z_n)$$

$$\vdots$$

$$z_n = P_n(z_1, \dots, z_n)$$

where the z_i are called *weight variables* (to distinguish them from variables in PERPL) and each P_i is a polynomial with nonnegative real coefficients. We write z for a vector in $[0, \infty]^n$ of assignments to the weight variables. If z, z' are such assignments, we write $z \le z'$ iff for all $i, z_i \le z'_i$. We write P(z) for the vector $[P_i(z_1, \ldots, z_n)]_{i=1}^n \in [0, \infty]^n$. Thus the MSPE can be written succinctly as z = P(z).

4.1 Constructing an MSPE

The first step of inference is to instantiate Figure 5 and Equation (8) for all the expressions of the program, which has the form of an MSPE.

PROPOSITION 4.2. The denotation of a program is a distribution whose weight values are components of the least solution (in $[0, \infty]^n$) of an MSPE.

PROOF. The MSPE has two kinds of weight variables:

- (1) For every subexpression Γ ; $\Delta \vdash e : \tau$, every $\delta \in \llbracket \Delta \rrbracket$, and every $v \in \llbracket \tau \rrbracket$, there is a weight variable $\llbracket e \rrbracket_{\eta}(\delta, v)$. Figure 5 gives an equation whose left-hand side is this weight variable and whose right-hand side is a polynomial in the weight variables with nonnegative coefficients. Examples of such equations are (6), (7), and (13).
- (2) For every global variable $x : \tau$ and every $v \in \llbracket \tau \rrbracket$, there is a weight variable $\eta(x)(v)$. Equation (8) gives an equation whose left-hand side is this weight variable and whose right-hand side is a weight variable of the first kind. Examples of such equations are (9) and (11).

The denotation of a program **define** $x_1 = e_1$; ...; **define** $x_n = e_n$; e with type τ is the distribution that assigns to each $v \in [\![\tau]\!]$ the weight $[\![e]\!](\emptyset; v)$, which is a weight variable of the first kind above.

All that needs to be shown is that the MSPE has finite size. Under Definition 3.6, $\llbracket \tau \rrbracket$ is finite for every τ , which can be shown by induction on the structure of τ . Also, $\llbracket \Delta \rrbracket$ is finite, being a finite Cartesian product of finite sets. So the number of weight variables is finite.

Even though the number of weight variables is finite, it can be exponential in the size of the program, because $[\![\Delta]\!]$ is the product of as many sets as Δ has variables. If only for this reason, PERPL inference is intractable in general. This blowup is not surprising given that, even without any recursive calls or data, PERPL can easily express all discrete Bayes nets [Cooper 1990] and conjunctive queries [Chandra and Merlin 1977], and inherits their intractability.

Even for a family of programs—such as the typical parser—for which the number of weight variables is polynomial in the size of the program, an unwisely chosen dependency can dramatically increase the degree of the polynomial; this is a concern as well for the polynomial-time dynamic-programming inference algorithms that we are trying to recover automatically. As with Bayes nets and conjunctive queries, finding the optimal inference strategy (tree decomposition) is NP-hard, but many heuristics help in practice. Also, in practice the weight variables can be packed into a tensor to take advantage of vectorized parallelism.

4.2 Solving the MSPE

The general strategy for solving an MSPE automatically is to decompose it into smaller MSPEs. If z_1, z_2 are weight variables, we write $z_1 < z_2$ if there is an equation whose left-hand side is z_2 and whose right-hand side contains z_1 . Then find the strongly connected components (SCCs) of <. In Example 3.11,

$$\eta(\text{flip})(\text{True}()) < \eta(\text{gen})(()) \quad \eta(\text{flip})(\text{False}()) < \eta(\text{gen})(()) \quad \eta(\text{gen})(()) < \eta(\text{gen})(())$$

so each of the three variables forms its own SCC. We visit each SCC in topological order, solving it and substituting its solution into the remaining equations.

The easiest, and most common, case is when an SCC has just one weight variable, whose equation (after substituting earlier weight variables) must be of the form z = w where w is a constant. Many existing PPLs with exact inference handle only this case, that is, when < is acyclic.

The next easiest case is when an SCC has more than one weight variable but (after substituting earlier weight variables) its equations are all linear. These can be solved *directly*—as opposed to iteratively—in the semiring $[0, \infty]$ by one of the algorithms of Lehmann [1977], such as Gaussian elimination. This case includes all loops (Example 2.2) and arises often in practice; for example, PCFG rules of the form $(X \to X)$ can be used an unbounded number of times. Most existing parsers limit how many times such rules can be applied [e.g., Collins 1999; Finkel et al. 2008; Taskar et al. 2004], but PERPL makes it easy to write code that handles such situations exactly.

The most general case is when the equations for the SCC are nonlinear. An example is parsing using a PCFG with rules of the form $(X \to \epsilon)$. Again, implementations usually resort to arbitrary limits on how such rules can be applied [e.g., Cai et al. 2011]. But these cases can be solved using a generalization of Newton's method to ω -continuous semirings [Esparza et al. 2007, 2010], which is guaranteed to converge to the least solution of an MSPE:

$$\begin{aligned} \mathbf{z}^{(0)} &= \mathbf{0} \\ \mathbf{z}^{(i+1)} &= \mathbf{z}^{(i)} + \left(\frac{\partial \mathbf{P}}{\partial \mathbf{z}}(\mathbf{z}^{(i)})\right)^* \left(\mathbf{P}(\mathbf{z}^{(i)}) - \mathbf{z}^{(i)}\right) \end{aligned}$$

where $\frac{\partial \mathbf{P}}{\partial \mathbf{z}}$ is the *formal* derivative of **P** (that is, defined only using the sum and product rules, not using limits). For any matrix **A**, the *closure* $\mathbf{A}^* = \sum_{i=0}^{\infty} \mathbf{A}^i$ can be computed by an algorithm of Lehmann [1977]. And $\infty - \infty$ can be defined to be anything (say, zero).

Although the iterates of Newton's method are in general only approximations, there is a known number of iterations that is guaranteed to reach any desired level of accuracy [Stewart et al. 2015]. This is also true of, for example, logarithms, and in the context of machine learning, such computations are generally considered "exact," in contrast to methods based on random sampling.

Example 4.3. In Equation (14), let $p = \frac{2}{3}$ and $q = \frac{1}{3}$. Then $P(z) = \frac{2}{3}z^2 + \frac{1}{3}$ and $\frac{\partial P}{\partial z}(z) = \frac{4}{3}z$, and using Newton's method to solve z = P(z) gives

$$z^{(0)} = 0$$

$$z^{(1)} = 0 + 0^* \left(\frac{1}{3} - 0\right) = \frac{1}{3} \approx 0.33333$$

$$z^{(2)} = \frac{1}{3} + \left(\frac{4}{9}\right)^* \left(\frac{11}{27} - \frac{1}{3}\right) = \frac{7}{15} \approx 0.46667$$

$$z^{(3)} = \frac{7}{15} + \left(\frac{28}{45}\right)^* \left(\frac{323}{675} - \frac{7}{15}\right) = \frac{127}{255} \approx 0.49804$$

$$z^{(4)} = \frac{127}{255} + \left(\frac{508}{765}\right)^* \left(\frac{97283}{195075} - \frac{127}{255}\right) = \frac{32767}{65536} \approx 0.49999$$

and so on. But suppose p=q=1. Then $P(z)=z^2+1$ and $\frac{\partial P}{\partial z}(z)=2z$, and Newton's method gives

$$z^{(0)} = 0$$

$$z^{(1)} = 0 + 0^*(1 - 0) = 1$$

$$z^{(2)} = 1 + 2^*(2 - 1) = \infty$$

$$z^{(3)} = \infty + \infty^*(\infty - \infty) = \infty$$

Accordingly, our implementation actually outputs inf.

In machine learning applications, it is often useful to optimize a quantity involving the distribution computed by a program, by adjusting parameters that the program depends on. In Example 2.4 for instance, having computed the likelihood L of some observed strings, we might want to maximize it by adjusting p. To adjust p by gradient descent, we can differentiate the MSPE denoted by the program with respect to p and then solve the resulting (linear) system of equations for $\frac{\mathrm{d}L}{\mathrm{d}p}$.

4.3 Factor Graph Grammars

Our implementation actually does not compile PERPL programs directly to MSPEs, but to *factor graph grammars* (FGGs) [Chiang and Riley 2020], whose inference algorithm constructs an MSPE in turn and solves it while taking advantage of vectorized parallelism. FGGs are a formalism for describing probability models on recursive structures that is more general than factor graphs, case–factor diagrams [McAllester et al. 2008] and sum–product networks [Poon and Domingos 2011]. Details of the translation to FGGs are in Appendix C.

5 RECURSIVE TYPES

In this section, we extend the language to allow recursive types (Section 5.1), and then, because our denotational semantics does not include recursive types and because known exact inference methods do not support recursive types, we show how to eliminate them (Sections 5.2 to 5.5).

5.1 Definitions

As shown in Figure 6, we add iso-recursive types [Pierce 2002] with introduction and elimination forms **fold** and **unfold** (also sometimes called **roll** and **unroll**; they do not mean catamorphism and anamorphism). Our definition of recursive types has two slightly unusual features. First, **unfold** expressions have a scope e', which will be used below (Section 5.2.2) when eliminating recursive types. Invariably, e' is a **case** expression, so we often write **case unfold** e **of** . . . as shorthand for **unfold** x = e **in case** x **of**

Second, the superscript t is a tag drawn from any infinite set such as \mathbb{N} . It is used to distinguish recursive types that would otherwise be considered equal, which is sometimes necessary for eliminating them (illustrated in Section 5.4 below). As with nominal typing, we consider two

Syntax
$$e := \mathbf{fold}_{\mu^t \alpha. \tau} \ e \ | \ \mathbf{unfold}_{\mu^t \alpha. \tau} \ x = e \ \mathbf{in} \ e' \qquad \tau := \mu^t \alpha. \tau \ | \ \alpha$$

Typing
$$\frac{\Gamma; \Delta \vdash e : \tau \{\alpha := \mu^t \alpha. \tau\}}{\Gamma; \Delta \vdash \mathbf{fold}_{\mu^t \alpha. \tau} \ e : \mu^t \alpha. \tau} \qquad \frac{\Gamma; \Delta \vdash e : \mu^t \alpha. \tau}{\Gamma; \Delta, \Delta' \vdash \mathbf{unfold}_{\mu^t \alpha. \tau} \ x = e \ \mathbf{in} \ e' : \tau'}{\Gamma; \Delta, \Delta' \vdash \mathbf{unfold}_{\mu^t \alpha. \tau} \ x = e \ \mathbf{in} \ e' : \tau'}$$

Reduction
$$\gamma \vdash \mathbf{unfold}_{\mu^t \alpha. \tau} \ x = (\mathbf{fold}_{\mu^t \alpha. \tau} \ v) \ \mathbf{in} \ e' \Longrightarrow \left\{ \left(1, e' \{x := v\} \right) \right\} \qquad (16)$$

Syntactic values
$$v := \mathbf{fold}_{\mu^t \alpha. \tau} \ v$$
Evaluation contexts
$$C := \mathbf{fold} \ C \mid \mathbf{unfold} \ x = C \ \mathbf{in} \ e \mid \mathbf{unfold} \ x = e \ \mathbf{in} \ C$$

Fig. 6. Syntax and semantics of recursive types.

recursive types that differ only in their tag to be two different types. Additionally, we require that different recursive types must have different tags; that is, if the program contains both $\mu^t \alpha. \tau_1$ and $\mu^t \alpha. \tau_2$, then $\tau_1 = \tau_2$. (This requirement is easy to satisfy, and used in Lemma A.12.) Tags can be inferred by using a different tag variable for each occurrence of μ , and unifying tag variables during type checking as necessary. (Our implementation of PERPL also infers tag polymorphism in global definitions and datatypes, by a straightforward generalization of Hindley–Milner–Damas type inference. It then eliminates this inferred polymorphism by monomorphization.)

Example 5.1. In Examples 5.2 and 5.3, we will transform the following simple program to eliminate the recursive type Nat $\equiv \mu^t \alpha$. (Zero Unit \oplus Succ α). This program samples a natural number n from an exponential distribution (by flipping a coin repeatedly and counting how many flips are **true** before the first **false**), and tests whether n is odd.

```
data Nat = Zero | Succ Nat
define sample : Nat = if flip then fold (Succ sample) else fold Zero
define odd (n: Nat) : Bool =
  unfold n' = n in case n' of Zero → false | Succ m → not (odd m)
odd sample
```

5.2 Eliminating Recursive Types

Recursive types such as Nat above cannot be handled by the conversion to an MSPE (Section 4) because they would give rise to an infinite number of weight variables. However, they can often be transformed away. In this section, we define these transformations and show how they work on Example 5.1. In Section 5.3 and Appendix A.3, we argue that they preserve program meaning, even in the presence of probabilistic effects. In Section 5.5, we provide a greedy algorithm that finds a successful sequence of transformations whenever one exists.

Our transformations are closely related to defunctionalization and refunctionalization, which are well-known. However, our formulation is slightly nonstandard, and our usage of them in a language with probabilistic effects is novel. Most notably, whereas traditional defunctionalization often introduces a recursive type, our usage of it eliminates a recursive type.

5.2.1 Defunctionalization. Suppose we want to eliminate a recursive type $\sigma = \mu^t \alpha$. $\bar{\sigma}$ from a given program. Let the n occurrences of \mathbf{fold}_{σ} in the program be \mathbf{fold}_{σ} $m_1, \ldots, \mathbf{fold}_{\sigma}$ m_n . For each i, let the nonglobal free variables of m_i form the tuple $\vec{y}_i : \vec{\varphi}_i$. We define a transformation $\mathcal{D}_{\sigma}[\![\cdot]\!]$ that changes terms of type σ to type $\varphi = \mathsf{Fold}_1 \ \vec{\varphi}_1 \oplus \cdots \oplus \mathsf{Fold}_n \ \vec{\varphi}_n$, where $\mathsf{Fold}_1, \ldots, \mathsf{Fold}_n$ are constructors. Because our purpose is to eliminate σ , it doesn't help to replace σ by another type that contains σ , and replacing σ in φ would cause infinite regress, so we require that φ not contain σ .

On types, our transformation $\mathcal{D}_{\sigma}[\![\cdot]\!]$ just changes occurrences of σ to φ , so it leaves $\vec{\varphi}_i$ unchanged. On programs, it adds a global function u_{σ} (u stands for unfold), of type $\varphi \multimap \bar{\sigma}\{\alpha := \varphi\}$:

define
$$u_{\sigma} = \lambda x : \varphi$$
. case x of Fold₁ $\vec{\eta}_1 \to \mathcal{D}_{\sigma} \llbracket m_1 \rrbracket \mid \cdots \mid \text{Fold}_n \vec{\eta}_n \to \mathcal{D}_{\sigma} \llbracket m_n \rrbracket$ (17)

On terms, it postpones from **fold**_{σ} to **unfold**_{σ} the work done by u_{σ} :

$$\mathcal{D}_{\sigma}[\![\mathsf{fold}_{\sigma} \ m_i]\!] = \mathsf{Fold}_i \ \vec{y}_i \quad \mathcal{D}_{\sigma}[\![\mathsf{unfold}_{\sigma} \ x = e \ \mathsf{in} \ e']\!] = \mathsf{let} \ x = u_{\sigma} \ \mathcal{D}_{\sigma}[\![e]\!] \ \mathsf{in} \ \mathcal{D}_{\sigma}[\![e']\!] \quad (18)$$

The other cases are uninteresting.

If the m_i are all abstractions, and if the e' in each unfold_σ x = e in e' is an application of x, then the \mathcal{D}_σ transformation is equivalent to defunctionalization [Danvy and Nielsen 2001; Reynolds 1972], with u_σ usually called apply_σ . Although the original purpose of defunctionalization was to get rid of λ -abstractions by postponing them to their applications, in general it can be used to get rid of any introduction forms by postponing them to their elimination forms; here, we use it to get rid of **folds** by postponing them to their **unfolds**.

Example 5.2. In Example 5.1, there are two occurrences of **fold**, inside sample. Since neither of them has any free variables, $\varphi = \text{Fold}_1 \, \text{Unit} \oplus \, \text{Fold}_2 \, \text{Unit}$ and so we can defunctionalize Nat. Create new types NatF and NatU, which take the roles of φ and $\bar{\sigma}\{\alpha := \varphi\}$ above, respectively:

```
data NatF = Fold_1 | Fold_2
data NatU = Zero | Succ NatF
```

The bodies of the **fold** expressions move into a new function u_Nat, which plays the role of u_{σ} :

```
define u_Nat (n: NatF) : NatU = 
 case \ n \ of \ Fold_1 \rightarrow Succ \ sample \ | \ Fold_2 \rightarrow Zero
```

The rest of the program is transformed by Equation (18):

5.2.2 Refunctionalization. Again suppose we want to eliminate a recursive type $\sigma = \mu^t \alpha.\bar{\sigma}$ from a given program. Instead of moving the work done at **fold**ing, we can move the work done at **unfold**ing. Without loss of generality, assume that every **unfold**_ σ expression binds the same variable name x, and let the n occurrences of **unfold**_ σ in the program be **unfold**_ σ $x = \cdots in m_i$, where $1 \le i \le n$. For each i, let the nonglobal free variables of the body $m_i : \varphi_i$, except x, form the tuple $\vec{y}_i : \vec{\varphi}_i$. We define a transformation $\mathcal{R}_{\sigma}[\![\cdot]\!]$ that changes terms of type σ to type

$$\varphi = (\vec{\varphi}_1 \multimap \varphi_1) \& \cdots \& (\vec{\varphi}_n \multimap \varphi_n). \tag{19}$$

Again because our purpose is to eliminate σ , we require that φ not contain σ .

On types, the transformation $\mathcal{R}_{\sigma}[\![\cdot]\!]$ just changes occurrences of σ to φ , so it leaves $\vec{\varphi}_i$ and φ_i unchanged. On programs, it adds a global function f_{σ} (f stands for fold), of type $\bar{\sigma}\{\alpha := \varphi\} \multimap \varphi$:

define
$$f_{\sigma} = \lambda x : \bar{\sigma}\{\alpha := \varphi\}. \langle \lambda \vec{y}_1. \mathcal{R}_{\sigma}[\![m_1]\!], \dots, \lambda \vec{y}_n. \mathcal{R}_{\sigma}[\![m_n]\!] \rangle$$
 (20)

On terms, it moves from **unfold**_{σ} to **fold**_{σ} the work done by f_{σ} :

$$\mathcal{R}_{\sigma}[\![\mathsf{fold}_{\sigma} e]\!] = f_{\sigma} \mathcal{R}_{\sigma}[\![e]\!] \qquad \qquad \mathcal{R}_{\sigma}[\![\mathsf{unfold}_{\sigma} x = e \text{ in } m_i]\!] = (\mathcal{R}_{\sigma}[\![e]\!].i) \vec{y}_i \qquad (21)$$

This is essentially the same as refunctionalization [Danvy and Millikin 2009; Danvy and Nielsen 2001]; in Danvy and Millikin's terminology, we have *disentangled* the **unfold**_{σ} expressions into the

apply functions $\lambda \vec{y}_i$. $\mathcal{R}_{\sigma}[\![m_i]\!]$ and *merged* them using an additive tuple (which works even if they have different types). The original definition of refunctionalization assumed that $\bar{\sigma}$ was a union type and created a function for each case; we create a single function f_{σ} for the same purpose.

Example 5.3. In Example 5.1, there is just one **unfold** expression. It has type **Bool**, and its body has no free variables other than n'. So instead of defunctionalizing Nat, we can also refunctionalize it. Create new types NatF and NatU, which take the roles of φ and $\bar{\sigma}\{\alpha:=\varphi\}$ above, respectively:

```
type NatF = Unit → Bool
data NatU = Zero | Succ NatF
```

(According to Equation (19), NatF should be a unary additive tuple type, which has no notation in this paper. To keep the presentation simple, we just omit unary tuple types here.)

The body of the **unfold** expression moves into a new function f_N at, which plays the role of f_{σ} :

```
define f_Nat (n': NatU) : NatF = \lambda(): Unit. case n' of Zero \rightarrow false | Succ m \rightarrow not (odd m) The rest of the program is transformed by Equation (21): define sample : NatF = if flip then f_Nat (Succ sample) else f_Nat Zero define odd (n: NatF) : Bool = n () odd sample
```

5.3 Correctness

It is easy to show that the transforms \mathcal{D} and \mathcal{R} preserve types; see Appendix A.2. In Appendix A.3, we show that they furthermore preserve the overall meaning of a program (a distribution over Unit, without loss of generality). Intuitively, this is because they merely change where the work done by u_{σ} or f_{σ} is expressed: \mathcal{D}_{σ} moves the work done by u_{σ} from fold_{σ} to unfold_{σ}, whereas \mathcal{R}_{σ} moves the work done by f_{σ} from unfold_{σ} to fold_{σ}.

5.4 Example: CFG Parsing

We next show how defunctionalization and refunctionalization together enable PERPL to handle some programs that neither transformation alone can. Specifically, we eliminate recursive types from the CFG parser in Example 2.4. This larger example demonstrates how the classic CYK algorithm is recovered, as well as some additional aspects of the transformations.

We repeat Example 2.4 below, with a few more details spelled out. First, we write **folds** and **unfolds** explicitly. Second, as mentioned above, tag inference automatically distinguishes two uses of String, which we call String¹ and String². If it were not for tags, we would not be able to handle this program: it would not be possible to defunctionalize String because of the free variable acc in **fold** (Cons A acc), and it would not be possible to refunctionalize String either because of the free variable xs inside the second **case unfold** ys expression.

Finally, some variables are used nonlinearly. For the nonrecursive type Terminal, this is harmless; Appendix B.1 discusses how to handle it. But for recursive types, we need to add a discard function to turn non-use into linear use; Appendix B.2 shows how this can be done automatically instead.

```
data String^1 = Nil^1 \mid Cons^1 Terminal String^1
data String^2 = Nil^2 \mid Cons^2 Terminal String^2
define gen (lhs: Nonterminal) (acc: String^1) : String^1 =
case lhs of S \rightarrow if flip then gen S (gen S acc) else fold (Cons^1 \land acc) define equal (xs: String^1) (ys: String^2) : Bool =
```

Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.

There are four occurrences of **fold**_{String}², all on the last line. Since none of them have any free variables, we can defunctionalize String². Create new types String²F and String²U, which take the roles of φ and $\bar{\sigma}\{\alpha:=\varphi\}$ in Section 5.2.1:

```
data String^2F = Fold_1 | Fold_2 | Fold_3 | Fold_4
data String^2U = Nil^2 | Cons^2 Terminal String^2F
```

The four values of type String²F can be thought of as positions between the symbols of the input string *aaa* (including the beginning and end of the string). The new function u_String² can then be thought of as taking a string position and returning the input symbol at that position, together with its successor position.

```
define u_String<sup>2</sup> (ys: String<sup>2</sup>F) : String<sup>2</sup>U = 

case ys of Fold<sub>1</sub> \rightarrow Cons<sup>2</sup> A Fold<sub>2</sub> | Fold<sub>2</sub> \rightarrow Cons<sup>2</sup> A Fold<sub>3</sub>

Fold<sub>3</sub> \rightarrow Cons<sup>2</sup> A Fold<sub>4</sub> | Fold<sub>4</sub> \rightarrow Nil<sup>2</sup>
```

Functions gen and discard¹ remain the same, but equal, discard², and the last line become

As for String¹, it cannot be defunctionalized, because **fold** (Cons¹ A acc) has a free variable acc of type String¹. Can it be refunctionalized? There are two **case unfold**_{String¹} expressions, one in equal and one in discard¹. Both bodies have type **Bool**; the first has a free variable ys with type String²F, and the second has no free variables. So yes, we can refunctionalize String¹ by creating new types String¹F and String¹U, which take the roles of φ and $\bar{\sigma}\{\alpha:=\varphi\}$ in Section 5.2.2:

```
type String^1F = (String^2F \multimap Bool) & (Unit \multimap Bool) data String^1U = Nil^1 | Cons^1 Terminal String^1F
```

A String 1F can be thought of as a string accessed through two "methods" [cf. Rendel et al. 2015]: the first (with type $String^2F \multimap Bool$) compares it with a suffix of the input string, and the second (with type $Unit \multimap Bool$) always returns **false**. These methods are implemented in the new function f_String^1 :

```
define f_String^1 (xs: String^1U) : String^1F = \langle \lambda ys: String^2F. \ case \ xs \ of
Nil^1 \rightarrow case \ u_String^2 \ ys \ of \ Nil^2 \rightarrow true \ | \ Cons^2 \ y \ ys \rightarrow discard^2 \ ys
Cons^1 \ x \ xs \rightarrow case \ u_String^2 \ ys \ of \ Nil^2 \rightarrow discard^1 \ xs
Cons^2 \ y \ ys \rightarrow x = y \ and \ equal \ xs \ ys,
\lambda(): \ Unit. \ case \ xs \ of \ Nil^1 \rightarrow false \ | \ Cons^1 \ x \ xs \rightarrow discard^1 \ xs \rangle
define equal (xs: String^1F) (ys: String^2F) : Bool = xs.1 \ ys
define discard^1 (xs: String^1F) : Bool = xs.2 ()
And occurrences of fold_{String^1} are changed to calls to f_String^1F.
define gen (lhs: String^1F) (acc: String^1F) : String^1F = case \ lhs \ of \ S \rightarrow if \ flip \ then \ gen \ S \ (gen \ S \ acc) \ else \ f_String^1 \ (Cons^1 \ A \ acc) equal (gen S \ (f_String^1 \ Nil^1)) Fold_1
```

No recursive types remain, so this program converts to an MSPE with a finite number of weight variables. For an input string $w = w_1 \cdots w_n$, there are $O(n^2)$ equations in $O(n^2)$ weight variables, and each equation has O(n) terms. They can be ordered so that each weight variable depends only on earlier variables, and solving them is equivalent to the CYK algorithm.

To show this equivalence with less clutter, we define

$$\overline{\ell} = 1 : (\text{Fold}_{\ell}() \mapsto \text{True}()) \in \llbracket \text{String}^1 \text{F} \rrbracket \qquad 1 \le \ell \le n+1. \tag{22}$$

The set $[String^1F] = [(String^2F \rightarrow Bool) \& (Unit \rightarrow Bool)]$ has 2n+4 possible values (not $O(2^n)$, thanks to our treatment of functions as nondeterministic pairs). Of these, the $\overline{\ell}$ are the interesting ones; each corresponds to a generated suffix that is equal to the input suffix starting at position ℓ . Thus, it will be easiest to think of them as string positions.

Then the equations simplify to

$$\eta(\text{gen})(S() \mapsto \overline{k} \mapsto \overline{i}) = \sum_{j} p \cdot \eta(\text{gen})(S() \mapsto \overline{j} \mapsto \overline{i}) \cdot \eta(\text{gen})(S() \mapsto \overline{k} \mapsto \overline{j})$$

$$+ q \cdot \eta(\text{f_String}^1)(\text{Cons}(A(), \overline{k}) \mapsto \overline{i}) \tag{23}$$

$$\eta(f_{\text{String}}^1)(\text{Cons}(A(), \overline{k}) \mapsto \overline{i}) = \mathbb{I}[w_i = a \land i = k - 1]$$
(24)

and the probability of the whole program being true is $\eta(\text{gen})(S() \mapsto \overline{n+1} \mapsto \overline{1})$. Equation (24) corresponds to the initialization in CYK; it gives weight 1 to occurrences of terminal symbols. Equation (23) corresponds to the main triple loop in CYK; it computes the weight of all derivations $S \Rightarrow^* w_i \cdots w_{k-1}$.

5.5 Inferring the Sequence of Transformations

In the previous section, we saw that not all recursive types are amenable to both defunctionalization and refunctionalization. Furthermore, some transformations can preclude others. For example, even though $String^2$ could be refunctionalized, it would preclude refunctionalizing $String^1$ because, in function equal, it would change the type of variable ys from $String^2$ to $(Unit \multimap Bool)$ & $((Terminal \otimes String^1) \multimap Bool)$ & $(Unit \multimap Bool)$. This type in turn contains $String^1$, and ys occurs free within the **case unfold**_{String1}xs expression. Although it may seem that such dependencies between transformations would make it difficult to find a successful sequence of transformations that eliminates all recursive types, the good news is that there is a simple greedy algorithm for deciding whether a successful sequence of transformations exists, and if so, to find it.

While there are any recursive types remaining:

- If there is a recursive type σ such that $\mathcal{D}_{\sigma}[\![\sigma]\!]$ contains no recursive type, defunctionalize σ .
- If there is a recursive type σ such that $\mathcal{R}_{\sigma}[\![\sigma]\!]$ contains no recursive type, refunctionalize σ .
- Else, there is no successful sequence of transformations.

We give the formal statement of correctness here and defer its proof to the appendix.

Definition 5.4. A DR-sequence S for a program p is a sequence $\mathcal{F}^{(1)}, \ldots, \mathcal{F}^{(n)}$, where each $\mathcal{F}^{(i)}$ is either \mathcal{D}_{σ} or \mathcal{R}_{σ} for some recursive type σ . We say that S is *successful* if it is empty and p has no recursive types, or if $S = \mathcal{F} \cdot S'$, where $\mathcal{F}[\![p]\!]$ is well-defined (that is, $\mathcal{F}[\![\sigma]\!]$ does not contain σ) and S' is a successful DR-sequence for $\mathcal{F}[\![p]\!]$.

Proposition 5.5. Let p be a program with at least one recursive type and a successful DR-sequence.

- (1) There is a recursive type σ and a transformation $\mathcal{F} \in \{\mathcal{D}, \mathcal{R}\}$ such that $\mathcal{F}_{\sigma}[\![\sigma]\!]$ contains no recursive type.
- (2) For any such σ and \mathcal{F} , the program $\mathcal{F}_{\sigma}[\![p]\!]$ is well-defined and has a successful DR-sequence.

```
Proof. See Appendix A.4.
```

A note on complexity: The transformations \mathcal{D} and \mathcal{R} only move code around, so any DR-sequence will not make the program much bigger. However, some successful DR-sequences may lead to more efficient inference than others. Moreover, there exist programs whose size is blown up exponentially by the monomorphization that needs to take place before the DR-sequence.

6 FURTHER EXAMPLES: PUSHDOWN AUTOMATA

PERPL automates the derivation of nonobvious inference algorithms from natural probabilistic models expressed using unbounded recursion. We've seen already how Section 5.4 derived the CYK algorithm from code that generates strings from a PCFG and compares them with an input string. In this section, we present three further examples, all related to pushdown automata (PDAs), to illustrate what kinds of recursive data structures PERPL can and cannot handle.

6.1 One Stack

The parser of Section 5.4 could also be implemented by converting the weighted CFG to a weighted nondeterministic PDA and running the PDA. Define a type for stacks:

```
data Stack = StkNil | StkCons Nonterminal Stack
```

(Our implementation of PERPL supports polymorphic datatypes such as List and monomorphizes their uses such as List Nonterminal.) The stack is initially just one S. The run_pda function below runs a nondeterministic PDA on a string. This PDA keeps no state, and accepts by empty stack.

```
define run_pda (zs: Stack) (ws: String) : Bool =
  case unfold zs of
    StkNil \rightarrow case \ unfold \ ws \ of \ Nil \rightarrow true
                                                                                    -- accept
                                     Cons w ws \rightarrow discard ws
                                                                                    -- reject
    StkCons z zs \rightarrow z = S and
                                                                                    -- pop S
       if flip then
         run_pda (fold (StkCons S (fold (StkCons S zs)))) ws
                                                                                    -- push SS
       else
         case unfold ws of Nil → discard_stk zs
                                                                                    -- reject
                              Cons w ws \rightarrow w = A and run_pda zs ws
                                                                                    -- scan a
```

```
define discard (ws: String) =
   case unfold ws of Nil → false | Cons w ws → discard ws
define discard_stk (zs: Stack) =
   case unfold zs of StkNil → false | StkCons z zs → discard_stk zs
run_pda (fold (StkCons S (fold StkNil)))
        (fold (Cons A (fold (Cons A (fold Nil)))))))
```

For simplicity, we have hard-coded the PDA into this program; it would be easy to add more symbols and transitions.

Unlike our previous examples using recursive data, a program that uses a stack like this is not structured as a producer-consumer pipeline. Nevertheless, String can be defunctionalized as in Section 5.4, yielding the following, where StringF, StringU, u_String, and discard are as before:

```
define run_pda (zs: Stack) (ws: StringF) : Bool =
  case unfold zs of
    StkNil \rightarrow case u_String ws of Nil \rightarrow true
                                                                                     -- accept
                                       Cons w ws \rightarrow discard ws
                                                                                    -- reject
    StkCons z zs \rightarrow z = S and
                                                                                    -- pop S
       if flip then
         run_pda (fold (StkCons S (fold (StkCons S zs)))) ws
                                                                                    -- push SS
       else
         case u_String ws of Nil \rightarrow discard_stk zs
                                                                                    -- reject
                                Cons w ws \rightarrow w = A and run_pda zs ws
                                                                                    -- scan a
run_pda (fold (StkCons S (fold StkNil))) Fold<sub>1</sub>
```

Although Stack cannot be defunctionalized because there are two occurrences of **fold**_{Stack} with a free variable zs of type Stack, the two occurrences of **unfold**_{Stack} meet the criterion for refunctionalization. So we refunctionalize Stack by creating new types StackF and StackU:

```
type StackF = (StringF → Bool) & (Unit → Bool)
data StackU = StkNil | StkCons Nonterminal StackF
define f_Stack (zs: StackU) : StackF =
   \langle \lambda ws: StringF. case zs of
        StkNil \rightarrow case u\_String ws of Nil \rightarrow true
                                                                                    -- accept
                                           Cons w ws \rightarrow discard ws
                                                                                    -- reject
        StkCons z zs \rightarrow z = S and
                                                                                    -- pop S
          if flip then
            run_pda (f_Stack (StkCons S (f_Stack (StkCons S zs)))) ws
                                                                                    -- push SS
          else
            case u_String ws of Nil \rightarrow discard_stk zs
                                                                                    -- reject
                                    Cons w ws \rightarrow w = A and run_pda zs ws,
                                                                                    -- scan a
    \lambda(): Unit. case zs of StkNil \rightarrow false | StkCons z zs \rightarrow discard_stk zs\rangle
define run_pda (zs: StackF) (ws: StringF) : Bool = zs.1 ws
define discard_stk (zs: StackF) : Bool = zs.2 ()
run_pda (f_Stack (StkCons S (f_Stack StkNil))) Fold<sub>1</sub>
```

Solving the corresponding MSPE is similar to Lang's algorithm for PDA recognition [Lang 1974]; as a matter of fact, refunctionalization has rederived a faster version of the algorithm found only recently [Butoi et al. 2022]. As observed by Danvy and Millikin [2009], this final program resembles the final program in Section 5.4, but written in continuation-passing style: the Stack data structure has been replaced by StackF, the type of continuations expecting a StringF or a Unit.

6.2 Two Stacks

An example of a program where defunctionalization and refunctionalization fail would be one that processes strings while maintaining two stacks instead of one. Since a two-stack PDA is equivalent to a Turing machine, such a program cannot be tractable, and it is desirable for PERPL to reject it at compile time. We omit a full program listing, but it would have a form like:

As in Section 6.1, the push operations render both Stack types non-defunctionalizable. Moreover, here **case unfold** left has free variable right, so refunctionalizing Stack¹ would change it into a type containing Stack², while **case unfold** right has free variable left, so refunctionalizing Stack² would change it into a type containing Stack¹. Thus, by Proposition 5.5, neither type is refunctionalizable.

6.3 A Stack of Stacks

In contrast to the two-stack PDA in Section 6.2, if multiple stacks are nested, refunctionalization does succeed, and it turns out to convert nested stacks into nested continuations. *Embedded PDAs* (EPDAs) are a generalization of PDAs where memory is structured not as a stack of symbols but as a stack of *stacks* of symbols [Vijayashanker 1994]. They are equivalent to tree-adjoining grammars (TAGs), which generalize CFGs. Moreover, TAGs can be parsed using a CYK-style algorithm [Vijay-Shankar and Joshi 1985]. Both of these results were major achievements in 1980s computational linguistics. Yet we show below that PERPL automatically takes a program that runs an EPDA and converts it into equations for a CYK-style TAG parser [cf. Alonso et al. 2000].

In the following code, we append a * to names relating to the stack-of-stacks.

```
data Stack = StkNil | StkCons Nonterminal Stack
data Stack* = StkNil* | StkCons* Stack Stack*
```

Whenever the top stack is empty, it is automatically popped. A move of the EPDA consists of popping a symbol z from the stack and then doing one of the following:

We assume a function called transition with type Nonterminal → Action. The following function decides whether the EPDA accepts a string. To sidestep the complication of explicit discard functions, it simply returns () if the EPDA accepts and **fails** otherwise.

```
define run_epda (zs*: Stack*) (ws: String) : Unit = case zs* of
  StkNil* \rightarrow case ws of Nil \rightarrow () | Cons w ws \rightarrow fail
  StkCons* z* zs* \rightarrow case z* of
    StkNil → run_epda zs* ws
    StkCons z zs \rightarrow case transition z of
       Pop \rightarrow run_epda (StkCons* zs zs*) ws
       Scan a \rightarrow case ws of Nil \rightarrow fail
                                Cons w ws \rightarrow if w = a then run_epda (StkCons* zs zs*) ws
                                                          else fail
       Push x y \rightarrow run_epda (StkCons* (StkCons x (StkCons y zs)) zs*) ws
       PushAbove x y \rightarrow
         run_epda (StkCons* (StkCons x StkNil) (StkCons* (StkCons y zs) zs*)) ws
       PushBelow x y \rightarrow
         run_epda (StkCons* (StkCons x zs) (StkCons* (StkCons y StkNil) zs*)) ws
run_epda (StkCons* (StkCons Z StkNil) StkNil*) (Cons A (Cons A (Cons A Nil)))
Like in Section 6.1, we defunctionalize String and refunctionalize Stack and Stack*, then simplify:
type StackF = Stack*F → StringF → Unit
type Stack*F = StringF → Unit
define f_StkNil* : Stack*F = \lambdaws: StringF. case u_String ws of
  Nil \rightarrow () \mid Cons w ws \rightarrow fail
define f_StkNil : StackF = \lambdazs*: Stack*F. \lambdaws: StringF. zs* ws
define f_StkCons (z: Nonterminal) (zs: StackF) : StackF =
  \lambda zs*: Stack*F. \lambda ws: StringF. case transition z of
    Pop \rightarrow zs zs* ws
    Scan a \rightarrow case u_String ws of Nil \rightarrow fail
                                        Cons w ws \rightarrow if w = a then zs zs* ws else fail
    Push x y \rightarrow f_StkCons x (f_StkCons y zs) zs* ws
    PushAbove x y \rightarrow f_StkCons x f_StkNil (f_StkCons y zs zs*) ws
    PushBelow x y \rightarrow f_StkCons x zs (f_StkCons y f_StkNil zs*) ws
f_StkCons Z f_StkNil f_StkNil* Fold<sub>1</sub>
```

As before, StringF can be thought of as a position in the input string, as can Stack*F. Then StackF can be thought of as a pair of string positions. The most computationally expensive expression is f_StkCons x (f_StkCons y zs), which represents $O(n^4)$ weight variables because it and its free variable zs both have type StackF. Each of these weight variables is a sum over $O(n^2)$ possible values of the subexpression f_StkCons y zs. Thus the program denotes a system of equations of size $O(n^6)$. Readers familiar with TAG will recognize $O(n^6)$ as the time complexity of CYK-style TAG parsing. Indeed, as we detail in Appendix D, when this program is converted to an FGG, the FGG is essentially a TAG, and PERPL's inference amounts to CYK-style TAG parsing.

Above, we noted that the result of refunctionalizing Section 6.1 was written in continuation-passing style. The twice-refunctionalized program above is written in extended continuation-passing

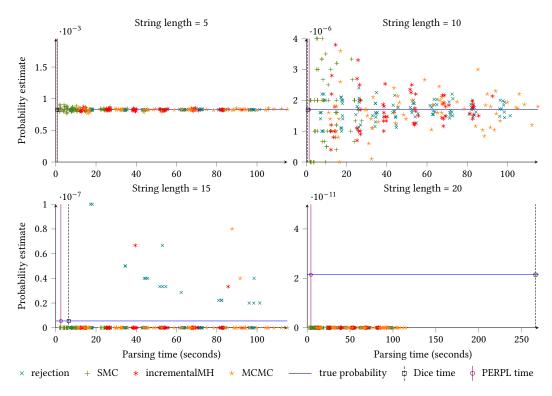


Fig. 7. The WebPPL version of Example 2.4 (PCFG parsing), using a variety of inference methods, displays high variance for all but the shortest strings. Dice and PERPL obtain exact probabilities much faster. For string length 15, one data point with probability estimate $3 \cdot 10^{-7}$ has been omitted. We ran rejection sampling up to 50 M samples, SMC up to 3 M particles, and incrementalMH and MCMC up to 30 M samples.

style [Danvy and Filinski 1990], with StackF being the type of continuations and Stack*F being the type of metacontinuations. Moreover, generalizations of EPDAs to automata with *k*-nested stacks (so-called *k*-EPDAs) can be expressed by similar PERPL programs and successfully compiled. We hence conjecture more generally that PERPL's program transformations relate the *control hierarchy* of context-free grammars [Weir 1992] to the (serendipitously same-named) *control hierarchy* of continuations [Danvy and Filinski 1990; Sitaram and Felleisen 1990].

7 BENCHMARK

We empirically demonstrate the improved speed and accuracy enabled by PERPL, by comparing against two existing probabilistic languages: WebPPL [Goodman and Stuhlmüller 2014], which supports unbounded recursive calls and data, but does not perform exact inference on them; and Dice [Holtzen et al. 2020], which performs exact inference, but only supports bounded loops. Because so much of PERPL's expressive power lies in its novel support for unbounded recursive data, it is tricky to find a benchmark that allows a quantitative comparison. We took the PCFG parsing problem in Example 2.4 and expressed it as idiomatically as we could (Appendix E): in WebPPL using unbounded recursion, and in Dice by unrolling a loop that builds up long parses from shorter ones. All experiments were performed on a 2.8 GHz CPU with 16 GB of RAM.

Comparison with approximate inference. Figure 7 shows that general-purpose approximate inference algorithms in WebPPL are inadequate for PCFG parsing, especially as the string parsed gets longer. Each scatterplot depicts the result of trying to parse a different string length, and each point represents one inference run. The horizontal axis shows the time taken; variation is caused primarily by different inference methods and different sample sizes. The vertical axis shows the probability estimated; variation is caused by random sampling, which has trouble explaining low-probability events: when the string length is 15, all probability estimates are either 0 or a wild overestimate, and when the string length exceeds 19, all estimates are 0. In contrast, Dice and PERPL produce answers within 0.000001% of the truth (shown by horizontal lines), in a fraction of the time taken by WebPPL (shown by vertical lines; also see below).

Comparison with exact inference. Figure 8 shows that PERPL scales better than Dice to parsing longer strings. The vertical axis shows time in log scale. Along the horizontal axis, we varied the string length between 1 and 100, but Dice ran out of memory at length 22. The plot shows that the running time of Dice (not to mention WebPPL) grows much more quickly than that of PERPL.

8 RELATED WORK

Exact inference and recursion. PERPL appears to be the first PPL that performs exact inference while allowing both unbounded recursive calls and *unbounded recursive data*. Here we consider other PPLs that support exact inference.

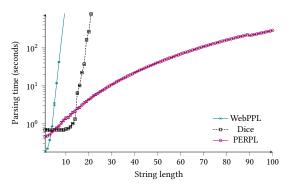


Fig. 8. The PCFG parser of Section 5.4 scales to longer strings much better than its Dice equivalent. Parsing time (*y*-axis) is on a logarithmic scale. The WebPPL curve shows how long it takes rejection sampling just to get estimates within 5% of the truth with probability 95%.

The "stochastic Lisp" of Koller et al. [1997] computes exact distributions by a generalization of variable elimination. Many PPLs that followed enable exact inference via graph representations. IBAL [Pfeffer 2001] compiles to a factor graph; Fun [Borgström et al. 2011] compiles to factor graphs with gates [Minka and Winn 2008]. Dice [Holtzen et al. 2020] and Bernoulliprob [Claret et al. 2013] compile to binary/algebraic decision diagrams [Bahar et al. 1997; Bryant 1986], which are also used by the probabilistic model checker PRISM [Kwiatkowska et al. 2011]. SPPL [Saad et al. 2021] compiles to sum-product networks [Poon and Domingos 2011], and FSPN [Stuhlmüller and Goodman 2012] generalizes sum-product networks to represent recursive dependencies. PSI [Gehr et al. 2016, 2020] and ProbZelus [Atkinson et al. 2022] perform exact inference symbolically. Finally, the expectations—and more generally *moments*—of distributions can be computed automatically by representing them with polynomials, interval arithmetic, and linear recurrences [Bartocci et al. 2019; Bouissou et al. 2016; Moosbrugger et al. 2022; Sankaranarayanan et al. 2020].

Most of these PPLs, to our knowledge, do not support unbounded recursive calls or loops, in the sense of allowing programs where for any N, there is a branch of computation with recursion/iteration depth greater than N and nonzero weight. There are four exceptions. BernoulliProb handles loops using fixed-point iteration (whereas PERPL always compiles loops to linear equations then solves them directly). FSPN handles recursive calls as in PERPL, by solving a system of equations. ProbZelus handles a stream of observations by maintaining a symbolic state, which can represent and update a distribution exactly. The methods for computing moments just mentioned are focused on unbounded loops, but do not handle other recursive calls such as those in a PCFG.

Most of these PPLs, to our knowledge, do not allow unbounded recursive data such as strings, trees, or stacks. Stochastic Lisp does, evaluating them lazily so they can be of unbounded or infinite size without necessarily causing nontermination. However, Pfeffer [2005] later noted that the interaction between lazy evaluation and memoization was problematic. This line of research continued in IBAL with a new evaluation strategy for lazy memoization, but later languages in this line, such as Figaro [Pfeffer 2016], supported exact inference only for programs with bounded recursive calls and data. Our use of defunctionalization can also be seen as lazy evaluation and covers similar cases to those that stochastic Lisp and IBAL intended to. But because it requires delayed computations to depend only on finite types, memoization is straightforward. Also, lazy evaluation would not turn programs using stacks into polynomial-time algorithms; refunctionalization does.

Probabilistic function semantics. The semantics of first-class probabilistic functions has long been known to require more than endowing a domain of functions with measurable-space structure [Aumann 1961]. Compared to recent semantics of probabilistic functions and recursion using quasi-Borel spaces [Heunen et al. 2017] and logical relations [Wand et al. 2018], our elementary treatment is limited to models where inference without sampling is possible, by solving MSPEs.

Linear logic. We build on linear logic [Girard 1987; Walker 2005] in several ways. First, the contrast between "eager" evaluation of \otimes and "lazy" evaluation of & draws on the typical computational interpretation of linear logic [Abramsky 1993]. Second, polarity [Girard 1991], which we use to allow local nonlinear bindings (Appendix B.1), has been used to structure the type system of a PPL [Ehrhard and Tasson 2019]. Finally, our treatment of linearly used functions as nondeterministic pairs, which dates back to encoding λ -calculus variables as Prolog variables [Pereira and Shieber 1987], is reminiscent of the symmetric monoidal closed category \mathcal{R}^{Π} of Laird et al. [2013], whose tensor product \otimes as well as exponential \multimap are just the Cartesian product of sets. Whereas the PPLs in this literature [Ehrhard et al. 2018; Ehrhard and Tasson 2019; Laird 2020; Laird et al. 2013] allow functions to be used any number of times and interpret them by infinitary means, we use linearity to *limit* expressivity to finitely many weight variables, reducing inference to solving an MSPE.

De- and refunctionalization. De- and refunctionalization [Danvy and Millikin 2009; Danvy and Nielsen 2001; Rendel et al. 2015; Reynolds 1972] are not new, but using defunctionalization to eliminate non-functions appears novel. Defunctionalization has been proven correct before [Banerjee et al. 2001; Nielsen 2000; Pottier and Gauthier 2006], but not for a language with effects other than nontermination. Refunctionalization is defined as the (left) inverse of defunctionalization, so the correctness of one follows from that of the other. However, our $\mathcal R$ transform incorporates additional steps of disentangling and merging apply functions presented by Danvy and Millikin [2009].

9 CONCLUSION

PERPL allows programmers to write probabilistic code using unbounded recursive calls and data and still maintain exact inference. We have seen that the compilation to MSPEs can automatically derive inference algorithms that originally took significant intellectual effort to discover.

We have implemented a PERPL compiler in Haskell, including nonlinear bindings for positive types (Appendix B.1), affine bindings for all types (Appendix B.2), and defunctionalization and refunctionalization (Sections 5.2 and 5.5). Rather than compile directly to an MSPE, it compiles to an FGG (Appendix C.2). Our implementation of FGGs, which employs the inference methods of Section 4.2 and makes them automatically differentiable using PyTorch [Paszke et al. 2019], will be the subject of a future paper. Both implementations are released as open-source software.¹

¹PERPL: https://github.com/diprism/perpl FGGs: https://github.com/diprism/fggs

REFERENCES

- Samson Abramsky. 1993. Computational Interpretations of Linear Logic. *Theoretical Computer Science* 111, 1–2 (1993), 3–57. https://doi.org/10.1016/0304-3975(93)90181-R
- Miguel A. Alonso, Éric Villemonte de la Clergerie, and Manuel Vilares. 2000. A redefinition of Embedded Push-Down Automata. In *Proc. International Workshop on Tree Adjoining Grammar and Related Frameworks (TAG+5).* 19–26. https://aclanthology.org/W00-2002
- Eric Atkinson, Charles Yuan, Guillaume Baudart, Louis Mandel, and Michael Carbin. 2022. Semi-Symbolic Inference for Efficient Streaming Probabilistic Programming. PACM on Programming Languages 6, OOPSLA2, Article 184 (2022), 29 pages. https://doi.org/10.1145/3563347
- Robert J. Aumann. 1961. Borel Structures for Function Spaces. *Illinois Journal of Mathematics* 5, 4 (1961), 614–630. https://doi.org/10.1215/ijm/1255631584
- R. I[ris] Bahar, E[rica] A. Frohm, C[harles] M. Gaona, G[ary] D. Hachtel, E[nrico] Macii, A[belardo] Pardo, and F[abio] Somenzi. 1997. Algebraic Decision Diagrams and Their Applications. Formal Methods in System Design 10 (1997), 171–206. https://doi.org/10.1023/A:1008699807402
- Anindya Banerjee, Nevin Heintze, and Jon G. Riecke. 2001. Design and Correctness of Program Transformations Based on Control-Flow Analysis. In *Theoretical Aspects of Computer Software: TACS 2001 (LNCS, 2215)*. Springer, 420–447. https://doi.org/10.1007/3-540-45500-0_21
- Andrew Barber. 1996. *Dual Intuitionistic Linear Logic*. Technical Report ECS-LFCS-96-347. Laboratory for Foundations of Computer Science, Department of Computer Science, University of Edinburgh. https://www.lfcs.inf.ed.ac.uk/reports/96/ECS-LFCS-96-347/
- Ezio Bartocci, Laura Kovács, and Miroslav Stankovič. 2019. Automatic Generation of Moment-Based Invariants for Prob-Solvable Loops. In *Automated Technology for Verification and Analysis: ATVA 2019 (LNCS, Vol. 11781).* Springer, 255–276. https://doi.org/10.1007/978-3-030-31784-3_15
- Michel Bauderon and Bruno Courcelle. 1987. Graph Expressions and Graph Rewritings. *Mathematical Systems Theory* 20 (1987), 83–127. https://doi.org/10.1007/BF01692060
- Taylor L. Booth and Richard A. Thompson. 1973. Applying Probability Measures to Abstract Languages. *IEEE Trans. Comput.* C-22, 5 (1973), 442–450. https://doi.org/10.1109/T-C.1973.223746
- Johannes Borgström, Andrew D. Gordon, Michael Greenberg, James Margetson, and Jurgen Van Gael. 2011. Measure Transformer Semantics for Bayesian Machine Learning. In Programming Languages and Systems: ESOP 2011 (LNCS, Vol. 6602). Springer, 77–96. https://doi.org/10.1007/978-3-642-19718-5_5
- Olivier Bouissou, Eric Goubault, Sylvie Putot, Aleksandar Chakarov, and Sriram Sankaranarayanan. 2016. Uncertainty Propagation Using Probabilistic Affine Forms and Concentration of Measure Inequalities. In *Tools and Algorithms for the Construction and Analysis of Systems: TACAS 2016 (LNCS, Vol. 9636)*. Springer, 225–243. https://doi.org/10.1007/978-3-662-49674-9_13
- Randal E. Bryant. 1986. Graph-Based Algorithms for Boolean Function Manipulation. *IEEE Trans. Comput.* C-35, 8 (1986), 677–691. https://doi.org/10.1109/TC.1986.1676819
- Wray L. Buntine. 1994. Operations for Learning with Graphical Models. J. Artificial Intelligence Research 2 (1994), 159–225. Alexandra Butoi, Brian DuSell, Tim Vieira, Ryan Cotterell, and David Chiang. 2022. Algorithms for Weighted Pushdown Automata. In Proc. EMNLP. https://aclanthology.org/2022.emnlp-main.656
- Shu Cai, David Chiang, and Yoav Goldberg. 2011. Language-Independent Parsing with Empty Elements. In *Proc. NAACL HLT*. 212–216. https://aclanthology.org/P11-2037
- Ashok K. Chandra and Philip M. Merlin. 1977. Optimal Implementation of Conjunctive Queries in Relational Data Bases. In *Proc. STOC.* 77–90. https://doi.org/10.1145/800105.803397
- Stanley F. Chen and Joshua Goodman. 1999. An empirical study of smoothing techniques for language modeling. *Computer Speech & Language* 13, 4 (1999), 359–394. https://doi.org/10.1006/csla.1999.0128
- David Chiang and Darcey Riley. 2020. Factor Graph Grammars. In *Proc. NeurIPS*. 6648–6658. https://proceedings.neurips.cc/paper/2020/hash/49ca03822497d26a3943d5084ed59130-Abstract.html
- Guillaume Claret, Sriram K. Rajamani, Aditya V. Nori, Andrew D. Gordon, and Johannes Borgström. 2013. Bayesian Inference Using Data Flow Analysis. In *Proc. ESEC/FSE 2013*. 92–102. https://doi.org/10.1145/2491411.2491423
- Michael Collins. 1999. *Head-Driven Statistical Models for Natural Language Parsing*. Ph. D. Dissertation. University of Pennsylvania.
- Gregory F. Cooper. 1990. The Computational Complexity of Probabilistic Inference Using Bayesian Belief Networks. *Artificial Intelligence* 42, 2 (1990), 393–405. https://doi.org/10.1016/0004-3702(90)90060-D
- Olivier Danvy and Andrzej Filinski. 1990. Abstracting Control. In *Proc. ACM Conference on Lisp and Functional Programming (LFP)*. 151–160. https://doi.org/10.1145/91556.91622
- Olivier Danvy and Kevin Millikin. 2009. Refunctionalization at Work. Science of Computer Programming 74, 8 (2009), 534–549. https://doi.org/10.1016/j.scico.2007.10.007

- Olivier Danvy and Lasse R. Nielsen. 2001. Defunctionalization at Work. In Proc. 3rd International Conference on Principles and Practice of Declarative Programming. 162–174. https://doi.org/10.1145/773184.773202
- Frank Drewes, Hans-Jörg Kreowski, and Annegret Habel. 1997. Hyperedge Replacement Graph Grammars. In *Handbook of Graph Grammars and Computing by Graph Transformation*, Grzegorz Rozenberg (Ed.). World Scientific, 95–162. https://doi.org/10.1142/9789812384720 0002
- Thomas Ehrhard, Michele Pagani, and Christine Tasson. 2018. Full Abstraction for Probabilistic PCF. J. ACM 65, 4, Article 23 (2018), 44 pages. https://doi.org/10.1145/3164540
- Thomas Ehrhard and Christine Tasson. 2019. Probabilistic Call by Push Value. Logical Methods in Computer Science 15, 1, Article 3 (2019), 46 pages. https://doi.org/10.23638/LMCS-15(1:3)2019
- Javier Esparza, Stefan Kiefer, and Michael Luttenberger. 2007. On Fixed Point Equations over Commutative Semirings. In *Proc. STACS*. 296–307. https://doi.org/10.1007/978-3-540-70918-3 26
- Javier Esparza, Stefan Kiefer, and Michael Luttenberger. 2008. Solving Monotone Polynomial Equations. In Proc. IFIP International Conference on Theoretical Computer Science (TCS). 285–298. https://doi.org/10.1007/978-0-387-09680-3_20 Invited paper.
- Javier Esparza, Stefan Kiefer, and Michael Luttenberger. 2010. Newtonian Program Analysis. J. ACM 57, 6, Article 33 (2010), 47 pages. https://doi.org/10.1145/1857914.1857917
- Kousha Etessami and Mihalis Yannakakis. 2009. Recursive Markov Chains, Stochastic Grammars, and Monotone Systems of Nonlinear Equations. J. ACM 56, 1, Article 1 (Feb. 2009), 66 pages. https://doi.org/10.1145/1462153.1462154
- Jenny Rose Finkel, Alex Kleeman, and Christopher D. Manning. 2008. Efficient, Feature-based, Conditional Random Field Parsing. In *Proc. ACL: HLT.* 959–967. https://aclanthology.org/P08-1109
- Timon Gehr, Sasa Misailovic, and Martin Vechev. 2016. PSI: Exact Symbolic Inference for Probabilistic Programs. In Computer Aided Verification: CAV 2016 (LNCS, Vol. 9779). 62–83. https://doi.org/10.1007/978-3-319-41528-4_4
- Timon Gehr, Samuel Steffen, and Martin Vechev. 2020. λPSI: Exact Inference for Higher-Order Probabilistic Programs. In *Proc. PLDI*. 883–897. https://doi.org/10.1145/3385412.3386006
- Jean-Yves Girard. 1987. Linear Logic. Theoretical Computer Science 50 (1987), 1–101. https://doi.org/10.1016/0304-3975(87)90045-4
- Jean-Yves Girard. 1991. A New Constructive Logic: Classic Logic. Mathematical Structures in Computer Science 1, 3 (1991), 255–296. https://doi.org/10.1017/S0960129500001328
- Noah D. Goodman and Andreas Stuhlmüller. 2014. The Design and Implementation of Probabilistic Programming Languages. http://dippl.org
- Alex Graves, Santiago Fernández, and Faustino Gomez. 2006. Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks. In *Proc. ICML*. 369–376.
- Annegret Habel and Hans-Jörg Kreowski. 1987. May we introduce to you: Hyperedge replacement. In *Graph Grammars and Their Application to Computer Science: Graph Grammars 1986 (LNCS, Vol. 291)*. Springer, 15–26. https://doi.org/10.1007/3-540-18771-5_41
- $Chris Heunen, Ohad Kammar, Sam Staton, and Hongseok Yang. 2017. \ A Convenient Category for Higher-Order Probability Theory. In {\it Proc. LICS. 1-12. https://dl.acm.org/doi/10.5555/3329995.3330072}$
- Steven Holtzen, Guy Van den Broeck, and Todd Millstein. 2020. Scaling Exact Inference for Discrete Probabilistic Programs. *PACM on Programming Languages* 4, OOPSLA, Article 140 (2020), 31 pages. https://doi.org/10.1145/3428208
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional LSTM-CRF Models for Sequence Tagging. https://arxiv.org/abs/1508.01991 arXiv:1508.01991v1.
- Daphne Koller and Nir Friedman. 2009. Probabilistic Graphical Models: Principles and Techniques. MIT Press.
- Daphne Koller, David McAllester, and Avi Pfeffer. 1997. Effective Bayesian Inference for Stochastic Programs. In *Proc. AAAI*. 740–747. https://www.aaai.org/Library/AAAI/1997/aaai97-115.php
- Dexter Kozen. 1981. Semantics of Probabilistic Programs. J. Comput. System Sci. 22, 3 (1981), 328–350. https://doi.org/10. 1016/0022-0000(81)90036-2
- Frank R. Kschischang, Brendan J. Frey, and Hans-Andrea Loeliger. 2001. Factor Graphs and the Sum-Product Algorithm. *IEEE Trans. Information Theory* 47, 2 (2001), 498–519. https://doi.org/10.1109/18.910572
- Marta Z. Kwiatkowska, Gethin Norman, and David Parker. 2011. PRISM 4.0: Verification of Probabilistic Real-Time Systems. In Computer Aided Verification (CAV) (LNCS, Vol. 6806). Springer, 585–591. https://doi.org/10.1007/978-3-642-22110-1_47
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In *Proc. ICML*. 282–289.
- James Laird. 2020. Weighted Models for Higher-Order Computation. *Information and Computation* 275, Article 104645 (2020), 34 pages. https://doi.org/10.1016/j.ic.2020.104645
- Jim Laird, Giulio Manzonetto, Guy McCusker, and Michele Pagani. 2013. Weighted Relational Models of Typed Lambda-Calculi. In *Proc. LICS*. 301–310. https://doi.org/10.1109/LICS.2013.36

Bernard Lang. 1974. Deterministic Techniques for Efficient Non-Deterministic Parsers. In *Proc. Colloquium on Automata, Languages, and Programming*. 255–269. https://doi.org/10.1007/978-3-662-21545-6_18

Bernard Lang. 1994. Recognition Can Be Harder Than Parsing. Computational Intelligence 10, 4 (1994), 486–494. https://doi.org/10.1111/j.1467-8640.1994.tb00011.x

Daniel J. Lehmann. 1977. Algebraic Structures for Transitive Closure. *Theoretical Computer Science* 4, 1 (1977), 59–76. https://doi.org/10.1016/0304-3975(77)90056-1

Jonathan May and Kevin Knight. 2006. Tiburon: A Weighted Tree Automata Toolkit. In *Proc. Conference on Implementation and Application of Automata (CIAA)*. 102–113. https://doi.org/10.1007/11812128_11

David McAllester, Michael Collins, and Fernando Pereira. 2008. Case-Factor Diagrams for Structured Probabilistic Modeling. J. Comput. System Sci. 74, 1 (2008), 84–96. https://doi.org/10.1016/j.jcss.2007.04.015

Tom Minka and John Winn. 2008. Gates. In Proc. NeurIPS. 1073-1080. https://papers.nips.cc/paper/3379-gates

Mehryar Mohri. 1997. Finite-State Transducers in Language and Speech Processing. *Computational Linguistics* 23, 2 (1997), 269–311. https://aclanthology.org/J97-2003

Marcel Moosbrugger, Miroslav Stankovič, Ezio Bartocci, and Laura Kovács. 2022. This is the Moment for Probabilistic Loops. *PACM on Programming Languages* 6, OOPSLA2, Article 178 (2022), 29 pages. https://doi.org/10.1145/3563341

Lasse R. Nielsen. 2000. A Denotational Investigation of Defunctionalization. Report RS-00-47. BRICS. https://www.brics.dk/RS/00/47/

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Proc. NeurIPS*. 8024–8035. http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library

Fernando C. N. Pereira and Stuart M. Shieber. 1987. *Prolog and Natural-Language Analysis*. Center for the Study of Language and Information.

Avi Pfeffer. 2001. IBAL: A Probabilistic Rational Programming Language. In Proc. IJCAI. 733-740.

Avi Pfeffer. 2005. The Design and Implementation of IBAL: A General-Purpose Probabilistic Language. Technical Report TR-12-05. Harvard Computer Science Group. http://nrs.harvard.edu/urn-3:HUL.InstRepos:25105000

Avi Pfeffer. 2016. Practical Probabilistic Programming. Manning.

Benjamin C. Pierce. 2002. Types and Programming Languages. MIT Press.

Hoifung Poon and Pedro Domingos. 2011. Sum-Product Networks: A New Deep Architecture. In Proc. UAI. 337–346. https://arxiv.org/abs/1202.3732

François Pottier and Nadji Gauthier. 2006. Polymorphic Typed Defunctionalization and Concretization. *Higher-Order and Symbolic Computation* 19, 1 (March 2006), 125–162. https://doi.org/10.1007/s10990-006-8611-7

Lawrence R. Rabiner. 1989. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proc. IEEE* 77, 2 (Feb. 1989), 257–286. https://doi.org/10.1109/5.18626

Tillmann Rendel, Julia Trieflinger, and Klaus Ostermann. 2015. Automatic Refunctionalization to a Language with Copattern Matching: With Applications to the Expression Problem. In *Proc. ICFP*. 269–279. https://doi.org/10.1145/2784731.2784763 John C. Reynolds. 1972. Definitional Interpreters for Higher-Order Programming Languages. In *Proc. ACM*. 717–740. https://doi.org/10.1145/800194.805852

Michael Riley, Cyril Allauzen, and Martin Jansche. 2009. OpenFst: An Open-Source, Weighted Finite-State Transducer Library and its Applications to Speech and Language. In *Proc. NAACL HLT, Companion Volume: Tutorial Abstracts.* 9–10. https://aclanthology.org/N09-4005

Feras A. Saad, Martin C. Rinard, and Vikash K. Mansinghka. 2021. SPPL: Probabilistic Programming with Fast Exact Symbolic Inference. In *Proc. PLDI.* 804–819. https://doi.org/10.1145/3453483.3454078

Sriram Sankaranarayanan, Yi Chou, Eric Goubault, and Sylvie Putot. 2020. Reasoning about Uncertainties in Discrete-Time Dynamical Systems using Polynomial Forms. In *Proc. NeurIPS*, Vol. 33. 17502–17513. https://proceedings.neurips.cc/paper/2020/hash/ca886eb9edb61a42256192745c72cd79-Abstract.html

Dorai Sitaram and Matthias Felleisen. 1990. Control Delimiters and Their Hierarchies. *Lisp and Symbolic Computation* 3, 1 (Jan. 1990), 67–99. https://doi.org/10.1007/BF01806126

Mitchell Stern, Jacob Andreas, and Dan Klein. 2017. A Minimal Span-Based Neural Constituency Parser. In *Proc. ACL (Volume 1: Long Papers).* 818–827. https://doi.org/10.18653/v1/P17-1076

Alistair Stewart, Kousha Etessami, and Mihalis Yannakakis. 2015. Upper Bounds for Newton's Method on Monotone Polynomial Systems, and P-Time Model Checking of Probabilistic One-Counter Automata. J. ACM 62, 4, Article 30 (2015), 33 pages. https://doi.org/10.1145/2789208

Andreas Stuhlmüller and Noah Goodman. 2012. A Dynamic Programming Algorithm for Inference in Recursive Probabilistic Programs. In *Proc. Workshop on Statistical Relational AI (StaRAI)*. https://arxiv.org/abs/1206.3555

- Ben Taskar, Dan Klein, Mike Collins, Daphne Koller, and Christopher Manning. 2004. Max-Margin Parsing. In *Proc. EMNLP*. 1–8. https://www.aclweb.org/anthology/W04-3201
- Jan-Willem van de Meent, Brooks Paige, Hongseok Yang, and Frank Wood. 2018. An Introduction to Probabilistic Programming. https://arxiv.org/abs/1809.10756 arXiv:1809.10756.
- K. Vijay-Shankar and Aravind K. Joshi. 1985. Some Computational Properties of Tree Adjoining Grammars. In *Proc. ACL*. 82–93. https://doi.org/10.3115/981210.981221
- K. Vijayashanker. 1994. A Study of Tree Adjoining Grammars. Ph. D. Dissertation. Univ. of Pennsylvania.
- John von Neumann. 1951. Various Techniques Used in Connection With Random Digits. *National Bureau of Standards Applied Mathematics Series* 12 (1951), 36–38.
- Rajan Walia, Praveen Narayanan, Jacques Carette, Sam Tobin-Hochstadt, and Chung-chieh Shan. 2019. From High-Level Inference Algorithms to Efficient Code. PACM on Programming Languages 3, ICFP, Article 98 (2019), 30 pages. https://doi.org/10.1145/3352468
- David Walker. 2005. Substructural Type Systems. In *Advanced Topics in Types and Programming Languages*, Benjamin C. Pierce (Ed.). MIT Press, Chapter 1, 3–43.
- Mitchell Wand, Ryan Culpepper, Theophilos Giannakopoulos, and Andrew Cobb. 2018. Contextual Equivalence for a Probabilistic Language with Continuous Random Variables and Recursion. *PACM on Programming Languages* 2, ICFP, Article 87 (2018), 38 pages. https://doi.org/10.1145/3236782
- David J. Weir. 1992. A Geometric Hierarchy Beyond Context-Free Languages. *Theoretical Computer Science* 104, 2 (1992), 235–261. https://doi.org/10.1016/0304-3975(92)90124-X

Received 2022-10-28; accepted 2023-02-25

A ADDITIONAL PROOFS

A.1 Proof of Lemma 3.13

By induction on the typing derivation of e'. We show a few core cases.

Case x_1 : Because Δ_0 is empty, δ_0 is empty as well.

Case $x \neq x_1$: Ruled out by inversion.

Case $\lambda x.e$: By Barendregt's variable convention, $x \neq x_1$ and $x \notin FV(e_1)$.

$$\begin{split} \llbracket (\lambda x.e)\{x_1 := e_1\} \rrbracket (\delta_0 \cup \delta_1, v_x \mapsto v_e) &= \llbracket \lambda x.(e\{x_1 := e_1\}) \rrbracket (\delta_0 \cup \delta_1, v_x \mapsto v_e) \\ &= \llbracket e\{x_1 := e_1\} \rrbracket (\delta_0 \cup \delta_1 \cup \{(x, v_x)\}, v_e) \\ &= \sum_{v_1} \llbracket e \rrbracket (\delta_0 \cup \{(x, v_x), (x_1, v_1)\}, v_e) \cdot \llbracket e_1 \rrbracket (\delta_1, v_1) \\ &= \sum_{v_1} \llbracket \lambda x.e \rrbracket (\delta_0 \cup \{(x_1, v_1)\}, v_x \mapsto v_e) \cdot \llbracket e_1 \rrbracket (\delta_1, v_1). \end{split}$$

Case e'_0 e'_1 : Thanks to linearity, x_1 occurs free in either e'_0 or e'_1 but not both. If $x_1 \in FV(e'_0)$,

$$\begin{split} & \llbracket (e'_0 \ e'_1) \{ x_1 := e_1 \} \rrbracket (\delta'_0 \cup \delta'_1 \cup \delta_1, v') \\ &= \llbracket (e'_0 \{ x_1 := e_1 \}) \ e'_1 \rrbracket (\delta'_0 \cup \delta'_1 \cup \delta_1, v') \\ &= \sum_{v'_1} \llbracket e'_0 \{ x_1 := e_1 \} \rrbracket (\delta'_0 \cup \delta_1, v'_1 \mapsto v') \cdot \llbracket e'_1 \rrbracket (\delta'_1, v'_1) \\ &= \sum_{v'_1} \sum_{v_1} \llbracket e'_0 \rrbracket (\delta'_0 \cup \{ (x_1, v_1) \}, v'_1 \mapsto v') \cdot \llbracket e_1 \rrbracket (\delta_1, v_1) \cdot \llbracket e'_1 \rrbracket (\delta'_1, v'_1) \\ &= \sum_{v_1} \llbracket e'_0 \ e'_1 \rrbracket (\delta'_0 \cup \delta'_1 \cup \{ (x_1, v_1) \}, v') \cdot \llbracket e_1 \rrbracket (\delta_1, v_1). \end{split}$$

Similarly if $x_1 \in FV(e'_1)$.

A.2 Proofs of Type Preservation

PROPOSITION A.1. If $e:\tau$, then $\mathcal{D}_{\sigma}[\![e]\!]:\mathcal{D}_{\sigma}[\![\tau]\!]$.

PROOF. By induction on the typing derivation of $e : \tau$. The only interesting cases are **fold**_{σ} and **unfold**_{σ} expressions. The nonlinear typing context is unchanged and elided below.

Case $\mathbf{fold}_{\sigma} m_i$: The translation $\mathcal{D}_{\sigma}[\![\mathbf{fold}_{\sigma} m_i]\!] = \mathsf{Fold}_i \vec{y}_i$ has the following typing derivation.

$$rac{ec{y}_i:ec{arphi}_i}{\mathsf{Fold}_i\,ec{y}_i:arphi}$$

Case **unfold**_{σ} x = e **in** e': Because $m_i : \bar{\sigma}\{\alpha := \sigma\}$, thanks to the induction hypothesis, the function u_{σ} has the following typing derivation.

$$\frac{\vec{y}_i: \vec{\varphi}_i \vdash \mathcal{D}_\sigma[\![m_i]\!]: \bar{\sigma}\{\alpha := \varphi\} \quad \text{for each } i}{x: \varphi \vdash \mathbf{case} \ x \ \mathbf{of} \ \mathsf{Fold}_1 \ \vec{y}_1 \to \mathcal{D}_\sigma[\![m_1]\!] \mid \cdots \mid \mathsf{Fold}_n \ \vec{y}_n \to \mathcal{D}_\sigma[\![m_n]\!]: \bar{\sigma}\{\alpha := \varphi\}}$$

$$\lambda x. \ \mathbf{case} \ x \ \mathbf{of} \ \mathsf{Fold}_1 \ \vec{y}_1 \to \mathcal{D}_\sigma[\![m_1]\!] \mid \cdots \mid \mathsf{Fold}_n \ \vec{y}_n \to \mathcal{D}_\sigma[\![m_n]\!]: \varphi \multimap \bar{\sigma}\{\alpha := \varphi\}$$

Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.

Then the translation of **unfold**_{σ} x = e **in** $e' : \tau$ has the following typing derivation.

$$\frac{u_{\sigma}: \varphi \multimap \bar{\sigma}\{\alpha := \varphi\} \quad \mathcal{D}_{\sigma}[\![e]\!] : \varphi}{u_{\sigma} \, \mathcal{D}_{\sigma}[\![e]\!] : \bar{\sigma}\{\alpha := \varphi\} \quad x : \bar{\sigma}\{\alpha := \varphi\} \vdash \mathcal{D}_{\sigma}[\![e']\!] : \mathcal{D}_{\sigma}[\![\tau]\!]}$$

$$\mathbf{let} \, x = u_{\sigma} \, \mathcal{D}_{\sigma}[\![e]\!] \, \mathbf{in} \, \mathcal{D}_{\sigma}[\![e']\!] : \mathcal{D}_{\sigma}[\![\tau]\!]$$

PROPOSITION A.2. If $e:\tau$, then $\mathcal{R}_{\sigma}[\![e]\!]:\mathcal{R}_{\sigma}[\![\tau]\!]$.

PROOF. By induction on the typing derivation of $e:\tau$. The only interesting cases are **fold** $_{\sigma}$ and **unfold** $_{\sigma}$ expressions. The nonlinear typing context is unchanged and elided below.

Case **fold**_{σ} e: Because $m_i: \varphi_i$, thanks to the induction hypothesis, the function f_{σ} has the following typing derivation.

$$\frac{x : \bar{\sigma}\{\alpha := \varphi\}, \vec{y}_i : \vec{\varphi}_i \vdash \mathcal{R}_{\sigma}\llbracket m_i \rrbracket : \varphi_i}{x : \bar{\sigma}\{\alpha := \varphi\} \vdash \lambda \vec{y}_i. \mathcal{R}_{\sigma}\llbracket m_i \rrbracket : \vec{\varphi}_i \multimap \varphi_i}$$
$$\frac{x : \bar{\sigma}\{\alpha := \varphi\} \vdash \langle \lambda \vec{y}_1. \mathcal{R}_{\sigma}\llbracket m_1 \rrbracket, \dots, \lambda \vec{y}_n. \mathcal{R}_{\sigma}\llbracket m_n \rrbracket \rangle : \varphi}{\lambda x. \langle \lambda \vec{y}_1. \mathcal{R}_{\sigma}\llbracket m_1 \rrbracket, \dots, \lambda \vec{y}_n. \mathcal{R}_{\sigma}\llbracket m_n \rrbracket \rangle : \bar{\sigma}\{\alpha := \varphi\} \multimap \varphi}$$

Then the translation of **fold** $_{\sigma}$ $e:\varphi$ has the following typing derivation.

$$\frac{f_{\sigma}:\bar{\sigma}\{\alpha:=\varphi\}\multimap\varphi\quad\mathcal{R}_{\sigma}[\![e]\!]:\bar{\sigma}\{\alpha:=\varphi\}}{f_{\sigma}\,\mathcal{R}_{\sigma}[\![e]\!]:\varphi}$$

Case $\operatorname{unfold}_{\sigma} x = m_i \operatorname{in} m_i$: The translation $\mathcal{R}_{\sigma}[[\operatorname{unfold}_{\sigma} x = m_i \operatorname{in} m_i]] = (\mathcal{R}_{\sigma}[[m_i]].i) \vec{y}_i$ has the following typing derivation.

$$\frac{\mathcal{R}_{\sigma}\llbracket m_{i} \rrbracket : \varphi}{\mathcal{R}_{\sigma}\llbracket m_{i} \rrbracket . i : \vec{\varphi}_{i} \multimap \varphi_{i} \qquad \vec{y}_{i} : \vec{\varphi}_{i}} \qquad \qquad \Box$$

$$(\mathcal{R}_{\sigma}\llbracket m_{i} \rrbracket . i) \ \vec{y}_{i} : \varphi_{i}$$

A.3 Correctness of Transformations

Sometimes, the programs before and after transformation can be executed in lockstep, given that our operational semantics allows evaluation anywhere in an expression. In Section 5.4, having fixed the target type $\sigma = \text{String}^2$ and located the n = 4 occurrences of $\mathbf{fold}_{\text{String}^2}$, we can treat the transformation $\mathcal{D}_{\text{String}^2}[\![\cdot]\!]$ as a relation \sim on terms,

$$e^r \sim e^d \quad \text{iff} \quad \mathcal{D}_{\text{String}^2} \llbracket e^r \rrbracket = e^d.$$
 (25)

Figure 9 shows an instance of such lockstep execution.

We'd like to be able to show that if $e^r \sim e^d$, then for any reduction performed on one side, both sides can be reduced (in zero or more steps) to related distributions $E^r \sim E^d$, as in the above example. Unfortunately, the relation defined in (25) is too simple to keep track of reductions, especially probabilistic choices, made inside \mathbf{fold}_σ expressions. For example, instead of parsing a literal string in Section 5.4, it may be useful to parse a random string. When an expression, say, $e^r = \mathbf{fold}_\sigma$ (amb) reduces on the left side, the corresponding reduction(s) on the right side would have to be made inside the global definition of u_σ . Worse, different occurrences of e^r can make different choices, and the global definition of u_σ only has room to record one of them.

Figure 10 defines an enriched relation \sim that accounts for this complexity. Like the simplistic attempt at \sim above, this relation depends implicitly on a defunctionalization target (fixing σ , m_i , and \vec{y}_i). However, what's related are no longer just terms ($e^r \sim e^d$), but distributions over terms

Fig. 9. Lockstep execution in a simple case of defunctionalization. We write e for each distribution $\{(1, e)\}$.

Mixture
$$\frac{\gamma + E_{j}^{r} \sim E_{j}^{d} \text{ for all } j \in J \text{ finite}}{\gamma + \sum_{j \in J} w_{j} \cdot E_{j}^{r} \sim \sum_{j \in J} w_{j} \cdot E_{j}^{d}}$$

$$\gamma + E_{j}^{r} \sim E_{j}^{d} \text{ for } 1 \leq j \leq |\vec{y}_{i}|$$

$$\gamma + \{(w_{1} \cdots w_{|\vec{y}_{i}|}, \mathbf{fold}_{\sigma} m_{i} \{\vec{y}_{i} := (e_{1}^{r}, \dots, e_{|\vec{y}_{i}|}^{r})\}) \mid (w_{j}, e_{j}^{r}) \in E_{j}^{r} \text{ for } 1 \leq j \leq |\vec{y}_{i}|\} \Longrightarrow^{*} E^{r}$$

$$\gamma + E^{r} \sim \{(w_{1} \cdots w_{|\vec{y}_{i}|}, F_{ol}d(e_{1}^{d}, \dots, e_{|\vec{y}_{i}|}^{d})) \mid (w_{j}, e_{j}^{d}) \in E_{j}^{d} \text{ for } 1 \leq j \leq |\vec{y}_{i}|\}$$

$$\gamma + E^{r} \sim E^{d} \qquad \gamma + E^{rr} \sim E^{rd}$$

$$\gamma + \{(ww', \mathbf{unfold}_{\sigma} x = e^{r} \mathbf{in} e^{rr}) \mid (w, e^{r}) \in E^{r}, (w', e^{rr}) \in E^{rr}\}$$

$$\sim \{(ww', \mathbf{let} \ x = u_{\sigma} e^{d} \mathbf{in} e^{rd}) \mid (w, e^{d}) \in E^{d}, (w', e^{rd}) \in E^{rd}\}$$

$$\gamma + \{(1, x)\} \sim \{(1, x)\} \qquad \gamma + \{(1, \mathbf{fail})\} \sim \{(1, \mathbf{fail})\}$$

$$\gamma + \{(w, \lambda x. e^{r}) \mid (w, e^{r}) \in E^{r}\}$$

$$\gamma + \{(w_{0}w_{1}, e_{0}^{r} e_{1}^{r}) \mid (w_{0}, e_{0}^{r}) \in E_{0}^{r}, (w_{1}, e_{1}^{r}) \in E_{1}^{r}\}$$

$$\sim \{(w, \lambda x. e^{d}) \mid (w, e^{d}) \in E^{d}\}$$

$$\sim \{(w_{0}w_{1}, e_{0}^{d} e_{1}^{d}) \mid (w_{0}, e_{0}^{d}) \in E_{0}^{d}, (w_{1}, e_{1}^{d}) \in E_{1}^{d}\}$$

and so on for all other syntax constructions

Fig. 10. Relating distributions of expressions before and after defunctionalization. In the Fold and Unfold rules, the definitions of σ , m_i , and \vec{y}_i are as in Section 5.2.1.

 $(E^r \sim E^d)$. These distributions are built by the Mixture rule, by taking linear combinations of related distributions. Also, because the reduction relation \Longrightarrow^* is parameterized by the global definitions γ on the left side, the relation \sim is parameterized by γ too.

We no longer define the relation \sim in terms of the transform \mathcal{D}_{σ} as in (25). Still, the rest of the inference rules in Figure 10 give rise to the transform \mathcal{D}_{σ} as a consequence:

Proposition A.3.
$$\{(1, e^r)\} \sim \{(1, \mathcal{D}_{\sigma} \llbracket e^r \rrbracket)\}.$$

PROOF. By induction on e^r . The interesting cases—those shown in (18)—use the Fold and Unfold rules. The other, uninteresting cases use the Congruence rules.

The Fold rule's second premise allows m_i inside **fold**_{σ} to be reduced arbitrarily and probabilistically to an entire distribution of terms and yet stay related to the same $Fold_i$ expression. These reductions have executed already on the left side, in E^r , but are "pent up" on the right side. They must wait for the $fold_{\sigma}$ on the left side to meet its annihilating $unfold_{\sigma}$, whose translation on the right side provides the u_{σ} whose definition contains m_i .

Lemma A.4 (Substitution preserves \sim). Suppose $\gamma \vdash E^r \sim E^d$ and $\gamma \vdash E_1^r \sim E_1^d$. If x_1 is a linear variable, or if x_1 is a nonlinear variable but E_1^r and E_1^d use no linear variables, then

$$\begin{split} \gamma &\vdash \left\{ \left(ww_1, e^r \{x_1 := e_1^r\}\right) \mid (w, e^r) \in E^r, (w_1, e_1^r) \in E_1^r \right\} \\ &\sim \left\{ \left(ww_1, e^d \{x_1 := e_1^d\}\right) \mid (w, e^d) \in E^d, (w_1, e_1^d) \in E_1^d \right\}. \end{split}$$

PROOF. By induction on the derivation of $\gamma \vdash E^r \sim E^d$.

THEOREM A.5. Whenever $E^r \sim E^d$:

- (1) If $E^r \Longrightarrow E^{r\dagger}$, then there exist $E^{r\ddagger}$ and $E^{d\ddagger}$ such that $E^{r\dagger} \Longrightarrow^* E^{r\ddagger}$, $E^d \Longrightarrow^* E^{d\ddagger}$, and $E^{r\ddagger} \sim E^{d\ddagger}$. (2) If $E^d \Longrightarrow E^{d\dagger}$, then there exist $E^{r\ddagger}$ and $E^{d\ddagger}$ such that $E^r \Longrightarrow^* E^{r\ddagger}$, $E^{d\dagger} \Longrightarrow^* E^{d\ddagger}$, and $E^{r\ddagger} \sim E^{d\ddagger}$.

Proof. First observe that any use of the Mixture rule in a derivation of $E^r \sim E^d$ either occurs at the outermost (bottommost) level or can be commuted to occur there. And if $E^r \sim E^d$ is derived using the Mixture rule at the outermost level, then the result follows from the definition of \Longrightarrow on distributions (Definition 3.1). Thus, we may assume from here on that the Mixture rule is not used, and so the only source of nondeterminism (Definition 3.12) is the second premise of the Fold rule. We proceed by simultaneous induction on the given derivations of \sim and \Longrightarrow .

If $E^r \sim E^d$ is derived using a Congruence rule, say, the one for application (the last rule shown in Figure 10), then the given reduction ($E^r \Longrightarrow E^{r\dagger}$ or $E^d \Longrightarrow E^{d\dagger}$) is either a β -reduction at the top level or a reduction of a subexpression (e_0^r or e_1^q or e_0^d or e_1^d) inside an evaluation context. Handle the β -reduction case using Lemma A.4. Handle the subexpression case using the induction hypothesis. In both cases, if the given reduction operates on a nondeterministic distribution (more precisely, if E_0 in Definition 3.1 is nonzero), then use the flexibility afforded by $E^{r\dagger} \Longrightarrow^* E^{r\ddagger}$ or $E^{d\dagger} \Longrightarrow^* E^{d\ddagger}$ to mimic the given reduction on the rest of the distribution.

If $E^r \sim E^d$ is derived using the Unfold rule, then the given reduction is (16) or (4i). Handle these possibilities like with a Congruence rule. The only difference is that, after handling (16) using Lemma A.4, any pent-up reductions from the Fold rule need to be released into $E^d \Longrightarrow^* E^{d^{\frac{1}{+}}}$.

If $E^r \sim E^d$ is derived using the Fold rule, then a reduction $E^r \Longrightarrow E^{r\dagger}$ can just be pent up in the witness for the second premise.

Finally, we can show that the program (possibly surrounded with some code that converts the answer to ()) denotes the same distribution before and after \mathcal{D} .

Proposition A.6.
$$\{(w^r, ())\} \sim \{(w^d, ())\}$$
 just in case $w^r = w^d$.

PROOF. By induction on the derivation of \sim .

Refunctionalization is treated similarly, using the \sim relation shown in Figure 11.

Inferring the Sequence of Transformations

At first glance, DR-sequences are difficult to reason about, because each transformation potentially changes the types available to subsequent transformations. Instead, we define a graph structure on the recursive types occurring in p and show that finding a successful DR-sequence is equivalent to a finding a particular kind of subgraph.

Fig. 11. Relating distributions of expressions before and after refunctionalization. In the Fold and Unfold rules, the definitions of σ , m_i , and \vec{y}_i are as in Section 5.2.2. Mixture and Congruence rules are like in Figure 10.

Fig. 12. (a) The full DR-graph of the program of Section 5.4. (b) The DR-subgraph of refunctionalizing both String¹ and String². (c) The DR-subgraph of refunctionalizing String¹ and defunctionalizing String².

Definition A.7. The DR-graph of a program p is an edge-labeled, directed graph whose nodes are the recursive types in p. There is an edge $\sigma \xrightarrow{\mathcal{D}} \tau$ iff $\mathcal{D}_{\sigma}[\![\sigma]\!] = \varphi$ contains τ , and an edge $\sigma \xrightarrow{\mathcal{R}} \tau$ iff $\mathcal{R}_{\sigma}[\![\sigma]\!] = \varphi$ contains τ .

Example A.8. The DR-graph of the program of Section 5.4 is shown in Figure 12a. Observe that because type $String^1$ has a self-loop labeled \mathcal{D} , defunctionalization of $String^1$ is not allowed. The \mathcal{R} edges form a cycle, which does not immediately prevent either type from being refunctionalized, but if we did refunctionalize String[2], the cycle would become a self-loop. Then $String^1$ would be neither defunctionalizable nor refunctionalizable. Intuitively, then, we want to avoid choices of defunctionalization and refunctionalization that form cycles.

Definition A.9. A DR-subgraph G for a program p is a subgraph of p's DR-graph where there is a mapping from each recursive type σ to a transformation $G(\sigma) \in \{\mathcal{D}, \mathcal{R}\}$ such that G keeps just the edges $\sigma \xrightarrow{G(\sigma)} \tau$ in the DR-graph. Intuitively, $G(\sigma)$ is what G plans to do with σ .

Example A.10. In Section 5.4, the DR-subgraph of refunctionalizing both String¹ and String², shown in Figure 12b, is cyclic, while the DR-subgraph corresponding to refunctionalizing String¹ and defunctionalizing String², shown in Figure 12c, is acyclic.

Definition A.11. Let G be a DR-subgraph for a program p. For any recursive type σ in p and transformation $\mathcal{F} \in \{\mathcal{D}, \mathcal{R}\}$, we write $\mathcal{F}_{\sigma}[\![G]\!]$ for the unique DR-subgraph of $\mathcal{F}_{\sigma}[\![p]\!]$ such that $\mathcal{F}_{\sigma}[\![G]\!](\mathcal{F}_{\sigma}[\![\tau]\!]) = G(\tau)$ for all $\tau \neq \sigma$. Intuitively, $\mathcal{F}_{\sigma}[\![G]\!]$ is what G plans to do with all the other recursive types after σ is eliminated via \mathcal{F} .

Lemma A.12. Let G be a DR-subgraph for a program p, and let σ be a recursive type in p.

(1) If $G(\sigma)_{\sigma}[\![p]\!]$ is defined (that is, if G does not have a self-loop $\sigma \to \sigma$), then $G(\sigma)_{\sigma}[\![G]\!]$ is acyclic iff G is.

(2) For any $\mathcal{F} \in \{\mathcal{D}, \mathcal{R}\}$ such that $\mathcal{F}_{\sigma}[\![\sigma]\!]$ contains no recursive type, $\mathcal{F}_{\sigma}[\![\sigma]\!]$ is acyclic iff G is.

PROOF. For any $\mathcal{F} \in \{\mathcal{D}, \mathcal{R}\}$ in general, $\mathcal{F}_{\sigma}[\![G]\!]$ is acyclic iff G is. This follows from the observation that for any $\tau, \tau' \neq \sigma$, the graph $\mathcal{F}_{\sigma}[\![G]\!]$ has the edge $\mathcal{F}_{\sigma}[\![\tau]\!] \to \mathcal{F}_{\sigma}[\![\tau']\!]$ iff the graph G has either the edge $\tau \to \tau'$ or the edges $\tau \to \sigma \to \tau'$. This observation holds because \mathcal{F}_{σ} is injective on nodes, thanks to the requirement that different recursive types in p must have different tags. \square

Proposition A.13. A program p has a successful DR-sequence iff it has an acyclic DR-subgraph G.

PROOF. By induction on the number of recursive types in p.

- (\Rightarrow) If p has a successful DR-sequence $\mathcal{F}_{\sigma} \cdot S'$, that means $\mathcal{F}_{\sigma}[\![p]\!]$ exists and has a successful DR-sequence (namely S'), so it has an acyclic DR-subgraph G' by the induction hypothesis. Let G be the unique DR-subgraph for p such that $G(\sigma) = \mathcal{F}$ and $G' = \mathcal{F}_{\sigma}[\![G]\!]$. Lemma A.12.1(\Rightarrow) says that G must be acyclic.
- (\Leftarrow) Let σ be a recursive type in p, and let $\mathcal{F} = G(\sigma)$. Since G is acyclic, the program $p' = \mathcal{F}_{\sigma}[\![p]\!]$ exists and has a DR-subgraph $G' = \mathcal{F}_{\sigma}[\![G]\!]$, which is also acyclic by Lemma A.12.1(\Leftarrow). By the induction hypothesis, p' has a successful DR-sequence S', so $\mathcal{F}_{\sigma} \cdot S'$ is a successful DR-sequence for p.

PROOF OF PROPOSITION 5.5. Because p has a successful DR-sequence, it has an acyclic DR-subgraph G. (1) Because G is acyclic and finite, it must have a node σ with no outgoing edges. Let $\mathcal{F} = G(\sigma)$. Then $\mathcal{F}_{\sigma}[\![\sigma]\!]$ contains no recursive type. (2) For any such σ and \mathcal{F} , by Lemma A.12.2(\Leftarrow), $\mathcal{F}_{\sigma}[\![\sigma]\!]$ is acyclic. So $\mathcal{F}_{\sigma}[\![p]\!]$ has an acyclic DR-subgraph and therefore a successful DR-sequence. \square

B RELAXING LINEARITY

Three extensions to our linear type system make it easier to write practical programs without fundamentally affecting the expressivity of the language or the efficiency of the implementation. These extensions are

- (1) allowing local nonlinear bindings of *positive* type;
- (2) allowing linear variables to be discarded unused (in other words, used affinely); and
- (3) allowing local nonlinear recursive bindings that do not use linear variables.

The example programs above already use the first two extensions. We describe each of these extensions in turn.

B.1 Positive Types

If τ has positive polarity in the sense of linear logic [Ehrhard and Tasson 2019; Girard 1991], that is, it uses only \oplus and \otimes , intuitively it should be possible to copy values of type τ , because a contraction (copying) function $\tau \multimap (\tau \otimes \tau)$ can be easily defined by induction on τ . Here, we show how to extend the type system to allow linear bindings to be turned into nonlinear bindings if they have positive type.

We say that τ is positive if it does not contain \multimap , &, or μ . Thus, functions and additive pairs must still be used linearly. Moreover, we decree that recursive types are not positive (so a value of recursive type must be consumed (**unfold**ed) exactly once) in preparation for transforming recursive types to other types that might not be positive.

Let us write **let** !x = e **in** e' to use $x : \tau$ any number of times in e'.

Syntax

$$e ::= \mathbf{let} \ !x = e \ \mathbf{in} \ e$$

Typing

$$\frac{\Gamma; \Delta \vdash e : \tau \qquad \tau \text{ is positive} \qquad \Gamma, x : \tau; \Delta' \vdash e' : \tau'}{\Gamma; \Delta, \Delta' \vdash \mathbf{let} ! x = e \text{ in } e' : \tau'}$$

The operational semantics of the new construct is as follows.

Reduction

$$\gamma \vdash \mathbf{let} ! x = v \text{ in } e' \Longrightarrow \{(1, e'\{x := v\})\}$$
 (26)

Evaluation contexts

$$C ::= let !x = C in e | let !x = e in C$$

Lemma B.1 (Nonlinear substitution preserves typing). (cf. Lemma 3.2) Suppose $\Gamma, x_1 \colon \tau_1; \Delta_0 \vdash e' \colon \tau'$ and $\Gamma; \cdot \vdash e_1 \colon \tau_1$. Then $\Gamma; \Delta_0 \vdash e' \{x_1 \coloneqq e_1\} \colon \tau'$.

PROOF. By induction on the typing derivation of e'.

Hence the new reduction (26) preserves typing (extending Proposition 3.4): v uses no linear variables, because it is a syntactic value whose type is positive.

The denotational semantics of the new construct extends the local environment δ with the new nonlinear binding.

$$[\![let !x = e in e']\!](\delta \cup \delta', v') = \sum_{v} [\![e]\!](\delta, v) \cdot [\![e']\!](\delta' \cup \{(x, v)\}, v')$$
(27)

This extension makes the division of names between the nonlinear and linear typing contexts Γ and Δ no longer match the division of names between the global and local environments γ and δ . To maintain this match, the denotation could alternatively be defined by extending the global environment γ with a distribution that is always deterministic.

Lemma B.2 (Nonlinear substitution preserves denotation). (cf. Lemma 3.13)

Suppose $\Gamma, x_1 : \tau_1; \Delta_0 \vdash e' : \tau'$ and $\Gamma; \cdot \vdash e_1 : \tau_1$. Then

$$[\![e'\{x_1 := e_1\}]\!]_{\eta}(\delta_0, v') = \sum_{v_1} [\![e']\!]_{\eta}(\delta_0 \cup \{(x_1, v_1)\}, v') \cdot [\![e_1]\!]_{\eta}(\emptyset, v_1)$$
(28)

for all $\eta \in \llbracket \Gamma \rrbracket$, $\delta_0 \in \llbracket \Delta_0 \rrbracket$, and $v' \in \llbracket \tau' \rrbracket$ such that $\llbracket e_1 \rrbracket_{\eta}$ is deterministic.

PROOF. By induction on the typing derivation of e'. We show a few core cases.

Case x_1 : As in Appendix A.1, but let $\delta_1 = \emptyset$.

Case $x \neq x_1$: The assumption that $[e_1]$ is deterministic justifies the second step:

$$[\![x\{x_1:=e_1\}]\!](\delta_0,v')=[\![x]\!](\delta_0,v')=\sum_{v_1}[\![x]\!](\delta_0\cup\{(x_1,v_1)\},v')\cdot[\![e_1]\!](\emptyset,v_1).$$

Case $\lambda x.e$: As in Appendix A.1, but let $\delta_1 = \emptyset$.

Case e_0' e_1' : The assumption that $[e_1]$ is deterministic justifies the second-to-last step:

$$\begin{split} & \big[\!\big[(e'_0 \ e'_1) \big\{ x_1 := e_1 \big\} \big]\!\big] (\delta'_0 \cup \delta'_1, v') \\ &= \big[\!\big[(e'_0 \big\{ x_1 := e_1 \big\} \big] (e'_1 \big\{ x_1 := e_1 \big\} \big]\!\big] (\delta'_0 \cup \delta'_1, v') \\ &= \sum_{v'_1} \big[\!\big[e'_0 \big\{ x_1 := e_1 \big\} \big]\!\big] (\delta'_0, v'_1 \mapsto v') \cdot \big[\!\big[e'_1 \big\{ x_1 := e_1 \big\} \big]\!\big] (\delta'_1, v'_1) \\ &= \sum_{v'} \Big(\sum_{v_1} \big[\!\big[e'_0 \big]\!\big] (\delta'_0 \cup \big\{ (x_1, v_1) \big\}, v'_1 \mapsto v') \cdot \big[\!\big[e_1 \big]\!\big] (\emptyset, v_1) \Big) \Big(\sum_{v_1} \big[\!\big[e'_1 \big]\!\big] (\delta'_1 \cup \big\{ (x_1, v_1) \big\}, v'_1) \cdot \big[\!\big[e_1 \big]\!\big] (\emptyset, v_1) \Big) \end{split}$$

Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.

$$\begin{split} &= \sum_{v_1'} \sum_{v_1} \llbracket e_0' \rrbracket (\delta_0' \cup \{(x_1, v_1)\}, v_1' \mapsto v') \cdot \llbracket e_1' \rrbracket (\delta_1' \cup \{(x_1, v_1)\}, v_1') \cdot \llbracket e_1 \rrbracket (\emptyset, v_1) \\ &= \sum_{v_1} \llbracket e_0' \ e_1' \rrbracket (\delta_0' \cup \delta_1' \cup \{(x_1, v_1)\}, v') \cdot \llbracket e_1 \rrbracket (\emptyset, v_1). \end{split}$$

Hence the new reduction (26) preserves denotation (extending Theorem 3.14): $\llbracket v \rrbracket$ is deterministic, because v is a syntactic value whose type is positive.

B.2 Affine Types

The previous extension does not add any flexibility for nonpositive types (those containing \neg , &, or μ). But we can further extend the language to allow values of nonpositive type to be used affinely (zero times or once). It is easy to express this extension in the type system, by allowing the linear typing context Δ to be non-empty in the rules for variables:

$$\frac{}{\Gamma, x : \tau; \Delta \vdash x : \tau} \qquad \frac{}{\Gamma; \Delta, x : \tau \vdash x : \tau}$$

However, allowing functions to be unused would be problematic for our denotational semantics, because any effects (nondeterminism and probabilities) inside the function body will be incurred regardless of whether the function is called; similarly for other nonpositive types.

To resolve this problem in a modular manner, we define a source-to-source transformation that operates on the typing derivation of a program and produces an equivalent program that uses values of nonpositive type only linearly. The idea is that whenever we want to create an affinely used function $\lambda x.e$, we create a &-pair containing a linearly used function ($\lambda x.e$) or nothing (that is, ()). To apply the function, we use the former; to discard it, we use the latter. We do something similar with recursive types, which can now be unfolded zero times or once, and additive pairs (&), which can also be eliminated zero times or once.

Formally, we define the transformation \mathcal{L} on types:

$$\mathcal{L}\llbracket \tau_{1} \multimap \tau' \rrbracket = (\mathcal{L}\llbracket \tau_{1} \rrbracket) \multimap \mathcal{L}\llbracket \tau' \rrbracket) \& \mathsf{Unit}$$

$$\mathcal{L}\llbracket \tau_{1} \otimes \cdots \otimes \tau_{n} \rrbracket = \mathcal{L}\llbracket \tau_{1} \rrbracket \otimes \cdots \otimes \mathcal{L}\llbracket \tau_{n} \rrbracket$$

$$\mathcal{L}\llbracket \tau_{1} \& \cdots \& \tau_{n} \rrbracket = (\mathcal{L}\llbracket \tau_{1} \rrbracket) \& \cdots \& \mathcal{L}\llbracket \tau_{n} \rrbracket) \& \mathsf{Unit}$$

$$\mathcal{L}\llbracket c_{1} \tau_{1} \oplus \cdots \oplus c_{n} \tau_{n} \rrbracket = c_{1} \mathcal{L}\llbracket \tau_{1} \rrbracket \oplus \cdots \oplus c_{n} \mathcal{L}\llbracket \tau_{n} \rrbracket$$

$$\mathcal{L}\llbracket \mu\alpha. \tau \rrbracket = \mu\alpha. \mathcal{L}\llbracket \tau \rrbracket$$

$$\mathcal{L}\llbracket \alpha \rrbracket = \alpha$$

On terms, we show the interesting cases:

The $\mathbb{Z}[\![\cdot]\!]$ transformation discards variables. How this is done depends on the pre-transformation type of each variable:

$$\mathcal{Z}[\![\cdot]\!] = ()$$

$$\mathcal{Z}[\![\Delta, x : \tau]\!] = \mathbf{let}\;() = \mathcal{Z}[\![\Delta]\!] \; \mathbf{in} \; \mathcal{Z}[\![x : \tau]\!]$$

$$\mathcal{Z}[\![x:\tau_1 \multimap \tau']\!] = x.2$$

$$\mathcal{Z}[\![x:\tau_1 \& \cdots \& \tau_n]\!] = x.2$$

$$\mathcal{Z}[\![x:\mu^t \alpha.\tau]\!] = \mathbf{discard}_{\mu^t \alpha.\tau} x$$

$$\mathcal{Z}[\![x:\tau_1 \otimes \cdots \otimes \tau_n]\!] = \mathbf{let} (x_1,\ldots,x_n) = x \text{ in } \mathcal{Z}[\![x_1:\tau_1,\ldots,x_n:\tau_n]\!]$$

$$\mathcal{Z}[\![x:c_1\tau_1 \oplus \cdots \oplus c_n\tau_n]\!] = \mathbf{case} x \text{ of } c_1x_1 \to \mathcal{Z}[\![x_1:\tau_1]\!] | \cdots | c_nx_n \to \mathcal{Z}[\![x_n:\tau_n]\!]$$

The **discard** $_{\mu^t\alpha.\tau}$ function used above is defined globally for each recursive type $\mu^t\alpha.\tau$. It calls itself wherever α appears in τ .

define discard
$$_{\mu^t\alpha.\tau} = \lambda f : (\mu^t\alpha.\tau)$$
. unfold $u = f$ in $\mathbb{Z}[\![u : \tau\{\alpha := \mu^t\alpha.\tau\}]\!]$

B.3 Local Recursion

Another mild extension to the language is to allow nonlinear (recursive) bindings locally, not just at the top level of a program. The body of these bindings cannot contain linear variables free. This extension is compatible with the previous ones.

Syntax

$$e := \mathbf{fix} \ x. \ e$$

Typing

$$\frac{\Gamma, x : \tau; \cdot \vdash e : \tau}{\Gamma : \cdot \vdash \mathbf{fix} \ x. \ e : \tau}$$

Reduction

$$\gamma \vdash \mathbf{fix} \ x. \ e \Longrightarrow \{(1, e\{x := e\})\}$$
 (29)

Denotation

$$\llbracket \mathbf{fix} \ x. \ e \rrbracket_{\eta}(\emptyset, v) = \mathrm{lfp}(\lambda f. \lambda v. \llbracket e \rrbracket_{\eta, x = f}(\emptyset, v))(v) \tag{30}$$

Unlike the previous two extensions, this extension has not been implemented; but it can be implemented simply by lambda lifting.

C TRANSLATING TO FACTOR GRAPH GRAMMARS

Factor graph grammars (FGGs) [Chiang and Riley 2020] are hyperedge replacement grammars (HRGs) [Bauderon and Courcelle 1987; Drewes et al. 1997; Habel and Kreowski 1987] for generating sets of factor graphs. They have recently been proposed as a unified formalism that can describe both repeated and alternative substructures [Chiang and Riley 2020]. Moreover, inference can be performed on an FGG without enumerating all the factor graphs generated.

C.1 Factor Graph Grammars

In this section, we give a condensed definition of FGGs [Chiang and Riley 2020]. Fix a finite set L^V of *node labels*, a finite set L^E of *edge labels*, and a function $type: L^E \to (L^V)^*$, which says, for each edge label, how many nodes an edge must be attached to, and what their labels must be. For any function $f: A \to B$, we extend f to strings over $A: f^*(a_1 \cdots a_n) = f(a_1) \cdots f(a_n)$.

Definition C.1. A hypergraph (or graph for short) is a tuple (V, E, att, vlab, elab, ext), where

- *V* names a finite set of *nodes*.
- *E* names a finite set of *edges*.
- $att: E \to V^*$ assigns to each edge e a sequence of zero or more attachment nodes.
- $vlab: V \to L^V$ and $elab: E \to L^E$ assign to each node and edge a label such that, for all e, $type(elab(e)) = vlab^*(att(e))$.

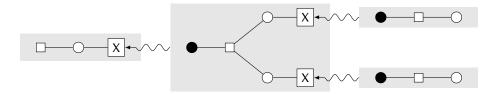
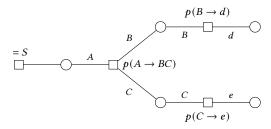


Fig. 13. An example derivation of the FGG in Example C.4. A wavy arrow means that the nonterminal at the head of the arrow is rewritten with the graph at the tail.

• $ext \in V^*$ is a sequence of zero or more *external* nodes.

A graph can be interpreted as a factor graph [Kschischang et al. 2001], as follows. Fix a mapping dom from L^V to sets; then each node v corresponds to a random variable over dom(vlab(v)). Extend dom to strings by defining $dom^*(\ell_1 \cdots \ell_n) = dom(\ell_1) \times \cdots \times dom(\ell_n)$. For each edge label $\ell \in L^E$, define a corresponding factor $fac_{\ell} : dom^*(type(\ell)) \to \mathbb{R}_{>0}$.

Example C.2. Below is a factor graph for trees of a certain shape generated by a PCFG with start symbol *S*. The variables range over nonterminal symbols or terminal symbols.



We draw an edge as a square with lines (called *tentacles*) to its attachment nodes. If there is more than one tentacle, we label each with a name. We write the edge's factor function next to it, possibly in terms of the tentacle names. A Boolean expression has value 1 if true and 0 if false.

Definition C.3. A hyperedge replacement graph grammar (HRG) is a tuple (N, T, P, S), where

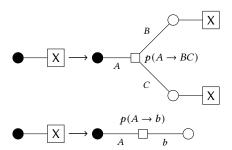
- $N \subseteq L^E$ is a finite set of *nonterminal symbols*.
- $T \subseteq L^E$ is a finite set of *terminal symbols*, disjoint from N.
- P is a finite set of *rules* of the form $(X \to R)$, where $X \in N$, and R is a hypergraph with node labels vlab, edge labels in $N \cup T$, and external nodes ext such that $type(X) = vlab^*(ext)$.
- $S \in N$ is a distinguished start nonterminal symbol.

An HRG generates a set of graphs; we omit a formal definition here. Most definitions require $type(S) = \epsilon$, including our original definition [Chiang and Riley 2020], but here we relax this requirement, to allow the graphs generated by an HRG to have external nodes.

Under *dom* and *fac*, an HRG can be interpreted as a way of generating factor graphs. We call such an HRG a *factor graph grammar* (FGG).

Example C.4. Below is an FGG for derivations of a PCFG in Chomsky normal form. The start symbol of the FGG is S'.

$$\begin{array}{c|c}
S' \longrightarrow & \square & X
\end{array}$$



This FGG generates an infinite number of factor graphs, including the one in Example C.2. The generation process can be pictured as in Figure 13.

We draw external nodes in black. If there is more than one, we write a name inside each. We draw an edge with nonterminal label *X* as a square with *X* inside and tentacles labeled with names.

Although the left-hand side is formally just a nonterminal symbol, we draw it like an edge, with replicas of the external nodes as attachment nodes.

The graphs generated by an FGG can be viewed as factor graphs, each of which defines a distribution over assignments. Recall that we allow the graphs generated by an FGG to have external nodes; for present purposes, we are interested in the marginal distribution over assignments to those external nodes.

Definition C.5. If V is a set of nodes, an assignment of V is a mapping ξ that takes each node $v \in V$ to a value in dom(vlab(v)). If $\vec{\ell}$ is a string of node labels, an assignment $\bar{\xi}$ of $\vec{\ell}$ is a sequence in $dom^*(\vec{\ell})$.

Let G = (N, T, P, S) be an FGG whose variables have finite domains. Write down the following system of equations: For each nonterminal $X \in N$ and assignment $\bar{\xi}$ of type(X),

$$w_X(\bar{\xi}) = \sum_{\substack{R\\(X \to R) \in P}} w_R(\bar{\xi}). \tag{31}$$

For each right-hand side R = (V, E, att, vlab, elab, ext) and assignment $\bar{\xi}$ of vlab(ext),

$$w_{R}(\bar{\xi}) = \sum_{\substack{\text{assignments } \xi \text{ of } V \text{ } e \in E \\ \xi^{*}(ext) = \bar{\xi}}} \prod_{e \in E} w_{elab(e)}(\xi^{*}(att(e))). \tag{32}$$

For each terminal $\ell \in T$ and assignment $\bar{\xi}$ of $type(\ell)$.

$$w_{\ell}(\bar{\xi}) = fac_{\ell}(\bar{\xi}). \tag{33}$$

Take the least fixed point of these equations (under the ordering $w \le w'$ iff $w_X(\bar{\xi}) \le w_X'(\bar{\xi})$ for all $X, \bar{\xi}$). Then w_X is the distribution over the external nodes of the graphs generated starting from X, and w_S is the distribution over the external nodes of the graphs generated by G. This is a monotone system of polynomial equations, which can be solved as in Section 4.2.

C.2 Translation Rules

Any PERPL program can be converted into an FGG that preserves semantics in the sense that the denotation of a subexpression e is equal to the weight function w_e of its translation (31).

Each judgement Γ ; $\Delta \vdash e : \tau$ in the typing derivation of a program translates to an FGG rule. The external nodes are:

- For each local binding $\tilde{x}: \tilde{\tau} \in \Delta$, there is an external node named $\tilde{x}: \tilde{\tau}$, which holds the value of \tilde{x} when e is evaluated.
- There is an external node named $v : \tau$, which holds the value of e.

We use plate notation [Buntine 1994; Koller and Friedman 2009] to depict multiple nodes: a rounded rectangle indicates that its contents are repeated, for each local variable \tilde{x} . We stress, however, that plates are only meta-notation; as will be clear in Example C.6 below, actual FGG rules do *not* use plates.

Programs. If the program has a global definition **define** x = e, add the following rule for occurrences of x:

$$x \longrightarrow e \stackrel{\mathsf{v}}{\longrightarrow} \mathbf{0} \tag{34a}$$

If the program is of the form \dots ; e, then the nonterminal e is the start symbol.

Variables. The rule for a local variable copies the node for that variable from the environment to the value:

Functions. A function becomes a node ranging over input-output pairs.

The rule for applications unpacks the function as an input-output pair, equates the input to the operand, and takes the output.

$$\begin{array}{c|c}
 & \tilde{x} & e_0 & v_0 \\
\hline
\tilde{x} & e_0 & e_1 & v_0
\end{array}$$

$$\begin{array}{c|c}
 & v_0 & v_0 \\
\hline
\tilde{x} & e_1 & v_0
\end{array}$$
(34d)

Random variables. The rules for amb and factor are very simple:

And there are no rules for fail, so that nonterminal fail has weight 0, as desired.

Unions. For a union type $c_1 \tau_1 \oplus \cdots \oplus c_n \tau_n$, the injections are straightforward to translate. For any constructor c_i :

$$\begin{array}{c|c}
\hline
\widetilde{\mathbf{x}} & c_i e \\
\hline
\widetilde{\mathbf{v}} & c_i e \\
\hline
\end{array}$$

$$\begin{array}{c|c}
v = c_i v_i \\
\hline
v \\
\hline
\end{array}$$

$$\begin{array}{c|c}
v = c_i v_i \\
\hline
\end{array}$$

A single **case**-expression translates to one rule for each constructor:

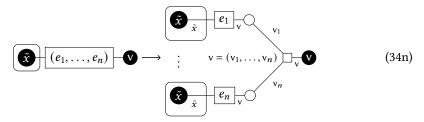
Since the desugaring of **if** leads to some superfluous nodes, we provide a direct translation as well:

This translation avoids a problem that van de Meent et al. [2018, Section 3.1] encounter when translating conditional expressions to factor graphs. In a factor graph, both arms of the conditional must be translated, and every assignment to the factor graph must assign values to variables in both arms, even though only one can be active at a time. Their translation requires some complicated machinery to work around this problem, but our translation to an FGG can simply generate a different graph for each arm.

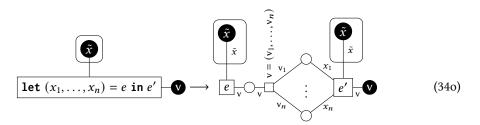
Tuples. A single additive tuple translates into multiple rules, each generating just one member. The elimination form *e.i* selects the factor graphs where the *i*-th member was generated.

$$\begin{array}{c|c}
\widehat{x} & \langle e_1, \dots, e_n \rangle & & & & \\
\widehat{x} & & & & \\
\widehat{x} & & & & \\
\hline
\widehat{x} & & & & \\
\hline$$

Multiplicative tuples are more straightforward:



Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.



Example C.6. Recall the program of Example 2.3. This program translates to the FGG in Figure 14. When we compute the sum-product of the FGG, we get min $\left(1, \frac{1-p}{p}\right)$ as desired.

C.3 Correctness

THEOREM C.7. Let $p = \text{define } x_1 = e_1 ; \dots ; \text{define } x_n = e_n ; e_0 \text{ be a program, and let } G \text{ be the } FGG \text{ translated from } p \text{ as defined above. Then } \eta \text{ as defined by } (8) \text{ exists iff } w \text{ as defined by } (31) \text{ exists, and if both exist, then for all } e \text{ in } p,$

$$w_e(\delta, \mathbf{v} = v) = [\![e]\!]_{\eta}(\delta, v). \tag{35}$$

In particular, $w_{e_0}(v = v) = [\![p]\!](v)$.

PROOF. We first show that (8) has a fixed point η iff (31) has a solution w such that (35) is satisfied; then we will show that (35) preserves ordering.

If η is a fixed point of (8), then set w according to (35). The proof that w is a solution to equation (31) is by induction on the typing derivation of e; we show just a few cases:

If $e = x_i$ is a global variable:

$$w_{x_i}(v = v) \stackrel{(35)}{=} [\![x_i]\!]_n(\emptyset, v) \stackrel{(5a)}{=} \eta(x_i)(v) \stackrel{(8)}{=} [\![e_i]\!]_n(\emptyset, v) \stackrel{(35)}{=} w_{e_i}(v = v).$$

which matches equation (31) for rule (34a). If e = x is a local variable:

$$w_x(x = v, v = v') \stackrel{\text{(35)}}{=} [\![x]\!]_n (\{(x, v)\}, v') \stackrel{\text{(5a)}}{=} [\![v = v']\!].$$

which matches equation (31) for rule (34b).

The most complex case is probably that of **case** expressions. For $\delta \in \llbracket \Delta \rrbracket$ and $\delta' \in \llbracket \Delta' \rrbracket$:

$$\begin{split} w_{e}(\delta \cup \delta', \mathbf{v} = v) & \stackrel{\text{(35)}}{=} \ \llbracket e \rrbracket \left(\delta \cup \delta', v \right) \\ & \stackrel{\text{(51)}}{=} \sum_{i=1}^{n} \sum_{v_{i}} \llbracket e_{0} \rrbracket \left(\delta, (i, v_{i}) \right) \cdot \llbracket e_{i} \rrbracket \left(\delta' \cup \{(x_{i}, v_{i})\}, v \right) \\ & \stackrel{\text{(35)}}{=} \sum_{i=1}^{n} \sum_{v_{i}} w_{e_{0}}(\delta, \mathbf{v} = (i, v_{i})) \cdot w_{e_{i}}(x_{i} = v_{i}, \delta', \mathbf{v} = v) \\ & = \sum_{i=1}^{n} \sum_{v_{0}} \sum_{v_{i}} w_{e_{0}}(\delta, \mathbf{v} = v_{0}) \cdot \mathbb{I}[v_{0} = (i, v_{i})] \cdot w_{e_{i}}(x_{i} = v_{i}, \delta', \mathbf{v} = v) \end{split}$$

which matches equation (31) for rules (34i).

Conversely, if w is a solution to (31), for all $x_i \in \Gamma$, set $\eta(x_i)(v) = w_{x_i}(v = v)$ by (35). We want to show that η is a fixed point of (8), which is very similar to the above.

Finally, consider two pairs of solutions related by (35):

$$w_e(\delta, \mathbf{v} = v) = [\![e]\!]_{\eta}(\delta, v)$$

$$w'_e(\delta, \mathbf{v} = v) = [\![e]\!]_{\eta'}(\delta, v).$$

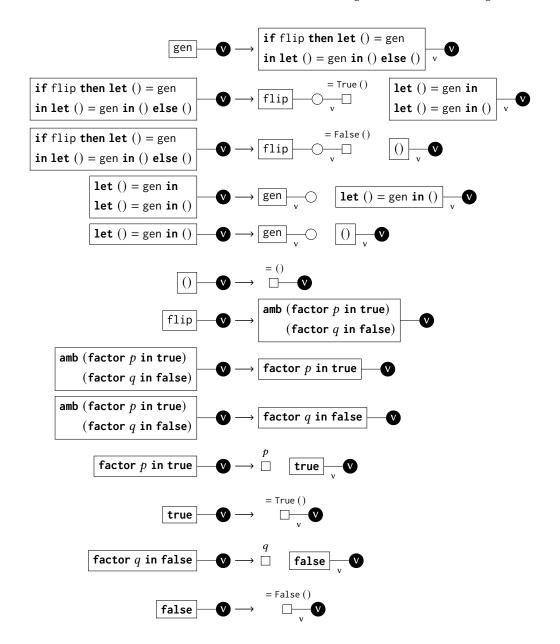


Fig. 14. FGG translated from the program in Example 2.3. The start symbol is gen.

If $w \le w'$, then it's easy to show that $\eta \le \eta'$. For each $x_i \in \Gamma$,

$$\eta(x_i)(v) = w_{x_i}(v = v) \le w'_{x_i}(v = v) = \eta'(x_i)(v).$$

Conversely, if $\eta \leq \eta'$, then for any x_i , $w_{x_i} \leq w'_{x_i}$. Since the equations defining w_e in terms of the w_{x_i} are polynomials with nonnegative coefficients, they are monotonic in the w_{x_i} . So for any e, $w_e \leq w'_e$.

Proc. ACM Program. Lang., Vol. 7, No. OOPSLA1, Article 98. Publication date: April 2023.

D DETAILS ABOUT EMBEDDED PUSHDOWN AUTOMATA

D.1 Definition and Normal Form

Definition D.1. An embedded pushdown automaton is a tuple $P = (Q, \Sigma, \Gamma, q_0, S, \delta)$, where

- *Q* is a finite set of states
- Σ is a finite input alphabet
- Γ is a finite stack alphabet
- $q_0 \in Q$ is the initial state
- $S \in \Gamma$ is the initial stack symbol
- δ is a set of transitions $q, A \xrightarrow{a} r, \beta$ where $q, r \in Q, A \in \Gamma, a \in \Sigma \cup \{\epsilon\}$, and $\beta \in (\Gamma^* \bot)^* (\Gamma^* \cdots) (\Gamma^* \bot)^*$, where \cdots stands for the rest of a stack, and \bot for the bottom of a stack.

A configuration of *P* is a tuple (q, α, w) where $\alpha \in (\Gamma^* \bot)^*$ and $w \in \Sigma^*$. If $(q, A \cdots \xrightarrow{a} r, \beta) \in \delta$, then we define

$$(q, A\alpha \perp \gamma, aw) \Rightarrow_P (r, \beta\{\cdots := \alpha \perp\} \gamma, w) \qquad \qquad \perp \notin \alpha$$
$$(q, \perp \alpha, w) \Rightarrow_P (q, \alpha, w)$$

We say that P accepts w (by empty stack) if $(q_0, S\perp, w) \Rightarrow_P^* (q, \epsilon, \epsilon)$ for any q.

PROPOSITION D.2. Every EPDA is equivalent to an EPDA whose transitions all have one of the following forms:

$$\begin{split} q, X & \cdots \xrightarrow{a} q', \cdots \\ q, X & \cdots \xrightarrow{\epsilon} q', YZ & \cdots \\ q, X & \cdots \xrightarrow{\epsilon} q', Y \bot Z & \cdots \\ q, X & \cdots \xrightarrow{\epsilon} q', Y & \cdots Z \bot \end{split}$$

PROPOSITION D.3. Every EPDA is equivalent to an EPDA with only one state.

PROOF. Create a new initial stack symbol S'. For each $X \in \Gamma$, create stack symbols X_{qrst} for all $q, r, s, t \in Q$. A stack symbol X_{qrst} on top of the top stack simulates being in state q. Then

$$\begin{array}{lll} \text{for each} & \text{create} \\ q, X \overset{a}{\longrightarrow} q', \cdots & X_{qq'rr} \overset{a}{\longrightarrow} \cdots \\ q, X \overset{a}{\longrightarrow} q', YZ \cdots & X_{qstr} \overset{a}{\longrightarrow} Y_{q'uvr}Z_{ustv} \cdots \\ q, X \overset{a}{\longrightarrow} q', Y\bot Z \cdots & X_{qtus} \overset{a}{\longrightarrow} Y_{q'rrr} \bot Z_{rtus} \cdots \\ q, X \overset{a}{\longrightarrow} q', Y \cdots Z\bot & X_{atus} \overset{a}{\longrightarrow} Y_{q'tur} \cdots Z_{rsss} \bot \end{array}$$

Finally, create transitions

$$S' \stackrel{\epsilon}{\longrightarrow} S_{q_0rrr} \cdots$$

The transitions maintain the following invariants:

- If X_{qrst} is immediately above $Y_{q'r's't'}$, then r = q' and s = t'.
- If X_{qrst} is the top symbol of a stack immediately above $Y_{q'r's't'}$, then t = q'.
- If X_{qrst} is the bottom symbol of a stack, then r = t.

When the pop transition pops a symbol $X_{qq'rr}$, if it is not at the bottom, the invariants ensure that the new top symbol is of the form $Y_{q'r's't'}$, simulating entering state q'. If $X_{qq'rr}$ is at the bottom, the invariants ensure that q' = r and therefore that the new top symbol is of the form $Y_{q'r's't'}$, again simulating entering state q'.

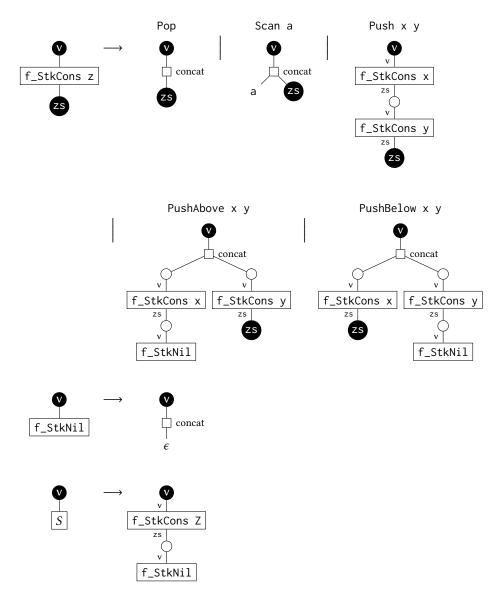


Fig. 15. Sketch of FGG compiled from final program of Section 6.3.

D.2 Relationship to TAG

Here we explain in more detail how the transformed EPDA program in Section 6.3 is related to a tree-adjoining grammar (TAG), and how inference on it amounts to TAG parsing. When the program is converted to an FGG (Appendix C), and smaller rules are fused into larger rules, the result looks like Figure 15. Each possible Action that transition z might take is converted to an FGG rule with f_StkCons z on the left. The external nodes and a have type StackF, which can be thought of as a pair of string positions, namely the starting position (usually called ws) and ending position (usually called zs*) of the substring that is scanned while the Stack is consumed.

For clarity, we've taken the liberty of drawing some strings directly into the FGG and writing factors labeled "concat" in place of more complicated logic that concatenates one or two substrings.

When arranged in this way, this FGG looks just like a TAG. A TAG could be defined as an HRG (Definition C.3) where all rule right-hand sides are ordered trees, which come in two kinds: *initial* trees have one external node, which is the root node, and *auxiliary* trees have exactly two external nodes, which are the root node and a node called the *foot* node. (TAGs with no auxiliary trees are called tree-substitution grammars and generate the same string languages as CFGs.) In Figure 15, the f_StkCons rules are the auxiliary trees, because f_StkCons z has type StackF \rightarrow StackF: the returned \bigcirc is the root node, and the argument \bigcirc is the foot node. On the other hand, the f_StkNil and S rules are the initial trees, because both f_StkNil and f_StkCons z f_StkNil have type StackF.

Any TAG can be converted to a normal form [Lang 1994] analogous to Chomsky normal form, in which every rule has at most two nonterminal symbols. In the two-nonterminal case, the two nonterminals can stand in three possible relationships: both dominating the foot node, one dominating the foot node and the other to its left, and one dominating the foot node and the other to its right. These possibilities correspond exactly to our FGG rules for Push, PushAbove, and PushBelow. There is also a rule that introduces a terminal symbol, which corresponds to our FGG rule for Scan.

E BENCHMARK PROGRAMS

The WebPPL benchmark results in Figure 7 were generated using programs like this:

```
var S = function() {
  if (flip(0.9))
    return ['A'];
    return S().concat(S());
}
var model = function() {
  return S().length;
}
display('50000000 rejection');
Infer({ model, method: 'rejection', samples: 50000000 }).getDist();
  The Dice benchmark results in Figure 8 were generated using programs like this:
fun A(i:int(3)) { if i < int(3,6) then i + int(3,1) else int(3,7) }
fun S0(i:int(3)) { int(3,7) }
fun S1(i:int(3)) { if flip 0.1 then S0(S0(i)) else A(i) }
fun S2(i:int(3)) { if flip 0.1 then S1(S1(i)) else A(i) }
fun S3(i:int(3)) { if flip 0.1 then S2(S2(i)) else A(i) }
fun S4(i:int(3)) { if flip 0.1 then S3(S3(i)) else A(i) }
fun S5(i:int(3)) { if flip 0.1 then S4(S4(i)) else A(i) }
fun S6(i:int(3)) { if flip 0.1 then S5(S5(i)) else A(i) }
S6(int(3,0)) == int(3,6)
```