

# Running Probabilistic Programs Backwards

Neil Toronto<sup>1</sup>      Jay McCarthy<sup>2</sup>      David Van Horn<sup>1</sup>  
neil.toronto@gmail.com    jay.mccarthy@gmail.com    dvanhorn@cs.umd.edu

<sup>1</sup>University of Maryland, College Park    <sup>2</sup>Vassar College

**Abstract.** Many probabilistic programming languages allow programs to be run under constraints in order to carry out Bayesian inference. Running programs under constraints *could* enable other uses such as rare event simulation and probabilistic verification—except that all such probabilistic languages are necessarily limited because they are defined or implemented in terms of an impoverished theory of probability. Measure-theoretic probability provides a more general foundation, but its generality makes finding computational content difficult.

We develop a measure-theoretic semantics for a first-order probabilistic language with recursion, which interprets programs as functions that compute preimages. Preimage functions are generally uncomputable, so we derive an approximating semantics. We implement the approximating semantics and use the implementation to carry out Bayesian inference, stochastic ray tracing (a rare event simulation), and probabilistic verification of floating-point error bounds.

**Keywords:** Probability, Semantics, Domain-Specific Languages

## 1 Introduction

One key feature usually distinguishes a probabilistic programming language from general-purpose languages: finding the probabilistic conditions under which stated constraints are satisfied. Often, a probabilistic program simulates a real-world random process and the constraints represent observed, real-world outcomes. Running the program under the constraints *infers causes from effects*.

Inferring probabilistic causes from observed outcomes is called **Bayesian inference**, a technique used widely in artificial intelligence. It has been successful in analyzing phenomena at all scales, from genomes to celestial bodies. Automating it is one of the primary drivers of probabilistic language development.

One of the simplest probabilistic programs that allows us to demonstrate Bayesian inference simulates the following process of flipping two coins.

1. Flip a fair coin; call the outcome  $x$ .
2. If  $x$  is heads, flip another fair coin. If  $x$  is tails, flip an unfair coin with heads probability 0.3 (tails probability 0.7). In either case, call the outcome  $y$ .

The following probabilistic program simulates this process.

```
let x := flip 0.5
    y := flip (if x = heads then 0.5 else 0.3)
in ⟨x, y⟩
```

(1)

Here, `flip q` returns `heads` with probability  $q$  and `tails` with probability  $1 - q$ .

The meaning of (1) is not the returned random value, but a **probability distribution** that describes the likelihoods of all possible returned random values. For discrete processes, this distribution can always be defined by a **probability mass function**: a mapping from possible values to their probabilities. These probabilities are computed by multiplying the probabilities of intermediate random values. For example, the probability of  $\langle \text{heads}, \text{heads} \rangle$  is  $0.5 \cdot 0.5 = 0.25$ , and the probability of  $\langle \text{tails}, \text{heads} \rangle$  (i.e. the second flip is unfair) is  $0.5 \cdot 0.3 = 0.15$ . The meaning of (1) is thus the probability mass function

$$\mathbf{p} := [\langle \text{heads}, \text{heads} \rangle \mapsto 0.25, \langle \text{heads}, \text{tails} \rangle \mapsto 0.25, \\ \langle \text{tails}, \text{heads} \rangle \mapsto 0.15, \langle \text{tails}, \text{tails} \rangle \mapsto 0.35] \quad (2)$$

Using  $\mathbf{p}$ , we can answer any question about the process under constraints. For example, if we do not know  $x$ , but constrain  $y$  to be `heads`, what is the probability that  $x$  is also `heads`? We compute the answer by dividing the probability of the outcome we are interested in (i.e.  $\langle x, y \rangle = \langle \text{heads}, \text{heads} \rangle$ ) by the total probability of outcomes in the constraint's corresponding subdomain  $\{\text{heads}, \text{tails}\} \times \{\text{heads}\}$ :

$$\frac{\mathbf{p} \langle \text{heads}, \text{heads} \rangle}{\sum_{z \in \{\text{heads}, \text{tails}\} \times \{\text{heads}\}} \mathbf{p} z} = \frac{0.25}{0.25 + 0.15} = 0.625 \quad (3)$$

Qualitatively,  $y$  being `heads` is a bit unusual if the second coin is unfair. Therefore, we infer that the second coin is most probably fair; i.e.  $x$  is most likely `heads`.

The time complexity of computing  $\mathbf{p}$  is generally exponential in the number of random choices, which is intractable for all but the simplest processes. One popular way to avoid this exponential explosion is to use advanced Monte Carlo algorithms to sample according to  $\mathbf{p}$  on the constraint's corresponding subdomain without explicitly enumerating that subdomain. The number of samples required is typically quadratic in the answer's desired accuracy [6, Sec. 12.2].

Probabilistic languages that are implemented using advanced Monte Carlo algorithms could be used not just for Bayesian inference, but for simulating **rare events** (i.e. very low-probability events) by encoding the events as constraints.

Stochastic ray tracing [29] is one such rare-event simulation task. As illustrated in Fig. 1, to carry out stochastic ray tracing, a probabilistic program simulates a light source emitting a single photon in a random direction, which is reflected or absorbed when it hits a wall. The program outputs the photon's path, which is constrained to pass through an aperture. Millions of paths that meet the constraint are sampled, then projected onto a simulated sensor array.

The program's main loop is a recursive function with two arguments: `path`, the photon's path so far as a list of points, and `dir`, the photon's current direction.

$$\begin{aligned} \text{simulate-photon path dir} &:= \\ \text{case } (\text{find-hit } (\text{fst path}) \text{ dir}) \text{ of} & \\ \text{absorb pt} &\longrightarrow \langle \text{pt}, \text{path} \rangle \\ \text{reflect pt norm} &\longrightarrow \text{simulate-photon } \langle \text{pt}, \text{path} \rangle (\text{random-half-dir norm}) \end{aligned} \quad (4)$$

Here, `find-hit (fst path) dir` finds the surface the photon hits. If the photon is absorbed, `find-hit` returns a data structure containing just the collision point

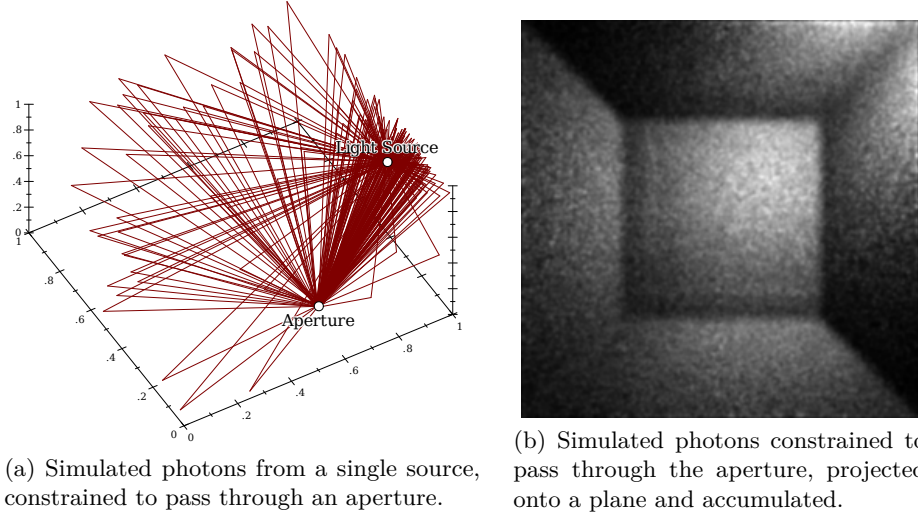


Fig. 1: Ray tracing by constraining the outputs of a probabilistic program.

`pt`. Otherwise, `find-hit` returns a data structure containing the collision point `pt` and surface normal `norm`, which `random-half-dir` uses to choose a new direction. Running `simulate-photon`  $\langle \text{pt}, \langle \rangle \rangle$  `dir`, where `pt` is the light source's location and `dir` is a random emission direction, generates a photon path. The head of the path (the last collision point) is constrained to be in the aperture. The remainder of the program is simple vector math that computes ray-plane intersections.

In contrast, hand-coded stochastic ray tracers, written in general-purpose languages, are much more complex and divorced from the physical processes they simulate, because they must interleave the advanced Monte Carlo algorithms that ensure the aperture constraint is met.

Unfortunately, while many probabilistic programming languages support random real numbers, none are capable of running a probabilistic program like (4) under constraints to carry out stochastic ray tracing. The reason is not lack of engineering or weak algorithms, but is theoretical at its core: they are all either defined or implemented using a naive theory of probability.

While probability mass functions cannot define distributions on  $\mathbb{R}$  that give positive probability to uncountably many values, there is a near-universal substitute that can: probability *density* functions. Density functions map single values to probability-like quantities, which makes them intuitively appealing and apparently simple. Unfortunately, density functions are not general enough to be used as probabilistic program meanings without imposing severe limitations on probabilistic languages. In particular, programs whose outputs are deterministic functions of random values and programs with recursion generally cannot denote density functions. The program in (4) exhibits both characteristics.

Measure-theoretic probability is a more powerful alternative to this naive probability theory based on probability mass and density functions. It not only

subsumes naive probability theory, but is capable of defining any computable probability distribution, and many uncomputable distributions. But while even the earliest work [14] on probabilistic languages is measure-theoretic, the theory’s generality has historically made finding useful computational content difficult.

We show that measure-theoretic probability can be made computational by

1. Using measure-theoretic probability to define a compositional, denotational semantics that gives a valid denotation to every program.
2. Deriving an approximating semantics, which allows computing answers to questions about probabilistic programs to arbitrary accuracy.
3. Implementing the approximating semantics and efficiently solving problems.

In fact, our primary implementation, *Dr. Bayes*, produced Fig. 1b by running a probabilistic program like (4) under an aperture constraint.

The rest of this report is organized as follows.

- Section 2 demonstrates why density functions are insufficient for interpreting probabilistic programs. It shows how measure-theoretic probability defines probability distributions using set-valued inverses, or *preimage functions*.
- Section 3 presents the categorical tools we use to derive many semantics from a single standard semantics in a way that makes them easy to prove correct.
- Section 4 defines the semantics of nonrecursive, nonprobabilistic programs, which interprets programs as preimage functions.
- Section 5 lifts this semantics to recursive, probabilistic programs.
- Section 6 derives a sound, implementable approximating semantics.
- Section 7 describes our implementations and gives examples, including probabilistic verification of floating-point error bounds.

In short, we show why and how to run probabilistic programs under constraints by computing preimage functions—that is, by running programs backwards.

## 2 Background

### 2.1 Probability Density Functions

Some distributions of real values can be defined by **probability density functions**: integrable functions  $\mathbf{p} : \mathbb{R}^n \rightarrow [0, \infty)$  that integrate to 1.

The simplest nontrivial probabilistic program is `random`, which returns a uniformly random value in the interval  $[0, 1]$ . The meaning of `random` is a probability distribution that can be defined by the density

$$\mathbf{p} : \mathbb{R} \rightarrow [0, \infty) \quad \mathbf{p} \, x := \begin{cases} 1 & \text{if } x \in [0, 1] \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Though  $\mathbf{p} \, x$  for any  $x$  indicates  $x$ ’s relative frequency,  $\mathbf{p} \, x$  is not a probability. Probabilities are obtained by integration. For example, the probability that `random` returns a value in  $[0, 0.5]$  is

$$\int_0^{0.5} (\mathbf{p} \, x) \, dx = \int_0^{0.5} 1 \, dx = \left[ x \right]_0^{0.5} = 0.5 - 0 = 0.5 \quad (6)$$

Similarly, the probability of  $[0.5, 0.5]$  or any other singleton set is zero. In fact, *every* probability density function integrates to zero on singleton sets.

This fact makes it trivial to write a probabilistic program whose distribution cannot be defined by a density. For example, consider

$$\text{max } \langle 0.5, \text{random} \rangle \quad (7)$$

where  $\text{max } \langle a, b \rangle$  returns the greater of the pair  $\langle a, b \rangle$ . This program evaluates to 0.5 whenever  $\text{random}$  returns a number in  $[0, 0.5]$ . In other words, the value of  $\text{max } \langle 0.5, \text{random} \rangle$  is in  $[0.5, 0.5]$  with probability 0.5. But if its distribution is defined by a density, then  $[0.5, 0.5]$  must have probability zero—not 0.5.

A probabilistic language without the  $\text{max}$  function can still be useful. It is fairly easy to compute densities for the outputs of single-argument functions that happen to have differentiable inverses, such as exponentiation and square root. But two-argument functions such as addition and multiplication require evaluating integrals, which generally do not have closed-form solutions.

Perhaps the most constricting limitation of probability density functions is that the number of dimensions must be finite and fixed. This limitation rules out recursive data types, and makes recursion so difficult that few probabilistic languages attempt to allow it.

## 2.2 Measures, and Measures of Preimages

Measure-theoretic probability gains its expressive power by mapping sets directly to probabilities. Functions that do so are called **probability measures**. For example, the distribution of  $\text{random}$  is defined by the probability measure

$$P : \mathcal{P} [0, 1] \rightarrow [0, 1] \quad P [a, b] = b - a \quad (8)$$

where  $\mathcal{P} [0, 1]$  is the powerset of  $[0, 1]$  and ‘ $\rightarrow$ ’ denotes a partial mapping. Though (8) apparently defines  $P$  only on intervals, it is regarded as defining  $P$  additionally on countable unions of intervals, their complements, countable unions of such, and so on. The resulting domain includes almost every subset of  $[0, 1]$  that can be written down.

Probability measures can be defined on any domain, including domains with variable and infinite dimension. They can also map singleton sets to nonzero probabilities, which we will demonstrate shortly by deriving a probability measure for  $\text{max } \langle 0.5, \text{random} \rangle$ .

Measure-theoretic probability takes great pains to separate random effects from the pure logic of mathematics. It does so in the same way Haskell and other purely functional programming languages allow random effects: by interpreting probabilistic processes as *deterministic functions* that operate on an assumed-random source. The probabilities of sets of outputs are uniquely determined by the probabilities of the corresponding sets of inputs.

Suppose we interpret  $\text{max } \langle 0.5, \text{random} \rangle$  as the deterministic function

$$f := \lambda r \in [0, 1]. \text{max } \langle 0.5, r \rangle \quad (9)$$

and assume that  $r$  is its uniform random source; i.e. that its distribution is  $P$  as defined in (8). To compute the probability that  $\max \langle 0.5, \text{random} \rangle$  evaluates to 0.5, we apply  $P$  to the set of all  $r$  for which  $f \ r \in [0.5, 0.5]$ , and get, as expected,

$$P \{r \in [0, 1] \mid f \ r \in [0.5, 0.5]\} = P [0, 0.5] = 0.5 - 0 = 0.5 \quad (10)$$

For any  $f$  and  $B$ , the set  $\{a \in \text{domain } f \mid f \ a \in B\}$  is called the **preimage of  $B$  under  $f$** . Functions that compute preimages are often denoted  $f^{-1}$  to emphasize that they are a sort of generalized inverse function. However, we find this notation confusing: inverse functions operate on *values* and may not be well-defined, whereas preimage functions operate on *sets* and are *always* well-defined.<sup>1</sup> Thus, we denote  $f$ 's preimage function by **preimage  $f$** . The probability that  $f$  outputs a value in  $B$  is therefore  $P ((\text{preimage } f) \ B)$ , or  $P (\text{preimage } f \ B)$ .

Though the distribution of  $\max \langle 0.5, \text{random} \rangle$ , or the output of  $f$ , has no probability density function, its probability measure is defined by

$$P_f : \mathcal{P} [0.5, 1] \rightarrow [0, 1] \quad P_f [a, b] = P (\text{preimage } f [a, b]) \quad (11)$$

An equivalent, more elegant definition is

$$P_f := P \circ (\text{preimage } f) \quad (12)$$

which clearly shows that  $P_f$  is factored into a part  $P$  that quantifies randomness, and a deterministic part **preimage  $f$**  that *runs  $f$  backwards on sets of outputs*.

This factorization confers the flexibility to interpret probabilistic programs by choosing any  $P$  and  $f$  for which  $P \circ (\text{preimage } f)$  is the correct measure. For  $P$ , we choose uniform measures on cartesian products of  $[0, 1]$  (e.g.  $[0, 1]^{\mathbb{N}}$ ) and interpret each **random** as a projection. Thus, for the remainder of this paper, we can concentrate solely on computing **preimage  $f$** .

Because **preimage  $f$**  is deterministic, techniques to compute it have applications outside of probabilistic programming; for example, constraint-functional languages, type inference, and verification. More immediately, its determinism means that, for the bulk of this paper, *readers do not need to know anything about probability, let alone measure theory*—only basic set theory.

### 2.3 Preimage Semantics

Several well-known identities suggest that preimages can be computed compositionally, which would make it possible to define a denotational semantics that interprets programs as preimage functions. For example, we have

$$\begin{aligned} \text{preimage id} &= \text{id} \\ \text{preimage } (f_2 \circ f_1) &= (\text{preimage } f_1) \circ (\text{preimage } f_2) \\ \text{preimage } \langle f_1, f_2 \rangle (B_1 \times B_2) &= (\text{preimage } f_1 \ B_1) \cap (\text{preimage } f_2 \ B_2) \end{aligned} \quad (13)$$

where  $\langle f_1, f_2 \rangle = \lambda a \in (\text{domain } f_1) \cap (\text{domain } f_2). \langle f_1 \ a, f_2 \ a \rangle$  constructs pairing functions and  $\text{id}$  is the identity function.

It might seem we can easily use identities like those in (13) directly to define a semantic function  $\llbracket \cdot \rrbracket_{\text{pre}}$  that interprets programs as preimage functions. Unfortunately, our task is not that simple, for the following reasons.

<sup>1</sup> If  $f^{-1} \ b$  is undefined, then the preimage of  $\{b\}$  under  $f$  is simply  $\emptyset$ .

1. The **preimage** function requires its argument to have an observable domain. This includes **extensional** functions, which are sets of input/output pairs (i.e. possibly infinite hash tables), but not **intensional** functions, which are syntactic rules for computing outputs from inputs (e.g. lambdas).<sup>2</sup>
2. We must ensure  $\text{preimage } f \ B$  is always in the domain of the chosen probability measure  $P$ . (Recall that probability measures are partial functions.) If this is true, we say  $f$  is **measurable**. Proving measurability is difficult, especially if  $f$  may not terminate.
3. The function  $\text{app} : (X \rightarrow Y) \times X \rightarrow Y$ , when restricted to measurable functions, is not generally measurable if we want good approximation properties [2]. This makes interpreting higher-order application difficult.

Implementing a language based on preimage semantics is complicated because

4. Ordinary set-based mathematics is unlike any implementation language.
5. It requires running programs written in a Turing-equivalent language backwards, efficiently, on possibly uncountable sets of outputs.

We address 1 and 4 by developing our semantics using  $\lambda_{\text{ZFC}}$  [28], an untyped, call-by-value  $\lambda$ -calculus with infinite sets, real numbers, extensional functions such as  $\lambda r \in [0, 1]. \max \langle 0.5, r \rangle$ , intensional functions such as  $\lambda r. \max \langle 0.5, r \rangle$ , a computable sublanguage, and an operational semantics. It is essentially ordinary mathematics extended with lambdas and general recursion, or equivalently a lambda calculus extended with uncountably infinite sets and set operations.

We have addressed difficulty 2 by proving that all programs' interpretations as functions are measurable if language primitives are measurable, including uncomputable primitives such as limits and real equality, regardless of nontermination. The proof interprets programs as extensional functions and applies well-known theorems from measure theory such as the identities in (13). Unfortunately, the required machinery does not fit in this report; see the first author's dissertation [27] for the entire development.

We avoid difficulty 3 for now by interpreting a language with *first-order* functions and recursion. We address 5 by deriving and implementing a *conservative approximation* of the preimage semantics, and using its approximations to compute measures of preimages with arbitrary accuracy.

## 2.4 Nonstandard Abstract Interpretation, Categorically

We interpret nonrecursive, nonprobabilistic programs three different ways, using

1. A **standard semantics**  $\llbracket \cdot \rrbracket_{\perp}$  that interprets programs that may raise errors (e.g. divide-by-zero) as functions.
2. A **nonstandard semantics**  $\llbracket \cdot \rrbracket_{\text{pre}}$  that interprets programs as preimage functions, which operate on uncountable sets, and are thus unimplementable.
3. An **approximating semantics**  $\llbracket \cdot \rrbracket'_{\text{pre}}$  that interprets programs as *approximating* preimage functions, which *are* implementable.

<sup>2</sup> The lambda  $\lambda r. \max \langle 0.5, r \rangle$  is intensional, but  $\lambda r \in [0, 1]. \max \langle 0.5, r \rangle$  constructs an extensional function by pairing every  $r \in [0, 1]$  with its corresponding  $\max \langle 0.5, r \rangle$ .

Of course, we must prove for any program  $p$ , that  $\llbracket p \rrbracket_{\text{pre}}$  correctly computes preimages under  $\llbracket \cdot \rrbracket_{\perp}$ , and that  $\llbracket p \rrbracket'_{\text{pre}}$  is sound with respect to  $\llbracket p \rrbracket_{\text{pre}}$ .

For recursive, probabilistic programs, we define three more semantic functions analogous to  $\llbracket \cdot \rrbracket_{\perp}$ ,  $\llbracket \cdot \rrbracket_{\text{pre}}$  and  $\llbracket \cdot \rrbracket'_{\text{pre}}$ , that have analogous proof obligations. We also prove that they correctly interpret nonrecursive, nonprobabilistic programs.

In the full development [27], two more semantic functions interpret programs as extensional functions, which are used to prove measurability. Another semantic function collects information needed for advanced Monte Carlo algorithms. In all, we have 9 related semantic functions, each defined by 11 or 12 rules, whose correctness and relationships must be proved by structural induction. Doing so is tedious and error-prone. We need a way to parameterize one semantic function on many meanings, where each “meaning” is simpler than a semantic function and ideally has exploitable properties.

Moggi [21] introduced monads as a categorical “metalanguage” for interpreting programs. Wadler [30] showed how to use monad categories in pure functional programming to encode and hide side effects such as mutation and randomness. Haskell programmers now primarily encode programs with side effects using **do-notation**, which is transformed into any monad. Essentially, Haskell has a built-in semantic function parameterized on a monad.

Other researchers have identified arrows [9] and idioms [18] as useful kinds of categories. Different kinds of categories are good for encoding different kinds of effects, and have different levels of expressiveness [15]. Arrows are good categories for interpreting first-order languages. We therefore interpret programs 9 different ways by parameterizing a semantic function on one of 9 arrow categories.

In our formulation, an arrow category consists of a type constructor and five combinators; each is thus half as complicated as the semantic function. Their categorical properties also allow two drastic simplifications. First, they allow proving the correctness of a semantic function  $\llbracket \cdot \rrbracket_{\text{b}}$  with respect to  $\llbracket \cdot \rrbracket_{\text{a}}$  by proving a simple theorem about arrows  $\text{a}$  and  $\text{b}$ . Second, they allow us to *derive* all the arrows for recursive, probabilistic programs from the arrows for nonrecursive, nonprobabilistic programs, by defining one function and proving one theorem.

## 2.5 Types and Notation

Because some arrows carry out uncountably infinite computations, we must define their combinators in a sufficiently powerful  $\lambda$ -calculus. We use  $\lambda_{\text{ZFC}}$  [28].

Though  $\lambda_{\text{ZFC}}$  is untyped, it helps to use a manually checked, auxiliary type system. For example, the types of some of  $\lambda_{\text{ZFC}}$ ’s primitives are those of membership  $(\in) : x \rightarrow \text{Set } x \rightarrow \text{Bool}$ , powerset  $\mathcal{P} : \text{Set } x \rightarrow \text{Set } (\text{Set } x)$ , big union  $\bigcup : \text{Set } (\text{Set } x) \rightarrow \text{Set } x$ , and the map-like image  $(x \rightarrow y) \rightarrow \text{Set } x \rightarrow \text{Set } y$ . We allow sets to be used as types, as in  $\max : \langle \mathbb{R}, \mathbb{R} \rangle \rightarrow \mathbb{R}$ .

More precisely, types are characterized by these rules:

- $x \rightarrow y$  is the type of intensional, partial functions from type  $x$  to type  $y$ .
- $\langle x, y \rangle$  is the type of pairs of values with types  $x$  and  $y$ .
- $\text{Set } x$  is the type of sets whose members have type  $x$ .



- An uppercase type variable such as  $X$  represents a set used as a type.
- $X^J$  is the type of extensional, total functions from set  $J$  to set  $X$ ; i.e. vectors of  $X$  indexed by  $J$ .

Because the inhabitants of the type  $\text{Set } X$  and  $\mathcal{P} X$  (i.e. subsets of the set  $X$ ) are the same, they are equivalent types. Similarly,  $\langle X, Y \rangle$  is equivalent to  $X \times Y$ .

We use adjacency (i.e.  $f \mathbf{a}$ ) to apply both intensional and extensional functions. For example, the first element of the infinite vector  $f : [0, 1]^{\mathbb{N}}$  is  $f 0$ .

Type constructors are defined using ‘ $::=$ ’; e.g.  $X \rightsquigarrow_{\perp} Y ::= X \rightarrow (Y \cup \{\perp\})$ .

Proofs, which we elide to save space, are in the first author’s dissertation [27].

### 3 Arrows and First-Order Semantics

Arrows [9], like monads [30], thread effects through computations in a way that imposes structure. But arrow computations are always

1. Function-like. The type constructor for arrow  $\mathbf{a}$  is written  $x \rightsquigarrow_{\mathbf{a}} y$  to connote this. In fact, the *function arrow*’s type constructor is  $x \rightsquigarrow y ::= x \rightarrow y$ .
2. First-order. There is no way to derive the higher-order application combinator  $\text{app} : (x \rightsquigarrow_{\mathbf{a}} y, x) \rightsquigarrow_{\mathbf{a}} y$  from the combinators that define arrow  $\mathbf{a}$ .

The first property makes arrows a good fit for a compositional translation from expressions to pure functions that operate on random sources. The second property makes arrows a good fit for the semantics of a first-order language.

#### 3.1 Arrow Combinators and Laws

Arrows factor computation into the following tasks.

1. Referring to pure, primitive functions.
2. Applying primitive or first-order functions.
3. Binding values to local variables and creating data structures.
4. Branching based on the results of prior computations.

The first four arrow combinators correspond respectively with each of these tasks. Another allows lazy branching in a call-by-value language such as  $\lambda_{\text{ZFC}}$ .

For laziness, we need a singleton type for thunks. We use the set  $1 := \{0\}$ .

**Definition 1** (**arrow**<sup>3</sup>). *A binary type constructor  $(\rightsquigarrow_{\mathbf{a}})$  and the combinators*

$\text{arr}_{\mathbf{a}} : (x \rightarrow y) \rightarrow (x \rightsquigarrow_{\mathbf{a}} y)$	lift
$(\ggg_{\mathbf{a}}) : (x \rightsquigarrow_{\mathbf{a}} y) \rightarrow (y \rightsquigarrow_{\mathbf{a}} z) \rightarrow (x \rightsquigarrow_{\mathbf{a}} z)$	compose
$(\&\&_{\mathbf{a}}) : (x \rightsquigarrow_{\mathbf{a}} y) \rightarrow (x \rightsquigarrow_{\mathbf{a}} z) \rightarrow (x \rightsquigarrow_{\mathbf{a}} \langle y, z \rangle)$	pair
$\text{ifte}_{\mathbf{a}} : (x \rightsquigarrow_{\mathbf{a}} \text{Bool}) \rightarrow (x \rightsquigarrow_{\mathbf{a}} y) \rightarrow (x \rightsquigarrow_{\mathbf{a}} y) \rightarrow (x \rightsquigarrow_{\mathbf{a}} y)$	if-then-else
$\text{lazy}_{\mathbf{a}} : (1 \rightarrow (x \rightsquigarrow_{\mathbf{a}} y)) \rightarrow (x \rightsquigarrow_{\mathbf{a}} y)$	laziness

define an **arrow** if certain monoid, homomorphism, and other laws hold [9].

<sup>3</sup> These are actually arrows *with choice*, which are typically defined using  $\text{first}_{\mathbf{a}}$  and  $\text{left}_{\mathbf{a}}$  instead of  $(\&\&_{\mathbf{a}})$  and  $\text{ifte}_{\mathbf{a}}$ . We find  $\text{ifte}_{\mathbf{a}}$  more natural for semantics than  $\text{left}_{\mathbf{a}}$ , and  $(\&\&_{\mathbf{a}})$  better matches the pairing preimage identity in (13).

For example, the **function arrow** is defined by the type constructor  $x \rightsquigarrow y$   $:= x \rightarrow y$  and the combinators

$$\begin{aligned}
\text{arr } f &:= f \\
(f_1 \ggg f_2) a &:= f_2 (f_1 a) \\
(f_1 \&\&\& f_2) a &:= \langle f_1 a, f_2 a \rangle \\
\text{ifte } f_1 f_2 f_3 a &:= \text{if } f_1 a \text{ then } f_2 a \text{ else } f_3 a \\
\text{lazy } f a &:= f 0 a
\end{aligned} \tag{14}$$

To demonstrate compositionally interpreting probabilistic programs as arrow computations, we interpret  $\text{max } \langle 0.5, \text{random} \rangle$  as a function arrow computation  $f : [0, 1] \rightsquigarrow \mathbb{R}$ . For any random source  $r \in [0, 1]$ , the interpretation of 0.5 should return 0.5, so 0.5 means  $\lambda r. 0.5$ , or  $\text{const } 0.5$  where  $\text{const } v := \lambda \_ . v$ . Assuming  $r \in [0, 1]$  is uniformly distributed,  $\text{random}$  means  $\lambda r. r$ , or  $\text{id}$ . We use  $(\&\&\&)$  to apply each of these interpretations to the random source to create a pair, and  $(\ggg)$  to send the pair to  $\text{max}$ . Thus,  $\text{max } \langle 0.5, \text{random} \rangle$ , interpreted as a function arrow computation, is  $f := ((\text{const } 0.5) \&\&\& \text{id}) \ggg \text{max}$ .

By substituting the definitions of  $\text{const}$ ,  $\text{id}$ ,  $(\&\&\&)$  and  $(\ggg)$ , we would find that  $f$  is equivalent to  $\lambda r. \text{max } \langle 0.5, r \rangle$ , similar to the interpretation in (9).

Only the function arrow can so cavalierly use pure functions as arrow computations. In any other arrow  $a$ , pure functions must be *lifted* using  $\text{arr}_a$ , to allow the arrow to manage any state or effects. Therefore, the interpretation of  $\text{max } \langle 0.5, \text{random} \rangle$  as an arrow  $a$  computation  $f_a : [0, 1] \rightsquigarrow_a \mathbb{R}$  is

$$f_a := (\text{arr}_a (\text{const } 0.5) \&\&\&_a \text{arr}_a \text{id}) \ggg_a \text{arr}_a \text{max} \tag{15}$$

So far, we have ignored the many arrow laws, which ensure that arrows are well-behaved (e.g. effects are correctly ordered) and are useful in proofs of theorems that quantify over arrows (i.e. nothing else is known about them). Fortunately, we can prove all the laws for an arrow  $b$  by defining it in terms of an arrow  $a$  for which the laws hold, and proving two properties about the lift from  $a$  to  $b$ . The first property is that the lift from  $a$  to  $b$  is distributive.

**Definition 2 (arrow homomorphism).**  $\text{lift}_b : (x \rightsquigarrow_a y) \rightarrow (x \rightsquigarrow_b y)$  is an *arrow homomorphism* from arrow  $a$  to arrow  $b$  if these distributive laws hold:

$$\text{lift}_b (\text{arr}_a f) \equiv \text{arr}_b f \tag{16}$$

$$\text{lift}_b (f_1 \ggg_a f_2) \equiv (\text{lift}_b f_1) \ggg_b (\text{lift}_b f_2) \tag{17}$$

$$\text{lift}_b (f_1 \&\&\&_a f_2) \equiv (\text{lift}_b f_1) \&\&\&_b (\text{lift}_b f_2) \tag{18}$$

$$\text{lift}_b (\text{ifte}_a f_1 f_2 f_3) \equiv \text{ifte}_b (\text{lift}_b f_1) (\text{lift}_b f_2) (\text{lift}_b f_3) \tag{19}$$

$$\text{lift}_b (\text{lazy}_a f) \equiv \text{lazy}_b \lambda 0. \text{lift}_b (f 0) \tag{20}$$

where “ $\equiv$ ” is an arrow-specific equivalence relation.

The second property is that the lift is right-invertible (i.e. surjective).

**Theorem 1 (right-invertible homomorphism implies arrow laws).** *If  $\text{lift}_b : (x \rightsquigarrow_a y) \rightarrow (x \rightsquigarrow_b y)$  is a right-invertible homomorphism from arrow  $a$  to arrow  $b$  and the arrow laws hold for  $a$ , then the arrow laws hold for  $b$ .*

$$\begin{aligned}
p &::= f := e; \dots; e \\
e &::= \text{let } e \text{ e} \mid \text{env } n \mid \text{if } e \text{ then } e \text{ else } e \mid \langle e, e \rangle \mid f \text{ e} \mid \delta \text{ e} \mid v \\
f &::= (\text{first-order function names}) \\
\delta &::= (\text{primitive function names}) \\
v &::= \langle v, v \rangle \mid \langle \rangle \mid \text{true} \mid \text{false} \mid (\text{other first-order constants}) \\
[f := e; \dots; e_b]_a &::= f := [e]_a; \dots; [e_b]_a & \llbracket \langle e_1, e_2 \rangle \rrbracket_a &::= \llbracket e_1 \rrbracket_a \ \&\&\&\_a \ \llbracket e_2 \rrbracket_a \\
[\text{let } e \text{ e}_b]_a &::= ([e]_a \ \&\&\&\_a \ \text{arr}_a \ \text{id}) \ggg_a [e_b]_a & \llbracket f \text{ e} \rrbracket_a &::= \llbracket (e, \langle \rangle) \rrbracket_a \ggg_a f \\
[\text{env } 0]_a &::= \text{arr}_a \ \text{fst} & \llbracket \delta \text{ e} \rrbracket_a &::= \llbracket e \rrbracket_a \ggg_a \text{arr}_a \ \delta \\
[\text{env } (n+1)]_a &::= \text{arr}_a \ \text{snd} \ggg_a [\text{env } n]_a & \llbracket v \rrbracket_a &::= \text{arr}_a \ (\text{const } v) \\
[\text{if } e_c \text{ then } e_t \text{ else } e_f]_a &::= \text{ifte}_a \llbracket e_c \rrbracket_a \ (\text{laz}_a \ \lambda 0. \llbracket e_t \rrbracket_a) \ (\text{laz}_a \ \lambda 0. \llbracket e_f \rrbracket_a)
\end{aligned}$$

where  $\text{const } v := \lambda \_ . v$       subject to  $\llbracket p \rrbracket_a : \langle \rangle \rightsquigarrow_a y$  for some  $y$

$\text{id} := \lambda v . v$

### 3.2 First-Order Let-Calculus Semantics

The result of applying  $[\cdot]_a$  is a  $\lambda_{\text{ZFC}}$  program in *environment-passing style* with the environment represented by a stack. The final expression has type  $\langle \rangle \rightsquigarrow_a y$ , where  $y$  is the type of the program’s output and  $\langle \rangle$  denotes the empty stack. A `let` expression uses pairing ( $\&&&_a$ ) to push a value onto the stack and composition ( $\ggg_a$ ) to pass the resulting stack to its body. First-order functions have type  $\langle x, \langle \rangle \rangle \rightsquigarrow_a y$  where  $x$  is the argument type and  $y$  is the return type. Application passes a stack with just an  $x$  using pairing and composition.

**Definition 3 (well-defined).** An expression (or program)  $e$  is **well-defined** under arrow  $\mathbf{a}$  if  $\llbracket e \rrbracket_{\mathbf{a}} : x \rightsquigarrow_{\mathbf{a}} y$  for some  $x$  and  $y$ , and  $\llbracket e \rrbracket_{\mathbf{a}}$  terminates.

**Theorem 2 (homomorphisms distribute over expressions).** *Let  $\text{lift}_b : (x \rightsquigarrow_a y) \rightarrow (x \rightsquigarrow_b y)$  be an arrow homomorphism. For all  $e$ ,  $\llbracket e \rrbracket_b \equiv \text{lift}_b \llbracket e \rrbracket_a$ .*

$$\begin{array}{ll}
X \rightsquigarrow_{\perp} Y ::= X \rightarrow Y_{\perp} & \text{ifte}_{\perp} f_1 f_2 f_3 a := \text{case } f_1 a \text{ of} \\
& \text{true} \rightarrow f_2 a \\
& \text{false} \rightarrow f_3 a \\
& \perp \rightarrow \perp \\
\text{arr}_{\perp} f a := f a & \text{lazy}_{\perp} f a := f 0 a \\
(f_1 \ggg_{\perp} f_2) a := \text{case } f_1 a \text{ of} & \\
& \perp \rightarrow \perp \\
& b \rightarrow f_2 b \\
(f_1 \&\&\&_{\perp} f_2) a := \text{case } \langle f_1 a, f_2 a \rangle \text{ of} & \\
& \langle \perp, \_ \rangle \rightarrow \perp \\
& \langle \_, \perp \rangle \rightarrow \perp \\
& \langle b_1, b_2 \rangle \rightarrow \langle b_1, b_2 \rangle
\end{array}$$

Fig. 3: Bottom arrow definitions.

Much of our development proceeds in the following way.

1. Define a simple arrow  $\mathbf{a}$  to interpret programs using  $\llbracket \cdot \rrbracket_{\mathbf{a}}$ .
  2. Define  $\text{lift}_{\mathbf{b}} : (x \rightsquigarrow_{\mathbf{a}} y) \rightarrow (x \rightsquigarrow_{\mathbf{b}} y)$  from arrow  $\mathbf{a}$  to arrow  $\mathbf{b}$  in such a way that  $\text{lift}_{\mathbf{b}} f$  for any  $f : x \rightsquigarrow_{\mathbf{a}} y$  is obviously correct.
  3. Prove  $\text{lift}_{\mathbf{b}}$  is a homomorphism; therefore  $\llbracket e \rrbracket_{\mathbf{b}}$  is correct (Theorem 2).
  4. Prove  $\text{lift}_{\mathbf{b}}$  is right-invertible; therefore  $\mathbf{b}$  obeys the arrow laws (Theorem 1).
- In shorter terms, *if arrow  $\mathbf{b}$  is defined in terms of a right-invertible homomorphism from  $\mathbf{a}$  to  $\mathbf{b}$ , then  $\llbracket \cdot \rrbracket_{\mathbf{b}}$  is correct with respect to  $\llbracket \cdot \rrbracket_{\mathbf{a}}$ .*

## 4 The Bottom and Preimage Arrows

The following commutative diagram shows the relationships between the exact arrows we define. We will use it as a sort of roadmap.

$$\begin{array}{ccc}
X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
\eta_{\perp}^* \downarrow & & \downarrow \eta_{\text{pre}^*} \\
X \rightsquigarrow_{\perp^*} Y & \xrightarrow{\text{lift}_{\text{pre}^*}} & X \rightsquigarrow_{\text{pre}^*} Y
\end{array} \tag{21}$$

In this section, we define the top row, which are the exact arrows for interpreting nonrecursive, nonprobabilistic programs.

### 4.1 The Bottom Arrow

To use Theorem 2 to prove correct the interpretations of expressions as preimage arrow computations, we need to define the preimage arrow in terms of a simpler arrow with easily understood behavior. The function arrow (14) is an obvious candidate. However, we will need to explicitly handle nontermination as an error value, so we need a slightly more complicated arrow.

Fig. 3 defines the **bottom arrow**, which is similar to the function arrow but propagates the error value  $\perp$ . Its computations have type  $X \rightsquigarrow_{\perp} Y ::= X \rightarrow Y_{\perp}$ , where  $Y_{\perp} ::= Y \cup \{\perp\}$ .

To prove the arrow laws, we need coarse enough notion of equivalence.

$X \xrightarrow{\text{pre}} Y ::= \langle \text{Set } Y, \text{Set } Y \rightarrow \text{Set } X \rangle$	$\langle \cdot, \cdot \rangle_{\text{pre}} : (X \xrightarrow{\text{pre}} Y_1) \rightarrow (X \xrightarrow{\text{pre}} Y_2) \rightarrow (X \xrightarrow{\text{pre}} Y_1 \times Y_2)$
$\text{pre} : (X \rightsquigarrow_{\perp} Y) \rightarrow (X \xrightarrow{\text{pre}} Y)$	$\langle \langle Y'_1, p_1 \rangle, \langle Y'_2, p_2 \rangle \rangle_{\text{pre}} :=$
$\text{pre } f \ A := \langle \text{image}_{\perp} f \ A, \text{preimage}_{\perp} f \ A \rangle$	$\text{let } Y' := Y'_1 \times Y'_2$ $p := \lambda B. \bigcup_{\langle b_1, b_2 \rangle \in B} (p_1 \{b_1\}) \cap (p_2 \{b_2\})$
$\emptyset_{\text{pre}} := \langle \emptyset, \lambda B. \emptyset \rangle$	in $\langle Y', p \rangle$
$\text{ap}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) \rightarrow \text{Set } Y \rightarrow \text{Set } X$	$(\circ_{\text{pre}}) : (Y \xrightarrow{\text{pre}} Z) \rightarrow (X \xrightarrow{\text{pre}} Y) \rightarrow (X \xrightarrow{\text{pre}} Z)$
$\text{ap}_{\text{pre}} \langle Y', p \rangle \ B := p \ (B \cap Y')$	$\langle Z', p_2 \rangle \circ_{\text{pre}} h_1 := \langle Z', \lambda C. \text{ap}_{\text{pre}} h_1 \ (p_2 \ C) \rangle$
$\text{range}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) \rightarrow \text{Set } Y$	$(\cup_{\text{pre}}) : (X \xrightarrow{\text{pre}} Y) \rightarrow (X \xrightarrow{\text{pre}} Y) \rightarrow (X \xrightarrow{\text{pre}} Y)$
$\text{range}_{\text{pre}} \langle Y', p \rangle := Y'$	$h_1 \cup_{\text{pre}} h_2 := \text{let } Y' := (\text{range}_{\text{pre}} h_1) \cup (\text{range}_{\text{pre}} h_2)$ $p := \lambda B. (\text{ap}_{\text{pre}} h_1 \ B) \cup (\text{ap}_{\text{pre}} h_2 \ B)$ in $\langle Y', p \rangle$
$\text{image}_{\perp} : (X \rightsquigarrow_{\perp} Y) \rightarrow \text{Set } X \rightarrow \text{Set } Y$	$\text{preimage}_{\perp} : (X \rightsquigarrow_{\perp} Y) \rightarrow \text{Set } X \rightarrow \text{Set } Y \rightarrow \text{Set } X$
$\text{image}_{\perp} f \ A := (\text{image } f \ A) \setminus \{\perp\}$	$\text{preimage}_{\perp} f \ A \ B := \{a \in A \mid f \ a \in B\}$

Fig. 4: Preimage functions and operations.

**Definition 4 (bottom arrow equivalence).** *Two computations  $f_1 : X \rightsquigarrow_{\perp} Y$  and  $f_2 : X \rightsquigarrow_{\perp} Y$  are equivalent, or  $f_1 \equiv f_2$ , when  $f_1 \ a = f_2 \ a$  for all  $a \in X$ .*

Using bottom arrow equivalence, it is easy to show that  $(\rightsquigarrow_{\perp})$  is isomorphic to the Maybe monad's Kleisli arrow. By Theorem 1, the arrow laws hold.

## 4.2 The Preimage Function Type and Operations

Before defining the preimage arrow, we need a type of preimage functions.  $\text{Set } Y \rightarrow \text{Set } X$  would be a good candidate, except that the  $(\ggg_{\text{pre}})$  combinator will require preimage functions to have observable domains, but instances of  $\text{Set } Y \rightarrow \text{Set } X$  may be intensional functions. We therefore define

$$X \xrightarrow{\text{pre}} Y ::= \langle \text{Set } Y, \text{Set } Y \rightarrow \text{Set } X \rangle \quad (22)$$

as the type of preimage functions. Fig. 4 defines the necessary operations on them. Operations  $\langle \cdot, \cdot \rangle_{\text{pre}}$  and  $(\circ_{\text{pre}})$  return preimage functions that compute preimages under pairing and composition, and are derived from the preimage identities in (13);  $(\cup_{\text{pre}})$  computes unions and is used to define  $\text{ifte}_{\text{pre}}$ .

Fig. 4 also defines  $\text{image}_{\perp}$  and  $\text{preimage}_{\perp}$  to operate on bottom arrow computations:  $\text{image}_{\perp} f \ A$  computes  $f$ 's range (with domain  $A$ ), and  $\text{preimage}_{\perp} f \ A$  returns a function that computes preimages under  $f$  restricted to  $A$ . Together, they can be used to convert bottom arrow computations to preimage functions:

$$\begin{aligned} \text{pre} : (X \rightsquigarrow_{\perp} Y) &\rightarrow \text{Set } X \rightarrow (X \xrightarrow{\text{pre}} Y) \\ \text{pre } f \ A &:= \langle \text{image}_{\perp} f \ A, \text{preimage}_{\perp} f \ A \rangle \end{aligned} \quad (23)$$

Lastly, the  $\text{ap}_{\text{pre}}$  function in Fig. 4 applies a preimage function to a set.

Preimage arrow correctness depends on  $\text{ap}_{\text{pre}}$  and  $\text{pre}$  behaving like  $\text{preimage}_{\perp}$ .

**Theorem 3 ( $\text{ap}_{\text{pre}}$  of  $\text{pre}$  computes preimages).** *Let  $f : X \rightsquigarrow_{\perp} Y$ . For all  $A \subseteq X$  and  $B \subseteq Y$ ,  $\text{ap}_{\text{pre}} (\text{pre } f \ A) \ B = \text{preimage}_{\perp} f \ A \ B$ .*

$$\begin{array}{ll}
X \rightsquigarrow_{\text{pre}} Y ::= \text{Set } X \rightarrow (X \rightrightarrows_{\text{pre}} Y) & \text{ifte}_{\text{pre}} h_1 h_2 h_3 A ::= \\
\text{arr}_{\text{pre}} ::= \text{lift}_{\text{pre}} \circ \text{arr}_{\perp} & \text{let } h'_1 := h_1 A \\
(h_1 \ggg_{\text{pre}} h_2) A ::= \text{let } h'_1 := h_1 A & h'_2 := h_2 (\text{ap}_{\text{pre}} h'_1 \{\text{true}\}) \\
& h'_3 := h_3 (\text{ap}_{\text{pre}} h'_1 \{\text{false}\}) \\
& \text{in } h'_2 \cup_{\text{pre}} h'_3 \\
& \text{in } h'_2 \circ_{\text{pre}} h'_1 \\
(h_1 \&\&_{\text{pre}} h_2) A ::= \langle h_1 A, h_2 A \rangle_{\text{pre}} & \text{lazy}_{\text{pre}} h A ::= \text{if } A = \emptyset \text{ then } \emptyset_{\text{pre}} \text{ else } h \ 0 \ A \\
& \text{lift}_{\text{pre}} ::= \text{pre}
\end{array}$$

Fig. 5: Preimage arrow definitions.

### 4.3 The Preimage Arrow

If we define the **preimage arrow** type constructor as

$$X \rightsquigarrow_{\text{pre}} Y ::= \text{Set } X \rightarrow (X \rightrightarrows_{\text{pre}} Y) \quad (24)$$

then we already have a lift  $\text{lift}_{\text{pre}} : (X \rightsquigarrow_{\perp} Y) \rightarrow (X \rightsquigarrow_{\text{pre}} Y)$  from the bottom arrow to the preimage arrow: **pre**. If  $\text{lift}_{\text{pre}}$  is **pre**, then by Theorem 3, lifted bottom arrow computations compute correct preimages, exactly as we should expect them to.

Fig. 5 defines the preimage arrow in terms of the preimage function operations in Fig. 4. For these definitions to make  $\text{lift}_{\text{pre}}$  a homomorphism, preimage arrow equivalence must mean “computes the same preimages.”

**Definition 5 (preimage arrow equivalence).** *Two preimage arrow computations  $h_1 : X \rightsquigarrow_{\text{pre}} Y$  and  $h_2 : X \rightsquigarrow_{\text{pre}} Y$  are equivalent, or  $h_1 \equiv h_2$ , when  $\text{ap}_{\text{pre}} (h_1 A) B = \text{ap}_{\text{pre}} (h_2 A) B$  for all  $A \subseteq X$  and  $B \subseteq Y$ .*

**Theorem 4 (preimage arrow correctness).**  *$\text{lift}_{\text{pre}}$  is a homomorphism.*

**Corollary 1 (semantic correctness).** *For all  $e$ ,  $\llbracket e \rrbracket_{\text{pre}} \equiv \text{lift}_{\text{pre}} \llbracket e \rrbracket_{\perp}$ .*

In other words,  $\llbracket e \rrbracket_{\text{pre}}$  always computes correct preimages under  $\llbracket e \rrbracket_{\perp}$ .

The type  $X \rightsquigarrow_{\text{pre}} Y$  does not constrain its inhabitants to behave intuitively; e.g.

$$\begin{array}{l}
\text{unruly} : \text{Bool} \rightsquigarrow_{\text{pre}} \text{Bool} \\
\text{unruly } A ::= \langle \text{Bool} \setminus A, \lambda B. B \rangle
\end{array} \quad (25)$$

So  $\text{ap}_{\text{pre}} (\text{unruly } \{\text{true}\}) \{\text{false}\} = \{\text{false}\} \cap (\text{Bool} \setminus \{\text{true}\}) = \{\text{false}\}$ —a “preimage” that does not even intersect the given domain  $\{\text{true}\}$ . Other examples show that preimage computations are not necessarily monotone, and lack other desirable properties. Those with desirable properties obey the following law.

**Definition 6 (preimage arrow law).** *Let  $h : X \rightsquigarrow_{\text{pre}} Y$ . If there exists an  $f : X \rightsquigarrow_{\perp} Y$  such that  $h \equiv \text{lift}_{\text{pre}} f$ , then  $h$  obeys the **preimage arrow law**.*

By homomorphism of  $\text{lift}_{\text{pre}}$ , preimage arrow combinators preserve the preimage arrow law. From here on, we assume all  $h : X \rightsquigarrow_{\text{pre}} Y$  obey it. By Definition 6,  $\text{lift}_{\text{pre}}$  has a right inverse; by Theorem 1, the arrow laws hold.

## 5 The Bottom\* and Preimage\* Arrows

We have defined the top of our roadmap:

$$\begin{array}{ccc}
 X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
 \eta_{\perp*} \downarrow & & \downarrow \eta_{\text{pre}*} \\
 X \rightsquigarrow_{\perp*} Y & \xrightarrow{\text{lift}_{\text{pre}*}} & X \rightsquigarrow_{\text{pre}*} Y
 \end{array} \tag{26}$$

so that  $\text{lift}_{\text{pre}}$  is a homomorphism. Now we move down each side and connect the bottom, in a way that makes every morphism a homomorphism.

Probabilistic functions that may not terminate, but terminate with probability 1, are common. For example, suppose `random` retrieves numbers in  $[0, 1]$  from an implicit random source. The following probabilistic function defines the well-known geometric distribution by counting the number of times `random`  $< p$ :

$$\text{geometric } p := \text{if } \text{random} < p \text{ then } 0 \text{ else } 1 + \text{geometric } p \tag{27}$$

For any  $p > 0$ , `geometric p` may not terminate, but the probability of never taking the “else” branch is  $(1 - p) \cdot (1 - p) \cdot (1 - p) \cdots = 0$ .

Suppose we interpret `geometric p` as  $h : R \rightsquigarrow_{\text{pre}} \mathbb{N}$ , a preimage arrow computation from random sources to naturals, and we have a probability measure  $P : \text{Set } R \rightarrow [0, 1]$ . The probability of  $N \subseteq \mathbb{N}$  is  $P(\text{ap}_{\text{pre}}(h \ R) \ N)$ . To compute this, we must

- Ensure  $\text{ap}_{\text{pre}}(h \ R) \ N$  terminates.
- Ensure each  $r \in R$  contains enough random numbers.
- Determine how `random` indexes numbers in  $r$ .

Ensuring  $\text{ap}_{\text{pre}}(h \ R) \ N$  terminates is the most difficult, but doing the other two will provide structure that makes it much easier.

### 5.1 Threading and Indexing

We need bottom and preimage arrows that thread a random source. To ensure random sources contain enough numbers, they should be infinite.

In a pure  $\lambda$ -calculus, random sources are typically infinite streams, threaded monadically: each computation receives and produces a random source. A little-used alternative is for the random source to be a tree [17], threaded applicatively: each computation receives, but does not produce, a random source. Combinators split the tree and pass subtrees to subcomputations.

With either alternative, for arrows, the resulting definitions are large, conceptually difficult, and hard to manipulate. Fortunately, it is relatively easy to assign each subcomputation a unique index into a tree-shaped random source and pass the random source unchanged. To do this, we need an indexing scheme.

**Definition 7 (binary indexing scheme).** *Let  $J$  be an index set,  $j_0 \in J$  a distinguished element, and  $\text{left} : J \rightarrow J$  and  $\text{right} : J \rightarrow J$  be total, injective functions. If for all  $j \in J$ ,  $j = \text{next } j_0$  for some finite composition  $\text{next}$  of  $\text{left}$  and  $\text{right}$ , then  $J$ ,  $j_0$ ,  $\text{left}$  and  $\text{right}$  define a **binary indexing scheme**.*

$$\begin{array}{ll}
x \rightsquigarrow_{a^*} y ::= J \rightarrow (\langle s, x \rangle \rightsquigarrow_a y) & \text{ifte}_{a^*} k_1 k_2 k_3 j ::= \\
\text{arr}_{a^*} ::= \eta_{a^*} \circ \text{arr}_a & \text{ifte}_a (k_1 (\text{left } j)) \\
(k_1 \ggg_{a^*} k_2) j ::= (\text{arr}_a \text{fst } \&\&_a k_1 (\text{left } j)) \ggg_a k_2 (\text{right } j) & (k_2 (\text{left } (\text{right } j))) \\
(k_1 \&\&_{a^*} k_2) j ::= k_1 (\text{left } j) \&\&_a k_2 (\text{right } j) & (k_3 (\text{right } (\text{right } j))) \\
& \text{lazy}_{a^*} k j ::= \text{lazy}_a \lambda 0. k 0 j \\
& \eta_{a^*} f j ::= \text{arr}_a \text{snd } \ggg_a f
\end{array}$$

Fig. 6: AStore (associative store) arrow transformer definitions.

For example, let  $J$  be the set of lists of  $\{0, 1\}$ ,  $j_0 := \langle \rangle$ , and  $\text{left } j := \langle 0, j \rangle$  and  $\text{right } j := \langle 1, j \rangle$ . In any case,  $J$  is countable, and can be thought of as a set of indexes into an infinite binary tree. Values of type  $A^J$  encode an infinite binary tree of  $A$  values as an infinite vector (i.e. total, extensional function).

We define random-source-threading variants of both the bottom and preimage arrows at once, by defining an **arrow transformer**: an arrow parameterized on another arrow. The AStore arrow transformer type constructor takes a store type  $s$  and an arrow  $x \rightsquigarrow_a y$ :

$$\text{AStore } s (x \rightsquigarrow_a y) ::= J \rightarrow (\langle s, x \rangle \rightsquigarrow_a y) \quad (28)$$

Reading the type, we see that computations receive an index  $j \in J$  and produce a computation that receives a store as well as an  $x$ . Lifting extracts the  $x$  from the input pair and sends it on to the original computation, ignoring  $j$ :

$$\begin{aligned}
\eta_{a^*} : (x \rightsquigarrow_a y) &\rightarrow \text{AStore } s (x \rightsquigarrow_a y) \\
\eta_{a^*} f j &:= \text{arr}_a \text{snd } \ggg_a f
\end{aligned} \quad (29)$$

Fig. 6 defines the remaining combinators. Each subcomputation receives  $\text{left } j$ ,  $\text{right } j$ , or some other unique binary index. We thus think of programs interpreted as AStore arrows as being completely unrolled into an infinite binary tree, with each expression labeled with its tree index.

## 5.2 Partial, Probabilistic Programs

To interpret probabilistic programs, we put infinite random trees in the store.

**Definition 8 (random source).** Define  $R := [0, 1]^J$ , the set of infinite binary trees whose node labels are in  $[0, 1]$ . A **random source** is any  $r \in R$ .

To interpret partial programs, we need to ensure termination. One ultimately implementable way is to have the store dictate which branch of each conditional, if any, is taken.

**Definition 9 (branch trace).** A **branch trace** is any  $t \in \text{Bool}_\perp^J$  such that  $t j = \text{true}$  or  $t j = \text{false}$  for no more than finitely many  $j \in J$ .

Let  $T \subset \text{Bool}_\perp^J$  be the largest set of branch traces.

Let  $X \rightsquigarrow_{a^*} Y ::= \text{AStore } (R \times T) (X \rightsquigarrow_a Y)$  be an AStore arrow type that threads both random stores and branch traces.



For probabilistic programs, we define a combinator  $\text{random}_{a^*}$  that returns the number at its tree index in the random source, and extend  $\llbracket \cdot \rrbracket_{a^*}$  for arrows  $a^*$  for which  $\text{random}_{a^*}$  is defined:

$$\begin{aligned} \text{random}_{a^*} &: X \rightsquigarrow_{a^*} [0, 1] \\ \text{random}_{a^*} j &:= \text{arr}_a (\text{fst} \ggg \text{fst} \ggg \pi j) \quad \llbracket \text{random} \rrbracket_{a^*} \equiv \text{random}_{a^*} \end{aligned} \quad (30)$$

where  $\pi : J \rightarrow X^J \rightarrow X$ , defined by  $\pi j f := f j$ , produces projection functions.

For partial programs, we define a combinator that reads branch traces, and an if-then-else combinator that ensures its test expression agrees with the trace:

$$\begin{aligned} \text{branch}_{a^*} &: X \rightsquigarrow_{a^*} \text{Bool} \\ \text{branch}_{a^*} j &:= \text{arr}_a (\text{fst} \ggg \text{snd} \ggg \pi j) \\ \text{ifte}_{a^*}^\downarrow &: (X \rightsquigarrow_{a^*} \text{Bool}) \rightarrow (X \rightsquigarrow_{a^*} Y) \rightarrow (X \rightsquigarrow_{a^*} Y) \rightarrow (X \rightsquigarrow_{a^*} Y) \\ \text{ifte}_{a^*}^\downarrow k_1 k_2 k_3 j &:= \text{ifte}_a ((k_1 (\text{left } j) \&\&_a \text{branch}_{a^*} j) \ggg_a \text{arr}_a \text{agrees}) \\ &\quad (k_2 (\text{left } (\text{right } j))) \\ &\quad (k_3 (\text{right } (\text{right } j))) \end{aligned} \quad (31)$$

where  $\text{agrees } \langle b_1, b_2 \rangle := \text{if } b_1 = b_2 \text{ then } b_1 \text{ else } \perp$ . Thus, if the branch trace does not agree with the test expression, it returns an error. We define a new semantic function  $\llbracket \cdot \rrbracket_{a^*}^\downarrow$  by replacing the if rule in  $\llbracket \cdot \rrbracket_{a^*}$ :

$$\llbracket \text{if } e_c \text{ then } e_t \text{ else } e_f \rrbracket_{a^*}^\downarrow \equiv \text{ifte}_{a^*}^\downarrow \llbracket e_c \rrbracket_{a^*}^\downarrow (\text{lazy}_{a^*} \lambda 0. \llbracket e_t \rrbracket_{a^*}^\downarrow) (\text{lazy}_{a^*} \lambda 0. \llbracket e_f \rrbracket_{a^*}^\downarrow) \quad (32)$$

For an AStore computation  $k$ , we obviously must run  $k$  on every branch trace in  $T$  and filter out  $\perp$ , or somehow find inputs  $\langle \langle r, t \rangle, a \rangle$  for which  $\text{agrees}$  never returns  $\perp$ . Preimage AStore arrows do the former by first computing an image, and the latter by computing preimages of sets that cannot contain  $\perp$ .

**Definition 10 (terminating, probabilistic arrows).** *Define*

$$\begin{aligned} X \rightsquigarrow_{\perp}^* Y &::= \text{AStore } (R \times T) (X \rightsquigarrow_{\perp} Y) \\ X \rightsquigarrow_{\text{pre}}^* Y &::= \text{AStore } (R \times T) (X \rightsquigarrow_{\text{pre}} Y) \end{aligned} \quad (33)$$

as the type constructors for the **bottom\*** and **preimage\*** arrows.

Suppose  $f := \llbracket e \rrbracket_{\perp}^\downarrow : X \rightsquigarrow_{\perp}^* Y$ . Its domain is  $X' := (R \times T) \times X$ . We assume each  $r \in R$  is random, but not  $t \in T$  nor  $a \in X$ ; therefore, neither  $T$  nor  $X$  should affect the probabilities of output sets. The probability of  $B \subseteq Y$  is therefore

$$\begin{aligned} &P (\text{image } (\text{fst} \ggg \text{fst}) (\text{preimage}_{\perp} f X' B)) \\ &= P (\text{image } (\text{fst} \ggg \text{fst}) (\text{ap}_{\text{pre}} (h X') B)) \end{aligned} \quad (34)$$

where  $h := \llbracket e \rrbracket_{\text{pre}}^\downarrow$ , if  $f$  and  $h$  always terminate and  $\llbracket \cdot \rrbracket_{\text{pre}}^\downarrow$  is correct.

### 5.3 Correctness and Termination

For correctness, we have two arrow lifts to prove homomorphic: one from pure computations to effectful (i.e. from those that do not access the store to those that do), and one from effectful computations to effectful. For both, we need AStore arrow equivalence to be a little coarser.

**Definition 11 (AStore arrow equivalence).** *Two AStore arrow computations  $k_1$  and  $k_2$  are equivalent, or  $k_1 \equiv k_2$ , when  $k_1 j \equiv k_2 j$  for all  $j \in J$ .*

**Theorem 5 (pure AStore arrow correctness).**  $\eta_{a^*}$  is a homomorphism.

**Corollary 2 (pure semantic correctness).** *For all pure  $e$ ,  $\llbracket e \rrbracket_{a^*} \equiv \eta_{a^*} \llbracket e \rrbracket_a$ .*

We need a lift between AStore arrows. Let  $x \rightsquigarrow_{a^*} y ::= \text{AStore } s \ (x \rightsquigarrow_a y)$ ,  $x \rightsquigarrow_{b^*} y ::= \text{AStore } s \ (x \rightsquigarrow_b y)$ , and  $\text{lift}_b : (x \rightsquigarrow_a y) \rightarrow (x \rightsquigarrow_b y)$ . Define

$$\begin{aligned} \text{lift}_{b^*} : (x \rightsquigarrow_{a^*} y) &\rightarrow (x \rightsquigarrow_{b^*} y) \\ \text{lift}_{b^*} f j &:= \text{lift}_b (f j) \end{aligned} \tag{35}$$

**Theorem 6 (effectful AStore arrow correctness).** *If  $\text{lift}_b$  is an arrow homomorphism from  $a$  to  $b$ , then  $\text{lift}_{b^*}$  is an arrow homomorphism from  $a^*$  to  $b^*$ .*

**Corollary 3 (preimage\* arrow correctness).**  $\text{lift}_{\text{pre}^*}$  is a homomorphism.

**Corollary 4 (effectful semantic correctness).** *For all expressions  $e$ ,  $\llbracket e \rrbracket_{\text{pre}^*} \equiv \text{lift}_{\text{pre}^*} \llbracket e \rrbracket_{\perp^*}$ .*

For termination, we need to define the largest domain on which  $\llbracket e \rrbracket_{a^*}^\downarrow$  and  $\llbracket e \rrbracket_{a^*}$  computations should agree.

**Definition 12 (maximal domain).** *Let  $f : X \rightsquigarrow_{\perp^*} Y$ . Its **maximal domain** is the largest  $A^* \subseteq (R \times T) \times X$  for which  $A^* = \{a \in A^* \mid f j_0 a \neq \perp\}$ .*

Because  $f j_0 a \neq \perp$  implies termination, all inputs in  $A^*$  are terminating.

**Theorem 7 (correct termination everywhere).** *Let  $\llbracket e \rrbracket_{\perp^*}^\downarrow : X \rightsquigarrow_{\perp^*} Y$  have maximal domain  $A^*$ , and  $X' := (R \times T) \times X$ . For all  $a \in X'$ ,  $A \subseteq X'$  and  $B \subseteq Y$ ,*

$$\begin{aligned} \llbracket e \rrbracket_{\perp^*}^\downarrow j_0 a &= \text{if } a \in A^* \text{ then } \llbracket e \rrbracket_{\perp^*} j_0 a \text{ else } \perp \\ \text{ap}_{\text{pre}} (\llbracket e \rrbracket_{\text{pre}^*}^\downarrow j_0 A) B &= \text{ap}_{\text{pre}} (\llbracket e \rrbracket_{\text{pre}^*} j_0 (A \cap A^*)) B \end{aligned} \tag{36}$$

In other words, preimages computed using  $\llbracket \cdot \rrbracket_{\text{pre}^*}^\downarrow$  always terminate, never include inputs that give rise to errors or nontermination, and are correct.

## 6 Approximating Semantics

We would like to compute preimages of uncountable sets, such as real intervals—but  $\text{preimage}_\perp f A B$  is uncomputable for most uncountable sets  $A$  and  $B$  no matter how cleverly they are represented. Further, because  $\text{pre}$ ,  $\text{lift}_{\text{pre}}$  and  $\text{arr}_{\text{pre}}$  are defined in terms of  $\text{preimage}_\perp$ , we cannot implement them.

Fortunately, we need only certain lifts. Fig. 7 gives explicit definitions for  $\text{arr}_{\text{pre}} \text{id}$ ,  $\text{arr}_{\text{pre}} \text{fst}$ ,  $\text{arr}_{\text{pre}} \text{snd}$ ,  $\text{arr}_{\text{pre}} (\text{const } b)$  and  $\text{arr}_{\text{pre}} (\pi j)$ . To implement them,

$\text{id}_{\text{pre}} A := \langle A, \lambda B. B \rangle$	$\text{const}_{\text{pre}} b A := \langle \{b\}, \lambda B. \text{if } B = \emptyset \text{ then } \emptyset \text{ else } A \rangle$
$\text{fst}_{\text{pre}} A := \langle \text{proj}_1 A, \text{unproj}_1 A \rangle$	$\pi_{\text{pre}} j A := \langle \text{proj } j A, \text{unproj } j A \rangle$
$\text{snd}_{\text{pre}} A := \langle \text{proj}_2 A, \text{unproj}_2 A \rangle$	
$\text{proj}_1 := \text{image fst}$	$\text{proj} : J \rightarrow \text{Set } X^J \rightarrow \text{Set } X$
$\text{proj}_2 := \text{image snd}$	$\text{proj } j A := \text{image } (\pi j) A$
$\text{unproj}_1 A B := A \cap (B \times \text{proj}_2 A)$	$\text{unproj} : J \rightarrow \text{Set } X^J \rightarrow \text{Set } X \rightarrow \text{Set } X^J$
$\text{unproj}_2 A B := A \cap (\text{proj}_1 A \times B)$	$\text{unproj } j A B := A \cap \prod_{i \in J} \text{if } j = i \text{ then } B \text{ else } \text{proj } j A$

Fig. 7: Preimage arrow lifts needed to interpret probabilistic programs.

we must model sets in a way that makes  $A = \emptyset$  decidable, and the following representable and finitely computable:

- $A \cap B$ ,  $\emptyset$ ,  $\{\text{true}\}$ ,  $\{\text{false}\}$  and  $\{b\}$  for every  $\text{const } b$
- $A_1 \times A_2$ ,  $\text{proj}_1 A$  and  $\text{proj}_2 A$
- $X^J$ ,  $\text{proj } j A$  and  $\text{unproj } j A B$

We first need to define families of sets under which these operations are closed.

**Definition 13 (rectangular family).** *Rect  $X$  denotes the **rectangular family** of subsets of  $X$ . Rect  $X$  must contain  $\emptyset$  and  $X$ , and be closed under finite intersections. Products must satisfy the following rules:*

$$\text{Rect } \langle X_1, X_2 \rangle = (\text{Rect } X_1) \boxtimes (\text{Rect } X_2) \quad (38)$$

$$\text{Rect } X^J = (\text{Rect } X)^{\boxtimes J} \quad (39)$$

where the following operations lift cartesian products to sets of sets:

$$\mathcal{A}_1 \boxtimes \mathcal{A}_2 := \{A_1 \times A_2 \mid A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2\} \quad (40)$$

$$\mathcal{A}^{\boxtimes J} := \bigcup_{J' \subset J \text{ finite}} \left\{ \prod_{j \in J'} A_j \mid A_j \in \mathcal{A}, j \in J' \iff A_j \subset \bigcup \mathcal{A} \right\} \quad (41)$$

We additionally define  $\text{Rect Bool} ::= \mathcal{P} \text{ Bool}$ . It is easy to show the collection of all rectangular families is closed under products, projections, and  $\text{unproj}$ .

Further, all of the operations in (37) can be exactly implemented if finite sets are modeled directly, sets in ordered spaces (such as  $\mathbb{R}$ ) are modeled by intervals, and sets in  $\text{Rect } \langle X_1, X_2 \rangle$  are modeled by pairs of type  $\langle \text{Rect } X_1, \text{Rect } X_2 \rangle$ . By (41), sets in  $\text{Rect } X^J$  have no more than finitely many axes that are proper subsets of  $X$ . They can be modeled by *finite* binary trees, where unrepresented axes are implicitly  $X$ .

The set of branch traces  $T$  is nonrectangular, but we can model  $T$  subsets by  $\text{Bool}_\perp^J$  rectangles, implicitly intersected with  $T$ .

**Theorem 8 (T model).** *If  $T' \in \text{Rect Bool}_\perp^J$  and  $j \in J$ , then  $\text{proj } j T' = \text{proj } j (T' \cap T)$ . If  $B \subseteq \text{Bool}_\perp$ , then  $\text{unproj } j (T' \cap T) B = \text{unproj } j T' B \cap T$ .*

$$\begin{array}{ll}
X \xrightarrow{\text{pre}}' Y ::= \langle \text{Rect } Y, \text{Rect } Y \rightarrow \text{Rect } X \rangle & \langle \cdot, \cdot \rangle'_{\text{pre}} : (X \xrightarrow{\text{pre}}' Y_1) \rightarrow (X \xrightarrow{\text{pre}}' Y_2) \rightarrow (X \xrightarrow{\text{pre}}' Y_1 \times Y_2) \\
\emptyset'_{\text{pre}} ::= \langle \emptyset, \lambda B. \emptyset \rangle & \langle \langle Y'_1, p_1 \rangle, \langle Y'_2, p_2 \rangle \rangle'_{\text{pre}} ::= \\
& \langle Y'_1 \times Y'_2, \lambda B. p_1 (\text{proj}_1 B) \cap p_2 (\text{proj}_2 B) \rangle \\
\text{ap}'_{\text{pre}} : (X \xrightarrow{\text{pre}}' Y) \rightarrow \text{Rect } Y \rightarrow \text{Rect } X & (\cup'_{\text{pre}}) : (X \xrightarrow{\text{pre}}' Y) \rightarrow (X \xrightarrow{\text{pre}}' Y) \rightarrow (X \xrightarrow{\text{pre}}' Y) \\
\text{ap}'_{\text{pre}} \langle Y', p \rangle B := p (B \cap Y') & \langle Y'_1, p_1 \rangle \cup'_{\text{pre}} \langle Y'_2, p_2 \rangle ::= \\
(\circ'_{\text{pre}}) : (Y \xrightarrow{\text{pre}}' Z) \rightarrow (X \xrightarrow{\text{pre}}' Y) \rightarrow (X \xrightarrow{\text{pre}}' Z) & \langle Y'_1 \vee Y'_2, \lambda B. \text{ap}'_{\text{pre}} \langle Y'_1, p_1 \rangle B \vee \text{ap}'_{\text{pre}} \langle Y'_2, p_2 \rangle B \rangle \\
\langle Z', p_2 \rangle \circ'_{\text{pre}} h_1 := \langle Z', \lambda C. \text{ap}'_{\text{pre}} h_1 (p_2 C) \rangle &
\end{array}$$

(a) Definitions for preimage functions that compute rectangular covers.

$$\begin{array}{ll}
X \xrightarrow{\text{pre}}' Y ::= \text{Rect } X \rightarrow (X \xrightarrow{\text{pre}}' Y) & \text{ifte}'_{\text{pre}} h_1 h_2 h_3 A := \\
& \text{let } h'_1 := h_1 A \\
(h_1 \ggg'_{\text{pre}} h_2) A := \text{let } h'_1 := h_1 A & h'_2 := h_2 (\text{ap}'_{\text{pre}} h'_1 \{\text{true}\}) \\
& h'_2 := h_2 (\text{range}'_{\text{pre}} h'_1) \\
& \text{in } h'_2 \circ'_{\text{pre}} h'_1 \\
(h_1 \&\&'_{\text{pre}} h_2) A := \langle h_1 A, h_2 A \rangle'_{\text{pre}} & \text{in } h'_2 \cup'_{\text{pre}} h'_3 \\
& \text{lazy}'_{\text{pre}} h A := \text{if } A = \emptyset \text{ then } \emptyset'_{\text{pre}} \text{ else } h \circ A
\end{array}$$

(b) Approximating preimage arrow, defined using approximating preimage functions.

$$\begin{array}{ll}
X \xrightarrow{\text{pre}^*} Y ::= \text{AStore } (R \times T) (X \xrightarrow{\text{pre}^*} Y) & \text{ifte}^{\text{ll}}_{\text{pre}^*} : (X \xrightarrow{\text{pre}^*} \text{Bool}) \rightarrow (X \xrightarrow{\text{pre}^*} Y) \rightarrow (X \xrightarrow{\text{pre}^*} Y) \rightarrow (X \xrightarrow{\text{pre}^*} Y) \\
\text{random}'_{\text{pre}^*} : X \xrightarrow{\text{pre}^*} [0, 1] & \text{ifte}^{\text{ll}}_{\text{pre}^*} k_1 k_2 k_3 j := \\
\text{random}'_{\text{pre}^*} j := & \text{let } \langle C_k, p_k \rangle := k_1 (\text{left } j) A \\
\text{fst}_{\text{pre}} \ggg'_{\text{pre}} \text{fst}_{\text{pre}} \ggg'_{\text{pre}} \pi_{\text{pre}} j & \langle C_b, p_b \rangle := \text{branch}_{\text{pre}^*} j A \\
& C_2 := C_k \cap C_b \cap \{\text{true}\} \\
& C_3 := C_k \cap C_b \cap \{\text{false}\} \\
\text{branch}'_{\text{pre}^*} : X \xrightarrow{\text{pre}^*} \text{Bool} & A_2 := p_k C_2 \cap p_b C_2 \\
& A_3 := p_k C_3 \cap p_b C_3 \\
\text{branch}'_{\text{pre}^*} j := & \text{in if } C_b = \{\text{true}, \text{false}\} \\
\text{fst}_{\text{pre}} \ggg'_{\text{pre}} \text{snd}_{\text{pre}} \ggg'_{\text{pre}} \pi_{\text{pre}} j & \text{then } \langle T, \lambda B. A_2 \vee A_3 \rangle \\
& \text{else } \langle k_2 (\text{left } (\text{right } j)) A_2 \cup'_{\text{pre}} k_3 (\text{right } (\text{right } j)) A_3 \rangle \\
\text{fst}'_{\text{pre}^*} := \eta'_{\text{pre}^*} \text{fst}_{\text{pre}}; \dots &
\end{array}$$

(c) Preimage\* arrow combinators for probabilistic choice and guaranteed termination. Fig. 6 (AStore arrow transformer) defines  $\eta'_{\text{pre}^*}$ ,  $(\ggg'_{\text{pre}^*})$ ,  $(\&\&'_{\text{pre}^*})$ ,  $\text{ifte}'_{\text{pre}^*}$  and  $\text{lazy}'_{\text{pre}^*}$ .

Fig. 8: Implementable arrows that approximate preimage arrows. Specific lifts such as  $\text{fst}_{\text{pre}} := \text{arr}_{\text{pre}} \text{fst}$  are computable (see Fig. 7), but  $\text{arr}'_{\text{pre}}$  is not.

Rectangular families are not closed under  $(\cup)$ . For conditionals, then, we need a lattice join  $(\vee)$  with respect to  $(\subseteq)$  with the following additional properties:

$$\begin{aligned}
(A_1 \times A_2) \vee (B_1 \times B_2) &= (A_1 \vee B_1) \times (A_2 \vee B_2) \\
(\prod_{j \in J} A_j) \vee (\prod_{j \in J} B_j) &= \prod_{j \in J} A_j \vee B_j
\end{aligned} \tag{42}$$

If for every nonproduct type  $X$ ,  $\text{Rect } X$  is closed under  $(\vee)$ , then rectangular families are clearly closed under  $(\vee)$ . Further, for any  $A$  and  $B$ ,  $A \cup B \subseteq A \vee B$ .

Fig. 8 defines approximating preimage arrows. Approximating preimage function operations (Fig. 8a) are defined in terms of lattice operations on rectangular families. Every approximating preimage arrow combinator (Fig. 8b) is defined the same way as its corresponding exact preimage arrow combinator, but using approximating preimage function operations instead of exact. Fig. 8c defines

$\text{random}'_{\text{pre}^*}$  and  $\text{branch}'_{\text{pre}^*}$  without using the uncomputable  $\text{arr}'_{\text{pre}^*}$ , and  $\text{ifte}'_{\text{pre}^*}$ , for interpreting expressions using  $\llbracket \cdot \rrbracket'_{\text{pre}^*}$  for guaranteed termination.

Let  $h := \llbracket e \rrbracket'_{\text{pre}^*} : X \rightsquigarrow_{\text{pre}^*} Y$  and  $h' := \llbracket e \rrbracket'_{\text{pre}^*} : X \rightsquigarrow_{\text{pre}^*}' Y$  for some expression  $e$ .

**Theorem 9 (monotone).**  $\lambda A B. \text{ap}'_{\text{pre}} (h' j_0 A) B$  is monotone in  $A$  and  $B$ .

**Theorem 10 (sound, terminating, decreasing).** For all  $A \in \text{Rect } \langle \langle R, T \rangle, X \rangle$  and  $B \in \text{Rect } Y$ ,  $\text{ap}_{\text{pre}} (h j_0 A) B \subseteq \text{ap}'_{\text{pre}} (h' j_0 A) B \subseteq A$ .

Given these properties, we might try to compute preimages of  $B$  by computing preimages with respect to increasingly fine discretizations of  $A$ .

**Definition 14 (preimage refinement algorithm).** Let  $B \in \text{Rect } Y$  and

$$\begin{aligned} \text{refine} &: \text{Rect } \langle \langle R, T \rangle, X \rangle \rightarrow \text{Rect } \langle \langle R, T \rangle, X \rangle \\ \text{refine } A &:= \text{ap}'_{\text{pre}} (h' j_0 A) B \end{aligned} \tag{43}$$

Define  $\text{partition} : \text{Rect } \langle \langle R, T \rangle, X \rangle \rightarrow \text{Set } (\text{Rect } \langle \langle R, T \rangle, X \rangle)$  to produce positive-measure, disjoint rectangles, and define

$$\begin{aligned} \text{refine}^* &: \text{Set } (\text{Rect } \langle \langle R, T \rangle, X \rangle) \rightarrow \text{Set } (\text{Rect } \langle \langle R, T \rangle, X \rangle) \\ \text{refine}^* \mathcal{A} &:= \text{image refine } (\bigcup_{A \in \mathcal{A}} \text{partition } A) \end{aligned} \tag{44}$$

For any  $A \in \text{Rect } \langle \langle R, T \rangle, X \rangle$ , iterate  $\text{refine}^*$  on  $\{A\}$ .

Theorem 9 guarantees that refining a partition of  $A$  never does worse than refining  $A$  itself. Theorem 10 guarantees that  $\text{refine } A \subseteq A$ , and that the algorithm is **sound**: the preimage of  $B$  is always contained in the covering partition  $\text{refine}^*$  returns. Ideally, it would be complete in the sense that covering partitions converge to a set that overapproximates by a measure-zero subset. Unfortunately, convergence fails on some examples that terminate with probability less than one. We leave completeness conditions for future work, and for now, use Monte Carlo algorithms, which depend only on soundness.

## 7 Implementations and Examples

We have three implementations: two direct implementations of the approximating semantics, and a less direct but more efficient one called **Dr. Bayes**. All of them can be found at <https://github.com/ntoronto/drbytes>.

Given a library for operating on rectangular sets, the approximating preimage arrows defined in Figs. 7 and 8 can be implemented with few changes in any practical  $\lambda$ -calculus. We have done so in Typed Racket [26] and Haskell [1]. Both implementations are almost line-for-line transliterations from the figures.

Dr. Bayes is written in Typed Racket. It includes  $\llbracket \cdot \rrbracket_{a^*}$  (Fig. 2), its extension  $\llbracket \cdot \rrbracket_{a^*}'$ , the bottom\* arrow (Figs. 3 and 6), the approximating preimage and preimage\* arrows (Figs. 7 and 8), and manual lifts such as  $\text{arr}'_{\text{pre}} (+)$  to compute approximate preimages under arithmetic. The preimage arrows operate on

a monomorphic rectangular set data type, which includes tagged rectangles and disjoint unions for ad-hoc polymorphism, and floating-point intervals to overapproximate real intervals.

Definition 14 outlines preimage refinement, a discretization algorithm that repeatedly shrinks and repartitions a program’s domain. *Dr. Bayes does not use this algorithm directly* because it is inefficient: good accuracy requires fine discretization, which is exponential in the number of discretized axes. Instead of *enumerating* covering partitions of the random source, Dr. Bayes *samples parts* from the covering partitions and then *samples a point* from each sampled part, with time complexity linear in the number of samples and discretized axes. It applies bottom\* arrow computations to the random source samples to get output samples, rejecting those outside the requested output set. Thus, it relies only on preimage refinement’s soundness.

We have tested Dr. Bayes on a variety of Bayesian inference tasks, including Bayesian regression and model selection [27]. Some of our Bayesian inference tests use recursion and constrain the outputs of deterministic functions, suggesting that Dr. Bayes and future probabilistic languages like it will allow practitioners to model real-world processes more expressively and precisely.

Recent work in probabilistic verification recasts it as a probabilistic inference task [8]. Given that Dr. Bayes’s runtime is designed to sample efficiently under low-probability constraints, using it to probabilistically verify that a program does not exhibit certain errors is fairly natural. To do so, we

1. Encode the program in a way that propagates and returns errors.
2. Run the program with the constraint that the output is an error.

Sometimes, Dr. Bayes can determine that the preimage of the constrained output set is  $\emptyset$ , which is a proof that the program never exhibits an error. Otherwise, the longer the program runs without returning samples, the likelier it is that the preimage has zero probability or is empty; i.e. that an error does not occur.

As an extended example, we consider verifying floating-point error bounds.

While Dr. Bayes’s numbers are implemented by floating-point intervals, semantically, they are real numbers. We therefore cannot represent floating-point numbers directly in Dr. Bayes—but we do not want to. We want *abstract* floating-point numbers, each consisting of an exact, real number and a bound on the relative error with which it is approximated. We define the following two structures to represent abstract floats.

```
(struct/drbytes float-any ())
(struct/drbytes float (value error))
```

An abstract value `(float v e)` represents every float between `(* v (- 1 e))` and `(* v (+ 1 e))` inclusive, while `(float-any)` represents NaN and other catastrophic error conditions. Abstract floating-point functions such as `flsqrt` compute exact results and use input error to compute bounds on output error:

```

(define/drbytes (flsqrt x)
  (if (float-any? x)
      x
      (let ([v (float-value x)]
            [e (float-error x)])
        (cond [(negative? v) (float-any)] ; NaN
              [(zero? v) (float 0 0)] ; exact case
              [else ; v is positive
               (float (sqrt v) ; exact square root
                      (+ (- 1 (sqrt (- 1 e))) ; relative error
                         (* 1/2 epsilon)))))) ; rounding error

```

We have similarly implemented abstract floating-point arithmetic, comparison, exponentials, and logarithms in Dr. Bayes.

Suppose we define an abstract floating-point implementation of the geometric distribution's inverse CDF using the formula  $(\log u)/(\log (1 - p))$ :

```

(define/drbytes (flgeometric-inv-cdf u p)
  (fl/ (fllog u) (fllog (fl- (float 1 0) p))))

```

We want the distribution of  $\langle u, p \rangle$  in  $(0, 1) \times (0, 1)$  with the value of

```

(float-error (flgeometric-inv-cdf (float u 0) (float p 0)))

```

constrained to  $(3 \cdot \varepsilon, \infty)$ , where  $\varepsilon \approx 2.22 \cdot 10^{-16}$  is floating-point epsilon for 64-bit floats. That is, we want the distribution of inputs for which the floating-point output may be more than 3 epsilons away from the exact output.

Dr. Bayes returns samples of  $\langle u, p \rangle$  within about  $(0, 1) \times (\varepsilon, 0.284)$ , a fairly large domain on which error is greater than 3 epsilons. Realizing that the rounding error in  $1 - p$  is magnified by  $\log$ 's relative error when  $p$  is small, we define

```

(define/drbytes (flgeometric-inv-cdf u p)
  (fl/ (fllog u) (fllog1p (flneg p))))

```

where `fllog1p` (abstractly) computes  $\log (1 + x)$  with high accuracy. Dr. Bayes reports that the preimage of  $(3 \cdot \varepsilon, \infty)$  is  $\emptyset$ . In fact, the preimage of  $(1.51 \cdot \varepsilon, \infty)$  is  $\emptyset$ , so `flgeometric-inv-cdf` introduces error of no more than 1.51 epsilons.

We have used this technique to verify error bounds on the implementations of `hypot`, `sqrt1pm1` and `sinh` in Racket's `math` library.

## 8 Related Work

Probabilistic languages can be approximately placed into two groups: those defined by a semantics, and those defined by an implementation.

Kozen's seminal work [14] on probabilistic semantics defines two measure-theoretic, denotational semantics, in two different styles: a **random-world semantics**,<sup>4</sup> that interprets programs as deterministic functions that operate on

<sup>4</sup> The term *random-world semantics* was coined by McAllester et al [17], and alludes to the *possible-world semantics* used to explain modal logic.

a random source, and a **distributional semantics** that interprets programs as probability measures. It seems that all semantics work thereafter is in one of these styles. For example, Hurd [10] develops a random-world semantics in HOL and uses it to formally verify randomized algorithms such as the Miller-Rabin primality test. Ours is also a random-world semantics.

Jones [11] defines the probability monad as a categorical metatheory for interpreting probabilistic programs as distributions. Ramsey and Pfeffer [24] reformulate it in terms of Haskell’s `return` and `>>=`, and use it to define a distributional semantics for a probabilistic lambda calculus. They implement the probability monad using probability mass functions, show that computing certain queries is inefficient, and devise an equivalent semantics that is more amenable to efficient implementation, for programs with finite probabilistic choice.

To put Infer.NET [20] on solid footing, Borgström et al [4] define a distributional semantics for a first-order probabilistic language with bounded loops and constraints, by transforming terms into arrow-like combinators that produce measures. But Infer.NET interprets programs as probability density functions,<sup>5</sup>so they develop a semantics that does the same and prove equivalence.

Borgström et al’s work and Ramsey and Pfeffer’s work are just two examples of a larger trend: while *defining* probabilistic languages can be done using measure theory, *implementing* them to support more than just straightforward evaluation (e.g. allowing constraints) has seemed hopeless enough to necessitate using a less explanatory theory of probability that has more obvious computational content. Indeed, the distributional semantics of Pfeffer’s IBAL [23] and Nori et al’s R2 [22] are defined in terms of probability mass and density functions in the first place. R2 lifts some of the resulting restrictions and speeds up sampling by propagating constraints toward the random values they refer to.

Some languages defined by an implementation are probabilistic Scheme [13], BUGS [16], BLOG [19], BLAISE [3], Church [7], and Kiselyov’s embedded language for OCaml [12]. Recently, Wingate et al [31] define nonstandard semantics that enable efficient inference, but do not define the languages. All of these languages are implemented in terms of probability mass or density functions.

Our work similar in structure to monadic abstract interpretation [25], which also parameterizes a semantics on categorical meanings. It is superficially similar to probabilistic abstract interpretation [5], which is a general framework for lifting static analyses into probabilistic settings. In contrast, we build a specific probabilistic language and—in contrast to most other abstract interpretation as well—use our abstract semantics to *run* programs rather than *analyze* them.

## 9 Conclusions and Future Work

To allow arbitrary constraints and recursion in probabilistic programs, we combined the power of measure theory with the unifying elegance of arrows. We

1. Defined a transformation from first-order programs to arbitrary arrows.

---

<sup>5</sup> More precisely, as factor graphs, which represent probability density functions.



2. Defined the bottom arrow as a standard translation target.
3. Derived the uncomputable preimage arrow as an alternative target.
4. Derived a sound, computable approximation of the preimage arrow, and enough computable lifts to transform programs.

We implemented the approximating semantics and carried out Bayesian inference, stochastic ray tracing, and probabilistic verification.

In the future, we intend to extend our first-order language’s expressiveness by adding lambdas, explore ways to use static or dynamic analyses to speed up Monte Carlo algorithms, and explore preimage computation’s connections to type checking and type inference. More broadly, we hope to advance probabilistic inference by providing a rich modeling language with an efficient, correct implementation, which allows general recursion and arbitrary constraints.

**Acknowledgments.** Special thanks to Mitch Wand for an additional careful review and helpful feedback.

## References

1. Haskell 98 language and libraries, the revised report (December 2002), <http://www.haskell.org/onlinereport/>
2. Aumann, R.J.: Borel structures for function spaces. *Illinois Journal of Mathematics* 5, 614–630 (1961)
3. Bonawitz, K.A.: *Composable Probabilistic Inference with Blaise*. Ph.D. thesis, Massachusetts Institute of Technology (2008)
4. Borgström, J., Gordon, A.D., Greenberg, M., Margetson, J., Gael, J.V.: Measure transformer semantics for Bayesian machine learning. In: *European Symposium on Programming*. pp. 77–96 (2011)
5. Cousot, P., Monerau, M.: Probabilistic abstract interpretation. In: *European Symposium on Programming*. pp. 166–190 (2012)
6. DeGroot, M., Schervish, M.: *Probability and Statistics*. Addison Wesley Publishing Company, Inc. (2012)
7. Goodman, N., Mansinghka, V., Roy, D., Bonawitz, K., Tenenbaum, J.: Church: a language for generative models. In: *Uncertainty in Artificial Intelligence* (2008)
8. Gulwani, S., Jojic, N.: Program verification as probabilistic inference. In: *Principles of Programming Languages*. pp. 277–289 (2007)
9. Hughes, J.: Generalizing monads to arrows. In: *Science of Computer Programming*. vol. 37, pp. 67–111 (2000)
10. Hurd, J.: *Formal Verification of Probabilistic Algorithms*. Ph.D. thesis, University of Cambridge (2002)
11. Jones, C.: *Probabilistic Non-Determinism*. Ph.D. thesis, Univ. of Edinburgh (1990)
12. Kiselyov, O., Shan, C.: Monolingual probabilistic programming using generalized coroutines. In: *Uncertainty in Artificial Intelligence* (2008)
13. Koller, D., McAllester, D., Pfeffer, A.: Effective Bayesian inference for stochastic programs. In: *14th National Conference on Artificial Intelligence* (August 1997)
14. Kozen, D.: Semantics of probabilistic programs. In: *Foundations of Computer Science* (1979)

15. Lindley, S., Wadler, P., Yallop, J.: Idioms are oblivious, arrows are meticulous, monads are promiscuous. *Electronic Notes in Theoretical Computer Science* (2008)
16. Lunn, D.J., Thomas, A., Best, N., Spiegelhalter, D.: WinBUGS – a Bayesian modelling framework. *Statistics and Computing* 10(4) (2000)
17. Mcallester, D., Milch, B., Goodman, N.D.: Random-world semantics and syntactic independence for expressive languages. Tech. rep., Massachusetts Institute of Technology (2008)
18. McBride, C., Paterson, R.: Applicative programming with effects. *Journal of Functional Programming* 18(1) (2008)
19. Milch, B., Marthi, B., Russell, S., Sontag, D., Ong, D., Kolobov, A.: BLOG: Probabilistic models with unknown objects. In: *International Joint Conference on Artificial Intelligence* (2005)
20. Minka, T., Winn, J., Guiver, J., Webster, S., Zaykov, Y., Yangel, B., Spengler, A., Bronskill, J.: Infer.NET 2.6 (2014), <http://research.microsoft.com/infernet>
21. Moggi, E.: Computational lambda-calculus and monads. In: *IEEE Symposium on Logic in Computer Science*. pp. 14–23 (1989)
22. Nori, A.V., Hur, C.K., Rajamani, S.K., Samuel, S.: R2: An efficient MCMC sampler for probabilistic programs. In: *AAAI Conference on Artificial Intelligence* (2014)
23. Pfeffer, A.: The design and implementation of IBAL: A general-purpose probabilistic language. In: *Statistical Relational Learning*. MIT Press (2007)
24. Ramsey, N., Pfeffer, A.: Stochastic lambda calculus and monads of probability distributions. In: *Principles of Programming Languages* (2002)
25. Sergey, I., Devriese, D., Might, M., Midtgaard, J., Darais, D., Clarke, D., Piessens, F.: Monadic abstract interpreters. In: *Programming Language Design and Implementation*. pp. 399–410 (2013)
26. Tobin-Hochstadt, S., Felleisen, M.: The design and implementation of typed Scheme. In: *Principles of Programming Languages*. pp. 395–406 (2008)
27. Toronto, N.: Trustworthy, Useful Languages for Probabilistic Modeling and Inference. Ph.D. thesis, Brigham Young University (2014), <http://students.cs.byu.edu/~ntoronto/dissertation.pdf>
28. Toronto, N., McCarthy, J.: Computing in Cantor’s paradise with  $\lambda_{ZFC}$ . In: *Functional and Logic Programming Symposium*. pp. 290–306 (2012)
29. Veach, E., Guibas, L.J.: Metropolis light transport. In: *ACM SIGGRAPH*. pp. 65–76 (1997)
30. Wadler, P.: Monads for functional programming. In: Jeuring, J., Meijer, E. (eds.) *Advanced Functional Programming* (2001)
31. Wingate, D., Goodman, N.D., Stuhlmüller, A., Siskind, J.M.: Nonstandard interpretations of probabilistic programs for efficient inference. In: *Neural Information Processing Systems*. pp. 1152–1160 (2011)