

## DERIVING PROBABILITY DENSITY FUNCTIONS FROM PROBABILISTIC FUNCTIONAL PROGRAMS

SOORAJ BHAT<sup>a</sup>, JOHANNES BORGSTRÖM<sup>b</sup>, ANDREW D. GORDON<sup>c</sup>, AND CLAUDIO RUSSO<sup>d</sup>

<sup>a</sup> Georgia Institute of Technology

<sup>b</sup> Uppsala University  
*e-mail address:* johannes.borgstrom@it.uu.se

<sup>c</sup> Microsoft Research and University of Edinburgh  
*e-mail address:* adg@microsoft.com

<sup>d</sup> Microsoft Research  
*e-mail address:* crusso@microsoft.com

---

**ABSTRACT.** The *probability density function* of a probability distribution is a fundamental concept in probability theory and a key ingredient in various widely used machine learning methods. However, the necessary framework for compiling probabilistic functional programs to density functions has only recently been developed. In this work, we present a density compiler for a probabilistic language with failure and both discrete and continuous distributions, and provide a proof of its soundness. The compiler greatly reduces the development effort of domain experts, which we demonstrate by solving inference problems from various scientific applications, such as modelling the global carbon cycle, using a standard Markov chain Monte Carlo framework.

### 1. INTRODUCTION

*Probabilistic programming* promises to arm data scientists with declarative languages for specifying their probabilistic models, while leaving the details of how to translate those models to efficient sampling or inference algorithms to a compiler. Many widely used machine learning techniques that might be employed by such a compiler use the *probability density function* (PDF) of the model as input. Such techniques include *maximum likelihood* or *maximum a posteriori estimation*, *L2 estimation*, *importance sampling*, and *Markov chain Monte Carlo* (MCMC) methods (Scott, 2001; Bishop, 2006).

However, despite their utility, density functions have been largely absent from the literature on probabilistic functional programming (Ramsey and Pfeffer, 2002; Goodman et al., 2008; Kiselyov and Shan, 2009). This is because the relationship between programs and their density functions is not straightforward: for a given program, the PDF may not exist or may be non-trivial to calculate. Such programs are not merely infrequent pathological curiosities but in fact arise in many ordinary scenarios. In this paper, we define, prove correct, and implement an algorithm for automatically

---

*Key words and phrases:* probability density function, probabilistic programming, Markov chain Monte Carlo.

computing PDFs for a large class of programs written in a rich probabilistic programming language. An abridged version of this paper was published as (Bhat et al., 2013).

**Probability density functions.** In this work, probabilistic programs correspond directly to *probability distributions*, which are important because they are a powerful formalism for data analysis. However, many techniques we would like to use require the *probability density function* of a distribution instead of the distribution itself. Unfortunately, not every distribution has a density function.

**Distributions.** One interpretation of a probabilistic program is that it is a simulation that can be run to generate a random sample from some set  $\Omega$  of possible outcomes. The corresponding probability distribution  $\mathbb{P}$  characterizes the program by assigning probabilities to different subsets of  $\Omega$  (*events*). The probability  $\mathbb{P}(A)$  for a subset  $A$  of  $\Omega$  corresponds to the proportion of runs that generate an outcome in  $A$ , in the limit of an infinite number of repeated runs of the simulation.

Consider for example a simple *mixture of Gaussians*, here written in Fun (Borgström et al., 2011), a probabilistic functional language embedded within F# (Syme et al., 2007).

```
if flip 0.7 then random(Gaussian(0.0, 1.0)) else random(Gaussian(4.0, 1.0))
```

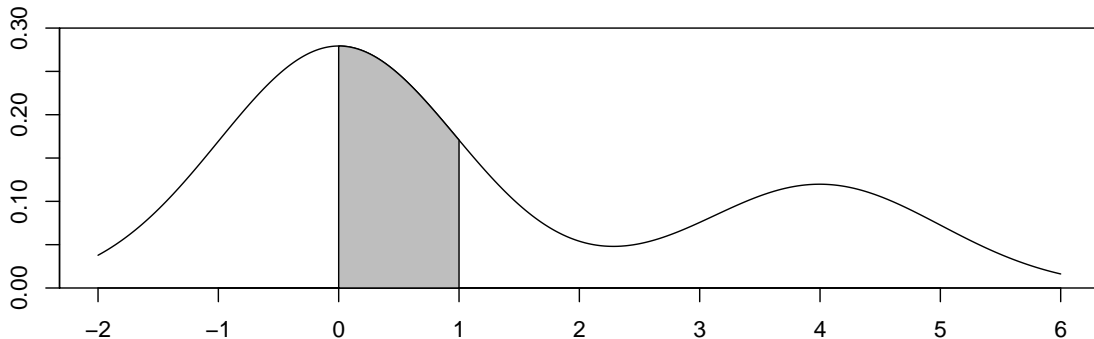
The program above specifies a distribution on the real line ( $\Omega$  is  $\mathbb{R}$ ) and corresponds to a generative process that flips a biased coin and then generates a number from one of two Gaussian distributions, both with standard deviation 1.0 but with mean either 0.0 or 4.0 depending on the result of the coin toss. In this example, we will be more likely to produce a value near 0.0 than near 4.0 because of the bias. The probability  $\mathbb{P}(A)$  for  $A = [0, 1]$ , for instance, is the proportion of runs that generate a number between 0 and 1.

**Densities.** A distribution  $\mathbb{P}$  is a function that takes subsets of  $\Omega$  as input, but for many purposes it turns out to be more convenient if we can find a function  $f$  that takes elements of  $\Omega$  *directly*, while still somehow capturing the same information as  $\mathbb{P}$ .

When  $\Omega$  is the real line, we are interested in a function  $f$  that satisfies  $\mathbb{P}(A) = \int_A f(x) dx$  for all intervals  $A$ , and we call  $f$  the *probability density function* (PDF) of the distribution  $\mathbb{P}$ . In other words,  $f$  is a function where the area under its curve on an interval  $A$  gives the probability  $\mathbb{P}(A)$  of generating an outcome falling in that interval. The PDF of our example (pictured below) is given by the function

$$f(x) = 0.7 \cdot \text{pdf.Gaussian}(0.0, 1.0, x) + 0.3 \cdot \text{pdf.Gaussian}(4.0, 1.0, x)$$

where `pdf.Gaussian(mean, sdev, ·)` is the PDF of a Gaussian distribution with mean `mean` and standard deviation `sdev` (the famous “bell curve” from statistics).



The PDF takes higher values where the generative process described above is more likely to generate an outcome. Now we see that the aforementioned probability  $\mathbb{P}(A)$  for  $A = [0, 1]$  is simply the area under this curve between 0 and 1. Note that, while there are some loose similarities, the expression

for the PDF is different from the expression comprising the source program. In more complicated programs, the correspondence with the PDF is even less obvious.

**Non-existence.** Sometimes a distribution does not have a PDF. For example, if we change the else-clause in our example to return 4.0 directly, instead of drawing from a Gaussian with mean 4.0, we get the following probabilistic program, which does not have a PDF:

```
if flip 0.7 then random(Gaussian(0.0, 1.0)) else 4.0
```

In short, the problem is that there is a non-zero amount of probability mass located on a zero-width interval (the process now returns 4.0 with probability 0.3), but integrals on such intervals yield zero, so we would never find a function that could satisfy the properties of being a PDF.

It is not always obvious which program modifications ruin the property of having a PDF, especially for multivariate distributions (thus far we have only given univariate distributions as examples). This can be a problem if one is innocently exploring different variations of the model. Details and examples are given by Bhat et al. (2012), who provide the theory for addressing this problem, which we extend and implement in this work.

**The task of data analysis.** So far we have detailed *what* PDFs are, but not *why* we want them. We motivate the desire by explaining one popular use-case of PDFs that arises when applying Bayesian learning to data analysis.

In the previous examples, the programs specify fully-known probability distributions. While there is indeed randomness in the probabilistic behavior of the samples they generate, the nature of this uncertainty is entirely known—we know the means and variances of the Gaussians, as well as the bias between the two, thus we can characterize the random behavior.

In real-world analysis tasks, we rarely have this luxury. Instead, we are often in the position of trying to figure out what the parameters should be, given the data we see. Thus, applications typically deal with *parameterized* models (families of distributions indexed by a parameter), and they try to learn something about which distributions in that family best explain the observed data. For example, in the following, `moG` is a parameterized model, indexed by the parameter `(mA,mB)`:

```
let moG (mA,mB) =
  if flip 0.7 then random(Gaussian(mA, 1.0)) else random(Gaussian(mB, 1.0))
```

It specifies an infinite number of distributions, one for each choice of `mA` and `mB`.

Now, as a data analyst, we may be presented a dataset that is a sequence of numeric values, and we may also have some domain-specific reason to believe that it can be well modelled as a biased mixture of Gaussians as specified by `moG`<sup>1</sup>. We now face the task of figuring out which choices of `mA` and `mB` are likely. Intuitively, if we see a clump of datapoints around 0.0, we might be inclined to believe that the mean of one of the Gaussians is 0.0. Note that this is a one-dimensional, probabilistic version of the problem of *clustering*. The means are the cluster centroids.

**Bayesian learning.** Bayesian modelling formalizes this task by requiring the modeller to provide as input a *prior distribution* over the parameters and a *generating distribution* over the data. These are used to construct a *posterior distribution* over the parameters as the output.

---

<sup>1</sup>Whether this is actually an appropriate modelling choice depends on whether it captures enough of the essence of the true, unknown data-generating process (*i.e.* Nature), and is an entirely separate discussion.

The prior distribution is a distribution over the possible values of the parameters and captures our belief about what values they are likely to take, before having seen any data (to a Bayesian, “belief” is synonymous with “probability distribution”). The following is one possible prior:

```
let prior () =
  let mA = random(Uniform(-10.0, 10.0)) in
  let mB = random(Uniform(-10.0, 10.0)) in
  (mA, mB)
```

This specifies that we are certain that the means lie between -10.0 and 10.0, but are otherwise uncertain, and our uncertainty about each mean is uniformly distributed between -10.0 and 10.0. This is of course a very particular assertion and is hopefully informed by domain knowledge. In practice, in the absence of domain knowledge, we can select a prior that reflects our open-mindedness about which values the means can take, such as a pair of Gaussians with a high standard deviation. The prior *produces* a distribution over pairs of means as output, rather than taking a pair of means as an input, as `moG` does.

The generating distribution is a model of how we believe Nature is generating a dataset given a specific choice of the parameter. In our example this is given by `moG` together with

```
let gen n (mA,mB) = [| for i in 1 .. n → moG (mA,mB) |]
```

This specifies a (parameterized) model for a dataset that is generated as an array of `n` independent and identically distributed (i.i.d.) values generated by `moG`.

The posterior distribution is a distribution over the possible values of the parameters and captures our belief about what values they are likely to take, *after* seeing the data. The posterior is related to the prior and generating distributions by Bayes’ rule, which gives us a way to describe how the prior is updated with the observed data to yield the posterior. Intuitively, this update represents the fact that our understanding about the world (our belief about the parameters) evolves based on the data we see. Bishop provides an excellent account of Bayesian learning (Bishop, 2006).

**Use-case of PDFs: Bayesian inference with MCMC.** Unfortunately, while prior distributions and generating distributions are often straightforward to work with (we have control over them as the modeller), the posterior distributions often end up intractable or unwieldy to work with (Bayes’ rule dictates their form).

*Markov chain Monte Carlo* (MCMC) methods are one class of techniques that let us actually do something productive with the posterior distribution—MCMC can be used to generate samples from the posterior distribution. The idea of MCMC is to construct a Markov chain in the parameter space of the model, whose equilibrium distribution is the posterior distribution over model parameters. Neal (1993) gives an excellent review of MCMC methods. Here we use Filzbach (Purves and Lyutsarev, 2012), an adaptive MCMC sampler based on the Metropolis-Hastings algorithm. All that is required for such algorithms is the ability to calculate the posterior density given a set of parameters, up to proportion. The posterior does not need to be from a mathematically convenient family of distributions. Samples from the posterior can then serve as its representation, or be used to calculate marginal distributions of parameters or other integrals under the posterior distribution.

The posterior density is a function of the PDFs of the various pieces of the model, so to perform inference using MCMC, we also need functions to compute the PDFs. Below, `pdf_moG` gives the PDF of a single data point, while `pdf_gen` gives the PDF of an array of independent data points drawn from the same distribution (iid).

```
let pdf_prior (mA,mB) = pdf_Uniform(-10.0, 10.0, mA) * pdf_Uniform(-10.0, 10.0, mB)
let pdf_moG (mA,mB) x = 0.7 * pdf_Gaussian(mA, 1.0, x) + 0.3 * pdf_Gaussian(mB, 1.0, x)
let pdf_gen (mA,mB) xs = product [| for x in xs → pdf_moG (mA,mB) x |]
```

(The `product` function multiplies together the elements of an array, returning 1.0 on the empty array.) Filzbach and other MCMC libraries require users to write these three functions<sup>2</sup>, in addition to the generative probabilistic functions `prior` and `gen` (which are used for model validation).

The goal of this paper is to instead compile these density functions from the generative code. This relieves domain experts from having to write the density code in the first place, as well as from the error-prone task of manually keeping their model code and their density code in synch. Instead, both the PDF and synthetic data are derived from the same declarative specification of the model.

**Contributions of this paper.** This work defines and applies automated techniques for computing densities to real inference problems from various scientific applications. The primary technical contribution is a *density compiler* that is correct, useful, and relatively simple and efficient. Specifically:

- We provide the first implementation of a density compiler based on the specification by Bhat et al. (2012). We compile programs in the probabilistic language Fun (described in Section 2.1) to their corresponding density functions (Section 3).
- We prove that the compilation algorithm is sound (Theorem 3.15). This is the first such proof for any variant of this compiler.
- We show that the compiler greatly reduces the development effort of domain experts by freeing them from writing tricky density code and that the produced code is comparable in performance to density functions hand-coded by experts. Our evaluation is based on textbook examples and on models from ecology (Section 4).

## 2. LANGUAGES

In order to describe the density compiler, we first specify its input (source) and output (target) language. Both languages are variants of a simple first-order functional language where the results of subcomputations can be bound to variables using a `let` construct.

**2.1. Fun: Probabilistic Expressions (Review).** Our source language is a version of the core calculus Fun (Borgström et al., 2011), without observation. To mark certain program points as impossible, we add a `fail` construct (Kiselyov and Shan, 2009). Fun is a first-order functional language without recursion that extends the language of Ramsey and Pfeffer (2002), and this version has a natural semantics in the sub-probability monad. Our implementation efficiently supports a richer language with records and fixed-size arrays and array comprehensions, which can be given a semantics in this core (records and arrays can be encoded as tuples, and comprehensions of fixed size as their unrolling).

<sup>2</sup>The actual implementation works with log-densities, as discussed in Section 4.

2.1.1. *Syntax and Types of Fun*: The language Fun has base types **int**, **real** and **unit**, product types (denoting pairs), and sum types (denoting tagged unions). A type is said to be *discrete* if it does not contain **real**. We let  $c$  range over constant data of base type,  $n$  over integers and  $r$  over real numbers. We write  $\text{ty}(c) = t$  to mean that constant  $c$  has type  $t$ .

### Types of Fun:

$$t, u ::= \mathbf{int} \mid \mathbf{real} \mid \mathbf{unit} \mid (t_1 * t_2) \mid (t_1 + t_2)$$

We take  $\mathbf{bool} \triangleq \mathbf{unit} + \mathbf{unit}$ , and let  $*$  associate to the right. We assume a collection of total deterministic functions on these types, including arithmetic and logical operators. Each operation  $\text{op}$  of arity  $n$  has a signature of the form  $\mathbf{val} \text{ op}: t_1 * \dots * t_n \rightarrow t_{n+1}$ . We also assume standard families of primitive probability distributions, including the following.

**Distributions:**  $\text{Dist}: (t_1 * \dots * t_n) \rightarrow t$

**Bernoulli**: (**real**)  $\rightarrow$  **bool**  
**Poisson**: (**real**)  $\rightarrow$  **int**  
**Gaussian**: (**real** \* **real**)  $\rightarrow$  **real**  
**Beta**: (**real** \* **real**)  $\rightarrow$  **real**  
**Gamma**: (**real** \* **real**)  $\rightarrow$  **real**

Above, the names  $x_i$  of the arguments to the distributions are present for documentation only. A **Bernoulli**(*bias*) distribution corresponds to a coin flip with probability *bias* to come up **true**. The **Poisson**(*rate*) distribution describes the number of occurrences of independent events that occur at the given average *rate*. The **Gaussian**(*mean, stdev*) distribution is also known as the normal distribution; its PDF has a symmetrical bell shape. The **Beta**( $a, b$ ) distribution is a suitable prior for the parameter of **Bernoulli** distributions, and intuitively means that  $a - 1$  counts of **true** and  $b - 1$  events of **false** have been observed. Similarly the **Gamma**(*shape, scale*) distribution is a suitable prior for the parameter of **Poisson**. The parameters of distributions only make sense within certain ranges (e.g., the bias of the **Bernoulli** distribution must be in the interval  $[0, 1]$ ). Outside these ranges, attempting to draw a value from the distribution (e.g., **Bernoulli**(2.0)) results in a failure (**fail** below).

### Expressions of Fun:

$V ::=$	value
$x$	variable
$c$	scalar constant
$(V, V)$	tuple constructor
$\mathbf{inl}_u V$	left sum constructor
$\mathbf{inr}_t V$	right sum constructor
$M, N ::=$	expression
$x \mid c$	variable and scalar constant
$(M, N) \mid \mathbf{fst} M \mid \mathbf{snd} M$	pairing and projections from a pair
$\mathbf{inl}_u M \mid \mathbf{inr}_t M$	sum constructors
$\mathbf{match} M \mathbf{with} \mathbf{inl} x_1 \rightarrow N_1 \mid \mathbf{inr} x_2 \rightarrow N_2$	matching (scope of $x_i$ is $N_i$ )
$\mathbf{let} x = M \mathbf{in} N$	let (scope of $x$ is $N$ )
$\text{op}(M)$	primitive operation (deterministic)
$\mathbf{random}(\text{Dist}(M))$	primitive distribution
$\mathbf{fail}_t$	failure

The **let** and **match** statements bind their variables  $(x, x_1, x_2)$ ; we identify expressions up to alpha-renaming of bound variables. Above, **inl** (resp. **inr**) generates a value corresponding to the left (resp. right) branch of a sum type. Values of sum type are deconstructed by the **match** construct, which behaves as either the left ( $N_1$ ) or the right ( $N_2$ ) branch depending on the result of  $M$ .

To ensure that a program has at most one type in a given typing environment, **inl** and **inr** are annotated with a type (see (FUN INL) below). The expression **fail** is annotated with the type it is used at. These types are included only for the convenience of our technical development, and can usually be inferred given a typable source program: we omit these types where they are not used.  $()$  is the **unit** constant.

A source language term  $M$  is pure, written “ $M$  pure”, iff  $M$  does not contain any occurrence of **random** or **fail**.

We write **Uniform** for **Beta(1.0,1.0)**. In the binders of **let** and **match** expressions, we let  $\_$  stand for a variable that does not appear free in the scope of the binder. We make use of standard sugar for **let**, such as writing  $M;N$  for **let**  $\_ = M$  **in**  $N$ . We write **if**  $M$  **then**  $N_1$  **else**  $N_2$  for **match**  $M$  **with inl**  $\_ \rightarrow N_1$  | **inr**  $\_ \rightarrow N_2$ ; this is most commonly used when  $M$  is Boolean. We let the tuple  $(M_1, M_2, \dots, M_n)$  stand for  $(M_1, (M_2, \dots, M_n))$ . Similarly, we write **let**  $x_1, x_2, \dots, x_n = V$  **in**  $N$  for **let**  $x_1 = \mathbf{fst} V$  **in let**  $z = \mathbf{snd} V$  **in let**  $x_2, \dots, x_n = z$  **in**  $N$  when  $z \not\# N$ .

When  $X$  is a term from some language (possibly with binders), we write  $x_1, \dots, x_n \not\# X$  if none of the  $x_i$  appear free in  $X$ .

We write  $\Gamma \vdash M : t$  to mean that in the type environment  $\Gamma = x_1 : t_1, \dots, x_n : t_n$  ( $x_i$  distinct) the expression  $M$  has type  $t$ . Apart from the following, the typing rules are standard. In (FUN INL), (FUN INR) (not shown) and (FUN FAIL), type annotations are used in order to obtain a unique type. In (FUN RANDOM), a random variable drawn from a distribution of type  $(x_1 : t_1 * \dots * x_n : t_n) \rightarrow t$  has type  $t$ .

#### Selected Typing Rules: $\Gamma \vdash M : t$

$\frac{\Gamma \vdash M : t}{\Gamma \vdash \mathbf{inl}_u M : t + u}$	$\frac{}{\Gamma \vdash \mathbf{fail}_t : t}$	$\frac{\text{Dist} : (t_1 * \dots * t_n) \rightarrow t \quad \Gamma \vdash M : (t_1 * \dots * t_n)}{\Gamma \vdash \mathbf{random}(\text{Dist}(M)) : t}$
--	--	---

Substitutions, ranged over by  $\sigma, \rho$ , are finite maps  $[x_1 \mapsto M_1, \dots, x_n \mapsto M_n]$  from variables to pure expressions. We write  $M\sigma$  for the result of substituting all free occurrences of variables  $x \in \text{dom}(\sigma)$  in  $M$  with  $\sigma(x)$ , avoiding capture of bound variables. To compose two substitutions with disjoint domains, we write  $[x_1 \mapsto M_1, \dots, x_n \mapsto M_n]\sigma$  for  $[x_1 \mapsto M_1\sigma, \dots, x_n \mapsto M_n\sigma] \cup \sigma$ . A substitution is called *closed* if the expressions in its range do not contain any free variables. A value substitution is a substitution where each expression in its range is a value. Below, we define what it means for a closed value substitution to be a valuation for a type environment.

#### Typing Rules for Closed Value Substitutions: $\Gamma \vdash \sigma$

$\frac{}{\varepsilon \vdash []}$	$\frac{\Gamma \vdash \sigma \quad \varepsilon \vdash V : t}{\Gamma, x : t \vdash \sigma[x \mapsto V]}$
----------------------------------	--

There is a default value at each type  $t$ , written  $0_t$ , that is returned from operations  $\text{op}$  where they otherwise would be undefined, e.g.  $r/0.0 = 0_{\text{real}} = \log(-1)$ .

**Default Value:**  $0_t$

$$\boxed{0_{\text{unit}} := () \quad 0_{\text{int}} := 0 \quad 0_{\text{real}} := 0.0 \quad 0_{t*u} := (0_t, 0_u) \quad 0_{t+u} := \text{inl } 0_t}$$

**2.1.2. Semantics of Fun.** As usual, for precision concerning probabilities over uncountable sets, we turn to measure theory. The interpretation of a type  $t$  is the set  $\mathbf{V}_t$  of closed values of type  $t$  (real numbers, integers etc.). Below we consider only Lebesgue-measurable sets of values, defined using the standard (Euclidian) metric, and ranged over by  $A, B$ . Indeed, the power of the axiom of choice is needed to construct a non-measurable set (Solovay, 1970).

A measure  $\mu$  over  $t$  is a function, from (measurable) subsets of  $\mathbf{V}_t$  to the non-negative real numbers extended with  $\infty$ , that is  $\sigma$ -additive, that is,  $\mu(\emptyset) = 0.0$  and  $\mu(\cup_i A_i) = \sum_i \mu(A_i)$  if  $A_1, A_2, \dots$  are pair-wise disjoint. We write  $|\mu|$  for  $\mu(\mathbf{V}_t)$ ; the measure  $\mu$  is called a probability measure if  $|\mu| = 1.0$ , and a sub-probability measure if  $|\mu| \leq 1.0$ .

We associate a default or *stock* measure to each type, inductively defined as the counting measure on  $\mathbb{Z}$  and  $\{\emptyset\}$ , the Lebesgue measure on  $\mathbb{R}$ , and the Lebesgue-completion of the product and disjoint sum, respectively, of the two measures for  $t * u$  and  $t + u$ . In particular, the counting measure on a discrete type assigns measure  $k$  to all sets of finite size  $k$ , and measure  $\infty$  to all infinite sets.

If  $f$  is a non-negative (measurable) function  $t \rightarrow \text{real}$ , we let  $\int_t f$  be the Lebesgue integral of  $f$  with respect to the stock measure on  $t$  if the integral is defined, and otherwise 0. This integral coincides with  $\sum_{x \in \mathbf{V}_t} f(x)$  for discrete types  $t$ , and with the standard Riemann integral (if it is defined) on  $t = \text{real}$ . We write  $\int_t f(x) dx$  for  $\int_t \lambda x. f(x)$ , and  $\int_t f(x) d\mu(x)$  for Lebesgue integration with respect to the measure  $\mu$  on  $t$ . Below, we often elide the index  $t$ ; indeed, we may consider any function  $t \rightarrow \text{real}$  as a function from the measurable space  $\uplus_u \mathbf{V}_u$  that is zero except on  $\mathbf{V}_t$ .

The Iverson brackets  $[p]$  are 1.0 if predicate  $p$  is true, and 0.0 otherwise. We write  $\int_A f$  for  $\int \lambda x. [x \in A] \cdot f(x)$  when  $A \subset \mathbf{V}_t$ . The function  $g$  is a *density* of  $\mu$  (with respect to the stock measure) if  $\int_A 1 d\mu(x) = \int_A g$  for all  $A$ . If  $\mu$  is a (sub-)probability measure, then we say that  $g$  as above is its PDF.

To turn expressions into density functions, we first need a way of interpreting a closed Fun expression  $M$  as a sub-probability measure  $\mathbb{P}_M$  over its return type. Open **fail**-free Fun expressions have a straightforward semantics (Ramsey and Pfeffer, 2002) as computations in the probability monad (Giry, 1982). In order to treat the **fail** primitive, we use an existing extension (Gordon et al., 2013) of the semantics of Ramsey and Pfeffer (2002) to a richer monad: the sub-probability monad (Panangaden, 1999)<sup>3</sup>. Compared to the operations of the probability monad, the sub-probability monad additionally admits a zero constant, yielding the zero measure. To accommodate the zero measure, the carrier set is extended from probability measures to sub-probability measures, i.e., admitting all  $\mu$  with  $|\mu| \leq 1$ .

Below we recapitulate the semantics of Fun by Gordon et al. (2013). Here  $\sigma$  is a closed value substitution whose domain contains all the free variables of  $M$ , and  $\text{detOp}(M)$  ranges over  $\text{op}(M)$ , **fst**  $M$ , **snd**  $M$ , **inl**  $M$  and **inr**  $M$ . We also let either  $f g (\text{inl } V) \triangleq f V$  and either  $f g (\text{inr } V) \triangleq g V$ .

<sup>3</sup>Sub-probabilities are also used in our compilation of **match** (and **if**) statements, where the probability that we have entered a particular branch may be less than 1.



**Monadic Semantics of Fun with fail:**  $\mathcal{P}[[M]] \sigma$ 

$(\mu \gg= f) A \triangleq \int f(x)(A) d\mu(x)$	Sub-probability monad's bind
$(\text{return } V) A \triangleq 1 \text{ if } V \in A, \text{ else } 0$	Sub-probability monad's return
$\text{zero } A \triangleq 0$	Sub-probability monad's zero

Below we assume that  $z \# N, N_1, N_2, \sigma$  and  $x, x_1, x_2 \# z, \sigma$ .

$\mathcal{P}[[x]] \sigma \triangleq \text{return } \sigma(x)$
$\mathcal{P}[[c]] \sigma \triangleq \text{return } c$
$\mathcal{P}[[\text{detOp}(M)]] \sigma \triangleq \mathcal{P}[[M]] \sigma \gg= \lambda x. \text{return detOp}(x)$
$\mathcal{P}[[ (M, N) ] ] \sigma \triangleq \mathcal{P}[[M]] \sigma \gg= \lambda z. \mathcal{P}[[N]] \sigma \gg= \lambda w. \text{return } (z, w)$
$\mathcal{P}[[\text{let } x = M \text{ in } N]] \sigma \triangleq \mathcal{P}[[M]] \sigma \gg= \lambda z. \mathcal{P}[[N]] (\sigma[x \mapsto z])$
$\mathcal{P}[[\text{match } M \text{ with inl } x_1 \rightarrow N_1 \mid \text{inr } x_2 \rightarrow N_2]] \sigma \triangleq$ $\mathcal{P}[[M]] \sigma \gg= \text{either } (\lambda z. \mathcal{P}[[N_1]] (\sigma[x_1 \mapsto z])) (\lambda z. \mathcal{P}[[N_2]] (\sigma[x_2 \mapsto z]))$
$\mathcal{P}[[\text{random}(\text{Dist}(M))]] \sigma \triangleq \mathcal{P}[[M]] \sigma \gg= \lambda z. \mu_{\text{Dist}(z)}$
$\mathcal{P}[[\text{fail}]] \sigma \triangleq \text{zero}$

We let the semantics of a closed expression  $M$  be  $\mathbb{P}_M \triangleq \mathcal{P}[[M]] \varepsilon$ , where  $\varepsilon$  denotes the empty substitution.

**Lemma 2.1.** *If  $\Gamma \vdash M : t$  and  $\Gamma \vdash \sigma$  then  $\mathcal{P}[[M]] \sigma$  is a sub-probability measure on type  $t$ .*

*Proof.* By induction on  $M$ . □

**2.2. Target Language for Density Computations.** For the target language of the density compiler, denoted  $\text{fun}$ , we use a pure version of Fun augmented with real-valued first-order functions and stock integration.

**Expressions of the Target Language:**  $E, F$ 

$E, F ::=$	target expression
$x \mid c$	variable and scalar constant
$(E, F) \mid \text{fst } E \mid \text{snd } E$	pairing and projections from a pair
$\text{inl}_u E \mid \text{inr}_t E$	sum constructors
$\text{match } E \text{ with inl } x_1 \rightarrow F_1 \mid \text{inr } x_2 \rightarrow F_2$	matching (scope of $x_i$ is $F_i$ )
$\text{let } x = E \text{ in } F$	let (scope of $x$ is $F$ )
$\text{op}(E)$	primitive operation
$\int_t \lambda(x_1, \dots, x_n). E$	stock integration

Above, the binders in **let** and **match** are as in Fun. Additionally, in  $\int_t \lambda(x_1, \dots, x_n). E$  the variables  $x_1, \dots, x_n$  bind into  $E$ .

If a Fun term  $M$  is pure then  $M$  is also an expression in the syntax of  $\text{fun}$ , and we silently treat it as such. In particular, a Fun substitution  $\sigma$  is also a valid  $\text{fun}$  substitution, and substitution application  $E\sigma$  for  $\text{fun}$  is defined in the same way as for Fun.

The typing rule involving integration is as follows. The other typing rules are as in Fun.

**Typing Rule for Integration:**  $\Gamma \vdash E : t$

(TARGET INT) $\frac{\Gamma, x_1 : t_1, \dots, x_n : t_n \vdash E : \mathbf{real}}{\Gamma \vdash \int_{t_1 * \dots * t_n} \lambda(x_1, \dots, x_n). E : \mathbf{real}}$
---

**Lemma 2.2** (Standard results for the type system of fun).

- (1) Substitution lemma: *if  $\Gamma, x : t, \Gamma' \vdash E : u$  and  $\Gamma \vdash F : t$ , then  $\Gamma, \Gamma' \vdash E[x \mapsto F] : u$ .*
- (2) Strengthening: *if  $\Gamma, x : t, \Gamma' \vdash E : u$  and  $x \# E$ , then  $\Gamma, \Gamma' \vdash E : u$ .*
- (3) Weakening: *if  $\Gamma, \Gamma' \vdash E : u$  and  $x \# \Gamma, \Gamma'$ , then  $\Gamma, x : t, \Gamma' \vdash E : u$ .*

2.2.1. *Semantics of fun.* The target language fun is equipped with a denotational semantics, written  $\mathcal{S}[[F]] \sigma$  where  $\sigma$  is a substitution of closed values for variables with  $\text{dom}(\sigma) \supseteq \text{fv}(F)$ . We define this semantics by re-interpreting the monadic semantics of Subsection 2.1.2 with respect to the identity monad: in this monad, `return` is the identity function, and `bind` ordinary function application. We rely on an auxiliary semantics  $\mathcal{S}[[\lambda(z_1, \dots, z_n). E]] \sigma$  that returns a function to be integrated.

**Identity Monad and Denotational Semantics of fun:**  $\mathcal{S}[[\lambda(z_1, \dots, z_n). E]] \sigma$  and  $\mathcal{S}[[E]] \sigma$

$V \gg= f \triangleq f(V)$	Identity monad's bind
$(\text{return } V) \triangleq V$	Identity monad's return
$\mathcal{S}[[\int_t \lambda(z_1, \dots, z_n). F]] \sigma \triangleq \begin{cases} \int_t \mathcal{S}[[\lambda(z_1, \dots, z_n). F]] \sigma & \text{if the integral is well-defined} \\ 0 & \text{otherwise} \end{cases}$	
(the other cases of $\mathcal{S}[[E]] \sigma$ are the same as the monadic semantics in Subsection 2.1.2)	
$\mathcal{S}[[\lambda(z_1, \dots, z_n). F]] \sigma \triangleq f$	
where $f(V_1, \dots, V_n) \triangleq \mathcal{S}[[F]] \sigma[z_1 := V_1, \dots, z_n := V_n]$ and $z_1, \dots, z_n \# \sigma$	

We write  $E \equiv F$  if there are  $\Gamma, t$  such that  $\Gamma \vdash E : t$  and  $\Gamma \vdash F : t$  and for all  $\sigma$  such that  $\Gamma \vdash \sigma$  we have  $\mathcal{S}[[E]] \sigma = \mathcal{S}[[F]] \sigma$ .

**Lemma 2.3.**

- (1) *If  $\Gamma \vdash E : t$  and  $\Gamma \vdash \sigma$  then  $\mathcal{S}[[E]] \sigma$  is a value of type  $t$ .*
- (2) *If  $\Gamma, x_1 : t_1, \dots, x_n : t_n \vdash E : \mathbf{real}$  and  $\Gamma \vdash \sigma$  and  $\text{dom}(\sigma) \supseteq \text{fv}(E)$  then  $\mathcal{S}[[\lambda(z_1, \dots, z_n). E]] \sigma$  is a function of type  $t_1 * \dots * t_n \rightarrow \mathbf{real}$ .*

*Proof.* (1) and (2) are proved jointly, by induction on  $E$ . □

### 3. THE DENSITY COMPILER

We compute the PDF of a Fun program by compilation into fun. Our compilation is based on that of Bhat et al. (2012), with modifications to treat **fail** statements, **match** (and general **if**) statements, pure (i.e., deterministic) **let** bindings, and integer arithmetic.

The compiler translates a well-typed Fun source term  $M$  into a function  $\lambda z. F$  computing the density (PDF) of  $M$ . Given an implementation of stock integration, the fun expression  $F[z \mapsto V]$  may be executed to evaluate the density of  $M$  at any value  $V$  of the type of  $M$ . Like traditional compilers, our compiler is compositional and deterministic, producing a unique translation if any

at all (Lemma 3.13). Unlike traditional compilers, our compiler is partial and will fail to produce a translation for some well-typed source terms. In particular, if  $M$  does not have a density function then the compiler will fail to produce an  $F$ . However, it may also fail if  $M$  has a PDF, but the compiler is just not complete enough to compute it. In particular, **let**-bound expressions must either be pure or have a PDF, even if their result is not used. The correctness statement for the compiler is given by Theorem 3.15.

We will use a version of the **moG** function from the introduction as a running example (Figure 1), with some expansion in order to make use of more of the translation rules.

```

1   let branch = random(Bernoulli(0.7)) in
2   let temp = random(Gaussian(0.0, 1.0)) in
3   match branch with
4     inl _ → random(Gaussian(mA, 1.0))
5     inr _ →
6       let result = temp + mB in
7       result

```

Figure 1: Expanded model for a mixture of two Gaussians

The structure of this section is as follows. In Section 3.1 we provide an intuitive outline of the compilation. We make preliminary definitions, such as the syntax of probability contexts  $\Upsilon$ , in Section 3.2. We define the compiler itself in Section 3.3 in terms of a couple of judgments. These judgments are inductively defined relations, but they in fact are partial functions and hence have a direct executable interpretation. Finally, in Section 3.4 we state and prove correctness of the compiler.

**3.1. Outline.** The simplest case in the density compilation is **fail**, which compiles to the function that always returns zero. The compilation works on the **let**-structure of the term: a sequence of random **lets**, as in **let**  $x_1 = \mathbf{random}(\text{Dist}_1(V_1))$  **in** ... **in**  $(x_1, \dots, x_n)$  is compiled to the product of the PDFs of the distributions  $\text{Dist}_1, \dots, \text{Dist}_n$ , following the chain rule of probability.

If the sequence of **lets** instead has a discrete deterministic return expression  $M$ , then  $M$  has a probability for each possible value  $V$ . This probability is computed by integrating the joint PDF of  $x_1, \dots, x_n$  over the set of values where  $M$  evaluates to  $V$ . A continuous deterministic return expression  $M$  is treated as a mathematical function  $f_M$ , using the change of variables rule of integration. In the one-dimensional case, if  $f_M(x)$  has inverse  $f^{-1}$ , the PDF of  $f_M$  at  $r$  is given by the PDF of  $x$  at  $f^{-1}(r)$ , multiplied with the derivative of  $f^{-1}(r)$ . Another simple case is projection  $M = x_n$ , where we simply integrate the joint PDF over the set of all possible values for the other variables  $x_1, \dots, x_{n-1}$ .

If a distribution  $\text{Dist}$  returns a sum type (e.g., **Bernoulli**) we can write  $\text{Dist}_l$  for the subdistribution yielding only the left part of the sum, and  $\text{Dist}_r$  for the right part. By additivity of probability, we can compile the **match** expression **match**  $\mathbf{random}(\text{Dist}(V))$  **with** **inl**  $y \rightarrow M$  **|** **inr**  $z \rightarrow N$  to the sum of the PDFs of **let**  $y = \mathbf{random}(\text{Dist}_l(V))$  **in**  $M$  and **let**  $z = \mathbf{random}(\text{Dist}_r(V))$  **in**  $N$ .

In a nested **let**, such as **let**  $x_i = (\mathbf{let} y_1 = \mathbf{random}(\text{Dist}(V))$  **in** ... **in**  $M_i)$  **in** ..., the expression bound to  $x_i$  denotes some subprobability distribution. We compute its PDF by recursively compiling the inner **let** sequence, holding  $x_1, \dots, x_{i-1}$  fixed. Pure **lets**, as in **let**  $x = M$  **in**  $N$  where  $M$  is pure, have the same PDF as  $N[x \mapsto M]$ . The compilation algorithm applies the substitution lazily to avoid introducing unnecessary copies of  $M$ .

**3.2. Probability contexts.** The density compilation is based on the let-structure of the expression. Variables that are bound in outer lets are referred to as parameters, and are treated as constants. A *probability context* gathers the variables that are bound in the current sequence of lets, together with the pure expressions defining the deterministic variables.

**Probability Context:  $\Upsilon$**

$\Upsilon ::=$	probability context
$\varepsilon$	empty context
$\Upsilon, x$	random variable
$\Upsilon, x = E$	deterministic variable

**Example 3.1.** The probability context at line 7 of Figure 1 is  $\Upsilon_7 := \text{branch}, \text{temp}, \text{result} = \text{temp} + \text{mB}$ , containing two random variables and one deterministic variable.

For a probability context to be well-formed, it has to be well-scoped and well-typed.

**Well-formed probability context:  $\Gamma \vdash \Upsilon \text{ wf}$**

(ENV EMPTY)	(ENV VAR)	(ENV CONST)
$\Gamma \vdash \varepsilon \text{ wf}$	$\frac{\Gamma \vdash \Upsilon \text{ wf} \quad \Gamma \vdash x : T \quad x \# \Upsilon}{\Gamma \vdash \Upsilon, x \text{ wf}}$	$\frac{\Gamma \vdash \Upsilon \text{ wf} \quad \Gamma \vdash x : T \quad x \# \Upsilon \quad \Gamma \vdash E : T}{\Gamma \vdash \Upsilon, x = E \text{ wf}}$

**Example 3.2.** The probability context  $\Upsilon_7$  is well-formed:  $\Gamma_7 \vdash \Upsilon_7 \text{ wf}$ , where the type context  $\Gamma_7 := \text{mA}:\text{real}, \text{mB}:\text{real}, \text{branch}:\text{bool}, \text{temp}:\text{real}, \text{result}:\text{real}$  also contains the types of the parameters  $\text{mA}$  and  $\text{mB}$ .

Given a well-formed context  $\Upsilon$ , we can extract the random variables  $\text{rands}(\Upsilon)$ , and an idempotent substitution  $\sigma_\Upsilon$  (i.e.,  $E\sigma_\Upsilon = (E\sigma_\Upsilon)\sigma_\Upsilon$  always) that gives values to the deterministic variables.

**Random variables  $\text{rands}(\Upsilon)$  and values of deterministic variables  $\sigma_\Upsilon$**

$\text{rands}(\varepsilon) \triangleq \varepsilon$	$\sigma_\varepsilon \triangleq []$
$\text{rands}(\Upsilon, x) \triangleq \text{rands}(\Upsilon), x$	$\sigma_{\Upsilon, x} \triangleq \sigma_\Upsilon$
$\text{rands}(\Upsilon, x = E) \triangleq \text{rands}(\Upsilon)$	$\sigma_{\Upsilon, x = E} \triangleq [x \mapsto E]\sigma_\Upsilon$

**Example 3.3.** We have  $\text{rands}(\Upsilon_7) = \text{branch}, \text{temp}$  and  $\sigma_{\Upsilon_7} = [\text{result} \mapsto \text{temp} + \text{mB}]$ .

**Lemma 3.4.** *If  $\Gamma \vdash \Upsilon \text{ wf}$  then  $\text{dom}(\sigma_\Upsilon) \# \text{range}(\sigma_\Upsilon)$*

*Proof.* By simultaneous induction on the derivations of  $\Gamma \vdash \Upsilon \text{ wf}$  and  $\sigma_\Upsilon$ . □

**3.3. Compilation rules.** A probability context  $\Upsilon$  is used together with a density expression ( $E$  below), which is an open term that expresses the joint density of the random variables in the context and the constraints that have been collected when choosing branches in **match** statements. Intuitively, the density expression  $E$  is the body of the PDF of the current branch. The main judgment of the compiler is  $\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$ , which associates a Fun term  $M$  with its density function  $z \mapsto F$ . Parameters may occur free in  $F$ , and  $z$  binds into  $F$ . The auxiliary judgment  $\Upsilon; E \vdash \text{marg}(\{x_1, \dots, x_k\}) \Rightarrow F$  yields a density expression  $F$  for the variables in its argument, marginalizing (i.e., integrating) out all other random variables in  $\Upsilon$  from  $E$ .

**Inductively Defined Judgments of the Compiler:**

$\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$	in $\Upsilon; E$ the function $\lambda z. F$ gives the PDF of $M$
$\Upsilon; E \vdash \text{marg}(X) \Rightarrow F$	in $\Upsilon; E$ expression $F$ gives the density of the variables in $X$

We begin the description of the compiler proper with the following judgment of marginal density, which computes an expression for the joint marginal PDF of the random variables in its argument. The variables in the argument are free in the computed expression. Below, we write  $x_1, \dots, x_n \setminus Y$  for the tuple of variables arising from  $x_1, \dots, x_n$  by deleting all instances of variables in  $Y$ , and dually for  $x_1, \dots, x_n \cap Y$ .

**Marginal Density:**  $\Upsilon; E \vdash \text{marg}(X) \Rightarrow F$ 

(MARGINAL)
$X \subseteq \text{rands}(\Upsilon) \quad \text{rands}(\Upsilon) \setminus X = y_1, \dots, y_n$
$\Upsilon; E \vdash \text{marg}(X) \Rightarrow \int \lambda(y_1, \dots, y_n). E \sigma_\Upsilon$

Here we first substitute in the deterministic **let**-bound variables, as given by  $\sigma_\Upsilon$ , and then integrate out the remaining random variables  $y_1, \dots, y_k$ . In the definition of the compiler,  $\text{marg}(X)$  is also used with  $X = \emptyset$ , to compute the probability of being in the current branch of the program.

**Example 3.5.** Here  $\Upsilon_7; E \vdash \text{marg}(\{\text{temp}\}) \Rightarrow F_{\text{temp}}$  where

$$F_{\text{temp}} := \int \lambda(\text{branch}). E[\text{result} \mapsto \text{temp} + \text{mB}],$$

which will be used when computing the PDF of the variable **result** on line 7 (cf. Examples 3.6, 3.11).

The main judgment of the compiler is the  $\text{dens}$  judgment  $\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$ , which gives the density  $z \mapsto F$  of  $M$  in the current context  $\Upsilon$ , where  $E$  is the accumulated body of the density function so far. In this judgment,  $z$  is binding into  $F$ . We introduce fresh variables during the compilation: in the rules below we assume that  $z, w \# \Upsilon, E, M$ .

**Density Compiler, base cases:**  $\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$ 

(VAR DET)	(VAR RND)
$\frac{(x = E') \in \Upsilon \quad \Upsilon; E \vdash \text{dens}(E') \Rightarrow \lambda z. F}{\Upsilon; E \vdash \text{dens}(x) \Rightarrow \lambda z. F}$	$\frac{x \in \text{rands}(\Upsilon) \quad \Upsilon; E \vdash \text{marg}(\{x\}) \Rightarrow F}{\Upsilon; E \vdash \text{dens}(x) \Rightarrow \lambda x. F}$
(CONSTANT)	(FAIL)
$\frac{\varepsilon \vdash V : t \quad t \text{ discrete} \quad \Upsilon; E \vdash \text{marg}(\emptyset) \Rightarrow F}{\Upsilon; E \vdash \text{dens}(V) \Rightarrow \lambda z. F \cdot [z = V]}$	$\frac{}{\Upsilon; E \vdash \text{dens}(\text{fail}) \Rightarrow \lambda z. 0.0}$

For a deterministic variable, (VAR DET) recurses into its definition. The rule (VAR RND) computes the marginal density of a random variable using the  $\text{marg}$  judgment. The (CONSTANT) rule states that the probability density of a discrete constant  $V$  (built from sums and products of integers and units) is the probability of being in the current branch at  $V$ , and 0 elsewhere. Note the absence of a rule for real number constants, since they do not possess a density. The (FAIL) rule gives that the density of **fail** is zero.

**Example 3.6.** By (VAR RND) we get  $\Upsilon_7; E \vdash \text{dens}(\text{temp}) \Rightarrow \lambda \text{temp}. F_{\text{temp}}$  as computed in Example 3.5.

To compute the PDF at line 7, (VAR DET) yields that  $\Upsilon_7; E \vdash \text{dens}(\text{result}) \Rightarrow \lambda z. F_7$  where  $\Upsilon_7; E \vdash \text{dens}(\text{temp} + \text{mB}) \Rightarrow \lambda z. F_7$  is computed in Example 3.11.

**Density Compiler, random variables :**  $\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$

<p>(RANDOM CONST)</p> $\frac{M \text{ pure} \quad \text{rands}(\Upsilon) \# (M\sigma_\Upsilon) \quad \Upsilon; E \vdash \text{marg}(\emptyset) \Rightarrow F}{\Upsilon; E \vdash \text{dens}(\mathbf{random}(\text{Dist}(M))) \Rightarrow \lambda z. F \cdot \text{PDF}_{\text{Dist}(M\sigma_\Upsilon)}(z)}$
<p>(RANDOM RND)</p> $\frac{\neg(M \text{ pure} \wedge \text{rands}(\Upsilon) \# (M\sigma_\Upsilon)) \quad \Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\mathbf{random}(\text{Dist}(M))) \Rightarrow \lambda z. \int \lambda w. F \cdot \text{PDF}_{\text{Dist}(w)}(z)}$

In (RANDOM CONST), a random variable drawn from a primitive distribution with a constant argument has the expected PDF (multiplied with the probability that we are in the current branch). Its precondition that  $M \text{ pure}$  and  $\text{rands}(\Upsilon) \# (M\sigma_\Upsilon)$  intuitively means that  $M$  is constant under  $\Upsilon$ . (RANDOM RND) treats calls to **random** with a random argument by marginalizing over the argument to the distribution. We here require that each primitive distribution  $\text{Dist}$  has a PDF for each value  $w$  of its arguments, denoted  $\text{PDF}_{\text{Dist}(w)}$ .

**Example 3.7.** Using rule (RANDOM CONST), we can compute the density at line 4 as  $\Upsilon; E \vdash \text{dens}(\mathbf{random}(\mathbf{Gaussian}(\text{mA}, 1.0))) \Rightarrow \lambda z. F_{\text{then}} \cdot \text{PDF}_{\mathbf{Gaussian}(\text{mA}, 1.0)}(z)$  where  $F_{\text{then}}$  intuitively yields the probability of being in the current branch, and is computed using (MARGINAL) as  $\int \lambda(\text{temp}, \text{branch}). E$ .

**Density Compiler, rules for tuples:**  $\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$

<p>(TUPLE VAR)</p> $\frac{\Upsilon; E \vdash \text{marg}(\{x_1, \dots, x_k\}) \Rightarrow F \quad k \geq 2 \quad x_1, \dots, x_k \text{ distinct}}{\Upsilon; E \vdash \text{dens}((x_1, \dots, x_k)) \Rightarrow \lambda z. \mathbf{let} \ x_1, \dots, x_k = z \ \mathbf{in} \ F}$
<p>(TUPLE PROJ L)</p> $\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\mathbf{fst} \ M) \Rightarrow \lambda z_1. \int \lambda z_2. \mathbf{let} \ w = (z_1, z_2) \ \mathbf{in} \ F}$
<p>(TUPLE PROJ R)</p> $\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\mathbf{snd} \ M) \Rightarrow \lambda z_2. \int \lambda z_1. \mathbf{let} \ w = (z_1, z_2) \ \mathbf{in} \ F}$

The rule (TUPLE VAR) computes the joint marginal density of a tuple of variables<sup>4</sup>. The rules (TUPLE PROJ L) and (TUPLE PROJ R) integrate out the right or the left component of a pair, respectively.

**Density Compiler, let:**  $\Upsilon; E \vdash \text{dens}(\mathbf{let} \ x = M \ \mathbf{in} \ N) \Rightarrow F$

<sup>4</sup>Joint marginal densities for tuples of expressions can be computed if those expressions are conditionally independent (Bhat et al., 2012). As an example,  $(x, y + 3)$  has a PDF whenever  $(x, y)$  does. However, the rules in this paper do not support such expressions, to avoid additional complexity.

$$\begin{array}{c}
\text{(LET DET)} \\
\frac{M \text{ pure} \quad \Upsilon, x = M; E \vdash \text{dens}(N) \Rightarrow \lambda z. F}{\Upsilon; E \vdash \text{dens}(\mathbf{let} \ x = M \ \mathbf{in} \ N) \Rightarrow \lambda z. F} \\
\hline
\end{array}
\qquad
\begin{array}{c}
\text{(LET RND)} \\
\frac{\neg(M \text{ pure}) \quad \varepsilon; 1 \vdash \text{dens}(M) \Rightarrow \lambda x. F_1 \quad \Upsilon, x; E \cdot F_1 \vdash \text{dens}(N) \Rightarrow \lambda z. F_2}{\Upsilon; E \vdash \text{dens}(\mathbf{let} \ x = M \ \mathbf{in} \ N) \Rightarrow \lambda z. F_2} \\
\hline
\end{array}$$

The rule (LET DET) simply adds a pure let-binding to the context. In (LET RND), we compute the density of the let-bound variable in an empty context, and multiply it into the current accumulated density when computing the density of the body.

**Example 3.8.** The density expression for the program fragment in Figure 1 is computed using (LET RND) as  $F_2$  where  $\mathbf{branch}; 1 \cdot F_{\mathbf{branch}} \vdash \text{dens}(N) \Rightarrow \lambda z. F_2$ , the expression  $N$  is lines 2-7, and  $F_{\mathbf{branch}} = (\int \lambda().1) \cdot \text{PDF}_{\text{Bernoulli}(0.7)}(z)$  is computed by  $\varepsilon; 1 \vdash \text{dens}(\mathbf{random}(\text{Bernoulli}(0.7))) \Rightarrow \lambda z. F_{\mathbf{branch}}$  using previously seen rules. Here  $F_{\mathbf{branch}} \equiv \text{PDF}_{\text{Bernoulli}(0.7)}$ , since  $\mathcal{S}[\int \lambda().1] \sigma = 1$ .

The **let** expression on line 2 of Figure 1 is also handled by (LET RND), while the one on line 6 is handled by (LET DET) since  $\mathbf{temp} + \mathbf{mB}$  pure.

For deterministic matches we use four deterministic operations, which we assume do not occur in source programs. We let  $\mathbf{isL} : t + u \rightarrow \mathbf{real}$  be the indicator function for the left branch defined as  $\mathbf{isL}(V) := \text{match } V \text{ with } \mathbf{inl} \ x : 1.0 \mid \mathbf{inr} \ x : 0.0$ , and dually for  $\mathbf{isR}$ . To destruct a value of sum type we use  $\mathbf{fromL} : t + u \rightarrow t$  defined as  $\mathbf{fromL}(V) := \text{match } V \text{ with } \mathbf{inl} \ x : x \mid \mathbf{inr} \ x : 0_t$ , and its dual  $\mathbf{fromR}$ .

**Density Compiler, rules for sums and match:**  $\Upsilon; E \vdash \text{dens}(\mathbf{match} \ M \ \mathbf{with} \ \dots) \Rightarrow \lambda z. F$

$$\begin{array}{c}
\text{(SUM CON L)} \\
\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F}{\Upsilon; E \vdash \text{dens}(\mathbf{inl} \ M) \Rightarrow \lambda w. \mathbf{match} \ w \ \mathbf{with} \ \mathbf{inl} \ z \rightarrow F \mid \mathbf{inr} \ _ \rightarrow 0} \\
\text{(SUM CON R)} \\
\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda z. F}{\Upsilon; E \vdash \text{dens}(\mathbf{inr} \ M) \Rightarrow \lambda w. \mathbf{match} \ w \ \mathbf{with} \ \mathbf{inl} \ _ \rightarrow 0 \mid \mathbf{inr} \ z \rightarrow F} \\
\text{(MATCH DET)} \\
\frac{M \text{ pure} \quad \Upsilon, x_1 = \mathbf{fromL}(M); E \cdot \mathbf{isL}(M) \vdash \text{dens}(N_1) \Rightarrow \lambda z. F_1 \quad \Upsilon, x_2 = \mathbf{fromR}(M); E \cdot \mathbf{isR}(M) \vdash \text{dens}(N_2) \Rightarrow \lambda z. F_2}{\Upsilon; E \vdash \text{dens}(\mathbf{match} \ M \ \mathbf{with} \ \mathbf{inl} \ x_1 \rightarrow N_1 \mid \mathbf{inr} \ x_2 \rightarrow N_2) \Rightarrow \lambda z. F_1 + F_2} \\
\text{(MATCH RND)} \\
\frac{\neg(M \text{ pure}) \quad \varepsilon; 1 \vdash \text{dens}(M) \Rightarrow \lambda w. F \quad \Upsilon, x_1; E \cdot \mathbf{let} \ w = \mathbf{inl} \ x_1 \ \mathbf{in} \ F \vdash \text{dens}(N_1) \Rightarrow \lambda z. F_1 \quad \Upsilon, x_2; E \cdot \mathbf{let} \ w = \mathbf{inr} \ x_2 \ \mathbf{in} \ F \vdash \text{dens}(N_2) \Rightarrow \lambda z. F_2}{\Upsilon; E \vdash \text{dens}(\mathbf{match} \ M \ \mathbf{with} \ \mathbf{inl} \ x_1 \rightarrow N_1 \mid \mathbf{inr} \ x_2 \rightarrow N_2) \Rightarrow \lambda z. F_1 + F_2} \\
\text{(FROML)} \\
\frac{(x = \mathbf{fromL}(M)) \in \Upsilon \quad \Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\mathbf{fromL}(M)) \Rightarrow \lambda z. \mathbf{let} \ w = \mathbf{inl} \ z \ \mathbf{in} \ F} \\
\text{(FROMR)} \\
\frac{(x = \mathbf{fromR}(M)) \in \Upsilon \quad \Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\mathbf{fromR}(M)) \Rightarrow \lambda z. \mathbf{let} \ w = \mathbf{inr} \ z \ \mathbf{in} \ F} \\
\hline
\end{array}$$

(SUM CON L) states that the density of **inl**  $M$  is the density of  $M$  in the left branch of a sum, and 0 in the right. Its dual is (FROML). (MATCH DET) is based on (LET DET), and we additionally multiply the constraint that we are in the correct branch (**isL**( $M$ ) or **isR**( $M$ )) with the joint density expression. We employ the functions **fromL** and **fromR** and their associated rules (FROML) and (FROMR) to avoid additional calls to (MATCH DET) arising from (VAR DET) if the compilation of the density of  $N_i$  requires computing the density of the match-bound variable  $y_i$ , as in **match fst**  $z$  **with inl**  $y_1 \rightarrow y_1$  **| inr**  $y_2 \rightarrow y_2$ . Since we assume that **fromL** and **fromR** do not appear in source programs, these rules are only ever used in the case described above. The (MATCH RND) rule is based on (LET RND), and we again multiply in the constraint that we are in the left (or right) branch of the **match**.

**Example 3.9.** The **match** selector on line 3 is a pure expression, so rule (MATCH DET) applies. For the left branch, we let  $E_4 \equiv \text{PDF}_{\text{Bernoulli}(0.7)}(\text{branch}) \cdot \text{PDF}_{\text{Gaussian}(0,1.0)}(\text{temp}) \cdot \text{isL}(\text{branch})$  and  $\Upsilon_4 = \text{branch}, \text{temp}, \_ = \text{fromL}(\text{branch})$  and compute

$$\Upsilon_4; E_4 \vdash \text{dens}(\text{random}(\text{Gaussian}(\text{mA}, 1.0))) \Rightarrow \lambda z. F_4$$

where  $F_4$  is computed using (RANDOM CONST) and (MARGINAL) as

$$\left( \int \lambda(\text{branch}, \text{temp}). E_4 \right) \cdot \text{PDF}_{\text{Gaussian}(\text{mA}, 1.0)}(z)$$

Here  $\mathcal{I}[\int \lambda(\text{branch}, \text{temp}). E_4] \sigma = 0.7$ , so the contribution of the left branch to the PDF of the **match** is the PDF of the branch scaled by the probability of entering the left branch. In general, this holds when the branch expression is independent from the body of the branch.

For the right branch, see Example 3.11. We then obtain the PDF of the match as the sum of the PDFs of the two branches.

Our implementation of the compiler uses the following derived rules for **if** statements where the branching expression is of type **bool**, and does not treat other sum types nor matches.

#### Derived rules for **if** statements

(IF DET)	
$M$ pure	$\Upsilon; E \cdot \text{isL}(M) \vdash \text{dens}(N_1) \Rightarrow \lambda z. F_1$ $\Upsilon; E \cdot \text{isR}(M) \vdash \text{dens}(N_2) \Rightarrow \lambda z. F_2$
$\Upsilon; E \vdash \text{dens}(\text{if } M \text{ then } N_1 \text{ else } N_2) \Rightarrow \lambda z. F_1 + F_2$	
(IF RND)	
$\neg(M \text{ pure})$	$\Upsilon; E \cdot \text{let } w = \text{true in } F \vdash \text{dens}(N_1) \Rightarrow \lambda z. F_1$ $\varepsilon; 1 \vdash \text{dens}(M) \Rightarrow \lambda w. F$ $\Upsilon; E \cdot \text{let } w = \text{false in } F \vdash \text{dens}(N_2) \Rightarrow \lambda z. F_2$
$\Upsilon; E \vdash \text{dens}(\text{if } M \text{ then } N_1 \text{ else } N_2) \Rightarrow \lambda z. F_1 + F_2$	

**Example 3.10.** Since the match-bound variables  $\_$  (lines 4 and 5) do not appear in the bodies of the match branches, we can instead use rule (IF DET) to avoid adding them to the probability context when computing the PDF of the body (cf. Example 3.1).

#### Density compiler, discrete operations : $\Upsilon; E \vdash \text{dens}(f(M)) \Rightarrow F$

(DISCRETE)	
$M$ pure	$\Upsilon; E \vdash \text{marg}(\{x_1, \dots, x_n\}) \Rightarrow F$ $f \notin \{\text{fromL}, \text{fromR}\}$ $f : t \rightarrow u$ $u$ discrete $\text{rands}(\Upsilon) \cap \text{fv}(M\sigma_\Upsilon) = x_1, \dots, x_n$
$\Upsilon; E \vdash \text{dens}(f(M)) \Rightarrow \lambda w. \int \lambda(x_1, \dots, x_n). F \cdot [w = f(M\sigma_\Upsilon)]$	



The (DISCRETE) rule for discrete operations, such as logical and comparison operations and integer arithmetic, computes the expectation of an indicator function over the joint distribution of the random variables occurring in the expression.

For numeric operations on real numbers we mimic the change of variable rule of integration (often summarized as “ $dx = \frac{dx}{dy} dy$ ”), multiplying the density of the argument with the derivative of the inverse of the operation. For operations of more than one argument (e.g., (PLUS RND) below), we instead use the matrix volume of the Jacobian matrix of the inverse operation (Ben-Israel, 1999). We only require that the operation is invertible on a restricted domain, namely where the PDF of its argument is non-zero. This is exemplified by the following rules.

**Density compiler, numeric operations on reals :**  $\Upsilon; E \vdash \text{dens}(f(M)) \Rightarrow F$

---

(INVERSE)	$\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(1/M) \Rightarrow \lambda z. (\text{let } w = 1/z \text{ in } F) \cdot (1/z^2)}$
(EXP)	$\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(\exp(M)) \Rightarrow \lambda z. \text{if } z > 0.0 \text{ then } (\text{let } w = \log(z) \text{ in } F) \cdot (1/z) \text{ else } 0.0}$
(LOG)	$\frac{\Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F \quad \neg \exists \sigma, r < 0, c \neq 0. \mathcal{S}[\text{let } w = r \text{ in } F] \sigma = c}{\Upsilon; E \vdash \text{dens}(\log(M)) \Rightarrow \lambda z. (\text{let } w = \exp(z) \text{ in } F) \cdot \exp(z)}$
(SCALE)	$\frac{c \neq 0 \quad \Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(c \cdot M) \Rightarrow \lambda z. (\text{let } w = z/c \text{ in } F) \cdot (1/\text{abs}(c))}$
(PLUS DET)	$\frac{N \text{ pure} \quad \text{rands}(\Upsilon) \# (N\sigma_\Upsilon) \quad \Upsilon; E \vdash \text{dens}(M) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(M + N) \Rightarrow \lambda z. \text{let } w = z - N\sigma_\Upsilon \text{ in } F}$
(PLUS RND)	$\frac{\neg(N \text{ pure} \wedge \text{rands}(\Upsilon) \# (N\sigma_\Upsilon)) \quad \Upsilon; E \vdash \text{dens}((M, N)) \Rightarrow \lambda w. F}{\Upsilon; E \vdash \text{dens}(M + N) \Rightarrow \lambda z. \int \lambda w_1. \text{let } w = (w_1, z - w_1) \text{ in } F}$

---

For the (LOG) rule above, we require that negative arguments to **log** have zero density.

**Example 3.11.** Letting  $E_7 \equiv \text{PDF}_{\text{Bernoulli}(0.7)}(\text{branch}) \cdot \text{PDF}_{\text{Gaussian}(0,1.0)}(\text{temp}) \cdot \text{isR}(\text{branch})$ , the sum on line 6 is evaluated using rule (PLUS DET) as

$$\Upsilon_7; E_7 \vdash \text{dens}(\text{temp} + \text{mB}) \Rightarrow \lambda z. F_7$$

where

$$F_7 = \text{let temp} = z - \text{mB in } \int \lambda \text{branch}. E_7 \equiv \lambda z. 0.3 \cdot \text{PDF}_{\text{Gaussian}(\text{mB}, 1.0)}(z)$$

since  $\text{PDF}_{\text{Gaussian}(0,1.0)}(z - r) = \text{PDF}_{\text{Gaussian}(r,1.0)}(z)$  for all  $r$ . Thus we obtain that the PDF of the program fragment in Figure 1 is given by

$$\lambda z. F_4 + F_7 \equiv \lambda z. 0.7 \cdot \text{PDF}_{\text{Gaussian}(\text{mA}, 1.0)}(z) + 0.3 \cdot \text{PDF}_{\text{Gaussian}(\text{mB}, 1.0)}(z).$$

Finally, as another example of compilation, the **if** statement in the program  
**let**  $p = \mathbf{random}(\mathbf{Beta}(1.0,1.0))$  **in** // the uniform distribution on the unit interval  
**let**  $b = \mathbf{random}(\mathbf{Bernoulli}(p))$  **in**  
**if**  $b$  **then**  $p+1.0$  **else**  $p$

is handled by (IF DET), yielding a density function that (modulo trivial integrals) is equivalent to

$$\begin{aligned} & \lambda z. \quad \mathbf{let} \ p = z - 1 \ \mathbf{in} \ \int \lambda b. [0 \leq p \leq 1] \cdot (\mathbf{if} \ b \ \mathbf{then} \ p \ \mathbf{else} \ 1 - p) \cdot \mathbf{isL}(b) \\ & \quad + \int \lambda b. [0 \leq z \leq 1] \cdot (\mathbf{if} \ b \ \mathbf{then} \ z \ \mathbf{else} \ 1 - z) \cdot \mathbf{isR}(b) \\ & \equiv \lambda z. \quad [1 \leq z \leq 2] \cdot (z - 1) + [0 \leq z \leq 1] \cdot (1 - z) \end{aligned}$$

**3.4. Compiler Correctness.** These derived judgments relate the types of the various terms occurring in the marg and dens judgments.

**Lemma 3.12** (Derived Judgments).

If  $\Gamma, \Gamma_Y \vdash Y$  wf and  $\text{dom}(\Gamma_Y) = \text{rands}(Y) \cup \text{dom}(\sigma_Y)$  and  $\Gamma, \Gamma_Y \vdash E : \mathbf{real}$  then

- (1) If  $Y; E \vdash \text{marg}(X) \Rightarrow F$  and  $X = \{x_1, \dots, x_n\}$  and  $\Gamma_Y \vdash (x_1, \dots, x_n) : (t_1 * \dots * t_n)$  then  $\Gamma, x_1 : t_1, \dots, x_n : t_n \vdash F : \mathbf{real}$ .
- (2) If  $Y; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$  and  $\Gamma, \Gamma_Y \vdash M : t$  then  $\Gamma, z : t \vdash F : \mathbf{real}$ .

*Proof.*

- (1) By inversion of (MARGINAL),  $F = \int \lambda(y_1, \dots, y_k). E \sigma_Y$  where  $X \subseteq \text{rands}(Y)$  and  $(\text{rands}(Y) \setminus X) = y_1, \dots, y_k$ . By  $\Gamma, \Gamma_Y \vdash Y$  wf we have that  $\Gamma, \Gamma_Y \vdash \sigma_Y(x) : \Gamma_Y(x)$  for all  $x \in \text{dom} \sigma_Y$ . Without loss of generality, let  $\Gamma_Y \equiv \Gamma_X, \Gamma_Y, \Gamma_{\sigma_Y}$  where  $\text{dom} \Gamma_X = X$ ,  $\text{dom} \Gamma_Y = \{y_1, \dots, y_k\}$  and  $\text{dom} \Gamma_{\sigma_Y} = \text{dom} \sigma_Y$ . By repeated application of Lemma 2.2.1 we obtain  $\Gamma, \Gamma_X, \Gamma_Y \vdash E \sigma_Y : \mathbf{real}$ . By (TARGET INT) we then derive  $\Gamma, \Gamma_X \vdash \int \lambda(y_1, \dots, y_n). E \sigma_Y : \mathbf{real}$ . By repeated inversion of  $\Gamma_Y \vdash (x_1, \dots, x_n) : (t_1 * \dots * t_n)$  we can conclude that  $\Gamma_X = x_1 : t_1, \dots, x_n : t_n$ , which gives us the result  $\Gamma, x_1 : t_1, \dots, x_n : t_n \vdash F : \mathbf{real}$ .
- (2) By induction on the derivation of  $Y; E \vdash \text{dens}(M) \Rightarrow \lambda z. F$ , using (1). □

The density compiler is deterministic.

**Lemma 3.13** (Determinism).

If  $\Gamma \vdash Y$  wf and  $\Gamma \vdash E : \mathbf{real}$  then

- (1) If  $Y; E \vdash \text{marg}(X) \Rightarrow F_1$  and  $Y; E \vdash \text{marg}(X) \Rightarrow F_2$  then  $F_1 = F_2$ .
- (2) If  $Y; E \vdash \text{dens}(M) \Rightarrow \lambda z. F_1$  and  $Y; E \vdash \text{dens}(M) \Rightarrow \lambda z. F_2$  then  $F_1 = F_2$ .

*Proof.*

- (1) The (MARGINAL) rule is deterministic.
- (2) By induction on  $M$ , using (1). In every case, at most one compilation rule applies. □

We also need a technical lemma relating pure Fun expressions with their fun counterparts.

**Lemma 3.14.** If  $M$  is pure and  $\Gamma \vdash M : t$  and  $\Gamma \vdash \sigma$  then  $\mathcal{P}[[M]] \sigma = \mathbf{return} \ \mathcal{S}[[M]] \sigma$ .

*Proof.* By induction on  $M$ . □

The soundness theorem asserts that, for all closed expressions  $M$ , the density function given by the density compiler indeed characterizes (via stock integration) the distribution of  $M$  given by the monadic semantics.

**Theorem 3.15** (Soundness). *If  $\varepsilon; 1 \vdash \text{dens}(M) \Rightarrow \lambda z. F$  and  $\varepsilon \vdash M : t$  then*

$$(\mathcal{P}[[M]] \varepsilon) A = \int_A \mathcal{S}[[\lambda z. F]] \varepsilon$$

*Proof.* We let  $\bar{y} := y_1, \dots, y_n := \text{rands}(\Upsilon)$ , and otherwise use the meta-variables from the derivation rules.

We prove the theorem by joint induction on the derivations of  $\text{dens}(M)$  and  $M : t$ , using the following induction hypothesis (IH):

$$\text{if } \Gamma, \Gamma_Y \vdash Y \text{ wf and } \text{dom}(\Gamma_Y) = \text{rands}(\Upsilon) \cup \text{dom}(\sigma_Y) \quad (3.1)$$

$$\text{and } \Gamma, \Gamma_Y \vdash M : t \text{ and } \Gamma, \Gamma_Y \vdash E : \mathbf{real} \text{ and } Y; E \vdash \text{dens}(M) \Rightarrow \lambda z. F \quad (3.2)$$

$$\text{and } \Gamma \vdash \rho \text{ and } \mu(B) = \int_B \mathcal{S}[[\lambda \bar{y}. E \sigma_Y]] \rho \text{ and } |\mu| \leq 1 \quad (3.3)$$

$$\text{and } (\forall \rho'. \Gamma_Y \vdash \rho' \wedge \mathcal{S}[[E]] \rho \rho' \neq 0.0 \Rightarrow \quad (3.4)$$

$$((\sigma_Y(x) = \mathbf{fromL}(M) \Rightarrow \exists V. \mathcal{S}[[M \sigma_Y]] \rho \rho' = \mathbf{inl} V) \quad (3.5)$$

$$\wedge (\sigma_Y(x) = \mathbf{fromR}(M) \Rightarrow \exists V. \mathcal{S}[[M \sigma_Y]] \rho \rho' = \mathbf{inr} V))) \quad (3.6)$$

$$\text{then } (\mu \gg = \lambda \bar{y}. \mathcal{P}[[M \sigma_Y]] \rho) A = \int_A \mathcal{S}[[\lambda z. F]] \rho. \quad (3.7)$$

We first note that since the density expression is only ever modified by multiplication with other real-valued expressions, the conjunct at 3.4-3.6 of IH can only be invalidated when a deterministic variable  $x = \mathbf{fromL}(M)$  or  $x = \mathbf{fromR}(M)$  is added to  $Y$ , which only can occur in rule (MATCH DET). In the left branch of the **match**,  $\mathcal{S}[[E \cdot \mathbf{isL}(M \sigma_Y)]] \rho \rho' \neq 0.0$  implies that  $\mathcal{S}[[\mathbf{isL}(M \sigma_Y)]] \rho \rho' \neq 0.0$ , so  $\mathcal{S}[[M \sigma_Y]] \rho \rho' = \mathbf{inl} V$  for some  $V$ . A symmetric argument applies to the right branch of the **match**.

We proceed with the induction.

(VAR DET)

$$\begin{aligned} \text{LHS} &= (\mu \gg = \lambda \bar{y}. \mathbf{return} z \sigma_Y \rho) A \\ (\text{by Lemma 3.14}) &= (\mu \gg = \lambda \bar{y}. \mathcal{P}[[E' \sigma_Y]] \rho) A \\ (\text{by IH}) &= \text{RHS} \end{aligned}$$

(VAR RND) Assume that  $x = y_i$ , and let  $\bar{y}' = \bar{y} \setminus x$ .

$$\text{LHS} = (\mu \gg = \lambda \bar{y}. \mathbf{return} y_i) A = \int \lambda \bar{y}. (\mathcal{S}[[E \sigma_Y]] \rho) \cdot [y_i \in A] = \int_A \lambda y_i. \int \mathcal{S}[[\lambda \bar{y}'. E \sigma_Y]] \rho = \text{RHS}$$

(CONSTANT)

$$\text{LHS} = (\mu \gg = \lambda \bar{y}. \mathbf{return} V) A = \int \lambda \bar{y}. (\mathcal{S}[[E \sigma_Y]] \rho) \cdot [V \in A] = \left( \int \mathcal{S}[[\lambda \bar{y}. E \sigma_Y]] \rho \right) \cdot \sum_{x \in A} [x = V] = \text{RHS}$$

(FAIL) LHS = 0.0 = RHS.

(RANDOM CONST)

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. (\mathcal{P} \llbracket M\sigma_Y \rrbracket \rho \gg= \lambda x. \mu_{\text{Dist}(x)}) A \\
(\text{by Lemma 3.14}) &= (\mu \gg= \lambda \bar{y}. (\text{return } (\mathcal{S} \llbracket M\sigma_Y \rrbracket \rho) \gg= \lambda x. \mu_{\text{Dist}(x)})) A \\
(\text{by monad laws}) &= (\mu \gg= \lambda \_ . \mu_{\text{Dist}(\mathcal{S} \llbracket M\sigma_Y \rrbracket \rho)}) A \\
&= \mu_{\text{Dist}(\mathcal{S} \llbracket M\sigma_Y \rrbracket \rho)}(A) \cdot \int \mathcal{S} \llbracket \lambda \bar{y}. E\sigma_Y \rrbracket \rho \\
&= \int_A \mathcal{S} \llbracket \lambda z. \left( \int \lambda \bar{y}. E\sigma_Y \right) \cdot \text{PDF}_{\text{Dist}(M\sigma_Y)}(z) \rrbracket \rho = \text{RHS}
\end{aligned}$$

(RANDOM RND) Let  $\nu A = \int_A \mathcal{S} \llbracket \lambda w. F \rrbracket \rho$ .

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. (\mathcal{P} \llbracket M\sigma_Y \rrbracket \rho \gg= \lambda z. \mu_{\text{Dist}(z)})) A \\
(\text{by monad laws}) &= ((\mu \gg= \lambda \bar{y}. \mathcal{P} \llbracket M\sigma_Y \rrbracket \rho) \gg= \lambda z. \mu_{\text{Dist}(z)}) A \\
(\text{by induction}) &= (\nu \gg= \lambda z. \mu_{\text{Dist}(z)}) A \\
&= \int \lambda z. \mu_{\text{Dist}(z)}(A) d\nu \\
&= \int \lambda z. ((\mathcal{S} \llbracket \lambda w. F \rrbracket \rho) z) \cdot \mu_{\text{Dist}(z)}(A) \\
&= \int \lambda w. (\mathcal{S} \llbracket F \rrbracket \rho) \cdot \mu_{\text{Dist}(w)}(A) \\
&= \int \lambda w. (\mathcal{S} \llbracket F \rrbracket \rho) \cdot \int_A \text{PDF}_{\text{Dist}(w)} \\
&= \int_A \lambda z. \int \lambda w. (\mathcal{S} \llbracket F \rrbracket \rho) \cdot \text{PDF}_{\text{Dist}(w)}(z) \\
&= \text{RHS}
\end{aligned}$$

(TUPLE VAR) Let  $\bar{x} = x_1, \dots, x_k$ , and  $\bar{y}' = \bar{y} \setminus \bar{x}$ . Let  $\nu A = \int_A \mathcal{S} \llbracket \lambda \bar{x}. F \rrbracket \rho$ .

$$\text{LHS} = (\mu \gg= \lambda \bar{y}. \text{return } (\bar{x})) A = \int \lambda \bar{y}. (\mathcal{S} \llbracket E\sigma_Y \rrbracket \rho \cdot [(\bar{x}) \in A]) = \int_A \lambda(\bar{z}). \int \mathcal{S} \llbracket \lambda \bar{y}'. E\sigma_Y \rrbracket \rho = \text{RHS}$$

(TUPLE PROJ L) Let  $\nu A = \int_A \mathcal{S} \llbracket \lambda w. F \rrbracket \rho$ .

$$\begin{aligned}
\text{LHS} &= (\nu \gg= \lambda w. \text{return } \mathbf{fst} w) A \\
&= \int \lambda w. [\mathbf{fst} w \in A] d\nu \\
&= \int \lambda w. [\mathbf{fst} w \in A] \cdot \mathcal{S} \llbracket F \rrbracket \rho \\
&= \int \lambda(z_1, z_2). \text{let } w = (z_1, z_2) \text{ in } [z_1 \in A] \cdot \mathcal{S} \llbracket F \rrbracket \rho \\
&= \int_A \lambda z_1. \int \lambda z_2. \mathcal{S} \llbracket \text{let } w = (z_1, z_2) \text{ in } F \rrbracket \rho = \text{RHS}
\end{aligned}$$

(TUPLE PROJ R) As (TUPLE PROJ L).

(LET DET)

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda v. \mathcal{P}[[N\sigma_Y]] \rho, x \mapsto v)) A \\
(\text{by Lemma 3.14}) &= (\mu \gg= \lambda \bar{y}. (\text{return } (\mathcal{S}[[M\sigma_Y]] \rho) \gg= \lambda v. \mathcal{P}[[N\sigma_Y]] \rho, x \mapsto v)) A \\
(\text{by monad laws}) &= (\mu \gg= \lambda \bar{y}. \mathcal{P}[[N\sigma_Y]] \rho, x \mapsto \mathcal{S}[[M\sigma_Y]] \rho) A \\
&= (\mu \gg= \lambda \bar{y}. \mathcal{P}[[N]] \sigma_{Y, x \mapsto M\rho}) A \\
&= \text{RHS}
\end{aligned}$$

(LET RND) Let  $v B := \int_B \mathcal{S}[[\lambda \bar{y}, x. (E \cdot F_1)\sigma_Y]] \rho$ . By induction  $(\mathcal{P}[[M]\rho') C = \int_C \mathcal{S}[[F_1]] \rho'$  whenever  $\Gamma, \Gamma_Y \vdash \rho'$ .

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda z. \mathcal{P}[[N\sigma_Y]] \rho, x \mapsto z)) A \\
(\text{by monad laws}) &= ((\mu \gg= \lambda \bar{y}. (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda z. \text{return } \bar{y}, z)) \gg= \lambda \bar{y}, z. \mathcal{P}[[N]] \sigma_Y \rho, x \mapsto z) A
\end{aligned}$$

Then

$$\begin{aligned}
&(\mu \gg= \lambda \bar{y}. (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda x. \text{return } \bar{y}, x)) B \\
&= \int \lambda \bar{y}. (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda x. \text{return } \bar{y}, x) B d\mu \\
&= \int \lambda \bar{y}. (\mathcal{S}[[E\sigma_Y]] \rho) \cdot (\mathcal{P}[[M\sigma_Y]] \rho \gg= \lambda x. \text{return } \bar{y}, x) B \\
&= \int \lambda \bar{y}. (\mathcal{S}[[E\sigma_Y]] \rho) \cdot \int \lambda x. [\bar{y}, x \in B] d\mathcal{P}[[M\sigma_Y]] \rho \\
(\text{induction}) &= \int \lambda \bar{y}. (\mathcal{S}[[E\sigma_Y]] \rho) \cdot \int \lambda x. \mathcal{S}[[F_1\sigma_Y]] \rho \cdot [\bar{y}, x \in B] \\
&= \int_B \lambda \bar{y}, x. \mathcal{S}[[E\sigma_Y \cdot F_1\sigma_Y]] \rho \\
&= v
\end{aligned}$$

By induction,  $(v \gg= \lambda \bar{y}, z. \mathcal{P}[[N\sigma_Y]] \rho, x \mapsto z) A = \text{RHS}$ .(SUM CON L) Let  $v A = \int_A \mathcal{S}[[\lambda z. F]] \rho$ .

$$\text{LHS} = (v \gg= \lambda z. \text{return } \mathbf{inl} z) A = \int \lambda z. [\mathbf{inl} z \in A] dv = \int_A \mathbf{either} \mathcal{S}[[F]] \rho \lambda \_ . 0 = \text{RHS}$$

(SUM CON R) As (SUM CON L).

(MATCH DET) Let

$$\begin{aligned}
L &= \{(V_1, \dots, V_n) \mid \exists W. \mathcal{S}[[M\sigma_Y]] \rho [y_1, \dots, y_n \mapsto V_1, \dots, V_n] = \mathbf{inl} W \\
&\text{and } \varepsilon \vdash V_i : \Gamma_Y(y_i) \text{ for all } i\}.
\end{aligned}$$

Then  $\mu = \mu_L + \mu_R$  where  $\mu_L(B) = \mu(B \cap L) = \int_B \lambda(\text{rands}(Y)). \mathcal{S}[[E\sigma_Y \cdot \mathbf{isL}(M\sigma_Y)]] \rho$  and  $\mu_R(B) = \mu(B \setminus L) = \int_B \lambda(\text{rands}(Y)). \mathcal{S}[[E\sigma_Y \cdot \mathbf{isR}(M\sigma_Y)]] \rho$ . By additivity, we have

$$\begin{aligned}
\text{LHS} &= (\mu_L + \mu_R \gg= \mathcal{P}[[M'\sigma_Y]] \rho) A \\
&= (\mu_L \gg= \mathcal{P}[[M'\sigma_Y]] \rho) A + (\mu_R \gg= \mathcal{P}[[M'\sigma_Y]] \rho) A
\end{aligned}$$

We let  $E_N := \text{match } \mathcal{S} \llbracket M\sigma_\Upsilon \rrbracket \rho$  with  $\text{inl } z : \mathcal{P} \llbracket N_1 \rrbracket (\sigma, x_1 \mapsto z) \mid \text{inr } z : \mathcal{P} \llbracket N_2 \rrbracket (\sigma, x_2 \mapsto z)$ .

Then

$$\begin{aligned}
(\mu_L \gg= \mathcal{P} \llbracket M'\sigma_\Upsilon \rrbracket \rho) A &= (\mu_L \gg= \lambda \bar{y}. (\mathcal{P} \llbracket M\sigma_\Upsilon \rrbracket \rho \gg= \text{either} \\
&\quad (\lambda z. \mathcal{P} \llbracket N_1 \rrbracket (\sigma, x_1 \mapsto z)) (\lambda z. \mathcal{P} \llbracket N_2 \rrbracket (\sigma, x_2 \mapsto z)))) A \\
\text{(by Lemma 3.14)} &= (\mu_L \gg= \lambda \bar{y}. (\text{return } (\mathcal{S} \llbracket M\sigma_\Upsilon \rrbracket \rho) \gg= \text{either } \dots)) A \\
\text{(by monad laws)} &= (\mu_L \gg= \lambda \bar{y}. E_N) A \\
&= \int_L \lambda \bar{y}. E_N(A) d\mu_L + \int_L \lambda \bar{y}. E_N(A) d\mu_L \\
&= \int_L \lambda \bar{y}. E_N(A) d\mu_L \\
(\forall \bar{V} \in L \exists W. \mathcal{S} \llbracket M\sigma_\Upsilon \rrbracket \rho [\bar{y} := \bar{V}] = \mathbf{inl } W) &= \int_L \lambda \bar{y}. \mathcal{P} \llbracket N_1\sigma_\Upsilon \rrbracket (\rho, x_1 \mapsto \text{fromL } \mathcal{S} \llbracket M\sigma_\Upsilon \rrbracket \rho)(A) d\mu_L \\
&= (\mu_L \gg= \lambda \bar{y}. \mathcal{P} \llbracket N_1\sigma_\Upsilon \rrbracket (\rho, x_1 \mapsto \text{fromL } \mathcal{S} \llbracket M\sigma_\Upsilon \rrbracket \rho)) A \\
\text{(by induction)} &= \int_A \mathcal{S} \llbracket \lambda z. F_1 \rrbracket \rho
\end{aligned}$$

Symmetrically,  $(\mu_R \gg= \mathcal{P} \llbracket M'\sigma_\Upsilon \rrbracket \rho) A = \int_A \mathcal{S} \llbracket \lambda z. F_2 \rrbracket \rho$ , so  $\text{LHS} = \int_A \mathcal{S} \llbracket \lambda z. F_1 \rrbracket \rho + \int_A \mathcal{S} \llbracket \lambda z. F_2 \rrbracket \rho = \int_A \lambda z. \mathcal{S} \llbracket F_1 \rrbracket \rho + \mathcal{S} \llbracket F_2 \rrbracket \rho = \text{RHS}$ .

(MATCH RND) Write  $E_i := \mathcal{P} \llbracket N_i\sigma_\Upsilon \rrbracket (\rho, x_i \mapsto z)$ . Here  $\text{either } \lambda z. E_1 \lambda z. E_2 =_\beta \lambda v. \text{match } v \text{ with } \text{inl } z : N_1 \mid \text{inr } z : N_2$ . As in case (LET RND) we let  $v B := \int_B \mathcal{S} \llbracket \lambda \bar{y}, x_1. E\sigma_\Upsilon \cdot \mathbf{let } w = \mathbf{inl } x_1 \text{ in } F\sigma_\Upsilon \rrbracket \rho$ , and get

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. (\mathcal{P} \llbracket M\sigma_\Upsilon \rrbracket \rho \gg= \lambda v. \text{return } \bar{y}, v) \\
&\quad \gg= \lambda \bar{y}, v. \text{match } v \text{ with } \text{inl } z : N_1 \mid \text{inr } z : N_2) A \\
&= v \gg= \lambda \bar{y}, v. \text{match } v \text{ with } \text{inl } z : N_1 \mid \text{inr } z : N_2) A \quad (*)
\end{aligned}$$

We proceed as in case (MATCH DET) but with  $L := \{(V_1, \dots, V_n, \mathbf{inl } W) \mid \varepsilon \vdash V_i : \Gamma_\Upsilon(y_i) \text{ for all } i\}$ , yielding  $(*) = \int_A \lambda z. \mathcal{S} \llbracket F_1 + F_2 \rrbracket \rho = \text{RHS}$ .

(FROML) Let  $v B := \int_B \mathcal{S} \llbracket \lambda w. F \rrbracket \rho$ .

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. \text{return } \text{fromL}(M\sigma_\Upsilon \rho)) A \\
\text{(monad law)} &= ((\mu \gg= \lambda \bar{y}. \text{return } M\sigma_\Upsilon \rho) \gg= \text{return } \circ \text{fromL}) A \\
\text{(induction)} &= (v \gg= \text{return } \circ \text{fromL}) A \\
\text{(definition)} &= \int_{t+u} \lambda w. [\text{fromL}(w) \in A] \cdot \mathcal{S} \llbracket F \rrbracket \rho
\end{aligned}$$

By part 3.4-3.6 of the IH,  $\mathcal{S} \llbracket F \rrbracket \rho \rho' \neq 0.0$  implies  $\mathcal{S} \llbracket M\sigma_\Upsilon \rho \rho' \rrbracket = \mathbf{inl } V$ , so we have

$$\int_{t+u} [\text{fromL}(x) \in A] \cdot \mathcal{S} \llbracket F \rrbracket \rho = \int_A \mathcal{S} \llbracket F \rrbracket \rho [z := \mathbf{inl } x] dx = \text{RHS}$$

(FROMR) As (FROML).

(DISCRETE) Let  $\nu B := \int_B \mathcal{S}[\lambda \bar{x}. F] \rho$  and  $\bar{y} \setminus \bar{x} = \bar{z}$

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. \text{return } f(M\sigma_\Gamma \rho)) A \\
&= \int \lambda \bar{y}. [f(\mathcal{S}[M\sigma_\Gamma] \rho) \in A] d\mu \\
&= \int \lambda \bar{y}. [f(\mathcal{S}[M\sigma_\Gamma] \rho) \in A] \cdot \mathcal{S}[E\sigma_\Gamma] \rho \\
&= \sum_{w \in A} \int \mathcal{S}[\lambda \bar{z}. [w = f(M\sigma_\Gamma)]] \cdot \int \lambda \bar{x}. E\sigma_\Gamma \rho = \text{RHS}
\end{aligned}$$

(PLUS RND) Let  $\nu B = \int_B \mathcal{S}[\lambda w. F] \rho$ .

$$\begin{aligned}
\text{LHS} &= ((\mu \gg= \lambda \bar{y}. \text{return } (M, N)\sigma_\Gamma \rho) \gg= \lambda(x_1, x_2). \text{return } x_1 + x_2) A \\
(\text{Lemma 3.14}) &= ((\mu \gg= \lambda \bar{y}. \mathcal{S}[(M, N)\sigma_\Gamma] \rho) \gg= \lambda(x_1, x_2). \text{return } x_1 + x_2) A \\
(\text{induction}) &= (\nu \gg= \lambda(x_1, x_2). \text{return } x_1 + x_2) A \\
&= \int \lambda x, y. [x + y \in A] d\nu \\
&= \int \lambda x, y. [x + y \in A] \cdot (\text{let } w = (x, y) \text{ in } \mathcal{S}[F] \rho) \\
(z := x + y) &= \int_A \lambda z. \int \lambda x. \mathcal{S}[\text{let } w = (x, z - x) \text{ in } F] \rho \\
&= \text{RHS}
\end{aligned}$$

**Numeric operations on real:** Assume that  $f$  is strictly monotonic and  $g := f^{-1}$  has a continuous derivative for  $f(x) \in A$ . Let  $\nu B := \int_B \mathcal{S}[\lambda w. F] \rho$ .

$$\begin{aligned}
\text{LHS} &= (\mu \gg= \lambda \bar{y}. \text{return } f(\mathcal{S}[M\sigma_\Gamma] \rho)) A \\
(\text{monad law}) &= ((\mu \gg= \lambda \bar{y}. \text{return } \mathcal{S}[M\sigma_\Gamma] \rho) \gg= \lambda w. \text{return } f(w)) A \\
(\text{induction}) &= (\nu \gg= \lambda w. \text{return } f(w)) A \\
&= \int \lambda w. [f(w) \in A] \cdot \mathcal{S}[F] \rho \\
(\text{change of variables}) &= \int_A \lambda z. (\mathcal{S}[\text{let } w = g(z) \text{ in } F] \rho) \cdot g'(w) = \text{RHS}
\end{aligned}$$

For the base case of the induction, we have that  $E = 1$ ,  $\mu$  is the probability measure on the unit type, and all of  $\Gamma$ ,  $\Gamma_\Gamma$ ,  $\sigma_\Gamma$  and  $\rho$  are empty. Clearly, (IH) holds for the base case.  $\square$

Part 3.4-3.6 of the induction hypothesis above is used when attempting to evaluate match-bound variables (e.g.,  $x = \text{fromL}(M)$ ) for valuations that give the other branch (e.g.,  $\mathcal{S}[M] \sigma = \text{inr } V$ ). For such valuations the density is always zero (since, e.g.,  $\text{isL}(\text{inr } V) = 0.0$ ).

#### 4. EVALUATION

We evaluate the compiler on several synthetic textbook examples and several real examples from scientific applications. We wish to validate that the density compiler handles these examples, and understand how much the compiler reduces the developer burden, and its performance impact.

**4.1. Implementation.** Since Fun is a sublanguage of F#, we implement our models as F# programs, and use the quotation mechanism of F# to capture their syntax trees. Running the F# program corresponds to sampling data from the model. To compute the PDF, the compiler takes the syntax tree (of F# type `Expr`) of the model and produces another `Expr` corresponding to a deterministic F# program as output. We then use run-time code generation to compile the generated `Expr` to MSIL bytecode, which is just-in-time compiled to executable machine code when called, just as for statically compiled F# code. Our implementation supports immutable arrays and records, which are both translated using adaptations of the corresponding rules for tuples. For efficiency, the implementation must avoid introducing redundant computations, translating the use of substitution in the formal rules to more efficient `let`-bindings that share the values of expressions that would otherwise be re-computed. As is common practice, our implementation and Filzbach (Purves and Lyutsarev, 2012) both work with the *logarithm* of the density, which avoids products of densities in favor of sums of log-densities where possible, to avoid numerical underflow. It also performs some simple but effective peephole optimization to elide canceling applications of `Log` and `Exp` and additions of 0.

### Code Examples.

**Example 4.1** (Mixture Of Gaussians). To illustrate the implementation, here is the actual F# code expressing a mixture of Gaussians (a variant of our introductory example):

```

type W = { bias: double; mean: double[]; sd: double[] }

[<Fun>]
let prior () =
  { bias = random(Uniform(0.0, 1.0))
    mean = [| for i in 0..1 → random(Uniform(-1000.0, 1000.0)) |]
    sd = [| for i in 0..1 → random(Uniform(100.0, 500.0)) |] }

let xs = [| for i in 1..100 → () |]

[<Fun>]
let model w = [| for x in xs →
  if random(Bernoulli(w.bias))
  then random(Gaussian(w.mean.[0],w.sd.[0])
  else random(Gaussian(w.mean.[1],w.sd.[1]) |]

```

The code uses both records (to structure the prior `w` of record type `W`) and arrays; both are encoded as tuples in the core language. The model function receives multiple inputs in an array `xs` and return an array multiple outputs from the model. The outputs of `model` are constructed using an *array comprehension*. The [`<Fun>`] attributes declare that definitions `prior` and `model` should be made available as quoted expression trees (as well as executable functions) so their code can be inspected by the density compiler.



The probability density function compiled for function `model` is (after manual reformatting to match the notation in this article):

```
let logPdf =
  fun w (ys:real[]) →
    logprodBy((fun i →
      let x = xs.[i]
      Log(Exp(Log(pdf_Gaussian(w.mean.[0], w.sd.[0], ys.[i])) +
        let b=true in Log(pdf_Bernoulli(w.bias,b)))
      +
      Exp(Log(pdf_Gaussian(w.mean.[1], w.sd.[1], ys.[i])) +
        let b=false in Log(pdf_Gaussian(w.bias, b))))),
    xs.Length)
```

The helper function `logprodBy(f, n)` computes the sum of the log densities  $f(0) + \dots + f(n-1)$ . Notice the insertion of logarithms to avoid underflow.

The effect of disabling our simple peephole optimizer is to produce both less readable and less efficient code:

```
let logPdf =
  fun w (ys:real[]) →
    logprodBy(
      (fun i →
        let x = xs.[i]
        Log(Exp(Log(pdf_Gaussian(w.mean.[0],w.sd.[0], ys.[i]))+
          Log(Exp(Log(1.0)+
            let b2=true
            Log(pdf_Bernoulli(w.bias, b2))+
            Log(Exp(Log(1.0)))))))+
          Exp(Log(pdf_Gaussian(w.mean.[1],w.sd.[1], ys.[i1]))+
            Log(Exp(Log(1.0)+
              let b2=false
              Log(pdf_Bernoulli(w.bias, b2))+
              Log(Exp(Log(1.0)))))))))
      xs.Length)
```

**Example 4.2** (Linear Regression). For an example involving arithmetic we take *linear regression*. Given some noisy sample of points, the task is to estimate the parameters of a line fitting the points, yielding the line's slope `a`, intercept `b` and an estimate of the `noise`.

The generative model is expressed as the following F# code:

```
type W = { a: double; b: double; noise: double }
```

```
[<Fun>]
let prior () =
  { a = random(Uniform(-1000.0, 1000.0))
    b = random(Uniform(-1000.0, 1000.0))
```

```

noise = random(Uniform(0.001, 100.0))

let xs = [| -100.0 .. 100.0 |]

[<Fun>]
let model w = [| for x in xs →
                let m = w.a * x + w.b
                let d = w.noise
                random(Gaussian(m, d)) |]

```

The (log) probability density function compiled for function `model` is:

```

let logPdf =
  fun w ys →
    logprodBy((fun i →
               let x = xs.[i]
               Log(pdf_Gaussian(w.a*x+w.b, w.noise, ys.[i])),
               xs.Length)

```

**Example 4.3** (Mixture Of Regressions). Combining aspects of the previous two examples we construct a mixture of two linear regressions, in which the slope and intercept of the line is selected by a latent boolean indicator variable.

```

type W = {bias:double; a: double[]; b: double[]; noise: double}

[<Fun>]
let prior () =
  { bias = rand.Uniform(0.0, 1.0)
    a = [| for i in 0..1 → rand.Uniform(-1000.0, 1000.0) |]
    b = [| for i in 0..1 → rand.Uniform(-1000.0, 1000.0) |]
    noise = rand.Uniform(0.001, 100.0) }

let xs = [| -100.0 .. 100.0 |]

[<Fun>]
let model w =
  [| for x in xs →
    if rand.Bernoulli(w.Bias)
    then let m = w.a.[0] * x + w.b.[0]
          rand.Gaussian(m, w.noise)
    else let m = w.a.[1] * x + w.b.[1]
          rand.Gaussian(m, w.noise) |]

```

Example	orig	LOC, orig	LOC, Fun		time (s), orig	time (s), Fun	
mixture of Gaussians	F#	32	20	0.63x	0.74	1.28	1.7x
linear regression	F#	27	18	0.67x	0.21	0.55	2.6x
mixture of regressions	F#	43	28	0.65x	1.02	2.09	2.0x
species distribution	C#	173	37	0.21x	36	67	1.9x
net primary productivity	C#	82	39	0.48x	3.7	6.1	1.6x
global carbon cycle	C#	1532	402	0.26x	n/a	301	n/a

Table 1: Lines-of-code and running time comparisons of synthetic and scientific models.

The (log) probability density function compiled for function `model` is:

```

let logPdf =
  fun w ys →
    logprodBy(
      (fun i →
        let x = xs.[i]
        Log(
          Exp(Log(pdf.Gaussian(w.a.[0]*x+w.b.[0], w.noise, ys.[i])) +
            let b=true in Log(pdf.Bernoulli(w.bias, b)))
          +
          Exp(Log(pdf.Gaussian(w.a.[1]*x+w.b.[1], w.noise, ys.[i])) +
            let b=false in Log(pdf.Bernoulli(w.bias, b))))),
      xs.Length);

```

**4.2. Metrics.** We consider scientific models with existing implementations for MCMC-based inference, written by domain experts. We are interested in how the modelling and inference experience would change, in terms of developer effort and performance impact, when adopting the Fun-based solution.

We assess the reduction in developer burden by measuring the code sizes (in lines-of-code (LOC)) of the original implementations of model and density code, and of the corresponding Fun model. For the synthetic examples, we have written both the model and the density code. The original implementations of the scientific models contain helper code such as I/O code for reading and writing data files in an application-specific format. Our LOC counts do not consider such helper code, but only count the code for generating synthetic data from the model, code for computing the logarithm of the posterior density of the model, and model-related code for setting up and interacting with Filzbach itself. We also compare the running times of the original implementations versus the Fun versions for MCMC-based inference using Filzbach, not including data manipulation before and after running inference.

**Synthetic examples.** Our synthetic examples are three classic problems: the unsupervised learning task *mixture of Gaussians* (Example 4.1) the supervised learning task *linear regression* (Example 4.2), and a *mixture of regressions* (Example 4.3). Example 4.1 can be thought of as a probabilistic version of *k-means clustering*: inference is trying to determine the unknown mixing bias and the means and variances of the Gaussian components. In Example 4.2 inference is trying to determine

the coefficients of the line. In Example 4.3 inference is trying to determine the coefficients and mixing bias of two lines.

**Species distribution** (McInerny and Purves, 2011). The species distribution problem is to give the probability that certain species will be present at a given site, based on climate factors. It is a problem of long-standing interest in ecology and has taken on new relevance in light of the issue of climate change. The particular model that we consider is designed to mitigate *regression dilution* arising from uncertainty in the predictor variables, for example, measurement error in temperature data. Inference tries to determine various features of the species and the environment, such as the optimal temperature preferred by a species, or the true temperature at a site (see code and density function in Appendix A).

**Global carbon cycle.** (Smith et al., 2012). The dynamics of the Earth’s climate are intertwined with the terrestrial carbon cycle, and better carbon models (modelling how carbon in the air gets converted to biomass) enable better constrained projections about these systems. We consider a fully data-constrained terrestrial carbon model, which is composed of various submodels for smaller processes such as *net primary productivity*, the fine root mortality rate, or the fraction of trees that are evergreen versus deciduous. Inference tries to determine the different parameters of these submodels.

**Discussion.** Table 1 reports the metrics for each example. The LOC numbers show significant reduction in code size, with more significant savings as the size of the model grows. The larger models (where the Fun versions are  $\approx 25\%$  of the size of the original) are more indicative of the savings in developer and maintenance effort because Filzbach interaction code, which is roughly the same in all models, takes up a larger fraction of the smaller models. We find the running times encouraging: we have made little attempt to optimize the generated code, and preliminary testing indicates that much of the performance slow-down is due to constant factors.

The global carbon cycle model is composed of submodels, each with their own dataset. Unfortunately, it is unclear from the original source code how this composition translates to a run of inference, making it difficult to know what constitutes a fair comparison. Thus, we do not report a running time for the full model. However, we can measure the running time of individual submodels, such as net primary productivity, where the data and control flow are simpler.

## 5. RELATED WORK

We have presented the first algorithm for deriving density functions from generative processes, that is both proved correct and implemented. An abridged version of this paper appears as Bhat et al. (2013). The correctness proof (for a version of the source language without pure **let** and general **match**) was recently mechanized in Isabelle by Eberl et al. (2015).

This paper builds on work by Bhat et al. (2012) who develop a theoretical framework for computing PDFs, but describe no implementation nor correctness proof. The density compiler of Section 3 has a simpler presentation, with two judgments compared to five, and has rules for pure **lets** and operations on integers. Our paper also uses a richer language (Fun), which adds **fail**, **match** and general **if** (and for performance reasons, pure **let**).

Gordon et al. (2013) describe a naive density calculation routine for Fun without random **lets**; this sublanguage does not cover many useful classes of models such as hierarchical and mixture models.

The BUGS system computes densities from declaratively specified models to perform Gibbs sampling (Gilks et al., 1994). However, the models are not compositional as in this work, and only the

joint density over all variables is possible. The AutoBayes system also computes densities for deriving maximum likelihood and Bayesian estimators for a significant class of statistical models (Schumann et al., 2008). It is not formally specified and does not appear to be compositional. Neither system addresses the non-existence of PDFs, presumably restricting expressivity in order to avoid the issue.

Stan (Stan Development Team, 2014) is a probabilistic programming language that supports various forms of MCMC. The Stan language is a derivative of BUGS, and is compiled to efficient C++ sampling code, in part based on automatically derived density functions. For models with latent variables, including mixture models such as mixture of Gaussians, and for models that perform non-linear computations on random values, Stan’s users are required to manually manipulate the log probability density function, using a primitive operation `increment_log_prob`. Stan employs *automatic differentiation* (Griewank and Walther, 2008) of log posteriors in order to apply gradient-based Hamiltonian MCMC algorithms. This relieves the user from coding error-prone derivatives.

Inference for the Church language also uses MCMC, but works with distributions over runs of a program instead of over its return value (Wingate et al., 2011), circumventing the need for a PDF.

There are many other systems for probabilistic programming, some of which provide a way to compute density functions or an analogous object; however, they sacrifice some other feature to do so. Several languages only provide support for finite, discrete distributions, but provide access to the probability mass function (Ramsey and Pfeffer, 2002; Kiselyov and Shan, 2009). Like BUGS and Stan, systems like the Hierarchical Bayes Compiler (Daumé III, 2007) are not formally defined and require models to be specified in a monolithic way, whereas large models in Fun can be composed of smaller models. Probabilistic logic languages like Markov Logic (Domingos et al., 2008) do not have generative semantics but instead have semantics in *undirected graphical models* which are equipped with *potential functions* that are analogous to density functions. The language is constructed in a way that guarantees the existence of a potential function, but which eliminates the possibility to express models that require pure **let**, such as the global carbon cycle model.

## 6. CONCLUSIONS AND FUTURE WORK

We have described a compiler for automatically computing probability density functions for programs from a rich Bayesian probabilistic programming language, proven the algorithm correct, and shown its applicability to real-world scientific models.

The inclusion of **fail** in the language appears useful for scientific models, giving a simple facility to exclude branches that are scientifically impossible from consideration. However, more investigation is needed to settle this claim.

A drawback of the compiler is that terms of composite type are required either to have a PDF or to be pure, ruling out terms such as `(0.0, random(Uniform))`. One possibility for future work would be to refine the types of expressions with determinacy information, and make use of this additional information to admit more joint distributions (cf. `(TUPLE VAR)`).

## ACKNOWLEDGMENTS

We thank Manuel Eberl and Tobias Nipkow for helpful comments, in particular on the semantics of the integration operator.

## REFERENCES

- A. Ben-Israel. The change of variables formula using matrix volume. *SIAM Journal of Matrix Analysis*, 21:300–312, 1999.
- S. Bhat, A. Agarwal, R. W. Vuduc, and A. G. Gray. A type theory for probability density functions. In J. Field and M. Hicks, editors, *POPL*, pages 545–556. ACM, 2012.
- S. Bhat, J. Borgström, A. D. Gordon, and C. Russo. Deriving probability density functions from probabilistic functional programs. In *Tools and Algorithms for the Construction and Analysis of Systems, 19th International Conference: TACAS 2013*, number 7795 in LNCS, pages 508–522. Springer, 2013.
- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer Science + Business Media, LLC, New York, NY, USA, 2006.
- J. Borgström, A. D. Gordon, M. Greenberg, J. Margetson, and J. Van Gael. Measure transformer semantics for Bayesian machine learning. In *European Symposium on Programming (ESOP’11)*, volume 6602 of LNCS, pages 77–96. Springer, 2011. Download available at <https://www.microsoft.com/en-us/research/project/infer-net-fun/>.
- H. Daumé III. HBC: Hierarchical Bayes Compiler, 2007. URL <http://hal3.name/HBC>.
- P. Domingos, S. Kok, D. Lowd, H. Poon, M. Richardson, and P. Singla. Markov logic. In L. De Raedt, P. Frasconi, K. Kersting, and S. Muggleton, editors, *Probabilistic inductive logic programming*, pages 92–117. Springer-Verlag, Berlin, Heidelberg, 2008.
- M. Eberl, J. Hölzl, and T. Nipkow. A verified compiler for probability density functions. In J. Vitek, editor, *24th European Symposium on Programming: ESOP 2015*, number 9032 in LNCS, pages 80–104. Springer, 2015.
- W. R. Gilks, A. Thomas, and D. J. Spiegelhalter. A language and program for complex Bayesian modelling. *The Statistician*, 43:169–178, 1994.
- M. Giry. A categorical approach to probability theory. In B. Banaschewski, editor, *Categorical Aspects of Topology and Analysis*, volume 915 of *Lecture Notes in Mathematics*, pages 68–85. Springer Berlin / Heidelberg, 1982.
- N. Goodman, V. K. Mansinghka, D. M. Roy, K. Bonawitz, and J. B. Tenenbaum. Church: a language for generative models. In *Uncertainty in Artificial Intelligence (UAI’08)*, pages 220–229. AUAI Press, 2008.
- A. D. Gordon, M. Aizatulin, J. Borgström, G. Claret, T. Graepel, A. Nori, S. Rajamani, and C. Russo. A model-learner pattern for Bayesian reasoning. In *POPL*, 2013.
- A. Griewank and A. Walther. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. SIAM, 2nd edition, 2008.
- O. Kiselyov and C. Shan. Embedded probabilistic programming. In *Domain-Specific Languages*, pages 360–384, 2009.
- G. McInerny and D. Purves. Fine-scale environmental variation in species distribution modelling: regression dilution, latent variables and neighbourly advice. *Methods in Ecology and Evolution*, 2(3):248–257, 2011.
- R. M. Neal. Probabilistic inference using Markov chain Monte Carlo methods. Technical Report CRG-TR-93-1, Dept. of Computer Science, University of Toronto, September 1993.
- P. Panangaden. The category of Markov kernels. *Electronic Notes in Theoretical Computer Science*, 22:171–187, 1999.
- D. Purves and V. Lyutsarev. *Filzbach User Guide*, 2012. Available at <https://www.microsoft.com/en-us/download/details.aspx?id=52465>.

- N. Ramsey and A. Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *POPL*, pages 154–165, 2002.
- J. Schumann, T. Pressburger, E. Denney, W. Buntine, and B. Fischer. AutoBayes program synthesis system users manual. Technical Report NASA/TM–2008–215366, NASA Ames Research Center, 2008.
- D. Scott. Parametric Statistical Modeling by Minimum Integrated Square Error. *Technometrics*, 43(3):274–285, 2001.
- M. J. Smith, M. C. Vanderwel, V. Lyutsarev, S. Emmott, and D. W. Purves. The climate dependence of the terrestrial carbon cycle; including parameter and structural uncertainties. *Biogeosciences Discussions*, 9:13439–13496, 2012.
- R. M. Solovay. A model of set-theory in which every set of reals is Lebesgue measurable. *The Annals of Mathematics, Second Series*, 92(1):1–56, 1970.
- Stan Development Team. Stan: A C++ library for probability and sampling, version 2.2, 2014. URL <http://mc-stan.org/>.
- D. Syme, A. Granicz, and A. Cisternino. *Expert F#*. Apress, 2007.
- D. Wingate, A. Stuhlmueeller, and N. Goodman. Lightweight implementations of probabilistic programming languages via transformational compilation. In *Proceedings of the 14th Intl. Conf. on Artificial Intelligence and Statistics*, page 131, 2011.

## APPENDIX A. SPECIES DISTRIBUTION (MCINERNEY AND PURVES, 2011)

```

let Nspecies = 20
let Nsamples = 2000

type W =
{ Topt: double[]
  Tbreadth: double[]
  MaxProb: double[]
  Terr: double
  Yerr: double
  Ttrue: double[] }

type Y = {Tobs: double[]; Y: double[][]}

let CalcSpProb w t sp =
  let z = (t - w.Topt.[sp]) / w.Tbreadth.[sp]
  w.MaxProb.[sp] * exp (- z*z)

[<Fun>]
let prior () =
{ Topt = [| for j in 0..Nspecies-1 → random(Uniform(0.1, 50.0)) |]
  Tbreadth = [| for j in 0..Nspecies-1 → random(Uniform(0.1, 50.0)) |]
  MaxProb = [| for j in 0..Nspecies-1 → random(Uniform(0.1, 1.0)) |]
  Terr = random(Uniform(0.1, 10.0))
  Yerr = random(Uniform(0.01, 0.5))
}

```

```

Ttrue = [| for i in 0..Nsamples-1 → random(Uniform(5.0, 30.0))|] }

let samplesR = [|0..Nsamples-1|]
let speciesR = [|0..Nspecies-1|]

[<Fun>]
let model w =
  let tobs = [| for i in samplesR → random(Gaussian(w.Ttrue.[i], w.Terr)) |]
  let y = [| for i in samplesR →
            [| for j in speciesR →
              let p = CalcSpProb w w.Ttrue.[i] j
              random(Gaussian(p, w.Yerr)) |] |]

  { Tobs = tobs
    Y = y }

// the generated log probability density function for model
let logPdf =
  fun w ys →
    let y = ys.Y
    let tobs = ys.Tobs
    let i0 = samplesR
    (logprodBy((fun i1 →
                let i=samplesR.[i1]
                Log(pdf_Gaussian(w.Ttrue.[i],w.Terr,tobs.[i1])),
                samplesR.Length)) +
    (let i2 = samplesR
      logprodBy((fun i3 →
                let i=samplesR.[i3]
                let i4=speciesR
                logprodBy((fun i5 →
                        let j=speciesR.[i5]
                        Log(pdf_Gaussian(CalcSpProb w (w.Ttrue.[i] j,w.Yerr,y.[i3].[i5])),
                        speciesR.Length)),
                samplesR.Length))

```