

Symbolic conditioning of arrays in probabilistic programs

PRAVEEN NARAYANAN AND CHUNG-CHIEH SHAN

ACM Reference format:

Praveen Narayanan and Chung-chieh Shan. 2017. Symbolic conditioning of arrays in probabilistic programs. *Proc. ACM Program. Lang.* 1, 1, Article 1 (June 2017), 24 pages.
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

ABSTRACT

Probabilistic programming systems automate machine learning inference methods and make them reusable. A common component of inference is the step of *conditioning*, which is addressed most generally by the measure-theoretic technique of *disintegration*. Recent work by Shan and Ramsey [2017] automates disintegration as a symbolic program transformation, thereby making inference more modular. This algorithm, however, is limited to conditioning a scalar variable. As a step towards tackling realistic machine learning problems, we have extended disintegration to symbolically condition arrays in probabilistic programs. The extended algorithm implements *lifted* disintegration, where repetition is treated symbolically and without unrolling loops. The technique uses a language of *index variables* for tracking expressions at various array levels. We find that the method works well for arbitrarily-sized arrays of independent random choices, with the conditioning step taking time linear in the number of indices needed to select an element.

1 INTRODUCTION

The guiding principle of machine learning is that of improving from experience. By using data to enhance our knowledge of the world, we can perform more effectively at tasks such as predicting the topics of documents or calculating the chance of heart attacks. The task at hand often determines the formalism for mechanically learning from data, and *Bayesian* techniques in particular are useful across a large class of problems.

Bayesian analyses consist of *modeling* followed by *inference*. In the modeling step we use probability distributions to express our beliefs prior to observing any data. In the inference step we answer questions posterior to observing data.

For example, say we want to analyze the running time of an iterative algorithm. We might model the algorithm as taking some initialization time along with some constant time per iteration. A *Bayesian linear regression* model would suit this analysis.

Figure 1 shows a Bayesian linear regression model. The goal of linear regression is to fit a line—characterized by its slope a and intercept b —through a collection of points. In the Bayesian setting, we develop this by:

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2017 Association for Computing Machinery.

2475-1421/2017/6-ART1 \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

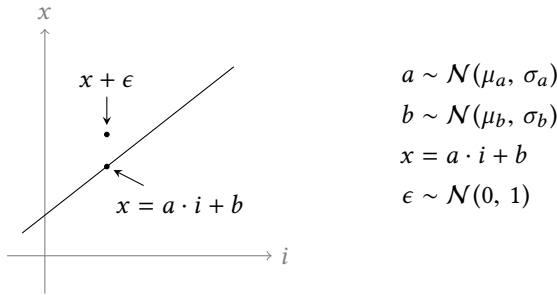


Fig. 1. A canonical sample (left) of a regression model (right) for some values of a , b , i , and ϵ . We observe $x + \epsilon$ and want to infer a distribution over a and b using this model.

- stating our prior beliefs about the slope and intercept—modeled here as Gaussian distributions with fixed means (μ_a, μ_b) and standard deviations (σ_a, σ_b) —and
- acknowledging the noise inherent in our measurements—modeled by ϵ here as a standard Gaussian distribution.

The variables a and b map to the time-per-iteration and the initialization time respectively. The variable i maps to the number of iterations for which we run our algorithm. Ideally we would run this experiment for various values of i —and this motivates our paper—but it is nevertheless instructive to understand how a Bayesian analysis is informed by a single datapoint.

Thus our analysis moves on to the inference step, where we want to answer questions about the line $(a \cdot i + b)$ in light of a noisy measurement of the running time $(x + \epsilon)$. Inference techniques require three quantities – the model, a data-introducing condition such as $x + \epsilon = t$, and the question of interest such as “What is the expected value of a ?”, or “What is the most likely value of b ?”. Inference consists of *conditioning*, which combines the model and the condition, followed by *querying*, which processes the question of interest.

Although modeling and inference are separate steps in theory, implementations of Bayesian techniques intertwine the two in practice. To classify documents by their topics, for example, the prevalent approach is to customize an inference algorithm towards a particular model [Resnik and Hardisty 2009]. Consequently, a small change to the task on paper (adjusting the model) can lead to a disproportionately large change to the implementation in code.

Probabilistic programming is about decoupling models and inference. With the advent of probabilistic programming languages, we are able to reuse an inference method on different models [Goodman et al. 2008; Milch et al. 2007], and try out different inference methods on a fixed model [Mansinghka et al. 2014; Narayanan et al. 2016; Wood et al. 2014]. This separation frees the implementor to explore trade-offs between intricate models and sophisticated inference methods.

Similarly, although conditioning and querying are separate steps of inference in theory, implementations of inference intertwine the two in practice. Shan and Ramsey [2017] show that we can decouple conditioning and querying using the measure-theoretic technique of *disintegration*, as advocated by Chang and Pollard [1997]. We adopt this approach and implement a modular inference technique where the first part of inference performs conditioning and returns a distribution that will be separately queried by the second part of inference.

Currently, modular inference techniques that separate out conditioning can only be used for small quantities of observational data. The existing methods simply unroll larger structures [Gehr

et al. 2016; Shan and Ramsey 2017], resulting in severe code growth. We desire an automatic conditioning method that observes multiply-indexed arrays of data symbolically and without unrolling. In this paper we achieve this goal via a technique that extends disintegration to handle arrays.

Before we describe array disintegration (in sections 4 and 5), we set up disintegration on a language without arrays (in sections 2 and 3). First, we use our running example (in the rest of this section) to describe modeling and inference in the context of probabilistic programming languages.

1.1 Models as probabilistic programs

Probabilistic programs express Bayesian models. Let us revisit the Bayesian linear regression model, which we can express as the following probabilistic program:

$$\begin{aligned} \text{blr} &: \mathbb{M}(\mathbb{R} \times (\mathbb{R} \times \mathbb{R})) & (1) \\ \text{blr} &= \text{do } \{a \leftarrow \text{normal } \mu_a \sigma_a; \\ & \quad b \leftarrow \text{normal } \mu_b \sigma_b; \\ & \quad x \leftarrow \text{return } (a \cdot i + b); \\ & \quad \epsilon \leftarrow \text{normal } 0 \ 1; \\ & \quad \text{return } (x + \epsilon, (a, b))\} \end{aligned}$$

This program is written in a probabilistic programming language that we call core Hakaru. Figure 2 defines its syntax.

Real numbers $r \in \mathbb{R}$	Variables x
Terms $e, m, M ::= x \mid r \mid -e \mid e + e \mid e \cdot e \mid \text{normalD } r \ r \ e \mid () \mid (e, e) \mid \text{fst } e \mid \text{snd } e$	
	$\mid \text{lebesgue} \mid \text{normal } r \ r \mid \text{return } e \mid \text{do } \{x \leftarrow m; M\} \mid \text{do } \{\text{factor } e; M\}$
Types $\alpha, \beta ::= \mathbb{1} \mid \mathbb{R} \mid \alpha \times \beta \mid \mathbb{M} \alpha$	

Fig. 2. Syntactic forms of core Hakaru

The program in equation (1) is composed of nested $\text{do } \{x \leftarrow e; e\}$ constructs, for which we use this shorthand notation:

$$\text{do } \{x_1 \leftarrow m_1; \text{do } \{x_2 \leftarrow m_2; \dots; \text{do } \{x_n \leftarrow m_n; M\} \dots\} \} \equiv \text{do } \{x_1 \leftarrow m_1; x_2 \leftarrow m_2; \dots; x_n \leftarrow m_n; M\}. \quad (2)$$

Core Hakaru programs are written in a *mochastic* (monadic-stochastic) style. Just as a monadic *bind* operator assigns a variable to the result of an effectful computation, the core Hakaru “ \leftarrow ” operator defines a variable distributed according to the expression on the right hand side. As a special case, the combination of \leftarrow and the **return** primitive corresponds to *let-binding* (written as = in the original model).

Terms that denote measures (which generalize distributions) have type \mathbb{M} in core Hakaru. Figure 3 describes the typing rules for terms of measure type; the omitted rules are standard.

$$\begin{array}{c} \frac{}{\Gamma \vdash \text{lebesgue} : \mathbb{M} \ \mathbb{R}} \quad \frac{}{\Gamma \vdash \text{normal } r \ r : \mathbb{M} \ \mathbb{R}} \quad \frac{\Gamma \vdash e : \alpha}{\Gamma \vdash \text{return } e : \mathbb{M} \ \alpha} \\ \frac{\Gamma \vdash m : \mathbb{M} \ \alpha \quad \Gamma, x : \alpha \vdash M : \mathbb{M} \ \beta}{\Gamma \vdash \text{do } \{x \leftarrow m; M\} : \mathbb{M} \ \beta} \quad \frac{\Gamma \vdash e : \mathbb{R} \quad \Gamma \vdash M : \mathbb{M} \ \alpha}{\Gamma \vdash \text{do } \{\text{factor } e; M\} : \mathbb{M} \ \alpha} \end{array}$$

Fig. 3. Typing rules for terms of measure type

The final line of the **blr** program in equation (1) returns the noisy reading of the runtime paired with the corresponding slope and intercept. This agrees with the type $\mathbb{M} (\mathbb{R} \times (\mathbb{R} \times \mathbb{R}))$, which states that the program denotes a measure over nested tuples of real numbers.

We can say that the program in equation (1) represents the *joint measure* $P(x + \epsilon, a, b)$. The joint decomposes into the *prior* $P(a, b)$, and the *likelihood* $P(x + \epsilon \mid a, b)$:

$$P(x + \epsilon, a, b) = P(a, b) \otimes P(x + \epsilon \mid a, b). \quad (3)$$

Let us define each element of the right-hand-side of equation (3).

The prior $P(a, b)$ represents what we believe, before observing data, about how a and b are distributed. In **blr** we express this as a program of type $\mathbb{M} (\mathbb{R} \times \mathbb{R})$:

$$P(a, b) = \text{do } \{a \leftarrow \text{normal } \mu_a \sigma_a; b \leftarrow \text{normal } \mu_b \sigma_b; \text{return } (a, b)\}. \quad (4)$$

The likelihood $P(x + \epsilon \mid a, b)$ is how we believe a measurement of the running time depends on a and b . It is a *conditional measure*, and we can express it as function from (a, b) to a program of type $\mathbb{M} \mathbb{R}$:

$$P(x + \epsilon \mid a, b) = \lambda(a, b). \text{do } \{x \leftarrow \text{return } (a \cdot i + b); \epsilon \leftarrow \text{normal } 0 \ 1; \text{return } (x + \epsilon)\}. \quad (5)$$

To create the joint measure, the factors in equation (3) must be linked together by an operation on measures that is notated \otimes and pronounced “bind x”:

$$\begin{aligned} \otimes : \mathbb{M} \alpha \rightarrow (\alpha \rightarrow \mathbb{M} \beta) \rightarrow \mathbb{M} (\beta \times \alpha) \\ m \otimes f = \text{do } \{a \leftarrow m; b \leftarrow f \ a; \text{return } (b, a)\}. \end{aligned}$$

As this definition illustrates, monadic bind is an important construct for composing larger programs out of smaller building blocks in general purpose languages, and core Hakaru applies the same principle in the probabilistic setting.

The grammar in figure 2 contains some constructs that we have not yet discussed. There are routine operators for arithmetic ($-$) and destructors for pair types (**fst**, **snd**). The **lebesgue** construct defines the Lebesgue measure—the unique-up-to-scale translation invariant measure on \mathbb{R} .

To understand **factor**, we note that measures can often be re-written in terms of **lebesgue** and **factor**. Formally, this means that for some measures m , there is a function $f : \mathbb{R} \rightarrow \mathbb{R}^+$ such that:

$$m \equiv \text{do } \{x \leftarrow \text{lebesgue}; \text{factor } (f \ x); \text{return } x\}. \quad (6)$$

This is what it means for m to have a *density f with respect to the Lebesgue measure*.

We can use **factor** in conjunction with **lebesgue** to express a large class of measures over the reals. For example, core Hakaru provides the **normalD** $\mu \sigma \ x$ construct, denoting the density of $\mathcal{N}(\mu, \sigma)$ at x with respect to **lebesgue**. This gives us an alternate definition for the Gaussian distribution:

$$\begin{aligned} \text{normal} : \mathbb{R} \rightarrow \mathbb{R} \rightarrow \mathbb{M} \mathbb{R} \\ \text{normal } \mu \sigma \equiv \text{do } \{x \leftarrow \text{lebesgue}; \text{factor } (\text{normalD } \mu \sigma \ x); \text{return } x\}. \end{aligned} \quad (7)$$

For more intuition about **factor** we can interpret a core Hakaru program as an importance sampler for its underlying measure. Now each sample from an importance sampler is paired with a non-negative weight. This weight is 1 for the primitive measures provided by the language, and is multiplied by the argument to **factor** in programs that use the construct.

When run as a sampler, the program in equation (7) produces real-valued samples, but how do we know that they follow the correct Gaussian distribution? The **factor** construct assures this by re-weighting the sample by the density of $\mathcal{N}(\mu, \sigma)$ at that value. While each real value is as likely to show up as any other (as they are all being roughly drawn from **lebesgue**), the importance weights will reshape the distribution of samples to approximate the required Gaussian distribution.

We note that sampling from programs that use **lebesgue** and **factor** is a good approach to intuition but a bad one for inference. This is largely because of the impossibility of sampling from the Lebesgue measure. We can handle this problem separately with a transformation that recognizes densities and converts a program that uses **lebesgue** to one that uses only primitive measures. The Hakaru probabilistic programming system provides such a simplification transformation [Carette and Shan 2016], typically used in the inference step to obtain more efficient samplers.

1.2 Inference on probabilistic programs

The modeling constructs shown thus far are largely consistent across probabilistic programming systems; the differences lie in the inference algorithms.

Typically we observe data for a subset of the variables in our model. While the **blr** model describes a measure over possible values of the variables a , b , and $x + \epsilon$, we observe data only for the variable $x + \epsilon$. The observed data restricts the space of possibilities, and we are interested in how a and b behave in this restricted space. In Bayesian terminology, we are interested in the *posterior distribution*.

Inference is about computing queries on the posterior. Now, many queries on a distribution can be approximately answered given a set of samples from that distribution. Thus we can intuitively perform inference by generating samples according to the model, throwing away those that are incompatible with our observations (by assigning to such samples a weight of 0), and computing our queries on the resultant set. This technique is also known as *rejection sampling*.

Although rejection sampling is useful for pedagogy, the inefficiency of wasted work makes it not so useful in practice. Additionally, the method cannot be used for observing continuous variables, when there is zero probability of matching our observations exactly.

In practice, probabilistic programming systems implement more efficient inference techniques to obtain and query posterior distributions. Systems have adapted approximate methods based on Markov Chain Monte Carlo (MCMC) [Gelman et al. 2015; Milch et al. 2007; Wingate et al. 2011] and variational inference [Kucukelbir et al. 2015; Wingate and Weber 2013], as well as exact techniques based on abstract interpretation, factor graphs, data flow analysis, and Bayes nets [Claret et al. 2013; McCallum et al. 2009; Minka et al. 2014].

Although there are various inference techniques, across them lies the same conceptual step of first obtaining the posterior distribution. This step is called *conditioning*. Every inference method first conditions the joint to obtain the posterior, and then converts the posterior into a unique representation suitable for querying.

Thus, a modular approach to inference would decouple the conditioning and the querying steps. Most systems however do not harness this modularity. Inference methods typically take three inputs—model, data, and query of interest—and repeat the work needed to perform conditioning. For example, consider a function **mh** that implements *Metropolis-Hastings* MCMC sampling (described in its most general form by Tierney [1998]) and has the following type:

$$\mathbf{mh} : \mathbb{M} (\alpha \times \beta) \rightarrow \alpha \rightarrow ([\beta] \rightarrow \gamma) \rightarrow \gamma. \quad (8)$$

This function takes a joint measure (the model) as its first argument and constrains its first dimension with an observation (the data) given by the second argument. The third argument is a function used for querying a set of samples (represented as $[\beta]$) and producing some output γ . Implicitly, this inference method performs conditioning, obtains the posterior *as a set of samples*, and then queries that set.

Suppose that in our running example we observe $x + \epsilon = t$, and that we want to compute the *expected value* of a given this observation. Given a function **expect** of type $[\mathbb{R} \times \mathbb{R}] \rightarrow \mathbb{R}$, we could

use Metropolis-Hastings sampling to approximate this value:

$$a_{\text{expected}} = \text{mh blr } t \text{ expect.} \quad (9)$$

Inference based on a slightly different method like *Gibbs sampling* [explained by Casella and George 1992; Resnik and Hardisty 2009] would have a type very similar to the one for **mh**. A **gibbs** method could even share the code used by **mh** to implement conditioning—after all both methods query a set of samples. Inference based on an exact method such as abstract interpretation, however, would need a different representation of the posterior, and might end up significantly duplicating the work done to perform conditioning.

We thus seek a reusable conditioning tool that produces the posterior in a format that can be manipulated by different inference methods. We choose this format to be a program in core Hakaru, and implement a conditioning program transformation that when given data takes us from a core Hakaru program representing the model to one representing the posterior. Let us specify this transformation formally and look at examples of what it would produce.

2 CONDITIONING VIA DISINTEGRATION

Let us specify a conditioning transformation on core Hakaru programs. To perform conditioning on our running example, we seek a function that, given the joint $P(x + \epsilon, a, b)$, produces the posterior $P(a, b \mid x + \epsilon)$. We could motivate our specification by applying a common intuition about conditioning in the form of Bayes' rule:

$$P(x + \epsilon) \cdot P(a, b \mid x + \epsilon) = P(x + \epsilon, a, b). \quad (10)$$

This rule, especially when rearranged so that the term $P(x + \epsilon)$ (called the *marginal* probability) appears on the right hand side as the denominator, tells us how to compute a posterior probability given the joint. Since we care about measures and not just probabilities, using \cdot does not make sense. We revise this equation to use \otimes :

$$P(x + \epsilon) \otimes P(a, b \mid x + \epsilon) = P(x + \epsilon, a, b). \quad (11)$$

We could say that equation (11) forms a *specification* for a conditioning transformation. If we build a tool that, when given the joint measure $P(x + \epsilon, a, b)$ produces the marginal $P(x + \epsilon)$ and conditional measure $P(a, b \mid x + \epsilon)$ such that equation (11) is satisfied, then the tool implements conditioning. When this relation holds, we say that $P(x + \epsilon)$ and $P(a, b \mid x + \epsilon)$ form a disintegration of $P(x + \epsilon, a, b)$. In such situations we can expand the definition of \otimes , and express the specification in core Hakaru:

$$\text{do } \{t \sim P(x + \epsilon); p \leftarrow P(a, b \mid x + \epsilon) t; \text{return } (t, p)\} \equiv P(x + \epsilon, a, b). \quad (12)$$

Now, if we are able to find a marginal and a conditional measure such that equation (11) (and thus equation (12)) is satisfied, and if we can assume that the marginal has a density f with respect to **lebesgue**, then we are better off building a transformation—let us call it **disintegrate**—that satisfies a different specification:

$$\text{do } \{t \sim \text{lebesgue}; p \leftarrow \text{disintegrate } m t; \text{return } (t, p)\} \equiv m \quad (13)$$

This seems strange, but equation (13) falls out from our assumptions:

$$P(x + \epsilon, a, b) \equiv [\text{equation (12)}] \quad (14)$$

$$\mathbf{do} \{t \sim P(x + \epsilon); p \sim P(a, b \mid x + \epsilon) t; \mathbf{return} (t, p)\} \equiv [\text{definition of density, equation (6)}] \quad (15)$$

$$\mathbf{do} \{t \sim \mathbf{do} \{t \sim \mathbf{lebesgue}; \mathbf{factor} (f t); \mathbf{return} t\}; p \sim P(a, b \mid x + \epsilon) t; \mathbf{return} (t, p)\} \equiv [\text{monad laws}] \quad (16)$$

$$\mathbf{do} \{t \sim \mathbf{lebesgue}; p \sim \mathbf{do} \{\mathbf{factor} (f t); P(a, b \mid x + \epsilon) t\}; \mathbf{return} (t, p)\}.$$

Our assumption that the marginal has a density with respect to **lebesgue** is not unrealistic—most marginals over \mathbb{R} that arise naturally in probabilistic programs have this property.

Additionally, in equation (16), the right hand side of the expression $p \sim \dots$ is not the true posterior, but rather something *proportional* to it. Fortunately, this is adequate for a large class of querying algorithms [Tierney 1998]. We thus require that **disintegrate** produces this term generally for all input measures m , giving us the specification in equation (13).

The specification in equation (13) also hints at the type for **disintegrate**:

$$\mathbf{disintegrate} : \mathbb{M} (\mathbb{R} \times \beta) \rightarrow \alpha \rightarrow \mathbb{M} \beta. \quad (17)$$

Disintegrations do not always exist, and when they do they are not unique. Ackerman et al. [2011] show that a disintegrator must fail to find solutions for a language that can express all and only *computable* distributions. Core Hakaru is not expressive enough to represent all computable distributions, so it is unclear what disintegrations can be expressed. In our system we express failure and non-uniqueness of disintegration by having the type of **disintegrate** return a *list* of possible answers. For ease of exposition we omit this detail in the paper.

For the **blr** example, our implemented **disintegrate** produces the following program:

$$\mathbf{disintegrate} \mathbf{blr} \equiv \lambda t. \mathbf{do} \{a \sim \mathbf{normal} \mu_a \sigma_a; b \sim \mathbf{normal} \mu_b \sigma_b; \mathbf{factor} (\mathbf{normalD} \ 0 \ 1 \ (t - (a \cdot i + b))); \mathbf{return} (a, b)\} \quad (18)$$

To understand this result, we can think of disintegration as “slicing” a measure. When we disintegrate a measure we slice it at the values of its observed variables (such as $x + \epsilon$), giving rise to a measure over the remaining variables (such as a and b).

Thus the result of calling **disintegrate** here is a function, and it represents the posterior measure over a and b given $x + \epsilon = t$. The input to the function is a real value t representing the data. The body of the function differs from the **blr** model in two ways. First, it is a distribution over a and b alone. Second, the bindings for x and ϵ are replaced by a **factor** expression that uses t .

In general, **disintegrate** replaces the binding for an observed variable with a **factor** involving the observed value. This factor comes from the probability mass assigned to the observed value by the likelihood (equation (5) in our running example). We can interpret this factor as re-weighting the priors over a and b , i.e., updating our beliefs after observing data.

Disintegrations that satisfy the specification in equation (13) represent posterior distributions [Chang and Pollard 1997; Shan and Ramsey 2017]. The **blr** example satisfies this specification:

```

do {t ~ lebesgue;           ≡ do {t ~ lebesgue;
  p ~ disintegrate blr t;    p ~ do {a ~ normal μa σa;
  return (t, p)}           b ~ normal μb σb;
                           factor (normalD 0 1 (t - (a * i + b)));
                           return (a, b);
                           return (t, p)}
                           ≡ do {t ~ lebesgue;
                              a ~ normal μa σa;
                              b ~ normal μb σb;
                              factor (normalD 0 1 (t - (a * i + b)));
                              return (t, (a, b))}
                           ≡ do {a ~ normal μa σa;
                              b ~ normal μb σb;
                              t ~ normal (a * i + b) 1;
                              return (t, (a, b))}
                           ≡ do {a ~ normal μa σa;
                              b ~ normal μb σb;
                              x ~ return (a * i + b);
                              ε ~ normal 0 1;
                              return (x + ε, (a, b))}
                           ≡ blr .

```

This paper is about extending disintegration to handle arrays. To motivate and describe this extension, we first implement our version of scalar disintegration based on the work of Shan and Ramsey [2017].

3 IMPLEMENTING DISINTEGRATION

We now detail our implementation of Shan and Ramsey’s disintegration technique. We implement disintegration as a program transformation whose input and output languages are both core Hakaru. Internally however, our disintegrator works with the extended language shown in figure 4.

Real numbers $r \in \mathbb{R}$	Variables x	Locations \bar{x}
Bindings	$b ::= \bar{x} \leftarrow m \mid \mathbf{factor} \ e$	
Heap	$h ::= [b; \dots; b]$	
Atomic terms	$u ::= x \mid -u \mid u + e \mid e + u \mid u * e \mid e * u \mid \mathbf{normalD} \ r \ r \ u \mid \mathbf{fst} \ u \mid \mathbf{snd} \ u$	
Head normal forms	$v ::= u \mid r \mid \mathbf{lebesgue} \mid \mathbf{normal} \ r \ r \mid \mathbf{return} \ e \mid \mathbf{do} \ \{x \leftarrow m; M\} \mid \mathbf{do} \ \{\mathbf{factor} \ e; M\} \mid () \mid (e, e)$	
Terms	$e, m, M ::= v \mid \bar{x} \mid -e \mid e + e \mid e * e \mid \mathbf{normalD} \ r \ r \ e \mid \mathbf{fst} \ e \mid \mathbf{snd} \ e$	
Types	$\alpha, \beta ::= \mathbb{1} \mid \mathbb{R} \mid \alpha \times \beta \mid \mathbb{M} \ \alpha$	

Fig. 4. The internal language of disintegration

The goal of the disintegrator is to constrain the values of observed variables in the input program. We cannot always satisfy the constraints, and in such cases **disintegrate** fails to produce a result. Thus we design our implementation around recording dependencies and propagating constraints on bound variables.

The internal language extends core Hakaru with *heaps*, *bindings*, and *locations*. A *heap* is an ordered collection of *bindings*, where each binding is either a **factor** term or a *location* bound (via \leftarrow)

to a measure expression. Locations (notated \bar{x}) are used by the disintegrator to track variables (notated x) originally bound in the input program. We order heaps with the oldest binding on the left, and the scope of a location follows this order.

To further streamline the recording and propagation of constraints, the internal language reorganizes terms into *atomic terms* and *head normal forms*. Atomic terms represent operations with at least one free variable. This makes observing atomic terms impossible, and disintegration (correctly) fails in these cases.

Disintegration is carried out mathematically by manipulating integrals [Chang and Pollard 1997; Shan and Ramsey 2017], and the disintegration of compound expressions (composed of multiple variables) is defined in terms of a change-of-variables operation. This operation involves differentiation with respect to one of the variables while keeping the others constant. To achieve this we need to evaluate the other variables to a *locally constant* head normal form. Thus we follow Shan and Ramsey [2017] in defining our disintegrator in terms of “forward” evaluating and “backward” constraining functions:

$$\triangleright\triangleright \text{ (“perform”)} \quad : \text{heap} \rightarrow [\mathbb{M} \alpha] \rightarrow (\text{emission} \times \text{heap} \times [\alpha]) \quad (19)$$

$$\triangleright \text{ (“evaluate”)} \quad : \text{heap} \rightarrow [\alpha] \rightarrow (\text{emission} \times \text{heap} \times [\alpha]) \quad (20)$$

$$\triangleleft\triangleleft \text{ (“constrain outcome”)} : \text{heap} \rightarrow [\mathbb{M} \alpha] \rightarrow [\alpha] \rightarrow (\text{emission} \times \text{heap}) \quad (21)$$

$$\triangleleft \text{ (“constrain value”)} \quad : \text{heap} \rightarrow [\alpha] \rightarrow [\alpha] \rightarrow (\text{emission} \times \text{heap}) \quad (22)$$

The four functions are defined in figure 5. The functions $\triangleright\triangleright$ and \triangleright both partially evaluate [Fischer et al. 2011, 2008] their input terms (notated $[\alpha]$) to head normal form (notated $[\alpha]$). Now the disintegrator is non-effectful in that no random numbers are generated while producing the posterior program. Thus for $\triangleright\triangleright$, evaluating measure terms to head normal form means simulating the act of making a random choice. We do this by *emitting* a binding for the measure to be evaluated, and returning a fresh variable that is atomic in the scope of the heap and the term during disintegration. The four functions of the disintegrator may add to but never look at this set of *emissions*.

Emissions occur in the $\triangleright\triangleright$ cases of **lebesgue**, **normal**, and an atomic measure term u . On a **return** measure $\triangleright\triangleright$ evaluates the point at which all probability mass is located.

Another role of $\triangleright\triangleright$ is to record bindings onto the heap. To track a variable x bound by $x \sim m$, $\triangleright\triangleright$ generates a location \bar{x} , stores the binding $\bar{x} \sim m$ on the heap, and replaces x with \bar{x} within its scope. The heap is also used for storing **factor** expressions. In both cases, the disintegrator is *lazy* [Launchbury 1993] in that these expressions are stored without inspection or evaluation.

For most cases, we define \triangleright as a standard partial evaluator. The cases for pair accessors **fst** and **snd** work by evaluating the argument and either recursively evaluating the sub-term or building a piece of code when the argument is atomic. For some operations that define atomic terms in figure 4 ($+$, $-$, and \cdot), we use *smart constructors* that can simplify their arguments in special cases.

When the input term is a location \bar{x} , \triangleright retrieves its binding from the heap and calls $\triangleright\triangleright$ on the associated measure term. The monadic binding for \bar{x} is modified into a deterministic **let**-binding that uses **return**, signalling the act of having made a random choice.

The functions $\triangleleft\triangleleft$ and \triangleleft constrain a term to be the observation t . Constraining a measure translates to computing the probability density (with respect to **lebesgue**) of that measure. Thus for a **normal** $r_1 r_2$ expression $\triangleleft\triangleleft$ stores **factor** (**normalD** $r_1 r_2 t$) on the heap, and for **lebesgue** it returns the heap unchanged. Just like $\triangleright\triangleright$, $\triangleleft\triangleleft$ also records monadic bind and **factor** expressions.

For a **return** term $\triangleleft\triangleleft$ calls \triangleleft on its point of mass. Many cases fail since **return** has no density with respect to **lebesgue** when the point mass is a constant. For the $-$ case we constrain the value of the \triangleleft to be $-t$, as we would when changing the variable of integration. In the case of $+$ we could

$\triangleright\triangleright$ (“perform”) : $heap \rightarrow [\mathbb{M} \alpha] \rightarrow (emission \times heap \times \lfloor \alpha \rfloor)$		
$\triangleright\triangleright h \llbracket m \rrbracket = \text{case } \llbracket m \rrbracket \text{ of}$		
$\llbracket u \rrbracket$	$\rightarrow (\bar{z} \leftarrow u, h, \llbracket \bar{z} \rrbracket)$	where u is atomic, \bar{z} is fresh (23)
$\llbracket \text{lebesgue} \rrbracket$	$\rightarrow (\bar{z} \leftarrow \text{lebesgue}, h, \llbracket \bar{z} \rrbracket)$	where \bar{z} is fresh (24)
$\llbracket \text{normal } r_1 r_2 \rrbracket$	$\rightarrow (\bar{z} \leftarrow \text{normal } r_1 r_2, h, \llbracket \bar{z} \rrbracket)$	where \bar{z} is fresh (25)
$\llbracket \text{return } e \rrbracket$	$\rightarrow \triangleright h \llbracket e \rrbracket$	(26)
$\llbracket \text{do } \{x \leftarrow e_1; e_2\} \rrbracket$	$\rightarrow \triangleright\triangleright [h; \bar{x} \leftarrow e_1] \llbracket e_2 \rrbracket \{x \mapsto \bar{x}\}$	where \bar{x} is fresh (27)
$\llbracket \text{do } \{\text{factor } e_1; e_2\} \rrbracket$	$\rightarrow \triangleright\triangleright [h; \text{factor } e_1] \llbracket e_2 \rrbracket$	(28)
$\llbracket e \rrbracket$	$\rightarrow ([s_1; s_2], h_2, \llbracket v \rrbracket)$	where e is not in head normal form, (29)
$(s_1, h_1, \llbracket m_1 \rrbracket) = \triangleright h \llbracket e \rrbracket, (s_2, h_2, \llbracket v \rrbracket) = \triangleright\triangleright h_1 \llbracket m_1 \rrbracket$		
<hr/>		
\triangleright (“evaluate”) : $heap \rightarrow [\alpha] \rightarrow (emission \times heap \times \lfloor \alpha \rfloor)$		
$\triangleright h \llbracket e \rrbracket = \text{case } \llbracket e \rrbracket \text{ of}$		
$\llbracket v \rrbracket$	$\rightarrow ([], h, \llbracket v \rrbracket)$	where v is in head normal form (30)
$\llbracket \text{fst } e \rrbracket$	$\rightarrow \text{case } \llbracket p \rrbracket \text{ of}$	where $(s_1, h_1, \llbracket p \rrbracket) = \triangleright h \llbracket e \rrbracket,$ (31)
	$\llbracket (e_1, _) \rrbracket \rightarrow ([s_1; s_2], h_2, \llbracket v \rrbracket)$	$(s_2, h_2, \llbracket v \rrbracket) = \triangleright h_1 \llbracket e_1 \rrbracket,$
	$\llbracket u \rrbracket \rightarrow (s_1, h_1, \llbracket \text{fst } u \rrbracket)$	u is atomic
$\llbracket \text{snd } e \rrbracket$	\rightarrow similar to the fst case	(32)
$\llbracket -e \rrbracket$	$\rightarrow (s, h_1, -\llbracket v \rrbracket)$	where $(s, h_1, \llbracket v \rrbracket) = \triangleright h \llbracket e \rrbracket$ (33)
$\llbracket e_1 + e_2 \rrbracket$	$\rightarrow ([s_1; s_2], h_2, \llbracket v_1 \rrbracket + \llbracket v_2 \rrbracket)$	where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h \llbracket e_1 \rrbracket,$ (34)
		$(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h_1 \llbracket e_2 \rrbracket$
$\llbracket e_1 \cdot e_2 \rrbracket$	$\rightarrow ([s_1; s_2], h_2, \llbracket v_1 \rrbracket \cdot \llbracket v_2 \rrbracket)$	where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h \llbracket e_1 \rrbracket,$ (35)
		$(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h_1 \llbracket e_2 \rrbracket$
$\llbracket \text{normalD } r_1 r_2 e \rrbracket$	$\rightarrow (s, h_1, \llbracket \text{normalD } r_1 r_2 v \rrbracket)$	where $(s, h_1, \llbracket v \rrbracket) = \triangleright h \llbracket e \rrbracket$ (36)
$\llbracket \bar{x} \rrbracket$	$\rightarrow (s, [h'_1; \bar{x} \leftarrow \text{return } v; h_2], \llbracket v \rrbracket)$	where $[h_1; \bar{x} \leftarrow e; h_2] = h, (s, h'_1, \llbracket v \rrbracket) = \triangleright\triangleright h_1 \llbracket e \rrbracket$ (37)
<hr/>		
$\triangleleft\triangleleft$ (“constrain outcome”) : $heap \rightarrow [\mathbb{M} \alpha] \rightarrow \lfloor \alpha \rfloor \rightarrow (emission \times heap)$		
$\triangleleft\triangleleft h \llbracket m \rrbracket t = \text{case } \llbracket m \rrbracket \text{ of}$		
$\llbracket u \rrbracket$	$\rightarrow \perp$	where u is atomic (38)
$\llbracket \text{lebesgue} \rrbracket$	$\rightarrow ([], h)$	(39)
$\llbracket \text{normal } r_1 r_2 \rrbracket$	$\rightarrow ([], [h; \text{factor } (\text{normalD } r_1 r_2 t)])$	(40)
$\llbracket \text{return } e \rrbracket$	$\rightarrow \triangleleft h \llbracket e \rrbracket t$	(41)
$\llbracket \text{do } \{x \leftarrow e_1; e_2\} \rrbracket$	$\rightarrow \triangleleft\triangleleft [h; \bar{x} \leftarrow e_1] \llbracket e_2 \rrbracket \{x \mapsto \bar{x}\} t$	where \bar{x} is fresh (42)
$\llbracket \text{do } \{\text{factor } e_1; e_2\} \rrbracket$	$\rightarrow \triangleleft\triangleleft [h; \text{factor } e_1] \llbracket e_2 \rrbracket t$	(43)
$\llbracket e \rrbracket$	$\rightarrow ([s_1; s_2], h_2)$	where e is not in head normal form, (44)
$(s_1, h_1, \llbracket m_1 \rrbracket) = \triangleright h \llbracket e \rrbracket, (s_2, h_2) = \triangleleft\triangleleft h_1 \llbracket m_1 \rrbracket t$		
<hr/>		
\triangleleft (“constrain value”) : $heap \rightarrow [\alpha] \rightarrow \lfloor \alpha \rfloor \rightarrow (emission \times heap)$		
$\triangleleft h \llbracket e \rrbracket t = \text{case } \llbracket e \rrbracket \text{ of}$		
$\llbracket \text{fst } e \rrbracket$	$\rightarrow \text{case } \llbracket p \rrbracket \text{ of}$	where $(s_1, h_1, \llbracket p \rrbracket) = \triangleright h \llbracket e \rrbracket,$ (45)
	$\llbracket (e_1, _) \rrbracket \rightarrow ([s_1; s_2], h_2)$	$(s_2, h_2) = \triangleleft h_1 \llbracket e_1 \rrbracket t,$
	$\llbracket u \rrbracket \rightarrow \perp$	u is atomic
$\llbracket \text{snd } e \rrbracket$	\rightarrow similar to the fst case	(46)
$\llbracket (e_1, e_2) \rrbracket$	$\rightarrow ([s_1; s_2], h_2)$	where $(s_1, h_1) = \triangleleft h \llbracket e_1 \rrbracket (\text{fst } t), (s_2, h_2) = \triangleleft h_1 \llbracket e_2 \rrbracket (\text{snd } t)$ (47)
$\llbracket -e \rrbracket$	$\rightarrow \triangleleft h \llbracket e \rrbracket (-t)$	(48)
$\llbracket e_1 + e_2 \rrbracket$	$\rightarrow ([s_1; s'_1], h'_1)$	where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h \llbracket e_1 \rrbracket, (s'_1, h'_1) = \triangleleft h_1 \llbracket e_2 \rrbracket (t - \llbracket v_1 \rrbracket)$ (49)
	$\sqcup ([s_2; s'_2], h'_2)$	$(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h \llbracket e_2 \rrbracket, (s'_2, h'_2) = \triangleleft h_2 \llbracket e_1 \rrbracket (t - \llbracket v_2 \rrbracket)$
$\llbracket v \rrbracket$	$\rightarrow \perp$	where v is in head normal form (50)
$\llbracket \bar{x} \rrbracket$	$\rightarrow (s, [h'_1; \bar{x} \leftarrow \text{return } t; h_2])$	where $[h_1; \bar{x} \leftarrow e; h_2] = h, (s, h'_1) = \triangleleft\triangleleft h_1 \llbracket e \rrbracket t$ (51)

Fig. 5. The implementation of our disintegrator

go forward on the first operand and constrain the second, or vice versa. This choice leads the disintegrator to provide two possible answers using \sqcup . In general, disintegration is non-deterministic in that it can produce multiple correct posterior distributions [Shan and Ramsey 2017]. In this paper we choose which output to describe based on simplicity.

On a location \bar{x} , \triangleleft retrieves the heap binding and calls $\triangleleft\triangleleft$ on the associated measure term. The location \bar{x} is then let-bound to **return** t , signalling that this variable has been observed.

In summary, disintegration has three failure modes – when constraining an atomic term, when constraining a term that has been evaluated, and when constraining the same term multiple times. The implementation is designed to track these constraints.

3.1 Defining and tracing an algorithm

Given this implementation, we can describe an algorithm to perform conditioning:

$$\begin{aligned} \text{disintegrate} &: \mathbb{M} (\alpha \times \beta) \rightarrow (\alpha \rightarrow \mathbb{M} \beta) \\ \text{disintegrate } m \ t &= \text{let } (e_1, h_1, v) = \triangleright\triangleright [] \ m \\ &\quad (e_2, h_2) = \triangleleft h_1 \ (\text{fst } v) \ t \\ &\quad \text{in do } \{e_1; e_2; h_2; \text{return } (\text{snd } v)\} \end{aligned} \quad (52)$$

The algorithm consists of three steps. Let us disintegrate **blr** and trace these steps:

- (1) First we execute the input program m symbolically by calling $\triangleright\triangleright [] \ m$. This populates the heap with each monadic (and **factor**) binding that we encounter in the syntax tree until we reach the final line of the program (returning a pair). This is shown in figure 6a.
- (2) Next we condition on the first of the pair. Given the type of **disintegrate**, we expect that this is the term to be conditioned. We introduce the observation by invoking $\triangleleft h_1 \ (\text{fst } v) \ t$. This is shown in figure 6b.

We use the notation $[\text{---}||\text{---}]$ as a shorthand denoting the *previous* heap, i.e., we use this when the heap is unchanged between successive lines of the trace. The notation $[\text{---}||\text{---}; b]$ denotes the previous heap concatenated with the binding b .

In this step we observe that $x + \epsilon = t$, and we choose to first evaluate \bar{x} , as seen by the sub-call $\triangleright [\dots] [\bar{x}]$. This in turn triggers evaluation of $\bar{a} \cdot i + \bar{b}$, and we emit bindings for a_z and b_z . We also indicate that \bar{a} and \bar{b} have been evaluated by updating them to let-bindings on the heap. This is informative to any procedures that might try later to constrain either of these locations. The result of evaluating \bar{x} is thus $a_z \cdot i + b_z$.

We use this result to constrain the noise $\bar{\epsilon}$, as seen in the sub-call $\triangleleft [\dots] [\bar{\epsilon}] \ (t - [a_z \cdot i + b_z])$. This introduces a **factor** binding into the heap expressing the density of **normal** 0 1 at the observed value.

- (3) Finally, we construct the posterior distribution by stitching together the emissions, final heap, and the second element of the pair:

```
do { $a_z \leftarrow \text{normal } \mu_a \ \sigma_a;$ 
 $b_z \leftarrow \text{normal } \mu_b \ \sigma_b;$ 
 $\bar{a} \leftarrow \text{return } a_z;$ 
 $\bar{b} \leftarrow \text{return } b_z;$ 
 $\bar{x} \leftarrow \text{return } (a_z \cdot i + b_z);$ 
factor (normalD 0 1 ( $t - (a_z \cdot i + b_z)$ ));
 $\bar{\epsilon} \leftarrow \text{return } (t - (a_z \cdot i + b_z));$ 
return ( $\bar{a}, \bar{b}$ )}
```

In this last step we also convert locations back to variables, beta-reduce, and eliminate unused bindings to obtain the succinct core Hakaru output shown in equation (18).

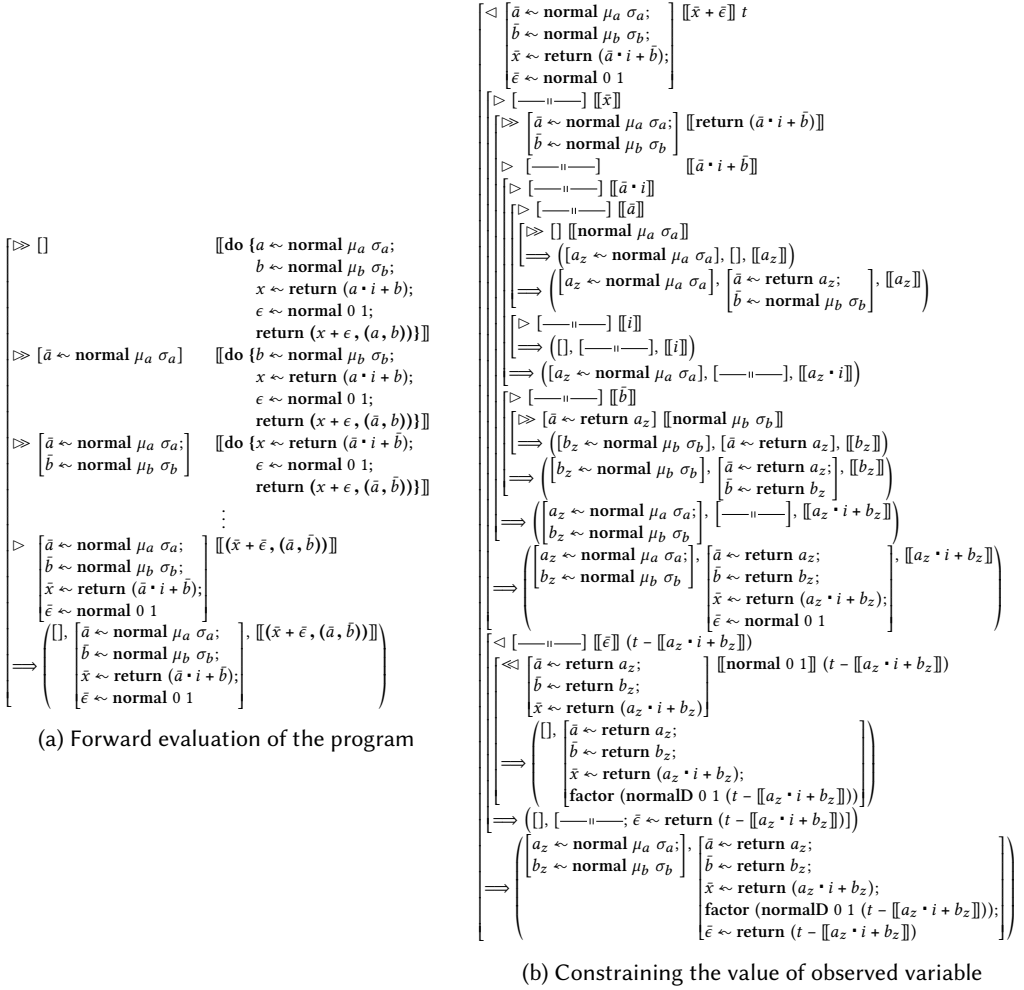


Fig. 6. The first two steps of disintegrating blr

These three steps form the template for every disintegration computed with the functions in figure 5. This includes cases involving tuples of data, which form an important intermediate step as we move towards more realistic examples involving arrays.

Making two measurements. We can develop the `blr` example so that we make *two* noisy measurements of the running time of our imagined algorithm. We can imagine running our experiment on *two* values of i_1 and i_2 which specify the number of iterations:

```

do { a ~ normal μa σa;
    b ~ normal μb σb;
    y1 ~ normal (a · i1 + b) 1;
    y2 ~ normal (a · i2 + b) 1;
    return ((y1, y2), (a, b)) }

```

(53)

Here we model the slope a and intercept b as we did in equation (1). For ease of exposition, we directly model the noisy measurements as y_1 and y_2 , making use of the fact that:

$$\text{do } \{x \leftarrow \text{return } (a \cdot i + b); \epsilon \leftarrow \text{normal } 0 \ 1; \text{return } (x + \epsilon)\} \equiv \text{normal } (a \cdot i + b) \ 1. \quad (54)$$

We can trace the algorithm as we did in the one-measurement case:

- (1) We populate the heap as we traverse the term. This is shown in figure 7a.
- (2) Observing a pair of measurements amounts to doing two observations— $\triangleleft [\dots][[\bar{y}_1]]$ (**fst** t) and $\triangleleft [\dots][[\bar{y}_2]]$ (**snd** t) in sequence.

This is depicted in figure 7b. Although this trace looks smaller than the one in figure 6b, it is only because we reformulated the model using the identity in equation (54). Without this reformulation we would be dealing with $(x_1 + \epsilon_1)$ and $(x_2 + \epsilon_2)$ instead of y_1 and y_2 , leading to roughly twice as many function calls when compared to figure 6b.

- (3) The final output follows the same structure as before, containing a **factor** term for each noisy measurement:

$$\begin{aligned} &\text{do } \{a \leftarrow \text{normal } \mu_a \ \sigma_a; && (55) \\ &\quad b \leftarrow \text{normal } \mu_b \ \sigma_b; \\ &\quad \text{factor } (\text{normalD } (a \cdot i_1 + b) \ 1 \ (\text{fst } t)); \\ &\quad \text{factor } (\text{normalD } (a \cdot i_2 + b) \ 1 \ (\text{snd } t)); \\ &\quad \text{return } (a, b)\} \end{aligned}$$

Observing a variable removes its binding and multiplies a weight into the measure. In this case we observe two variables, and the total weight is a product of the weights occurring in the two **factor** expressions.

3.2 The problem with unrolling

As touched upon by the trace in figure 7b, the disintegrator *unrolls* any structures that are observed in the input program. This creates inefficiencies that are exacerbated with increasing data.

In the model in equation (53) we construct a pair out of two separately bound noisy values; in a sense this input program is already unrolled. However, the output of the disintegrator would not change if we instead generated a pair of noisy measurements using a function (such as \otimes):

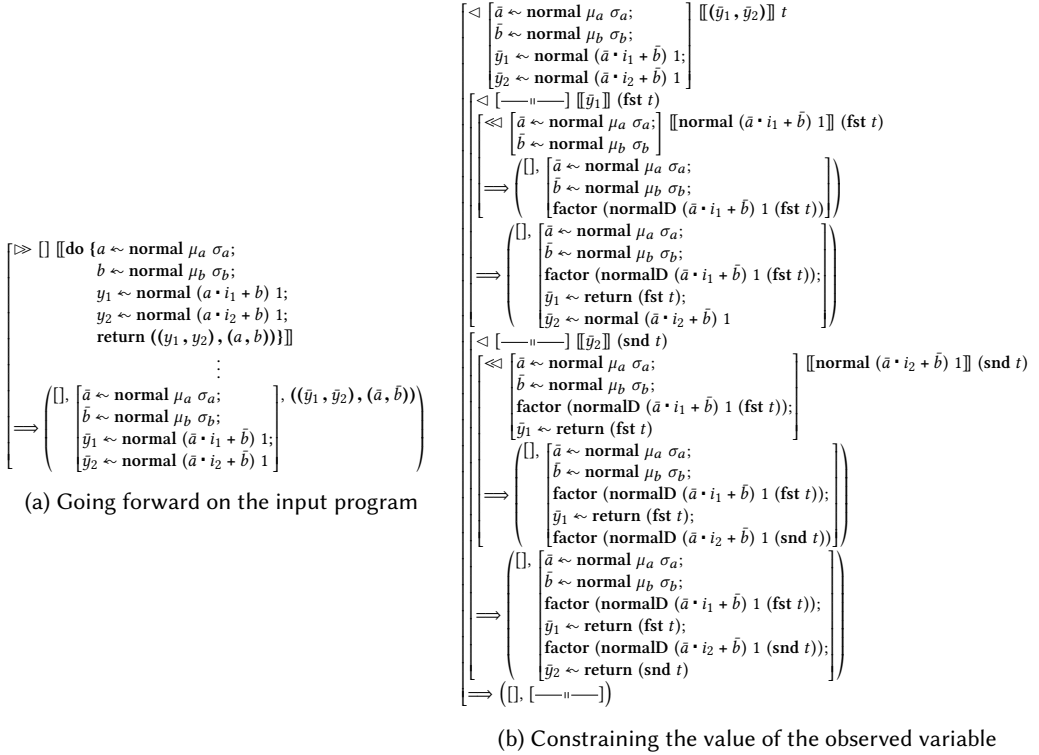
$$\begin{aligned} &\text{do } \{a \leftarrow \text{normal } \mu_a \ \sigma_a; && (56) \\ &\quad b \leftarrow \text{normal } \mu_b \ \sigma_b; \\ &\quad y \leftarrow (\text{normal } (a \cdot i_1 + b) \ 1) \otimes (\lambda_. \text{normal } (a \cdot i_2 + b) \ 1); \\ &\quad \text{return } (y, (a, b))\} \end{aligned}$$

The disintegrator in figure 5 would produce a result identical to the one seen in equation (55)—the pair of noisy values would be unrolled and the output program would be flattened in relation to the input model.

Ideally we'd make several measurements to estimate the slope and intercept in **blr**. We might use tuples or arrays of arbitrary length to store these measurements, and the earlier method of unrolling the pair structure will have to be generalized to these arbitrary-length structures. Such an approach is conceptually sound; systems such as PSI [Gehr et al. 2016] unroll their arrays.

We could encode arrays in core Hakaru using nested tuples: $(x + \epsilon_1, (x + \epsilon_2, (x + \epsilon_3, \dots)))$. The disintegrator in figure 5 can already handle such nested tuples, but producing a disintegration will take time and space linear in the size of the input, i.e., in the number of scalar elements in the array. This is not desirable, considering that we expect to condition our models on large quantities of data.

The problem seems especially grim when we consider deeper structures, i.e., when we move from singly-indexed to multiply-indexed arrays. Whereas an array of n numbers takes $O(n)$ time

Fig. 7. The first two steps of disintegrating two-measurement **blr**

and space to disintegrate, an array of n arrays each of n numbers takes $O(n^2)$ time and space. In a modular approach to inference, this increase in the size of the output program would adversely affect any transformation that we may wish to use downstream from **disintegrate**.

This motivates a new technique for disintegrating probabilistic programs that contain arrays of stochastic variables, one that manipulates arrays symbolically and without unrolling. First, we extend the language of programs to allow the expression and manipulation of arrays.

4 AN ARRAY PROBABILISTIC LANGUAGE

Real numbers $r \in \mathbb{R}$	Natural numbers $n \in \mathbb{N}$	Variables x	
Terms $e, l, m, M ::= x \mid r \mid n \mid -e \mid e + e \mid e \cdot e \mid \text{normalD } r \ r \ e \mid e = e \mid () \mid (e, e) \mid \text{fst } e \mid \text{snd } e$			
$\mid \text{counting } \mid \text{lebesgue } \mid \text{normal } r \ r \ \mid \text{return } e \mid \text{do } \{ x \leftarrow m; M \} \mid \text{do } \{ \text{factor } e; M \}$			
$\mid \text{array } l \ x. \ e \mid \text{idx } e \ e \mid \text{plate } l \ x. \ m \mid \text{size } e$			
Types	$\alpha, \beta ::= \mathbb{1} \mid \mathbb{R} \mid \mathbb{N} \mid \alpha \times \beta \mid \langle \alpha \rangle \mid \mathbb{M} \ \alpha$		

Fig. 8. Core Hakaru extended with array, plate, idx, and size

We extend core Hakaru with constructs for writing array probabilistic programs. Figure 8 describes this extended language, while figure 9 shows the typing rules for array constructs. We can

$$\begin{array}{c}
\frac{}{\Gamma \vdash \mathbf{counting} : \mathbb{M} \mathbb{N}} \quad \frac{\Gamma \vdash l : \mathbb{N} \quad \Gamma, x:\mathbb{N} \vdash m : \mathbb{M} \alpha}{\Gamma \vdash \mathbf{plate} \ l \ x. \ m : \mathbb{M} \langle \alpha \rangle} \quad \frac{\Gamma \vdash l : \mathbb{N} \quad \Gamma, x:\mathbb{N} \vdash e : \alpha}{\Gamma \vdash \mathbf{array} \ l \ x. \ e : \langle \alpha \rangle} \\
\frac{\Gamma \vdash e_1 : \langle \alpha \rangle \quad \Gamma \vdash e_2 : \mathbb{N}}{\Gamma \vdash \mathbf{idx} \ e_1 \ e_2 : \alpha} \quad \frac{\Gamma \vdash e : \langle \alpha \rangle}{\Gamma \vdash \mathbf{size} \ e : \mathbb{N}} \quad \frac{\Gamma \vdash e_1 : \mathbb{N} \quad \Gamma \vdash e_2 : \mathbb{N}}{\Gamma \vdash e_1 = e_2 : \mathbb{N}}
\end{array}$$

Fig. 9. Typing rules for the extensions to core Hakaru

understand this language by example, where we update **blr** to make k noisy measurements of the running time:

$$\begin{array}{l}
\mathbf{do} \{ a \sim \mathbf{normal} \ \mu_a \ \sigma_a; \\
\quad b \sim \mathbf{normal} \ \mu_b \ \sigma_b; \\
\quad y \sim \mathbf{plate} \ k \ i. \ (\mathbf{normal} \ (a \cdot 100 \cdot i + b) \ 1); \\
\mathbf{return} \ (y, (a, b)) \}
\end{array} \tag{57}$$

In this model we use a new construct: **plate**. This is a primitive used for expressing loops containing stochastic computation. The construct derives its name from “plate notation”, a shorthand used in the literature on graphical models and Bayesian inference to represent variables that repeat.

The output of **plate** is an array that collects the result of each loop iteration. The first argument is the size of this array, i.e., the number of loop iterations. The second argument is a variable (of type \mathbb{N}) representing an index for each iteration. This variable takes scope over the third argument, which represents the body of the loop. The body has measure type, indicating that each iteration is a stochastic computation.

In the model above, we use **plate** to draw a noisy measurement of each of the running times from 100 iterations to $k \times 100$ iterations. These k measurements are stored into the array y that is bound as the result of calling **plate**.

The language in figure 8 also introduces the **array** constructor for generating an array. The first two arguments are identical to those for the **plate** primitive. The third argument is an expression that describes each element of the array. For example, **array** 100 x . x defines an array of the first 100 natural numbers.

A few more additions to core Hakaru, as introduced in figure 8, remain to be discussed. The set of types now includes the natural numbers \mathbb{N} and the type of arrays, notated $\langle \alpha \rangle$. Terms can contain natural numbers n , while the addition and negation operations are overloaded to work on both real and natural numbers. Furthermore, we assume that natural numbers can be implicitly promoted to reals when necessary. The **counting** primitive denotes the unique translation-invariant measure over \mathbb{N} . While the **size** primitive obtains the length of an array, the **idx** primitive indexes into an array at a given expression of type \mathbb{N} . Finally, the $=$ primitive returns 1 if two expressions represent the same natural number, and 0 otherwise.

5 DISINTEGRATING PROGRAMS WITH ARRAYS

Given a probabilistic program written in the language of figure 8, we want to extend **disintegrate** to automatically produce the posterior distribution in the same language. For the model in equation (57), we want to produce:

$$\begin{aligned}
 \lambda t. \text{ do } \{ & a \leftarrow \text{normal } \mu_a \sigma_a; \\
 & b \leftarrow \text{normal } \mu_b \sigma_b; \\
 & _ \leftarrow \text{plate } k \ i. \text{ do } \{ \text{factor } (\text{normalD } (a \cdot 100 \cdot i + b) \ 1 \ (\text{idx } t \ i)); \\
 & \quad \text{return } (); \\
 & \text{factor } (k = \text{size } t); \\
 & \text{return } (a, b) \}
 \end{aligned} \tag{58}$$

This is a conditional measure where the input t is an array of real values. This output generalizes the result in equation (55) from a having a product of 2 to a product of k factors, each of which is derived from the density of the Gaussian distribution representing a noisy measurement. To set up **disintegrate** to produce this symbolic product, we extend the language of figure 8.

5.1 A language of repeated bindings

Real numbers $r \in \mathbb{R}$	Natural numbers $n \in \mathbb{N}$	Variables x	Locations \bar{x}	Indices \hat{x}
Bindings	$b ::= \bar{x} \leftarrow [\hat{x} \cdot l \dots]. m \mid \text{factor } [\hat{x} \cdot l \dots]. e$			
Heap	$h ::= [b; \dots; b]$			
Atomic terms	$u ::= x \mid -u \mid u + e \mid e + u \mid u \cdot e \mid e \cdot u \mid \text{normalD } r \ r \ u$ $\mid \text{fst } u \mid \text{snd } u \mid \text{idx } u \ e \mid \text{idx } e \ u \mid \text{size } u \mid u = e \mid e = u$			
Head normal forms $v ::=$	$u \mid r \mid n \mid () \mid (e, e) \mid \text{counting} \mid \text{lebesgue} \mid \text{normal } r \ r \mid \text{return } e$ $\mid \text{do } \{x \leftarrow m; M\} \mid \text{do } \{\text{factor } e; M\} \mid \text{plate } l \ x. m \mid \text{array } l \ x. e$			
Terms	$e, l, m, M ::= v \mid \bar{x}[\hat{x} \dots] \mid \langle \bar{x}[\hat{x} \dots] \rangle \mid \hat{x} \mid -e \mid e + e \mid e \cdot e \mid \text{normalD } r \ r \ e \mid e = e$ $\mid \text{fst } e \mid \text{snd } e \mid \text{idx } e \ e \mid \text{size } e$			
Types	$\alpha, \beta ::= \mathbb{1} \mid \mathbb{R} \mid \mathbb{N} \mid \alpha \times \beta \mid \langle \alpha \rangle \mid \mathbb{M} \ \alpha$			

Fig. 10. Internal language for disintegrating array programs

Figure 10 shows the internal language used for disintegrating array probabilistic programs. As in the case of the previous internal language (of figure 4), we track bound variables by extending the input language with a heap of bindings. When a binding originates under an array, it would in principle be repeated many times on the heap. To express this repetition succinctly, the internal language introduces *index variables*.

Index variables, or *indices* for short, are variables of type \mathbb{N} that are used for tracking repeated expressions. Notated \hat{x} , they are created by the disintegrator when it goes under the binder of **array** and **plate** expressions. Just as a location \bar{x} represents the variable bound from $x \leftarrow e$, an index \hat{x} represents the bound variable in **array** $l \ x. e$ or **plate** $l \ x. e$.

While a location is bound on the heap, an index is bound by a list notated $[\hat{x} \cdot l \dots]$. Here each element of the list is an index variable paired with an expression of type \mathbb{N} that denotes the length of the corresponding array. Since the language permits arbitrary nesting of arrays, the list can contain zero or more elements. The ordering of the elements matters here, and the list—being a

binding form—contains no repetitions. We label this structure a *pairwise-distinct list*, and refer to its instances as *binding-lists*.

The binding-list represents repetition. An expression together with a binding-list such as the pairwise-distinct list $[\hat{x}_1 \cdot l_1, \dots, \hat{x}_n \cdot l_n]$ compactly represents that expression repeated $l_1 \cdot \dots \cdot l_n$ times, with any uses of the variables $\hat{x}_1, \dots, \hat{x}_n$ updated accordingly. When an expression is augmented with a binding-list we deem the expression *lifted*.

Every binding is now lifted. This means that locations and weights are (conceptually) repeated zero or more times on the heap. When a location is used, it must be “applied” to a list of indices to select from all the symbolic bindings under that location. This list—notated $[\hat{x}\dots]$ —contains indices alone, and is akin to an array *accessor* or *selector*. We refer to its instances as *selector-lists*.

The scoping rules for locations on the heap remain unchanged from the previous disintegrator—older bindings are to the left and take scope over younger bindings to the right. Binding-lists, on the other hand, take scope only over the rest of the *current* binding. This is better illustrated via an example heap:

$$[\bar{a} \sim [], \mathbf{normal} \ 0 \ 1; \bar{b} \sim [\hat{x}\cdot 50]. \mathbf{normal} \ \bar{a}[] \ 1; \mathbf{factor} \ [\hat{x}\cdot 30]. (\bar{b}[\hat{x}] + \hat{x})] \quad (59)$$

The \bar{b} in $\bar{b}[\hat{x}] + \hat{x}$ refers to the binding $\bar{b} \sim [\hat{x}\cdot 50]$. $\mathbf{normal} \ \bar{a}[] \ 1$. On the other hand, the \hat{x} in $\bar{b}[\hat{x}] + \hat{x}$ refers to $\hat{x}\cdot 30$ and not $\hat{x}\cdot 50$.

Just as conditioning on a pair of terms involves disintegrating each subterm, conditioning on an array involves successive disintegration on each array element. Thus, the four functions of disintegration are themselves lifted:

$$\triangleright\triangleright \text{ (“perform”)} \quad : \text{heap} \rightarrow \text{inds} \rightarrow \llbracket \mathbb{M} \alpha \rrbracket \rightarrow (\text{emission} \times \text{heap} \times \llbracket \alpha \rrbracket) \quad (60)$$

$$\triangleright \text{ (“evaluate”)} \quad : \text{heap} \rightarrow \text{inds} \rightarrow \llbracket \alpha \rrbracket \rightarrow (\text{emission} \times \text{heap} \times \llbracket \alpha \rrbracket) \quad (61)$$

$$\triangleleft\triangleleft \text{ (“constrain outcome”)} : \text{heap} \rightarrow \text{inds} \rightarrow \llbracket \mathbb{M} \alpha \rrbracket \rightarrow \llbracket \alpha \rrbracket \rightarrow (\text{emission} \times \text{heap}) \quad (62)$$

$$\triangleleft \text{ (“constrain value”)} \quad : \text{heap} \rightarrow \text{inds} \rightarrow \llbracket \alpha \rrbracket \rightarrow \llbracket \alpha \rrbracket \rightarrow (\text{emission} \times \text{heap}) \quad (63)$$

The binding-lists here take scope over the rest of the arguments—the term being disintegrated and the observation (when present). This captures any indices and the selector-lists of locations used within the term or the observation.

5.2 Extending the previous disintegrator

Figures 11 and 12 illustrate the four functions implementing lifted disintegration, which we describe in more detail here. Most of the extensions are straightforward. Where earlier we may have emitted a binding—as in the case of $\triangleright\triangleright [\dots] \llbracket \llbracket \text{lebesgue} \rrbracket \rrbracket$ —we now emit a lifted binding, and where earlier we may have stored a binding on the heap—as in the case of $\triangleright\triangleright [\dots] \llbracket \llbracket \mathbf{do} \{x \sim e_1; e_2\} \rrbracket \rrbracket$ —we now store a lifted binding. In both cases the bindings are lifted with the binding-list of the current function. In all cases where a fresh location is generated, all bindings of that location are selected for use, i.e. we apply each use of that location to a list of all indices in the current binding-list.

Thus, in most cases the lifting binding-list remains unchanged. We can now consider one-by-one the cases that use the binding-list non-trivially.

The **plate** constructor is a measure operation that encapsulates repeated measure terms. For this case $\triangleright\triangleright$ and $\triangleleft\triangleleft$ both store on the heap a binding for e in **plate** $l \ x. \ e$. Since **plate** itself repeats $e \ l$ times, the binding is lifted by the current binding-list *extended* with a fresh element $\hat{x}\cdot l$. Since **plate** returns an array, $\triangleright\triangleright$ and $\triangleleft\triangleleft$ need to recur on all l bindings in an array form. For this purpose we use the *multiloc* constructor $\langle \bar{x}[\hat{x}\dots] \rangle$, which denotes an array of locations. It is applied to a selector-list whose length is one less than the size of the binding-list for the corresponding location bound on the heap.

\triangleright (“perform”) : $heap \rightarrow inds \rightarrow \llbracket \mathbb{M} \alpha \rrbracket \rightarrow (emission \times heap \times \llbracket \alpha \rrbracket)$	
$\triangleright h [\hat{x} \cdot l_1 \dots] \llbracket m \rrbracket = \text{case } \llbracket m \rrbracket \text{ of}$	
$\llbracket u \rrbracket \rightarrow (\bar{z}[\hat{x} \cdot l_1 \dots] \leftarrow u, h, \llbracket \bar{z}[\hat{x} \dots] \rrbracket)$	where u is atomic, \bar{z} is fresh (64)
$\llbracket \text{lebesgue} \rrbracket \rightarrow (\bar{z}[\hat{x} \cdot l_1 \dots] \leftarrow \text{lebesgue}, h, \llbracket \bar{z}[\hat{x} \dots] \rrbracket)$	where \bar{z} is fresh (65)
$\llbracket \text{normal } r_1 r_2 \rrbracket \rightarrow (\bar{z}[\hat{x} \cdot l_1 \dots] \leftarrow \text{normal } r_1 r_2, h, \llbracket \bar{z}[\hat{x} \dots] \rrbracket)$	(66)
$\llbracket \text{return } e \rrbracket \rightarrow \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket$	(67)
$\llbracket \text{do } \{x \leftarrow e_1; e_2\} \rrbracket \rightarrow \triangleright [h; \bar{x}[\hat{x} \cdot l_1 \dots] \leftarrow e_1][\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket \{x \mapsto \bar{x}[\hat{x} \dots]\}$	where \bar{x} is fresh (68)
$\llbracket \text{do } \{\text{factor } e_1; e_2\} \rrbracket \rightarrow \triangleright [h; \text{factor } [\hat{x} \cdot l_1 \dots] e_1][\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$	(69)
$\llbracket \text{plate } l x. e \rrbracket \rightarrow \triangleright [h; \bar{p}[\hat{x} \cdot l_1 \dots \hat{y} \cdot l] \leftarrow e \{x \mapsto \hat{y}\}][\hat{x} \cdot l_1 \dots] \llbracket \bar{p}[\hat{x} \dots] \rrbracket$	where \bar{p}, \hat{y} are fresh (70)
$\llbracket e \rrbracket \rightarrow (s_1; s_2, h_2, \llbracket v \rrbracket)$	where e is not in head normal form, (71)
	$(s_1, h_1, \llbracket m_1 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket,$
	$(s_2, h_2, \llbracket v \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket m_1 \rrbracket$

\triangleright (“evaluate”) : $heap \rightarrow inds \rightarrow \llbracket \alpha \rrbracket \rightarrow (emission \times heap \times \llbracket \alpha \rrbracket)$	
$\triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket = \text{case } \llbracket e \rrbracket \text{ of}$	
$\llbracket v \rrbracket \rightarrow ([], h, \llbracket v \rrbracket)$	where v is in head normal form (72)
$\llbracket \text{fst } e \rrbracket \rightarrow \text{case } \llbracket p \rrbracket \text{ of}$	where $(s_1, h_1, \llbracket p \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket,$ (73)
$\llbracket (e_1, _) \rrbracket \rightarrow (s_1; s_2, h_2, \llbracket v \rrbracket)$	$(s_2, h_2, \llbracket v \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket,$
$\llbracket u \rrbracket \rightarrow (s_1, h_1, \llbracket \text{fst } u \rrbracket)$	u is atomic
$\llbracket \text{snd } e \rrbracket \rightarrow \text{similar to the } \text{fst} \text{ case}$	(74)
$\llbracket \text{idx } e_1 e_2 \rrbracket \rightarrow \text{if } a \text{ is atomic then } (s_1, h_1, \llbracket \text{idx } a e_2 \rrbracket) \text{ else}$	where $(s_1, h_1, \llbracket a \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket$ (75)
$\text{if } z \text{ is atomic then } (s_1; s_2, h_2, \llbracket \text{idx } a z \rrbracket) \text{ else}$	$(s_2, h_2, \llbracket z \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$
$\text{case } \llbracket (a, z) \rrbracket \text{ of}$	
$\llbracket (\text{array } l x. e, \hat{z}) \rrbracket \rightarrow (s_1; s_2; s_3, h_3, \llbracket v \rrbracket)$	$(s_3, h_3, \llbracket v \rrbracket) = \triangleright h_2 [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket \{x \mapsto \hat{z}\}$
$\llbracket ((\hat{x}[\hat{y} \dots]), \hat{z}) \rrbracket \rightarrow (s_1; s_2; s_3, h_3, \llbracket v \rrbracket)$	$(s_3, h_3, \llbracket v \rrbracket) = \triangleright h_2 [\hat{x} \cdot l_1 \dots] \llbracket \hat{x}[\hat{y} \dots \hat{z}] \rrbracket$
$\llbracket \text{size } e \rrbracket \rightarrow (s, h_1, \text{size } \llbracket v \rrbracket)$	where $(s, h_1, \llbracket v \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket$ (76)
$\llbracket -e \rrbracket \rightarrow (s, h_1, -\llbracket v \rrbracket)$	where $(s, h_1, \llbracket v \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket$ (77)
$\llbracket e_1 + e_2 \rrbracket \rightarrow (s_1; s_2, h_2, \llbracket v_1 \rrbracket + \llbracket v_2 \rrbracket)$	where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket,$ (78)
	$(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$
$\llbracket e_1 \cdot e_2 \rrbracket \rightarrow (s_1; s_2, h_2, \llbracket v_1 \rrbracket \cdot \llbracket v_2 \rrbracket)$	where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket,$ (79)
	$(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$
$\llbracket \text{normalD } r_1 r_2 e \rrbracket \rightarrow (s, h_1, \llbracket \text{normalD } r_1 r_2 v \rrbracket)$	where $(s, h_1, \llbracket v \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket$ (80)
$\llbracket \hat{x}[\hat{y} \dots] \rrbracket \rightarrow (s, [h'_1; \bar{x}[\hat{z} \cdot l_3 \dots] \leftarrow \text{return } v; h_2], \llbracket v \rrbracket \{[\hat{z} \dots] \mapsto [\hat{y} \dots]\})$	where $[h_1; \bar{x}[\hat{z} \cdot l_3 \dots] \leftarrow e; h_2] = h,$ (81)
	$ \llbracket \hat{y} \dots \rrbracket = \llbracket \hat{z} \cdot l_3 \dots \rrbracket ,$
	$(s, h'_1, \llbracket v \rrbracket) = \triangleright h_1 [\hat{z} \cdot l_3 \dots] \llbracket e \rrbracket$

Fig. 11. The forward functions handling arrays

The `idx` operator generalizes `fst` and `snd` to the array case. Thus when either operand of `idx` $e_1 e_2$ is atomic, the result looks similar to handling pair destructors on an atomic operand: \triangleright returns code while \triangleleft fails. When the index operand (e_2) is an index variable, we symbolically access an element of the array operand e_1 . For an `array` expression this involves substitution of the index variable into the body of the array. For a `multiloc` expression we extend its selector-list to select the corresponding location. Both \triangleright and \triangleleft recur on the newly accessed term. We note that neither function currently handles the case of indexing into an array at non-atomic and non-index operands such as literal numbers.

\llcorner (“constrain outcome”) : $heap \rightarrow inds \rightarrow [\mathbb{M} \alpha] \rightarrow [\alpha] \rightarrow (emission \times heap)$	
$\llcorner h [\hat{x} \cdot l_1 \dots] \llbracket m \rrbracket t = \text{case } \llbracket m \rrbracket \text{ of}$	
$\llbracket u \rrbracket$	$\rightarrow \perp$ where u is atomic (82)
$\llbracket \text{lebesgue} \rrbracket$	$\rightarrow ([], h)$ (83)
$\llbracket \text{normal } r_1 r_2 \rrbracket$	$\rightarrow ([], [h; \text{factor } [\hat{x} \cdot l_1 \dots] (\text{normalD } r_1 r_2 t)])$ (84)
$\llbracket \text{return } e \rrbracket$	$\rightarrow \llcorner h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket t$ (85)
$\llbracket \text{do } \{x \leftarrow e_1; e_2\} \rrbracket$	$\rightarrow \llcorner [h; \bar{x}[\hat{x} \cdot l_1 \dots] \leftarrow e_1] \llbracket e_2 \rrbracket \{x \mapsto \bar{x}[\hat{x} \dots]\} t$ where \bar{x} is fresh (86)
$\llbracket \text{do } \{\text{factor } e_1; e_2\} \rrbracket$	$\rightarrow \llcorner [h; \text{factor } [\hat{x} \cdot l_1 \dots] e_1] \llbracket e_2 \rrbracket t$ (87)
$\llbracket \text{plate } l x. e \rrbracket$	$\rightarrow \llcorner [h; \bar{p}[\hat{x} \cdot l_1 \dots \hat{y} \cdot l] \leftarrow e \{x \mapsto \hat{y}\}][\hat{x} \cdot l_1 \dots] \llbracket \langle \bar{p}[\hat{x} \dots] \rangle \rrbracket t$ where \bar{p}, \hat{y} are fresh (88)
$\llbracket e \rrbracket$	$\rightarrow ([s_1; s_2], h_2)$ where e is not in head normal form, (89)
	$(s_1, h_1, \llbracket m_1 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket,$
	$(s_2, h_2) = \llcorner h_1 [\hat{x} \cdot l_1 \dots] \llbracket m_1 \rrbracket t$
<hr/>	
\llcorner (“constrain value”) : $heap \rightarrow inds \rightarrow [\alpha] \rightarrow [\alpha] \rightarrow (emission \times heap)$	
$\llcorner h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket t = \text{case } \llbracket e \rrbracket \text{ of}$	
$\llbracket \text{fst } e \rrbracket$	$\rightarrow \text{case } \llbracket p \rrbracket \text{ of}$ where $(s_1, h_1, \llbracket p \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket,$ (90)
	$\llbracket (e_1, _) \rrbracket \rightarrow ([s_1; s_2], h_2)$ $(s_2, h_2) = \llcorner h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket t,$
	$\llbracket u \rrbracket \rightarrow \perp$ u is atomic
$\llbracket \text{snd } e \rrbracket$	\rightarrow similar to the fst case (91)
$\llbracket \text{idx } e_1 e_2 \rrbracket$	\rightarrow if a is atomic then \perp else where $(s_1, h_1, \llbracket a \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket$ (92)
	if z is atomic then \perp else case $\llbracket (a, z) \rrbracket$ of $(s_2, h_2, \llbracket z \rrbracket) = \triangleright h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$
	$\llbracket (\text{array } l x. e, \hat{z}) \rrbracket \rightarrow ([s_1; s_2; s_3], h_3)$ $(s_3, h_3) = \llcorner h_2 [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket \{x \mapsto \hat{z}\} t$
	$\llbracket (\langle \bar{x}[\hat{y} \dots] \rangle, \hat{z}) \rrbracket \rightarrow ([s_1; s_2; s_3], h_3)$ $(s_3, h_3) = \llcorner h_2 [\hat{x} \cdot l_1 \dots] \llbracket \bar{x}[\hat{y} \dots \hat{z}] \rrbracket t$
$\llbracket () \rrbracket$	$\rightarrow ([], h)$ (93)
$\llbracket (e_1, e_2) \rrbracket$	$\rightarrow ([s_1; s_2], h_2)$ where $(s_1, h_1) = \llcorner h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket$ (fst t), (94)
	$(s_2, h_2) = \llcorner h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket$ (snd t)
$\llbracket \text{array } l x. e \rrbracket$	$\rightarrow \llcorner [h; \text{factor } [\hat{x} \cdot l_1 \dots] (l = \text{size } t)] [\hat{x} \cdot l_1 \dots \hat{y} \cdot l] \llbracket e \rrbracket \{x \mapsto \hat{y}\} (\text{idx } t \hat{y})$ (95)
	where \hat{y} is fresh
$\llbracket \langle \bar{x}[\hat{y} \dots] \rangle \rrbracket$	$\rightarrow \llcorner [h; \text{factor } [\hat{x} \cdot l_1 \dots] (l = \text{size } t)] [\hat{x} \cdot l_1 \dots \hat{q} \cdot l] \llbracket \bar{x}[\hat{y} \dots \hat{q}] \rrbracket (\text{idx } t \hat{q})$ (96)
	where $[h_1; \bar{x}[\hat{z} \cdot l_3 \dots \hat{p} \cdot l] \leftarrow m; h_2], \hat{q}$ is fresh, $ \llbracket \hat{y} \dots \rrbracket = \llbracket \hat{z} \cdot l_3 \dots \hat{p} \cdot l \rrbracket - 1,$
$\llbracket \text{sin } e \rrbracket$	$\rightarrow \perp$ (97)
$\llbracket -e \rrbracket$	$\rightarrow \llcorner h [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket (-t)$ (98)
$\llbracket e_1 + e_2 \rrbracket$	$\rightarrow ([s_1; s'_1], h'_1)$ where $(s_1, h_1, \llbracket v_1 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket,$ (99)
	$(s'_1, h'_1) = \llcorner h_1 [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket (t - \llbracket v_1 \rrbracket)$
	$\sqcup ([s_2; s'_2], h'_2)$ where $(s_2, h_2, \llbracket v_2 \rrbracket) = \triangleright h [\hat{x} \cdot l_1 \dots] \llbracket e_2 \rrbracket,$
	$(s'_2, h'_2) = \llcorner h_2 [\hat{x} \cdot l_1 \dots] \llbracket e_1 \rrbracket (t - \llbracket v_2 \rrbracket)$
$\llbracket v \rrbracket$	$\rightarrow \perp$ where v is in head normal form (100)
$\llbracket \bar{x}[\hat{y} \dots] \rrbracket$	$\rightarrow (s, [h'_1; \bar{x}[\hat{x} \cdot l_1 \dots] \leftarrow \text{return } t; h_2])$ where $[h_1; \bar{x}[\hat{z} \cdot l_3 \dots] \leftarrow e; h_2] = h,$ (101)
	$ \llbracket \hat{y} \dots \rrbracket = \llbracket \hat{z} \cdot l_3 \dots \rrbracket ,$
	$\llbracket \hat{y} \dots \rrbracket$ is a permutation of $\llbracket \hat{x} \dots \rrbracket,$
	$(s, h'_1) = \llcorner h_1 [\hat{x} \cdot l_1 \dots] \llbracket e \rrbracket \{\llbracket \hat{z} \dots \rrbracket \mapsto \llbracket \hat{y} \dots \rrbracket\} t$

Fig. 12. The backward functions handling arrays

The **array** and **multiloc** constructors are head normal forms; \triangleright does not treat them specially. For \triangleleft the case for pairs is informative—we go under the structure and make successive calls to \triangleleft for each element. We repeat disintegration by extending the binding-list for \triangleleft with a fresh index variable. In the **array** case we substitute this index under the binder and call \triangleleft on the body. In the **multiloc** case we select a location by extending its selector-list, and we call \triangleleft on the resulting term.

Constraining the **array** and **multiloc** constructors repeatedly constrains each element to be the corresponding element of the observed variable t (obtained using **idx**). While t is guaranteed by the type-system to be an array, we introduce a guard that constrains the size of t . This size check appears in the output in equation (58). To see the need for this constraint we consider the new specification of **disintegrate**.

Specifying disintegration on arrays. The specification of disintegration depends on the type of the observed variable. Previously we specified disintegration when the input data is a real number, where **lebesgue** is used as the base measure. When observing pairs of real numbers, we update our specification and disintegrate with respect to **lebesgue** \otimes ($\lambda_.$ **lebesgue**). When observing arrays of real numbers, the specification becomes:

$$\text{do } \{n \leftarrow \text{counting}; t \leftarrow \text{plate } n _ . \text{lebesgue}; p \leftarrow \text{disintegrate } m \ t; \text{return } (t, p)\} \equiv m \quad (102)$$

Here the subexpression **do** $\{n \leftarrow \text{counting}; \text{plate } n _ . \text{lebesgue}\}$ actually defines the *disjoint union* of Lebesgue measures on \mathbb{R}^n for each value of n . To select the correct base measure from this disjoint union, the output of **disintegrate** must apply a guard checking the size of the input array t and the observed variable inside the model.

Now only the cases for locations remain to be described. When handling locations the disintegrator has to compare three quantities—the binding-list $[\hat{x} \cdot l_1 \dots]$ of the function call, the binding-list $[\hat{z} \cdot l_3 \dots]$ lifting this location’s binding on the heap, and the selector-list $[\hat{y} \dots]$ selecting which bindings we want to evaluate or constrain. Both \triangleright and \triangleleft check that $[\hat{y} \dots]$ and $[\hat{z} \cdot l_3 \dots]$ have the same number of elements.

When evaluating a location, we simply evaluate *all* the bindings defined for that location on the heap. Thus \triangleright performs the associated measure term as before, but lifts the call to $\triangleright\triangleright$ with $[\hat{z} \cdot l_3 \dots]$ as the binding-list. When we return from this call we use $[\hat{y} \dots]$ to select the required bindings in head normal form.

This is a departure from our previously lazy approach to evaluation. If we wanted to select a subset of the bindings, for example, if we only required the diagonal elements of an array (wherein $[\hat{y} \dots]$ would contain duplicates), we would be still be evaluating all the bound elements. This may cause the disintegrator to fail more often in cases that evaluate and constrain mutually exclusive sections of an array. Nevertheless, we take the current approach for ease of implementation and assign extensions for future work.

For constraining a location, we require that $[\hat{y} \dots]$ be a permutation of $[\hat{x} \cdot l_1 \dots]$. In other words, for each invocation of \triangleleft as per $[\hat{x} \cdot l_1 \dots]$, we want to constrain a distinct binding on the heap. This requirement extends to arrays the idea that we cannot observe the same term multiple times. Given this condition, we lift $\triangleleft\triangleleft$ on the associated measure term with $[\hat{x} \cdot l_1 \dots]$. Before we invoke $\triangleleft\triangleleft$ we preserve the order of observation by substituting $[\hat{y} \dots]$ in place of $[\hat{z} \cdot l_3 \dots]$ in the measure term.

5.3 Putting it all together

We can now describe a complete method of conditioning the k -measurement **blr** model from equation (57). The algorithm in equation (52) is still valid for the new disintegrator, consisting of three steps:

$$\begin{aligned} & \triangleright \square \square \llbracket \text{do } \{ a \sim \text{normal } \mu_a \sigma_a; \\ & \quad b \sim \text{normal } \mu_b \sigma_b; \\ & \quad y \sim \text{plate } k \text{ i. (normal } (a \cdot 100 \cdot i + b) \text{ 1);} \\ & \quad \text{return } (y, (a, b)) \rrbracket \\ & \quad \vdots \\ & \Rightarrow \left(\square, \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b; \\ \bar{y} \sim \square. \text{plate } k \text{ i. (normal } (\bar{a}[] \cdot 100 \cdot i + \bar{b}[]) \text{ 1)} \end{array} \right] \right) \llbracket (\bar{y}[], (\bar{a}[], \bar{b}[])) \rrbracket \end{aligned}$$

(a) Going forward on the input program

$$\begin{aligned} & \triangleleft \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b; \\ \bar{y} \sim \square. \text{plate } k \text{ i. (normal } (\bar{a}[] \cdot 100 \cdot i + \bar{b}[]) \text{ 1)} \end{array} \right] \square \llbracket \bar{y}[] \rrbracket t \\ & \triangleleft \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b \end{array} \right] \square \llbracket \text{plate } k \text{ i. (normal } (\bar{a}[] \cdot 100 \cdot i + \bar{b}[]) \text{ 1)} \rrbracket t \\ & \triangleleft \left[\text{---}; \bar{p} \sim [\hat{i}.k]. \text{normal } (\bar{a}[] \cdot 100 \cdot \hat{i} + \bar{b}[]) \text{ 1} \right] \square \llbracket \langle \bar{p}[] \rangle \rrbracket \quad t \\ & \triangleleft \left[\text{---}; \text{factor } [], (k = \text{size } t) \right] \quad [\hat{j}.k] \llbracket \bar{p}[\hat{j}] \rrbracket \quad (\text{idx } t \hat{j}) \\ & \triangleleft \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b \end{array} \right] [\hat{j}.k] \llbracket \text{normal } (\bar{a}[] \cdot 100 \cdot \hat{j} + \bar{b}[]) \text{ 1} \rrbracket (\text{idx } t \hat{j}) \\ & \Rightarrow \left(\square, \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b; \\ \text{factor } [\hat{j}.k]. (\text{normalD } (\bar{a}[] \cdot 100 \cdot \hat{j} + \bar{b}[]) \text{ 1 } (\text{idx } t \hat{j})) \end{array} \right] \right) \\ & \Rightarrow \left(\square, \left[\begin{array}{l} \bar{a} \sim \square. \text{normal } \mu_a \sigma_a; \\ \bar{b} \sim \square. \text{normal } \mu_b \sigma_b; \\ \text{factor } [\hat{j}.k]. (\text{normalD } (\bar{a}[] \cdot 100 \cdot \hat{j} + \bar{b}[]) \text{ 1 } (\text{idx } t \hat{j})); \\ \bar{p} \sim [\hat{j}.k]. \text{return } (\text{idx } t \hat{j}); \\ \text{factor } [], (k = \text{size } t) \end{array} \right] \right) \\ & \Rightarrow (\square, \left[\text{---}; \bar{y} \sim \square. \text{return } t \right]) \end{aligned}$$

(b) Constraining the observed array

Fig. 13. The first two steps of disintegrating k -measurement `blr`

- (1) We scan the input program for the observed term.
- (2) Figure 13b shows the second step of performing lifted disintegration on this model. Here we note that there is constant number of calls to \triangleleft . This number is independent of the size of the observed array, in contrast to the case of disintegrating (nested) pairs.
- (3) The only change is in the final step of converting the resultant program from the internal language to the input language. In the case of lifted disintegration, this conversion involves the following changes:

- Expressions that represent lifted bindings are converted to **plate** expressions:

$$\bar{p} \sim [\hat{x}_1.l_1, \dots, \hat{x}_n.l_n]. m \implies p \sim \text{plate } l_1 x_1. (\dots \text{plate } l_n x_n. m) \quad (103)$$

$$\text{factor } [\hat{x}_1.l_1, \dots, \hat{x}_n.l_n]. e \implies _ \sim \text{plate } l_1 x_1. (\dots \text{plate } l_n x_n. \text{do } \{ \text{factor } e; \text{return } () \}) \quad (104)$$

- Expressions that select lifted bindings are converted to **idx** expressions:

$$\bar{x}[\hat{x}_1, \dots, \hat{x}_n] \implies \text{idx } (\dots (\text{idx } x x_1) \dots) x_n \quad (105)$$

$$\langle \bar{x}[\hat{x}_1, \dots, \hat{x}_n] \rangle \implies \text{idx } (\dots (\text{idx } x x_1) \dots) x_n \quad (106)$$

6 EVALUATION

Our disintegrator was evaluated using Hakaru, a suite of program transformations for a language with additional types (such as booleans), operations (such as `exp`, `log`, and `sin`) and distributions (such as binomial, beta, bernoulli, categorical, uniform, and gamma). We looked at the benchmarks

distributed with the R2 inference system [Nori et al. 2014], which comprise 11 of the 21 benchmarks used by Gehr et al. [2016] to evaluate their simplifier for probabilistic programs.

Of these 11 benchmarks, 4 models do not use arrays:

- **Burglar alarm:** Probability of burglary given sounding of alarm
- **Grass:** Probability that it is raining given grass is wet
- **Noisy or:** Infer the values of boolean nodes related by noisy-or functions
- **Two coins:** Marginal distribution of two fair coins when at least one landed on tails

Lifted disintegration matches the behavior of Shan and Ramsey’s disintegrator in these cases. The machinery for symbolic conditioning of arrays developed in this paper does not hamper the performance of the disintegrator on models that do not observe arrays.

The remaining 7 models use arrays:

- **Clinical trial:** Evaluate a medical treatment based on control and experimental outcomes
- **Coin bias:** Estimate the bias of a coin
- **Digit recognition:** Recognize digits from handwriting
- **HIV:** Tuning the parameters of a linear HIV dynamical model
- **Linear regression:** Fit a line through observed points
- **Survey unbiased:** Estimate the gender bias of a Gaussian modeled population
- **True skill:** Use the outcomes of a series of games to infer the skill levels of players

Lifted disintegration correctly handles 6 out of these 7 models. The remaining benchmark—*True skill*—contains conditionals inside arrays, and it is a work in progress to condition such models symbolically and without unrolling arrays.

On these array benchmarks the original disintegration algorithm requires time linear in the number of scalar elements in the observed array. We expect the PSI system of Gehr et al. [2016] to exhibit similar performance as it unrolls loops as well. This is why in their evaluation Gehr et al. use data sets truncated to up to 784 data points.

In contrast, the current approach using lifted disintegration takes time linear in the number of indices used to select an element of the array. All of the benchmarks that we have discussed use either singly-indexed or doubly-indexed arrays, giving lifted disintegration a big performance edge over other symbolic conditioning approaches. In the case of PSI it is tricky to quantify the slowdown because PSI combines conditioning with simplification whereas Hakaru separates these two transformations.

In addition to the R2 benchmarks, lifted disintegration works for various microbenchmarks that manipulate arrays via mapping, transposing, and copying. The transformation produces posterior distributions in time independent of the number of scalar elements of the array.

7 DISCUSSION

We have extended the automatic disintegrator described by Shan and Ramsey [2017] to condition arrays of stochastic variables. Notably, the extended disintegrator handles arrays *symbolically*, i.e., without unrolling the loop. The time taken for our automatic disintegrator to produce the posterior distribution is independent of the length – and linearly proportional to the dimensionality – of the observed data array. This is only the beginning, as there are many foreseeable extensions to array disintegration.

The languages described in this paper lack conditionals because currently our system does not handle control-flow uncertainty under arrays. The disintegrator lazily flattens the input program, which means that bindings inside loops may end up lifted out into their own array forms. Lifting conditionals out of loops in the same manner can lead to incorrect results. For example, consider

the program:

```
do { $\mu \leftarrow$  normal 0 1;
   $p \leftarrow$  plate 50  $\_.$  do { $b \leftarrow$  bern 0.5;
    if  $b$  (normal  $\mu$  1) (normal ( $\mu + 5$ ) 1)};
  return ( $p, \mu$ )}
```

(107)

Here pulling the conditional out of the plate body changes the meaning of the program. The current system avoids this by failing upon encountering a conditional when disintegration is lifted, i.e., when the current binding-list is non-empty. While we can perceive a solution that duplicates code for this kind of problem, we seek a more symbolic solution for a variety of programs that mix arrays and conditionals.

One use case of conditionals is in expressing array accesses of the form “all elements except one”. This would allow us to symbolically produce the conditional proposal distributions of Gibbs samplers, which are useful in a wide variety of applications including document classification [Resnik and Hardisty 2009].

The `plate` primitive can only express loops over stochastic expressions such that element are drawn independently. Looking to the future, we desire symbolic conditioning of arrays where the elements have chain-like dependencies as well. This would be useful for disintegrating Hidden Markov Models and topic models [Wallach 2006].

We believe that the lifted technique described in this paper elegantly solves the problem of symbolically conditioning arrays. The approach of using indices to representing repetition shows promise in being extensible – we envision using boolean variables to represent conditionals, and generalizing from array-sizes to *shapes* that represent more exotic data structures. We look forward to the role of sophisticated symbolic conditioning in a world of increasingly rich probabilistic models.

REFERENCES

- Nathanael L. Ackerman, Cameron E. Freer, and Daniel M. Roy. 2011. Noncomputable Conditional Distributions. In *Proceedings of the 2011 IEEE 26th Annual Symposium on Logic in Computer Science (LICS '11)*. IEEE Computer Society, Washington, DC, USA, 107–116. <https://doi.org/10.1109/LICS.2011.49>
- Jacques Carette and Chung-Chieh Shan. 2016. *Simplifying Probabilistic Programs Using Computer Algebra*. Springer International Publishing, Cham, 135–152. https://doi.org/10.1007/978-3-319-28228-2_9
- George Casella and Edward I. George. 1992. Explaining the Gibbs Sampler. *The American Statistician* 46, 3 (1992), 167–174. <http://www.jstor.org/stable/2685208>
- Joseph T. Chang and David Pollard. 1997. Conditioning as Disintegration. *Statistica Neerlandica* 51, 3 (1997), 287–317.
- Guillaume Claret, Sriram K. Rajamani, Aditya V. Nori, Andrew D. Gordon, and Johannes Borgström. 2013. *Bayesian Inference Using Data Flow Analysis*. Technical Report MSR-TR-2013-27. Microsoft Research. <http://research.microsoft.com/apps/pubs/default.aspx?id=171611>
- Sebastian Fischer, Oleg Kiselyov, and Chung-chieh Shan. 2011. Purely Functional Lazy Nondeterministic Programming. *Journal of Functional Programming* 21, 4–5 (2011), 413–465.
- Sebastian Fischer, Josep Silva, Salvador Tamarit, and Germán Vidal. 2008. Preserving Sharing in the Partial Evaluation of Lazy Functional Programs. In *Revised Selected Papers from LOPSTR 2007: 17th International Symposium on Logic-Based Program Synthesis and Transformation (Lecture Notes in Computer Science)*. Springer, Berlin, 74–89.
- Timon Gehr, Sasa Misailovic, and Martin Vechev. 2016. *PSI: Exact Symbolic Inference for Probabilistic Programs*. Springer International Publishing, Cham, 62–83. https://doi.org/10.1007/978-3-319-41528-4_4
- Andrew Gelman, Daniel Lee, and Jiqiang Guo. 2015. Stan: A probabilistic programming language for Bayesian inference and optimization. (2015).
- Noah D. Goodman, Vikash K. Mansinghka, Daniel M. Roy, Keith Bonawitz, and Joshua B. Tenenbaum. 2008. Church: a language for generative models. In *Proc. of Uncertainty in Artificial Intelligence*. http://danroy.org/papers/church_GooManRoyBonTen-UAI-2008.pdf
- A. Kucukelbir, R. Ranganath, A. Gelman, and D. M. Blei. 2015. Automatic Variational Inference in Stan. *ArXiv e-prints* (June 2015). [arXiv:stat.ML/1506.03431](https://arxiv.org/abs/1506.03431)

- John Launchbury. 1993. A Natural Semantics for Lazy Evaluation. In *POPL '93: Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM Press, New York, 144–154.
- Vikash K. Mansinghka, Daniel Selsam, and Yura N. Perov. 2014. Venture: a higher-order probabilistic programming platform with programmable inference. *CoRR* abs/1404.0099 (2014). <http://arxiv.org/abs/1404.0099>
- Andrew McCallum, Karl Schultz, and Sameer Singh. 2009. FACTORIE: Probabilistic Programming via Imperatively Defined Factor Graphs. In *Neural Information Processing Systems (NIPS)*.
- Brian Milch, Bhaskara Marthi, Stuart Russell, David Sontag, Daniel L. Ong, and Andrey Kolobov. 2007. BLOG: Probabilistic Models with Unknown Objects. In *Statistical Relational Learning*, Lise Getoor and Ben Taskar (Eds.). MIT Press. <http://sites.google.com/site/bmilch/papers/blog-chapter.pdf>
- T. Minka, J. M. Winn, J. P. Guiver, S. Webster, Y. Zaykov, B. Yangel, A. Spengler, and J. Bronskill. 2014. Infer.NET 2.6. (2014). <http://research.microsoft.com/infernet> Microsoft Research Cambridge.
- Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. 2016. Probabilistic Inference by Program Transformation in Hakaru (System Description). In *Functional and Logic Programming: 13th International Symposium, FLOPS 2016 (Lecture Notes in Computer Science)*. Springer, Berlin, 62–79.
- Aditya V. Nori, Chung-Kil Hur, Sriram K. Rajamani, and Selva Samuel. 2014. R2: An Efficient MCMC Sampler for Probabilistic Programs. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*. AAAI Press, 2476–2482.
- Philip Resnik and Eric Hardisty. 2009. Gibbs Sampling for the Uninitiated. (2009).
- Chung-chieh Shan and Norman Ramsey. 2017. Exact Bayesian Inference by Symbolic Disintegration. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL 2017)*. ACM, New York, NY, USA, 130–144. <https://doi.org/10.1145/3009837.3009852>
- Luke Tierney. 1998. A Note on Metropolis-Hastings Kernels for General State Spaces. *The Annals of Applied Probability* 8, 1 (1998), 1–9.
- Hanna M. Wallach. 2006. Topic modeling: beyond bag-of-words. In *In Proceedings of the 23rd International Conference on Machine Learning*. 977–984.
- David Wingate, Andreas Stuhlmüller, and Noah D. Goodman. 2011. Lightweight Implementations of Probabilistic Programming Languages Via Transformational Compilation. In *Proceedings of AISTATS 2011: 14th International Conference on Artificial Intelligence and Statistics (JMLR Workshop and Conference Proceedings)*. MIT Press, Cambridge, 770–778.
- D. Wingate and T. Weber. 2013. Automated Variational Inference in Probabilistic Programming. *ArXiv e-prints* (Jan. 2013). [arXiv:stat.ML/1301.1299](https://arxiv.org/abs/1301.1299)
- Frank Wood, Jan Willem van de Meent, and Vikash Mansinghka. 2014. A New Approach to Probabilistic Programming Inference. In *Proceedings of the 17th International conference on Artificial Intelligence and Statistics*. 1024–1032.