

Denotational Validation of Higher-Order Bayesian Inference

ADAM ŚCIBIOR, University of Cambridge, England, UK and Max Planck Institute for Intelligent Systems, Germany

OHAD KAMMAR, University of Oxford, England, UK

MATTHIJS VÁKÁR, University of Oxford, England, UK

SAM STATON, University of Oxford, England, UK

HONGSEOK YANG, KAIST, South Korea

YUFEI CAI, Universität Tübingen, Germany

KLAUS OSTERMANN, Universität Tübingen, Germany

SEAN K. MOSS, University of Cambridge, England and University of Oxford, England, UK

CHRIS HEUNEN, University of Edinburgh, Scotland, UK

ZOUBIN GHAHRAMANI, University of Cambridge, England, UK and Uber AI Labs, California, USA

We present a modular semantic account of Bayesian inference algorithms for probabilistic programming languages, as used in data science and machine learning. Sophisticated inference algorithms are often explained in terms of composition of smaller parts. However, neither their theoretical justification nor their implementation reflects this modularity. We show how to conceptualise and analyse such inference algorithms as manipulating intermediate representations of probabilistic programs using higher-order functions and inductive types, and their denotational semantics.

Semantic accounts of continuous distributions use measurable spaces. However, our use of higher-order functions presents a substantial technical difficulty: it is impossible to define a measurable space structure over the collection of measurable functions between arbitrary measurable spaces that is compatible with standard operations on those functions, such as function application. We overcome this difficulty using quasi-Borel spaces, a recently proposed mathematical structure that supports both function spaces and continuous distributions.

We define a class of semantic structures for representing probabilistic programs, and semantic validity criteria for transformations of these representations in terms of distribution preservation. We develop a collection of building blocks for composing representations. We use these building blocks to validate common inference algorithms such as Sequential Monte Carlo and Markov Chain Monte Carlo. To emphasize the connection between the semantic manipulation and its traditional measure theoretic origins, we use Kock's synthetic measure theory. We demonstrate its usefulness by proving a quasi-Borel counterpart to the Metropolis-Hastings-Green theorem.

CCS Concepts: • **Mathematics of computing** → **Metropolis-Hastings algorithm; Sequential Monte Carlo methods; Theory of computation** → **Probabilistic computation; Bayesian analysis; Denotational semantics; Software and its engineering** → **Language types; Functional languages; Interpreters; Domain specific languages; Computing methodologies** → *Machine learning*;

Additional Key Words and Phrases: quasi-Borel spaces, synthetic measure theory, Bayesian inference, applied category theory, commutative monads, Kock integration, initial algebra semantics, sigma-monoids



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

© 2018 Copyright held by the owner/author(s).

2475-1421/2018/1-ART60

<https://doi.org/10.1145/3158148>

ACM Reference Format:

Adam Ścibior, Ohad Kammar, Matthijs Vákár, Sam Staton, Hongseok Yang, Yufei Cai, Klaus Ostermann, Sean K. Moss, Chris Heunen, and Zoubin Ghahramani. 2018. Denotational Validation of Higher-Order Bayesian Inference. *Proc. ACM Program. Lang.* 2, POPL, Article 60 (January 2018), 29 pages. <https://doi.org/10.1145/3158148>

1 INTRODUCTION

One of the key challenges in Bayesian data analysis is to develop or find an efficient algorithm for estimating the posterior distribution of a probabilistic model with respect to a given data set. This posterior distribution combines prior knowledge encoded in the model and information present in the data set consistently according to the rules of probability theory, but its mathematical definition often involves integration or summation over a large index set and does not yield to an efficient computation strategy immediately. A data scientist typically has to make one of the suboptimal decisions: she has to consult a large body of specialised research in order to pick an algorithm suitable for her model, or to change the model so that it falls into one of those cases with efficient known algorithms for posterior inference, or to face the challenge directly by developing a new algorithm for herself.

Recent probabilistic programming languages aim to resolve this dilemma. They include constructs for describing probability distributions and conditioning, and enable data scientists to express sophisticated probabilistic models as programs. More importantly, they come with the implementation of multiple algorithms for performing posterior inference for models and data sets expressed in the languages. The grand vision is that by using these languages, a data scientist no longer has to worry about the choice or design of such an inference algorithm but focuses on the design of an appropriate model, instead.

In this paper, we provide a denotational validation of inference algorithms for higher-order probabilistic programming languages, such as Church [Goodman et al. 2008], Anglican [Wood et al. 2014] and Venture [Mansinghka et al. 2014]. The correctness of these algorithms is subtle. The early version of the lightweight Metropolis-Hastings algorithm had a bug because of an incorrect acceptance ratio [Wingate et al. 2011]. The correctness often relies on intricate interplay between facts from probability theory and those from programming language theory. Moreover, correctness typically requires stronger results from probability theory than those used for the usual \mathbb{R}^n case in the machine-learning community (e.g., Green’s measure-theoretic justification of Markov Chain Monte Carlo rather than the usual one for \mathbb{R}^n based on density functions).

Our starting point is the body of existing results on validating inference algorithms for probabilistic programs [Borgström et al. 2016; Hur et al. 2015]. Those earlier results tend to be based on operational semantics, and often (not always) focus on first-order programs. By working in a modular way with monads, denotational semantics and higher-order functions, we are able to validate sophisticated inference algorithms, such as resample-move Sequential Monte Carlo [Doucet and Johansen 2011], that are complex yet modular, being composed of smaller reusable components, by combining our semantic analysis of these components.

The probabilistic programming language considered in the paper includes continuous distributions, which means that semantic accounts of them or their inference algorithms need to use measure theory and Lebesgue integration. Meanwhile, our semantic account uses a meta-language with higher-order functions for specifying and interpreting intermediate representations of probabilistic programs that are manipulated by components of inference algorithms. Such higher-order functions let us achieve modularity and handle higher-order functions in the target probabilistic programming language. These two features cause a tension because it is impossible to define a measurable space structure over the collection of measurable functions between arbitrary measurable

spaces that is compatible with standard operations on those functions, such as function application. We resolve the tension using quasi-Borel spaces [Heunen et al. 2017], a recently proposed mathematical structure that supports both function spaces and continuous distributions.

We define a semantic class of structures for various intermediate representations of probabilistic programs, and semantic validity criteria for transformations of these representations in terms of distribution preservation. We develop a collection of building blocks for composing representations. We use these building blocks to validate common inference algorithms such as Sequential Monte Carlo and Markov Chain Monte Carlo. To emphasize the connection between the semantic manipulation and its traditional measure theoretic origins, we use Kock’s synthetic measure theory. We demonstrate its usefulness by proving a quasi-Borel counterpart to the Metropolis-Hastings-Green theorem.

To ease the presentation, we proceed in two steps. First, we present our development in the discrete setting, where the set-theoretic account is simpler and more accessible. Then, after developing an appropriate mathematical toolbox, we transfer this account to the continuous case. Inference in the continuous setting, while conceptually very similar to the discrete case, is inseparable from our development. The semantic foundation for continuous distributions over higher-order functions has been very problematic in the past. The fact that our approach *does* generalise to the continuous case, and does so smoothly, is one of our significant contributions, only brought about through the careful combination of quasi-Borel spaces, synthetic measure theory, the metalanguage, and the inference building blocks.

The rest of the paper is structured as follows. Sec. 2 presents a core calculus, our metalanguage, with its type system and set-theoretic denotational semantics. Sec. 3 presents the core ideas of our development in a simpler discrete set-theoretic setting. Sec. 4 reviews the mathematical concepts required for dealing with continuous distributions. Sec. 5 presents representations and transformations for continuous distributions. Sec. 6 decomposes the common Sequential Monte Carlo inference algorithm into simpler inference representations and transformations. Sec. 7 similarly decomposes the general Trace Markov Chain Monte Carlo algorithm. Sec. 8 concludes. Basic results in synthetic measure theory are listed for the reader’s convenience in Appendix A.

2 THE CORE CALCULUS

We use a variant of the simply-typed λ -calculus with sums and inductive types, base types and constructors, primitives, and primitive recursion, but without effects. We also use monad-like constructs in the spirit of Moggi’s computational λ -calculus [1989]. The core calculus is very simple, and at places we need an inherently semantic treatment, which the core calculus alone cannot express. In those cases, we resort directly to the semantic structures, sets or spaces. However, the calculus still serves a very important purpose: every type and function expressed in it denote well-formed objects and well-formed morphisms. In the continuous case, using this calculus yields correct-by-construction quasi-Borel spaces and their morphisms, avoiding a tedious and error-prone manual verification. Using the core calculus also brings our theoretical development closer to potential implementations in functional languages.

2.1 Syntax

Fig. 1 (top) presents the types of our core calculus. To support inductive types, we include type variables, taken from a countable set ranged over by $\alpha, \beta, \gamma, \dots$. Our kind system will later ensure these type variables are *strictly positive*: they can only appear free covariantly – to the right of a function type. Variant types use constructor labels taken from a countable set ranged over by $\ell, \ell_1, \ell_2, \dots$. Variant types are in fact partial functions with a finite domain from the set of constructor labels to the set of types. When σ is a variant type, we write $(\ell \tau) \in \sigma$ for the assertion

$\tau, \sigma, \rho ::=$	types		
α	positive variable	$ \tau \rightarrow \sigma$	function
$\{ \ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n \}$	variant	$ A$	base
$\mathbf{1} \mid \tau * \sigma$	finite product	$ F \tau$	base constructors
$\mu\alpha.\tau$	inductive type	$\Gamma := x_1 : \tau_1, \dots, x_n : \tau_n$	variable contexts
$t, s, r ::=$	terms		
x	variable	$\mathbf{match} \ t$	binary products
$\tau.\ell t$	variant constructor	$\mathbf{with} \ (x, y) \rightarrow s$	
$() \mid (t, s)$	nullary and binary tuples	$\mathbf{match} \ t$	inductive types
$\tau.\mathbf{roll}$	iso-inductive constructor	$\mathbf{with} \ \mathbf{roll} \ x \rightarrow s$	
$\lambda x : \tau.t$	function abstraction	$\tau.\mathbf{fold} \ t$	inductive recursion
$\mathbf{match} \ t$	pattern matching: variants	$t \ s$	function application
$\mathbf{with} \ \{ \ell_1 x_1 \rightarrow s_1 \mid \dots \mid \ell_n x_n \rightarrow s_n \}$		φ	primitive

Fig. 1. Core calculus types (top) and terms (bottom)

that σ assigns the type τ to ℓ . We include the standard unit type, binary products, and function types. We include unary uninterpreted base types and constructors. While we use a list syntax for variable contexts Γ , they are in fact partial functions with a finite domain from the countable set of variables, ranged over by x, y, z, \dots , to the set of types.

We desugar stand-alone labels in a variant type $\{\dots \mid \ell \mid \dots\}$ to the unit type $\{\dots \mid \ell () \mid \dots\}$. We also desugar seemingly-recursive type declarations $\tau := \sigma[\alpha \mapsto \tau]$ to $\tau := \mu\alpha.\sigma$.

Example 2.1. The type of booleans is given by $\mathbf{bool} := \{\mathbf{True} \mid \mathbf{False}\}$. The type of natural numbers is given by $\mathbb{N} := \{\mathbf{Zero} \mid \mathbf{Succ} \ \mathbb{N}\}$ desugaring to $\mathbb{N} := \mu\alpha.\{\mathbf{Zero} \mid \mathbf{Succ} \ \alpha\}$. The type of α -lists is given by $\mathbf{List} \ \alpha := \{\mathbf{Nil} \mid \mathbf{Cons} \ \alpha * \mathbf{List} \ \alpha\}$, desugaring to $\mathbf{List} \ \alpha := \mu\beta.\{\mathbf{Nil} \mid \mathbf{Cons} \ \alpha * \beta\}$.

Base types and constructors allow us to include semantic type declarations into our calculus. For example, we will always include the following base types:

- \mathbb{I} : unit interval $[0, 1]$;
- $\overline{\mathbb{R}}$: extended real line $[-\infty, \infty]$;
- $\overline{\mathbb{R}}_+$: non-negative extended reals $[0, \infty]$;
- \mathbb{R} : real line $(-\infty, \infty)$;
- \mathbb{R}_+ : non-negative reals $[0, \infty)$; and

In addition, once we define a type constructor such as $\mathbf{List} \ \alpha$, we will later reuse it as a base type constructor $\mathbf{List} \ \tau$, effectively working in an extended calculus. Thus we are working with a family of calculi, extending the base signature with each type definition in our development.

Fig. 1 (bottom) presents the terms in our core calculus. Variant constructor terms $\tau.\ell t$ are annotated with their variant type τ to avoid label clashes. The tupling constructors are standard. We use *iso-inductive* types: construction of inductive types requires an explicit rolling of the inductive definition such as $\mathbb{N}.\mathbf{roll} \ (\mathbf{Zero}())$. Variable binding in function abstraction is *intrinsically typed* in standard Church-style. We include standard pattern matching constructs for variants, binary products, and inductive types. We include a structural recursion construct $\tau.\mathbf{fold}$ for every inductive type τ . Function application is standard, as is the inclusion of primitives.

To ease the construction of terms, we use the standard syntactic sugar (e.g. $\mathbf{let} \ x = t \ \mathbf{in} \ s$ for $(\lambda x.t)s$, $\mathbf{if} \ \mathbf{then} \ \mathbf{else}$ for pattern matching booleans), informally elide types from the terms, elide \mathbf{roll} ing/unrolling inductive types, and informally use nested pattern matching.

Example 2.2. For $\mathbf{List} \ \tau = \mu\alpha.\{\mathbf{Nil} \mid \mathbf{Cons} \ \tau * \alpha\}$, we can express standard list manipulation:

$$\begin{aligned}
 x :: x_s = \mathbf{Cons}(x, x_s) & \quad \mathbf{foldr} \ a \ f = \mathbf{List} \ \tau.\mathbf{fold} \ \lambda\{\mathbf{Nil} \rightarrow a \mid \mathbf{Cons}(x, b) \rightarrow f(x, b)\} \\
 x_s \# y_s = \mathbf{foldr} \ y_s \ (\::) \ x_s & \quad \mathbf{map} \ f \ x_s = \mathbf{foldr} \ [\] \ (\lambda\{(y, y_s) \rightarrow (f(y), y_s)\})
 \end{aligned}$$

$$\begin{array}{c}
\frac{}{\Delta \vdash_k \alpha : \text{type}} (\alpha \in \Delta) \quad \frac{\text{for all } 1 \leq i \leq n: \Delta \vdash_k \tau_i : \text{type}}{\Delta \vdash_k \{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\} : \text{type}} \quad \frac{}{\Delta \vdash_k \mathbf{1} : \text{type}} \\
\frac{\Delta \vdash_k \tau : \text{type} \quad \Delta \vdash_k \sigma : \text{type}}{\Delta \vdash_k \tau * \sigma : \text{type}} \quad \frac{\Delta \vdash_k \mu\alpha.\tau : \text{type}}{\Delta \vdash_k \tau : \text{type}} \quad \frac{\Delta \vdash_k \tau \rightarrow \sigma : \text{type}}{\text{for all } (x : \tau) \in \Gamma: \vdash_k \tau : \text{type}} \\
\frac{}{\Delta \vdash_k A : \text{type}} \quad \frac{}{\Delta \vdash_k F \tau : \text{type}} \quad \frac{}{\vdash_k \Gamma : \text{context}}
\end{array}$$

Fig. 2. Core calculus kind system

$$\begin{array}{c}
\frac{}{\Gamma \vdash t : \tau} ((x : \tau) \in \Gamma) \quad \frac{\Gamma \vdash t : \tau_i}{\Gamma \vdash \tau.\ell_i t : \tau} ((\ell_i \tau_i) \in \tau) \quad \frac{}{\Gamma \vdash () : \mathbf{1}} \\
\frac{\Gamma \vdash (t, s) : \tau * \sigma \quad \Gamma \vdash \tau.\text{roll} : (\sigma[\alpha \mapsto \tau]) \rightarrow \tau \quad \frac{}{\Gamma, x : \tau \vdash t : \sigma} (\tau = \mu\alpha.\sigma)}{\Gamma \vdash t : \{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\} \quad \text{for each } 1 \leq i \leq n: \Gamma, x_i : \tau_i \vdash s_i : \tau} \\
\frac{\Gamma \vdash \text{match } t \text{ with } \{\ell_1 x_1 \rightarrow s_1 \mid \dots \mid \ell_n x_n \rightarrow s_n\} : \tau}{\Gamma \vdash t : \sigma * \rho \quad \Gamma, x : \sigma, y : \rho \vdash s : \tau} \quad \frac{\Gamma \vdash t : \mu\alpha.\sigma \quad \Gamma, x : \sigma[\alpha \mapsto \mu\alpha.\sigma] \vdash s : \tau}{\Gamma \vdash \text{match } t \text{ with } (x, y) \rightarrow s : \tau} \\
\frac{\Gamma \vdash t : (\sigma[\alpha \mapsto \rho]) \rightarrow \rho \quad \frac{}{\Gamma \vdash t : \sigma \rightarrow \tau} (\tau = \mu\alpha.\sigma)}{\Gamma \vdash \tau.\text{fold } t : \tau \rightarrow \rho} \quad \frac{\Gamma \vdash \text{match } t \text{ with roll } x \rightarrow s : \tau}{\Gamma \vdash t s : \tau} \quad \frac{}{\Gamma \vdash \varphi : \tau_\varphi}
\end{array}$$

Fig. 3. Core calculus type system

where we abbreviate $[a_1, \dots, a_n]$ to $\text{Cons}(a_1, \dots, \text{Cons}(a_n, \text{Nil}) \dots)$.

2.2 Type System

To ensure the well-formedness of types, which involve type variables, we use a simple kind system, presented in Fig. 2. Each kinding judgement $\Delta \vdash_k \tau : \text{type}$ asserts that a given type τ is well-formed in the *type variable context* Δ , which is finite set of type variables.

The kinding judgements are standard. All type variables must be bound by the enclosing context, or by an inductive type binder. The contravariant position in the function type $\tau \rightarrow \sigma$ must contain a *closed* type, ensuring that free type variables can only appear in strictly positive positions. Variable contexts Γ must only assign closed types.

Example 2.3. The types from Ex. 2.1 are well-kinded: $\vdash_k \text{bool}, \mathbb{N}, \text{List } \alpha : \text{type}$.

We define capture avoiding substitution of types for type variables in the standard way, which obeys the usual structural properties. Henceforth we consider only well-formed types in context, leaving the context implicit wherever possible, and gloss over issues of alpha-convertibility of bound type variables.

To type terms, we assume each primitive φ has a well-formed type $\vdash_k \tau_\varphi : \text{type}$ associated with it. Fig. 3 presents the resulting type system. Each typing judgement $\Gamma \vdash t : \tau$ asserts that a given term t is well-typed with the well-formed closed type $\vdash_k \tau : \text{type}$ in the variable context $\vdash_k \Gamma : \text{context}$.

The rules are standard. By design, every term has at most one type in a given context.

Example 2.4. Once desugared, the list manipulation terms from Ex. 2.2 have types:

$$\begin{aligned} (::) &: \tau * \text{List } \tau \rightarrow \text{List } \tau & \text{foldr} &: \sigma * (\tau * \sigma \rightarrow \sigma) * \text{List } \tau \rightarrow \sigma \\ \text{map} &: (\tau \rightarrow \sigma) \rightarrow (\text{List } \tau \rightarrow \text{List } \sigma) & (+) &: (\text{List } \tau) * (\text{List } \tau) \rightarrow \text{List } \tau \end{aligned}$$

2.3 Primitive Recursion

As is well-known [Geuvers and Poll 2007; Hutton 1999], structural recursion on inductive types allows us to express primitive recursion. By ‘primitive recursion’, we mean recursing through values of an inductive type $\mu\alpha.\sigma$ using a term of the form: $\Gamma, k : \sigma[\alpha \mapsto (\mu\alpha.\sigma) * \rho] \vdash t : \rho$ with the intention that t can use either arbitrary (total) processing on the sub-structures of its input k , or make a primitive recursive call to itself with a sub-structure. In order to desugar such a term into a function of type $\tau * (\mu\alpha.\sigma) \rightarrow \rho$, we use terms of the following type, defined by induction on types:

$$\pi_{\alpha.\sigma,\rho} : \sigma[\alpha \mapsto (\mu\alpha.\sigma) * \rho] \rightarrow \sigma[\alpha \mapsto \mu\alpha.\sigma]$$

and interpret the primitive recursive declaration t embodied by:

$$\Gamma, x : \mu\alpha.\sigma \vdash \mathbf{match} (\mu\alpha.\sigma).\mathbf{fold} \left(\lambda k : \sigma[\alpha \mapsto (\mu\alpha.\sigma) * \rho]. (\mathbf{roll} \pi_{\alpha.\sigma,\rho} k, t) \right) x \\ \mathbf{with} (_, r) \rightarrow r : \sigma$$

This translation is global in nature: the structure of the term π depends on the type of t . Thus, it does not constitute a *macro* translation [Felleisen 1991]. With this point in mind, we will allow ourselves to use primitive recursive definitions.

Example 2.5. We define a function $\text{aggr} : \text{List}(\mathbb{R}_+ * X) \rightarrow \text{List}(\mathbb{R}_+ * X)$ which takes a list of weighted values and aggregates all the weights based on their values. We make use of the auxiliary function $\text{add} : (\mathbb{R}_+ * X) * \text{List}(\mathbb{R}_+ * X) \rightarrow \text{List}(\mathbb{R}_+ * X)$, which adds a weighted value to an already aggregated list. We define add by primitive recursion:

$$\begin{aligned} \text{add}((s, a), x_s) &:= \mathbf{match} x_s \mathbf{with} \{ [] && \rightarrow [(s, a)] && \text{-- new entry} \\ (r, x) &:: x_s \rightarrow \mathbf{if} \quad x = a && && \\ && \mathbf{then} (s + r, a) &:: x_s && \text{-- accumulate} \\ && \mathbf{else} (r, x) &:: \text{add}((s, a), x_s) \} && \text{-- recurse} \end{aligned}$$

and set $\text{aggr} := \text{foldr} [] \text{add}$. This example makes use of an equality predicate between X elements, restricting its applicability.

2.4 Denotational Semantics

We give a set-theoretic semantics to the calculus. In such set-theoretic semantics, types-in-context $\Delta \vdash_k \tau : \text{type}$ are interpreted as functors $\llbracket \tau \rrbracket : \mathbf{Set}^\Delta \rightarrow \mathbf{Set}$, i.e., $\llbracket \tau \rrbracket$ assigns a set $\llbracket \tau \rrbracket (X_\alpha)_{\alpha \in \Delta}$ for every Δ -indexed tuple of sets, and a function

$$\llbracket \tau \rrbracket (f_\alpha : X_\alpha \rightarrow Y_\alpha)_{\alpha \in \Delta} : \llbracket \tau \rrbracket (X_\alpha) \rightarrow \llbracket \tau \rrbracket (Y_\alpha)$$

for every Δ -indexed tuple of functions between the sets with corresponding index, and this assignment preserves composition and identities.

In order to interpret iso-inductive types $\mu\alpha.\tau$, we need canonical isomorphisms between the sets $\llbracket \tau \rrbracket (\llbracket \mu\alpha.\tau \rrbracket) \cong \llbracket \mu\alpha.\tau \rrbracket$. We will do this in a standard way, by interpreting $\llbracket \mu\alpha.\tau \rrbracket$ as the initial algebra for the functor $\llbracket \tau \rrbracket : [\mathbf{Set}^\Delta \rightarrow \mathbf{Set}] \rightarrow [\mathbf{Set}^\Delta \rightarrow \mathbf{Set}]$. This means that for every functor $A : \mathbf{Set}^\Delta \rightarrow \mathbf{Set}$ with a natural family of functions $\{a_X : (\llbracket \tau \rrbracket A)(X) \rightarrow A(X)\}_{X \in \mathbf{Set}^\Delta}$, there is a canonical natural family of functions $\{\text{fold}_X : \llbracket \mu\alpha.\tau \rrbracket (X) \rightarrow A(X)\}_{X \in \mathbf{Set}^\Delta}$.

A technical requirement is needed to ensure that this initial algebra exists: we fix a regular cardinal κ , and demand that each type denotes a κ -ranked functor (ranked functor for short), that is,

$$\begin{aligned}
\llbracket \alpha \rrbracket d &:= d(\alpha) & \llbracket \{\ell_1 \tau_1 \mid \dots \mid \ell_n \tau_n\} \rrbracket d &:= \sum_{i=1}^n \llbracket \tau_i \rrbracket d & \llbracket \mathbf{1} \rrbracket d &:= \mathbf{1} \\
\llbracket \tau * \sigma \rrbracket d &:= (\llbracket \tau \rrbracket d) \times (\llbracket \sigma \rrbracket d) & \llbracket A \rrbracket d &:= \mathcal{B}[A] \\
\llbracket F \tau \rrbracket d &:= \mathcal{B}[\llbracket F \rrbracket(\llbracket \tau \rrbracket d)] & \llbracket \mu \alpha. \tau \rrbracket d &:= \mu X. \llbracket \tau \rrbracket d[\alpha \mapsto X] & \llbracket \tau \rightarrow \sigma \rrbracket d &:= (\llbracket \sigma \rrbracket d)^{\llbracket \tau \rrbracket d}
\end{aligned}$$

Fig. 4. Core calculus type-level semantics

that it denotes a functor that preserves κ -filtered colimits¹. The κ -ranked functors are closed under composition, products, sums, and initial algebras. Initial algebras for κ -ranked functors on locally presentable categories always exist, because they can be built in an iterative way by transfinite induction (see e.g. [Kelly 1980]).

2.4.1 Set-Theoretic Interpretation. To interpret types, we assume a given interpretation $\mathcal{B}[-]$ of the base types A as sets $\mathcal{B}[A]$ and of base type constructors F as ranked functors $\mathcal{B}[F] : \mathbf{Set} \rightarrow \mathbf{Set}$. We then interpret each well-formed type in context $\Delta \vdash_{\kappa} \tau : \text{type}$ as a ranked functor $\llbracket \tau \rrbracket : \mathbf{Set}^{\Delta} \rightarrow \mathbf{Set}$, as depicted in Fig. 4.

In this definition, the parameter d may be either a tuple of sets or functions. When interpreting type variables, we write $d(\alpha)$ for the α -indexed component of d . The interpretation of simple types uses disjoint unions, singletons, finite products, and exponentials, i.e. the bi-cartesian closed structure of \mathbf{Set} . We interpret inductive types $\llbracket \mu \alpha. \tau \rrbracket d$ using the initial algebra for the ranked functor $\lambda X. \llbracket \tau \rrbracket d[\alpha \mapsto X] : \mathbf{Set} \rightarrow \mathbf{Set}$. In the semantics of the function type $\tau \rightarrow \sigma$, the exponential makes no use of the functor's arguments, and relies on the fact that all type variables are strictly positive. We use the given interpretation of base types and type constructors to interpret them.

LEMMA 2.6. *The semantics of types is well-defined: every well-formed type $\Delta \vdash_{\kappa} \tau : \text{type}$ denotes a ranked functor $\llbracket \tau \rrbracket : \mathbf{Set}^{\Delta} \rightarrow \mathbf{Set}$. In particular, every closed type denotes a set.*

The proof is by induction on the kinding judgements, using well-known properties of \mathbf{Set} . We will always interpret the base types \mathbb{I}, \mathbb{R} , etc. by the sets they represent.

Example 2.7. We calculate the denotations of the types from Ex. 2.1. Booleans denote a two-element set $\llbracket \text{bool} \rrbracket = \{\text{False}, \text{True}\}$, and the natural numbers denote the set of natural numbers $\llbracket \mathbb{N} \rrbracket = \mathbb{N}$. By Lemma 2.6, $\llbracket \text{List} \rrbracket$ denotes a ranked functor $\text{List} : \mathbf{Set} \rightarrow \mathbf{Set}$, and this functor is given by the set of sequences of X -elements $\text{List } X := \bigcup_{n \in \mathbb{N}} X^n$.

Beyond establishing the well-definedness of the semantic interpretation, Lemma 2.6 equips us with syntactic means to define ranked functors. Once defined, we can add these functors to our collection of base types (in an extended instance of the core calculus). In the sequel, we will often restrict a given ranked functor $F : \mathbf{Set} \rightarrow \mathbf{Set}$ by specifying a subset $G X \subseteq F X$. Doing so is analogous to imposing an *invariant* on a datatype. The subsets $G X$ form a *subfunctor* $G \subseteq F$ precisely if they are closed under the functorial action of F , i.e., for every function $f : X \rightarrow Y$ and $a \in G X$, $F f(a) \in G Y$.

LEMMA 2.8. *Subfunctors of ranked functors over \mathbf{Set} are ranked.*

We can prove this lemma directly, but it also follows from a higher-level argument using the commutation of finite limits and κ -directed colimits in \mathbf{Set} .

¹We do not use simpler classes of functors, such as *polynomial functors* or *containers*, as they are not closed under subfunctors, given by subsets in the discrete case and subspaces in the continuous case, which we need in the sequel.

<p>instance <i>Monad</i> (List) where return $x = [x]$ $x_s \gg= f = \text{foldr } [] (\lambda(x, y_s). f(x) \# y_s) x_s$</p> <p>(a) Declaring monadic interfaces</p>	<table border="0"> <thead> <tr> <th style="text-align: left; border-bottom: 1px solid black;">Sugar</th> <th style="text-align: left; border-bottom: 1px solid black;">Elaboration</th> </tr> </thead> <tbody> <tr> <td>• $x \leftarrow t; s$</td> <td>$t \gg= \lambda x. s$</td> </tr> <tr> <td>• return t</td> <td>$\text{return}^T t$</td> </tr> <tr> <td>• $t; s$</td> <td>$_ \leftarrow t; s$</td> </tr> </tbody> </table> <p>(b) Haskell's do-notation</p>	Sugar	Elaboration	• $x \leftarrow t; s$	$t \gg= \lambda x. s$	• return t	$\text{return}^T t$	• $t; s$	$_ \leftarrow t; s$
Sugar	Elaboration								
• $x \leftarrow t; s$	$t \gg= \lambda x. s$								
• return t	$\text{return}^T t$								
• $t; s$	$_ \leftarrow t; s$								

Fig. 5. Monadic programming notation

2.5 Monadic Programming

In the sequel, we will be working with types that support a monadic programming style. More precisely, a *monadic interface* \underline{T} consists of a triple $\underline{T} = (T, \text{return}^{\underline{T}}, \gg=^{\underline{T}})$ where: T assigns to each set X a set TX ; $\text{return}^{\underline{T}}$ assigns to each set X a function $\text{return}_X^{\underline{T}} : X \rightarrow TX$; and $\gg=^{\underline{T}}$ assigns to each pair of sets X and Y a function $\gg=_{X,Y}^{\underline{T}} : TX \times (TY)^X \rightarrow TY$. We borrow Haskell's type-class syntax to define such interfaces. As an example, Fig. 5a defines a monadic interface over List.

Each such monadic interface \underline{T} allows us to use standard do-notation summarised in Fig. 5b. Though simple in principle, we must take care when treating this notation as syntactic sugar, as choosing the appropriate function return_X or $\gg=_{X,Y}$ at each desugaring step must take typing information into account. When we use do-notation in the sequel, we ensure that such choices can be disambiguated. Finally, we will delimit our use of do-notation to within a *do-block* $\underline{T}.\text{do}\{\dots\}$, omitting the monadic interface \underline{T} or the entire delimiter when either is clear from the context.

Importantly, we do not insist that a monadic interface satisfies the monad associativity and unit laws: $(\text{return } x) \gg= f = f(x)$, $a \gg= \text{return} = a$, and $(a \gg= f) \gg= g = a \gg= (\lambda x. (f x \gg= g))$.

3 DISCRETE INFERENCE

We can now lay-out the core ideas in the simpler, set-theoretic case: a semantic structure for higher-order (discrete) probabilistic programs, intermediate representations of these programs for the purpose of inference, valid transformations between these representations, and modular building blocks for creating new representations and transformations from existing ones. For simplicity, we consider representations and transformations from simple rather naive inference algorithms only in this section. In Sec. 6 and Sec. 7, we show how the core ideas here apply to advanced algorithms when aided with further technical developments.

3.1 The Mass Function Monad

For our purposes, probabilistic programming languages contain standard control-flow mechanisms and data types, such as our core calculus, together with *probabilistic choice* and *conditioning* operations. In the discrete case, these are given by two effectful operations:

$$\frac{}{\Gamma \vdash_{\text{comp}} \mathbf{flip} : \text{bool}} \qquad \frac{\Gamma \vdash t : \mathbb{R}_+}{\Gamma \vdash_{\text{comp}} \mathbf{score } t : \mathbf{1}}$$

In Bayesian probabilistic programming, we think of **flip** as drawing from a (uniform) prior distribution on `bool`, and of **score** as recording a likelihood. Typically, one calls $\mathbf{score}(f(x))$ where f is a density function of a distribution, which records the likelihood of observing data x from the distribution f . The score might be zero, a hard constraint: this path is impossible. The score might be in the unit interval, the probability of a discrete observation; but in general a likelihood function can take any positive real value. The inference problem is to approximate the posterior distribution, from the unnormalized posterior defined by the program, combining a prior and likelihood.

To give a set-theoretic semantic structure to such a higher-order language with these two constructs, it suffices to give a monadic interface \underline{T} for which the associativity and unit laws hold, together with two functions:

$$\text{flip} : [\mathbf{1}] \rightarrow T [\text{bool}] \qquad \text{score} : [\mathbb{R}_+] \rightarrow T [\mathbf{1}]$$

For the purposes of the discrete development, the following monad fits the bill. A (finite) *mass function* over a set X is a function $\mu : X \rightarrow \mathbb{R}_+$ for which there exists a finite set $F \subseteq X$ such that μ is 0 outside F : in other words, the support set $\text{supp } \mu := \{x \in X \mid \mu(x) \neq 0\}$ is finite. For every set X , let $\text{Mass } X := \{\mu : X \rightarrow \mathbb{R}_+ \mid \mu \text{ is a mass function}\}$. The *mass function monad* is given by:

$$\begin{aligned} \underline{\text{Mass}} &:= \text{instance Monad (Mass) where} \\ &\quad \text{return } x_0 = \lambda x. \text{ if } (x = x_0) \text{ then } \mathbf{1} \text{ else } 0 \\ &\quad \mu \gg= f = \lambda y. \sum_{x \in \text{supp } \mu} \mu(x) \cdot (f(x)(y)) \end{aligned}$$

and we set $\text{flip} = \lambda_. \frac{1}{2}$ and $\text{score } r = \lambda\{ () \rightarrow r\}$. Intuitively, values of $\text{Mass } X$ represent unnormalized probabilistic computations of a result in X . From the Bayesian perspective, the meaning of a program is the unnormalized posterior.

LEMMA 3.1. *The monadic interface $\underline{\text{Mass}}$ defines a ranked monad over Set .*

This monad is also known as the *free positive cone monad*, as it constructs the ‘positive fragment’ of a vector space over the field of reals with basis X .

3.2 Inference Representations

The mass function semantics is accurate, but idealised: realistic implementations cannot be expected to compute mass functions at arbitrary types, and especially at higher-order types. Instead, probabilistic inference engines would manipulate some representation of the program, while maintaining its semantics.

Definition 3.2. A *discrete inference representation* \underline{T} is a sextuple

$$\underline{T} = \left(T, \text{return}^T, \gg=^T, \text{flip}^T, \text{score}^T, m^T \right)$$

consisting of:

- a monadic interface $(T, \text{return}^T, \gg=^T)$;
- two functions $\text{flip}^T : \mathbf{1} \rightarrow T \mathbf{2}$ and $\text{score}^T : \mathbb{R}_+ \rightarrow T \mathbf{1}$, where $\mathbf{1} := [\mathbf{1}]$, $\mathbf{2} := [\text{bool}]$; and
- an assignment of a *meaning* function $m_X^T : TX \rightarrow \text{Mass } X$ for every set X

such that the following laws hold for all sets X, Y , and $x \in X$, $a \in TX$, $r \in \mathbb{R}_+$, and $f : X \rightarrow TY$:

$$\text{return}^{\text{Mass}} x = m(\text{return}^T x) \qquad m(a \gg=^T f) = (m a) \gg=^{\text{Mass}} \lambda x. m(f x)$$

$$m(\text{flip}^T) = \text{flip}^{\text{Mass}} \qquad m(\text{score}^T r) = \text{score}^{\text{Mass}} r$$

As with monadic interfaces, we use a type-class notation for defining inference representations.

Example 3.3 (Discrete weighted sampler). Consider the type

$$\text{Term } \alpha := \{\text{Return } (\mathbb{R}_+ * \alpha) \mid \text{Flip } (\text{Term } \alpha * \text{Term } \alpha)\}$$

which induces a ranked functor Term . The elements of $\text{Term } X$ are binary trees, which we call terms, whose leaves contain weighted values of type X . Fig. 6a presents the inference representation structure of the functor Term . Flip represents a probabilistic choice while Return holds the final value and the total weight for the branch. Thus an immediately returning computation is

instance *Discrete Monad* (Term) **where**

return $x = \text{Return}(1, x)$

$a \gg= f = \text{let } (scale : \mathbb{R}_+ * \text{Term } X \rightarrow \text{Term } X) = \quad \text{-- uses primitive recursion}$
 $\lambda s. \lambda \{ \text{Return}(r, x) \rightarrow \text{Return}(s \cdot r, x)$
 $\quad | \text{Flip}(k_{\text{False}}, k_{\text{True}}) \rightarrow \text{Flip}(scale(s, k_{\text{False}}), scale(s, k_{\text{True}})) \}$
in match a **with** $\{$
 $\quad \text{Return}(r, x) \rightarrow scale(r, f \ x)$
 $\quad | \text{Flip}(k_{\text{False}}, k_{\text{True}}) \rightarrow \text{Flip}(k_{\text{False}} \gg= f, \text{-- uses primitive recursion}$
 $\quad \quad \quad k_{\text{True}} \gg= f) \}$

flip = Flip (Return (1, False), Return (1, True))

score $r = \text{Return}(r, ())$

$m \ a = \text{fold } \lambda \{ \text{Return}(r, x) \rightarrow \text{Mass} . \text{do} \{ \text{score } r; \text{return } x \}$
 $\quad | \text{Flip}(\mu_{\text{False}}, \mu_{\text{True}}) \rightarrow \text{Mass} . \text{do} \{ x \leftarrow \text{flip};$
 $\quad \quad \quad \text{if } x \text{ then } \mu_{\text{True}} \text{ else } \mu_{\text{False}} \}$

(a) Discrete weighted sampler representation

instance *Discrete Monad* (Enum) **where**

return $x = [(1, x)]$

$x_s \gg= f = \text{let } (scale : \mathbb{R}_+ * \text{Enum } X \rightarrow \text{Enum } X) =$
 $\lambda \{ (r, x_s) \rightarrow \text{map } \lambda \{ (s, y) \rightarrow (r \cdot s, y) \}$
 $\quad \quad \quad x_s \}$
in foldr $[]$
 $\lambda \{ ((r, x), y_s) \rightarrow scale(r, f \ x) \# y_s \}$
 $\quad \quad \quad x_s$

flip = $[(\frac{1}{2}, \text{False}), (\frac{1}{2}, \text{True})]$

score $r = [(r, ())]$

$m \ x_s = \lambda a. \quad \text{-- } m \ x_s \ a = \sum_{\substack{(r,x) \in x_s \\ x=a}} r$

foldr $0 \left(\lambda \{ ((r, x), s) \rightarrow \right.$
 $\quad \quad \quad \text{if } x = a \text{ then } r + s \text{ else } s \}$ $\left. \right) x_s$

(b) Discrete enumeration sampler

instance *Inf Trans* (W) **where**

lift $_T \ a = \underline{T} . \text{do} \{ x \leftarrow a;$

$\quad \quad \quad \text{return}(1, x) \}$

return $_{W \underline{T}} \ x = \text{return}^T(1, x)$

$a \gg=_{W \underline{T}} \ f = \underline{T} . \text{do} \{ (r, x) \leftarrow a;$

$\quad \quad \quad (s, y) \leftarrow f(x);$

$\quad \quad \quad \text{return}(r \cdot s, y) \}$

flip $_{W \underline{T}} = \text{lift } \text{flip}^T$

score $_{W \underline{T}} \ r = \text{return}^T(r, ())$

$m_{W \underline{T}} \ a = \lambda x. \sum_{(r,x) \in \text{supp } m^T(a)} m^T(a) \ r$

(tmap $t) \ x = \underline{t}_{\mathbb{R}_+ * X}$

(c) Discrete weighting transformer

Fig. 6. Example inference representations (a,b) and transformers (c)

represented by a leaf with weight 1. The auxiliary function *scale* in the definition of $\gg=$ scales the leaves of its input term by the input weight. The function $\gg=$ itself substitutes terms for the leaves according to its input function *f*, making sure the newly grafted terms are scaled appropriately. The probabilistic choice operation **flip** constructs a single node with each leaf recording the probabilistic choice *unweighted*. Conditioning records the input weight.

The meaning function recurses over the term, replacing each node representing a probabilistic choice by probabilistic choice of the mass function monad, and reweighting the end result appropriately.

The main step in validating the inference representation laws involves $\gg=$: first show that composing the meaning function with the auxiliary function *scale* scales the meaning of the input term appropriately, and then proceed by structural induction on terms.

The weighted sampler representation in fact forms a proper monad over **Set**: it is the free monad for an algebraic theory with a binary operation `flip` and unary operations `scorer`, subject to $\mathbf{flip}(\mathbf{score}_r(x), \mathbf{score}_r(y)) = \mathbf{score}_r(\mathbf{flip}(x, y))$. As the mass function monad also validates these equations, the meaning function is then the unique monad morphism from **Term** to **Mass** preserving the operations `flip` and `score`.

However, we emphasise that an inference representation need not form a proper monad, and that the meaning function need not be a monad morphism. Indeed, the Pop Sam representation introduced in Sec. 6 is not a monad and most of the non-trivial inference transformations we discuss are not monad morphisms.

The weighted sampler representation allows us to incorporate both intensional and operational aspects into our development. Bayesian inference ultimately reduces a representation into probabilistic simulation. The weighted sampler representation can thus act as an internal representation of this simulation. Moreover, its continuous analogue will allow us to manipulate traces when analysing the Trace Markov Chain Monte Carlo algorithm in Sec. 7.

Example 3.4 (Enumeration). The type $\mathbf{Enum} \alpha := \mathbf{List}(\mathbb{R}_+ * \alpha)$ induces a ranked functor **Enum**. Elements of $\mathbf{Enum} X$ form an enumeration of the mass function they represent, with the same value x potentially appearing multiple times with different weights. Values not appearing in the list at all have weight 0.

Fig. 6b presents an inference representation structure using **Enum**. Returning a value lists the unique non-zero point mass. The \gg operation applies the given function to each element listed, scales the list appropriately and accumulates all intermediate lists. The choice operation enumerates both branches with equal probability, and conditioning inserts a scaling factor. The meaning function assigns to an element the sum of its weights. This definition uses an equality predicate.

Establishing the inference representation laws is straightforward.

3.3 Inference Transformations

We can now define the central validity criterion in our development. We decompose Bayesian inference algorithms into smaller transformations between inference representations. To be correct, these transformations need to preserve the meaning of the representation they manipulate:

Definition 3.5. Let $\underline{T}, \underline{S}$ be two inference representations. A *discrete inference transformation* $\underline{t} : \underline{T} \rightarrow \underline{S}$ assigns to each set X a function $\underline{t}_X : TX \rightarrow SX$ satisfying $m^{\underline{T}}(a) = m^{\underline{S}}(\underline{t}_X(a))$ for every $a \in TX$.

This validity criterion guarantees nothing beyond the preservation of the overall mass function of our representation. The transformed representation may not be better for inference along any axis, such as better convergence properties or execution time. It is up to the inference algorithm designer to convince herself of such properties by other means: formal, empirical, or heuristic.

Some transformations change the representation type:

Example 3.6 (Enumeration). Define a transformation: $\underline{t} : \underline{\mathbf{Term}} \rightarrow \underline{\mathbf{Enum}}$ by:

$$\underline{t} := \lambda \{ \begin{array}{l} \mathbf{Return}(r, x) \quad \rightarrow \mathbf{Enum} . \mathbf{do} \{ \mathbf{score} \ r; \mathbf{return} \ x \} \\ \mathbf{Flip}(x_s^{\mathbf{False}}, x_s^{\mathbf{True}}) \rightarrow \mathbf{Enum} . \mathbf{do} \{ b \leftarrow \mathbf{flip}; \mathbf{if} \ b \ \mathbf{then} \ x_s^{\mathbf{True}} \ \mathbf{else} \ x_s^{\mathbf{False}} \} \end{array}$$

Straightforward calculation shows it preserves the meaning functions.

The last example is a special case: analogous functions form inference transformations $\underline{t}_{\underline{T}} : \underline{\mathbf{Term}} \rightarrow \underline{T}$ for every discrete inference representation \underline{T} . To establish meaning preservation, calculate that both $m^{\underline{\mathbf{Term}}}$ and $m^{\underline{T}} \circ \underline{t}_{\underline{T}}$ are monad morphisms that preserve probabilistic choice and conditioning and appeal to the initiality of $\underline{\mathbf{Term}}$.

An inference transformation need not be natural:

Example 3.7 (Aggregation). Recall the functions $\text{aggr}_X : \text{List}(\mathbb{R}_+ * X) \rightarrow \text{List}(\mathbb{R}_+ * X)$ from Ex. 2.5 which aggregate list elements according to their X component by summing their weights. It forms an inference transformation $\text{aggr} : \underline{\text{Enum}} \rightarrow \underline{\text{Enum}}$. The meaning preservation proof uses straightforward structural induction. Note that aggr is not a natural transformation.

3.4 Inference Transformers

We can decompose the weighted sampler representation $\underline{\text{Term}}$, which forms a monad, by transforming the *discrete sampler* representation $\text{DSam } X := \{\text{Return } X \mid \text{Sample}(\text{DSam } X * \text{DSam } X)\}$ with the following *writer monad transformer* $\text{WT } X := T(\mathbb{R}_+ * X)$, i.e. $\text{Term} = \text{WDSam}$. Such decompositions form basic building blocks for constructing and reasoning about more sophisticated representations.

Definition 3.8. An inference transformer \underline{F} is a triple $(F, \text{tmap}^F, \text{lift}^F)$ whose components assign:

- inference representation $F \underline{T}$ to every inference representation \underline{T} ;
- inference transformation $\text{tmap}^F \underline{t} : \underline{FT} \rightarrow \underline{FS}$ to every inference transformation $\underline{t} : \underline{T} \rightarrow \underline{S}$;
and
- inference transformation $\text{lift}_T : \underline{T} \rightarrow F \underline{T}$ to every inference representation \underline{T} .

We use type-class notation for defining inference transformers.

Example 3.9. The *weighting* inference transformer structure on $\text{WT } X := T(\mathbb{R}_+ * X)$ is given in Fig. 6c. We lift a representation in \underline{T} into $\text{WT } \underline{T}$ by assigning weight 1 to it. The monadic interface uses the standard writer monad for the multiplication structure on \mathbb{R}_+ , accumulating the weights as computation proceeds. We lift the probabilistic choice from \underline{T} , but crucially we reimplement a *new* conditioning operation using the explicitly given weights. The mass function meaning of a representation then accumulates the mass of all weights associated to a given value. We transform an inference transformation by picking the component of the appropriate type.

It is straightforward to show that $\text{WT } \underline{T}$ is an inference representation, using preservation of return and $\gg=$ by the meaning function to reduce the proof to manipulations of weighted sums over \mathbb{R}_+ . Establishing the validity of lift and tmap is straightforward.

The weighting transformer *augments* the representation with a new conditioning operation, but *transforms* its choice operation to the new representation. We will later see more examples of both kinds.

3.5 Summary

We have introduced our three core abstractions, inference representations, transformations, and transformers, in relation to a mathematical semantic structure, the mass function monad. The examples so far show that the higher-order structure in our core calculus acts as a useful glue for manipulating and defining these abstractions. In the continuous case, we will also use this higher-order structure to represent computations over the real numbers.

4 PRELIMINARIES

In order to generalise this higher-order treatment to the continuous case, we first need to review and develop the mathematical theory of quasi-Borel spaces. Our development uses Kock's synthetic measure theory [2012], which allows us to reason analogously to measure theory. In order to present the synthetic theory, we briefly review the required category theoretic concepts. These sections are aimed at readers who are interested in the categorical context of our development. Other readers may continue directly to § 4.3.

4.1 Category Theory

Basic Notions. We assume basic familiarity with categories \mathcal{C} , \mathcal{D} , functors $F, G : \mathcal{C} \rightarrow \mathcal{D}$, and natural transformations $\alpha, \beta : F \rightarrow G$, and their theory of limits, colimits, and adjunctions. To fix notation, a *cartesian closed category* is a category with finite products, denoted by $\mathbb{1}$, \times , $\prod_{i=1}^n$, and exponentials, denoted by X^Y . In this subsection, we use the fragment of our core calculus consisting of the simply-typed λ -calculus (with sums, if necessary) to more compactly review the relevant concepts.

Monads. A *strong monad* \underline{T} over a cartesian closed category is a triple $(T, \text{return}, \gg=)$ consisting of an assignment of an object TX and a morphism $\text{return}_X : X \rightarrow TX$ for every object X , and an assignment of a morphism $\gg=_{X,Y} : TX \times (TY)^X \rightarrow TY$, satisfying the monad laws from § 2.5. Given a monad \underline{T} , a \underline{T} -*algebra* A is a pair $(|A|, \gg=^A)$ consisting of an object $|A|$, called the *carrier*, and an assignment of a morphism $\gg=^A_X : |A|^X \rightarrow |A|^{TX}$ to every object X satisfying

$$(\text{return } x \gg=^A f) = f \ x \quad \text{and} \quad ((a \gg= f) \gg=^A g) = a \gg= (\lambda x. f(x) \gg=^A g).$$

The pair $(TX, \gg=)$ always forms a T -algebra called the *free* \underline{T} -algebra over X . The *Eilenberg-Moore* category $\mathcal{C}^{\underline{T}}$ for a monad \underline{T} consists of \underline{T} -algebras and their homomorphism. The *Kleisli* category $\mathcal{C}_{\underline{T}}$ consists of the same objects as \mathcal{C} , but morphisms from X to Y in $\mathcal{C}_{\underline{T}}$ are morphisms $X \rightarrow TY$ in \mathcal{C} . The Kleisli category $\mathcal{C}_{\underline{T}}$ inherits any coproducts \mathcal{C} has. A strong monad \underline{T} is *commutative* when, for every

$$a : \underline{TX}, b : \underline{TY} \vdash \underline{T}.do \{x \leftarrow a; y \leftarrow b; \text{return}(x, y)\} = \underline{T}.do \{y \leftarrow b; x \leftarrow a; \text{return}(x, y)\}$$

(The notion of strong/commutative monad is due to [Kock 1972]; our formulation of algebras also appears in [Marmolejo and Wood 2010].)

Biproducts. A zero object \mathbb{Z} is both initial and terminal. A category has (*finite, countable, etc.*) *biproducts* if it has a zero object (and hence zero morphisms $\mathbf{0}_{X,Y} : X \rightarrow \mathbb{Z} \rightarrow Y$) and the following canonical morphisms are invertible:

$$\left[(\delta_{i,j})_{j \in I} \right]_{i \in I} : \sum_{i \in I} X_i \rightarrow \prod_{j \in I} X_j \quad \text{where: } \delta_{i,i} := \text{id}_{X_i}, \delta_{i,j} := \mathbf{0}_{X_i, X_j} \text{ for } i \neq j.$$

Algebraic Structure. Recall the notion of a *commutative monoid* $(M, 1, \cdot)$ in a category with finite products. We extend it to countably many arguments. Let \mathcal{C} be a category with countable products. A σ -*monoid* (see also [Haghverdi and Scott 2006]) is a triple $(M, 0, \Sigma)$ consisting of: an object M ; a morphism $0 : \mathbb{1} \rightarrow M$; and a morphism $\Sigma : M^{\mathbb{N}} \rightarrow M$ such that:

- setting $\delta_0 := \text{id}_M : M \rightarrow M$ and $\delta_i := 0 \circ ! : M \rightarrow \mathbb{1} \rightarrow M$, $i > 0$, we have $\Sigma \circ (\delta_i)_{i \in \mathbb{N}} = \delta_0$; and
- for every bijection $\varphi : \mathbb{N} \cong \mathbb{N} \times \mathbb{N}$, $a_{(-,-)} : M^{\mathbb{N} \times \mathbb{N}} \vdash \Sigma \left(\left(a_{(i,j)} \right)_{j \in \mathbb{N}} \right)_{i \in \mathbb{N}} = \Sigma \left(a_{\varphi(k)} \right)_{k \in \mathbb{N}}$.

PROPOSITION 4.1. *In a category with countable biproducts, each object M is a σ -monoid via:*

$$\mathbf{0}_{\mathbb{1}, M} : \mathbb{1} \rightarrow M \quad \Sigma : \prod_{i \in \mathbb{N}} M \cong \sum_{i \in \mathbb{N}} M \xrightarrow{\nabla} M \text{ where } \nabla \text{ is the codiagonal.}$$

Every morphism is a σ -monoid homomorphism with respect to this structure.

A σ -*semiring* is a quintuple $(S, 1, \cdot, 0, \Sigma)$ consisting of: a commutative monoid $(S, 1, \cdot)$; and a σ -monoid $(S, 0, \Sigma)$, such that $a : S, b_- \in S^{\mathbb{N}} \vdash a \cdot \Sigma (b_i)_{i \in \mathbb{N}} = \Sigma (a \cdot b_i)_{i \in \mathbb{N}}$. Given a σ -semiring $(S, 1, \cdot, 0, \Sigma)$, an S -*module* is a pair (M, \odot) consisting of a σ -monoid M ; and a morphism $\odot : S \times M \rightarrow M$ satisfying: $x = x, 0_S \odot x = 0_M$, $(a \cdot b) \odot x = a \odot (b \odot x)$, $(\Sigma (a_n)_{n \in \mathbb{N}}) \odot x = \Sigma (a_n \odot x)_{n \in \mathbb{N}}$.

4.2 Synthetic Measure Theory

Synthetic mathematics identifies structure and axioms from which we can recover the main concepts and results of specific mathematical theories, and transport them to new settings. We now briefly recount the relevant parts of Kock’s [2012] development. (In the finite discrete case, this is also related to Jacobs’s [2017] work on effectuses.)

4.2.1 Axioms and Structure. Let C be a cartesian closed category with countable products and coproducts, and let \underline{M} be a commutative monad over C . If the morphism $! : M0 \rightarrow \mathbb{1}$ is invertible, then both the Eilenberg-Moore category $C^{\underline{M}}$ and the Kleisli category $C_{\underline{M}}$ have zero objects. As a consequence, we have a canonical \underline{M} -homomorphism $\gg= \left[(\delta_{i,j})_j \right]_i : M \sum_{i \in \mathbb{N}} X_i \rightarrow \prod_{j \in \mathbb{N}} M X_j$.

Definition 4.2. A *measure category* is a pair (C, \underline{M}) consisting of a cartesian closed category C with countable products and coproducts, and equalisers; and a commutative monad \underline{M} over C such that the morphisms $! : M0 \rightarrow \mathbb{1}$ and $\gg= \left[(\delta_{i,j})_j \right]_i : M \sum_{i \in \mathbb{N}} X_i \rightarrow \prod_{j \in \mathbb{N}} M X_j$ are invertible.

We fix a measure category (C, \underline{M}) for the remainder of this section. The intuition is that $M X$ is the object of distributions/measures over X . Kock shows that, while short, the above definition has surprisingly many consequences.

Both the Eilenberg-Moore and the Kleisli categories have countable biproducts, and as a consequence, all \underline{M} -algebras have a σ -monoid structure and all \underline{M} -homomorphisms are σ -monoid homomorphisms with respect to it. Moreover, this structure on the free algebra on the terminal object $R := M \mathbb{1}$ extends to a σ -semiring structure by setting: $1 := \mathbf{return}()$ and $r \cdot s := \underline{M}.\mathbf{do} \{r; s\}$. Kock calls this structure the σ -semiring of scalars. Each \underline{M} -algebra A has an R -module structure:

$$r : R, a : |A| \vdash r \odot a := \underline{M}.\mathbf{do} \{r; a\}$$

As C has equalisers, for each object X , we may form the equaliser $P X \xrightarrow{\text{sub}_X} M X \xrightarrow[\mathbb{1}]{M!} R$ because $R = M \mathbb{1}$. Each sub_X is monic, the monadic structure factors through sub turning P into a commutative monad \underline{P} , and $\text{sub} : \underline{P} \rightarrow \underline{M}$ into a strong monad monomorphism.

The morphism $M! : M X \rightarrow R$ is called the *total measure* morphism, and P is then the sub-object of all the measures with total measure 1, and so we think of it as the object of *probability measures* over X . For example, every \underline{P} -algebra is closed under *convex* linear combinations of scalars: if $r_- : \mathbb{N} \rightarrow R$ satisfies $\sum (r_i)_i = 1$ then $\underline{\mu}_- : (P X)^{\mathbb{N}} \vdash M!(\sum (r_i \odot \underline{\mu}_i)_i) = 1$.

4.2.2 Notation and Basic Properties. Kock’s theory shines brightly when we adopt a measure-theoretic notation, as in Fig. 7, by thinking of $M X$ as the object of measures over X , and R as the object of scalars these measures take values in. The functorial action of the monad allows us to push measures along morphisms, and pushing all the measure into the terminal object gives a scalar we think of as the total measure of an object. The monadic return acts as a dirac distribution. The main advantage is the *Kock integral*, synonymous to the monadic $\gg=$. The main difference between the Kock integral \int and the usual Lebesgue integral \int from measure theory is that the Kock integral evaluates to a *measure*, and not a scalar. Calculating with the Kock integral is analogous to using Lebesgue integrals with respect to a generic test function, and proceeding by algebraic manipulation. The scalar rescaling \odot allows us to rescale a distribution by an arbitrary weight function. A *kernel* is a morphism $k : X \rightarrow M Y$, and we use the usual notation for integration against a kernel and iterated integration. We define the product measure by iterated integration. Finally, the $\gg=$ operation of an \underline{M} -algebra A gives rise to an expectation operation. Here we will only make use of the scalars’ algebra structure, which generalises the usual Lebesgue integral.

	Notation	Meaning	Terminology
	R	$:= M\mathbb{1}$	Scalars
$f : Y^X, \underline{\mu} : MX$	$\vdash f_*\underline{\mu}$	$:= (Mf)(\underline{\mu})$	Push-forward
$\underline{\mu} : MX$	$\vdash \underline{\mu}(\overline{X})$	$:= !_*\underline{\mu}$	The total measure
$x : X$	$\vdash \underline{\delta}_x$	$:= \mathbf{return}(x)$	Dirac distribution
$\underline{\mu} : MX, f : (MY)^X$	$\vdash \int_X f(x)\underline{\mu}(dx)$	$:= \underline{\mu} \gg= f$	Kock integral
$w : R^X, \underline{\mu} : MX$	$\vdash w \odot \underline{\mu}$	$:= \int_X (w(x) \odot \underline{\delta}_x)\underline{\mu}(dx)$	Rescaling
$\left[\begin{array}{l} f : (TZ)^{X \times Y}, \\ x : X, k : (TY)^X \end{array} \right]$	$\vdash \int_Y f(x, y)k(x, dy)$	$:= \int_Y f(x, y)k(x)(dy)$	Kernel integration
$\left[\begin{array}{l} f : (MX)^{X \times Y}, \\ \underline{\mu} \in M(X \times Y) \end{array} \right]$	$\vdash \iint_{X \times Y} f(x, y)\underline{\mu}(dx, dy)$	$:= \int_{X \times Y} f(z)\underline{\mu}(dz)$	Iterated integrals
$\underline{\mu} : MX, \underline{\nu} : MY$	$\vdash \underline{\mu} \otimes \underline{\nu}$	$:= \int_X \left(\int_Y \underline{\delta}_{(x, y)}\underline{\nu}(dy) \right) \underline{\mu}(dx)$	Product measure
$\underline{\mu} : MX, f : A ^X$	$\vdash \mathbb{E}_{x \sim \underline{\mu}}^A[f(x)]$	$:= \underline{\mu} \gg= f$	Expectation
$f : R^X, \underline{\mu} : MX$	$\vdash \int_X f(x)\underline{\mu}(dx)$	$:= \mathbb{E}_{x \sim \underline{\mu}}^R[f(x)]$	Lebesgue integral

Fig. 7. Synthetic measure theory notation

The justification for this notation is that it obeys the expected properties, which we now survey. The commutativity of the monad lets us change the order of integration:

THEOREM 4.3 (FUBINI-TONELLI). *For every pair of objects X, Y in a measure category (C, \underline{M}) :*

$$\iint_{X \times Y} f(x, y)(\underline{\mu} \otimes \underline{\nu})(dx, dy) = \int_Y \underline{\nu}(dy) \int_X \underline{\mu}(dx) f(x, y) = \int_Y \underline{\nu}(dy) \int_X \underline{\mu}(dx) f(x, y)$$

Moreover, for every M -algebra A :

$$\underline{\mu} : MX, \underline{\nu} : MY, f : |A|^{X \times Y} \vdash \mathbb{E}_{\substack{x \sim \underline{\mu} \\ y \sim \underline{\nu}}}^A[f(x, y)] = \mathbb{E}_{x \sim \underline{\mu}}^A[\mathbb{E}_{y \sim \underline{\nu}}^A[f(x, y)]] = \mathbb{E}_{y \sim \underline{\nu}}^A[\mathbb{E}_{x \sim \underline{\mu}}^A[f(x, y)]]$$

As usual, we allow placing the binder $\underline{\mu}(dx)$ on either side of the integrand $f(x)$.

Integrals and expectation interact well with the R -module structure in the sense that they are homomorphisms in both arguments. The precise statement of this fact can be found in Appendix A.

The push-forward operation interacts with rescaling in the following way:

THEOREM 4.4 (FROBENIUS RECIPROCITY). *For all objects X, Y in a measure category (C, \underline{M}) :*

$$w : R^X, \underline{\mu} : MX, f : Y^X \vdash w \odot (f_*\underline{\mu}) = f_*((w \odot f) \odot \underline{\mu})$$

When calculating in this notation, we use the equations in Appendix A where we present a toolbox for synthetic measure theory. This toolbox includes most of the equations we come to expect from standard measure theory, like the change of variables law. To validate them, inline the definitions and proceed using the usual category-theoretic properties.

The following two sections contain relevant extensions to Kock's theory.

4.2.3 Radon-Nikodym Derivatives. The Radon-Nikodym Theorem is a powerful tool in measure theory, and we now phrase a synthetic counterpart. As usual in the synthetic setting, we set the definitions up such that the theorem will be true. Doing so highlights the difference between three measure-theoretic concepts that coincide in measure theory, but may differ in the synthetic setting.

Let $\underline{\mu}, \underline{\nu} \in MX$ be measures. We say that $\underline{\nu}$ is *absolutely continuous* with respect to $\underline{\mu}$, and write $\underline{\nu} \ll \underline{\mu}$, when there exists a morphism $w : X \rightarrow R$ such that $\underline{\nu} = w \odot \underline{\mu}$. Given two morphisms $w, v : X \rightarrow R$ and a measure $\underline{\mu} \in MX$, we say that w and v are *equal $\underline{\mu}$ -almost everywhere* ($\underline{\mu}$ -a.e.) when $w \odot \underline{\mu} = v \odot \underline{\mu}$. A *measurable property* over X is a morphism $P : X \rightarrow \text{bool}$. Given a measure $\underline{\mu} \in MX$ a measurable property P over X *holds $\underline{\mu}$ -a.e.*, when the morphism $[P] := \lambda x. \text{ if } P \ x \ \text{then } 1 \ \text{else } 0$ is equal $\underline{\mu}$ -a.e. to 1.

THEOREM 4.5 (RADON-NIKODYM). *Let (C, M) be a well-pointed measure category. For every $\underline{\nu} \ll \underline{\mu}$ in MX , there exists a $\underline{\mu}$ -a.e. unique morphism $\frac{d\underline{\nu}}{d\underline{\mu}} : X \rightarrow R$ satisfying $\frac{d\underline{\nu}}{d\underline{\mu}} \odot \underline{\mu} = \underline{\nu}$.*

4.2.4 Kernels. We say that a kernel $k : X \rightarrow MY$ is *Markov* when, for all x , $k(x, Y) = 1$, i.e., when k factors through the object of probability measures via $\text{sub} : P \rightarrow M$. We now restrict attention to kernels $k : X \rightarrow MX$ over the same object X . We say that such a kernel *preserves* a measure $\underline{\mu}$ when $\underline{\mu} \gg k = \underline{\mu}$. Recall the morphism swap $:= \lambda(x, y). (y, x) : X \times Y \rightarrow Y \times X$. Given a measure $\underline{\mu} \in MX$ and a kernel k , we define the *box product* by $\underline{\mu} \boxtimes k := \iint_{X \times X} \delta_{(x, y)} \underline{\mu}(dx) k(x, dy)$. A kernel k is *reversible* with respect to a measure $\underline{\mu} \in MX$ when $\text{swap}_*(\underline{\mu} \boxtimes k) = \underline{\mu} \boxtimes k$.

The following standard results on kernels transfer into the synthetic setting. If a *Markov* kernel k is reversible with respect to $\underline{\mu}$, then k preserves $\underline{\mu}$. Kernels obtained by rescaling the Dirac kernel, i.e., $\lambda x. w(x) \odot \delta_x$ are reversible w.r.t. all measures. Finally, linear combinations $\lambda x. \sum_{n \in \mathbb{N}} r_n \odot k_n(x)$ of reversible kernels w.r.t. $\underline{\mu}$ are also reversible w.r.t. $\underline{\mu}$.

4.3 Quasi-Borel Spaces

It remains to show that there is a concrete model of synthetic measure theory that contains the classical measure theoretic ideas that are central to probability theory and inference. This is novel because [Kock's work \[2012\]](#) is targeted at the geometric/topological setting, whereas probability theory is based around Borel sets rather than open sets. It is non-trivial because the traditional setting for measure theory does not support higher-order functions [[Aumann 1961](#)] and commutativity of integration is subtle in general. In this section we resolve these problems by combining some recent discoveries [[Heunen et al. 2017](#); [Staton 2017](#)], and exhibit a model of synthetic measure theory which contains classical measure theory, for instance:

- the σ -semiring over the morphisms $\mathbb{1} \rightarrow R$ is isomorphic to the usual σ -semiring over the extended non-negative reals, $\overline{\mathbb{R}}_+$;
- this isomorphism induces a bijective correspondence between the morphisms $R \rightarrow \mathbb{1} + \mathbb{1}$ and the Borel subsets of $\overline{\mathbb{R}}_+$, as characteristic functions, and also between the morphisms $R \rightarrow R$ and the measurable functions $\overline{\mathbb{R}}_+ \rightarrow \overline{\mathbb{R}}_+$;
- it also induces an injection of the morphisms $\mathbb{1} \rightarrow M(R)$ into the set of Borel measures on $\overline{\mathbb{R}}_+$, whose image contains all the probability measures; the morphisms $R \rightarrow M(R)$ include all the Borel probability kernels;
- the canonical morphism $R^R \times M(R) \rightarrow R$, $(f, \underline{\mu}) \mapsto \int f(x) \underline{\mu}(dx)$, corresponds to classical Lebesgue integration.

Moreover, each object X can be seen as a set $U(X) = C(\mathbb{1}, X)$ with structure, because the category is well-pointed, in the sense that the morphisms $X \rightarrow Y$ are a subset of the functions $U(X) \rightarrow U(Y)$.

4.3.1 Rudiments of Classical Measure Theory. Measurable spaces are the cornerstone of conventional measure theory, supporting a notion of measure.

Recall that a σ -algebra on a set X is a set Σ_X of subsets of X that is closed under countable unions and complements. A *measurable space* is a set together with a σ -algebra. A *measure* is a σ -additive

function $\Sigma_X \rightarrow \overline{\mathbb{R}}_+$. A function f between measurable spaces is *measurable* if the inverse image of every measurable set according to f is measurable.

For example, on a Euclidean space $\overline{\mathbb{R}}^n$ we can consider the Borel sets, which form the smallest σ -algebra containing the open cubes. There is a canonical measure on $\overline{\mathbb{R}}^n$, the *Lebesgue measure*, which assigns to each cube its volume, and thus to every measurable function $f: \overline{\mathbb{R}}^n \rightarrow \overline{\mathbb{R}}_+$ a Lebesgue integral $\int_{\overline{\mathbb{R}}^n} f \in \overline{\mathbb{R}}_+$. A slightly more general class of measures is the *σ -finite measures*, which include the Lebesgue measures and are closed under disjoint unions and product measures.

A measurable space that arises from the Borel sets of a Polish space is called a *standard Borel space*. In fact, every standard Borel space is either countable or isomorphic to $\overline{\mathbb{R}}$. Standard Borel spaces are closed under countable products and countable disjoint unions.

4.3.2 Quasi-Borel Spaces. In this section we fix an uncountable standard Borel space, \mathfrak{R} . For example, $\mathfrak{R} = \overline{\mathbb{R}}$. The basic idea of quasi-Borel spaces is that rather than focusing on measurable sets of a set X , as in classical measure theory, one should focus on the admissible random elements $\mathfrak{R} \rightarrow X$.

Definition 4.6 ([Heunen et al. 2017]). A *quasi-Borel space (QBS)* is a set X together with a set of functions $M_X \subseteq [\mathfrak{R}, X]$ such that (i) all the constant functions are in M_X , (ii) M_X is closed under precomposition with measurable functions on \mathfrak{R} , and (iii) M_X satisfies the piecewise condition: if $\mathfrak{R} = \biguplus_{i=1}^{\infty} U_i$, where U_i is Borel measurable and $\alpha_i \in M_X$ for all i , then $\biguplus_{i=1}^{\infty} \alpha_i \cap (U_i \times X)$ is in M_X .

A *morphism* $f: X \rightarrow Y$ is a function that respects the structure, i.e. if $\alpha \in M_X$ then $(f \circ \alpha) \in M_Y$. Morphisms compose as functions, and we have a category **QBS**.

A QBS X is a *subspace* of a QBS Y if $X \subseteq Y$ and $M_X = \{\alpha: \mathfrak{R} \rightarrow X \mid \alpha \in M_Y\}$.

A measurable space X can be turned into a QBS when given the set of measurable functions $\mathfrak{R} \rightarrow X$ as M_X . When X and Y are standard Borel spaces considered as QBSes this way, **QBS**(X, Y) comprises the measurable functions, so **QBS** can be thought of as a conservative extension of the universe of standard Borel spaces. The three conditions on quasi-Borel spaces ensure that coproducts and products of standard Borel spaces retain their universal properties in **QBS**. In fact, the category of QBSs has all limits and colimits. It is also cartesian closed; e.g., $\mathbb{R}^{\mathbb{R}} := \mathbf{QBS}(\mathbb{R}, \mathbb{R})$, and $M_{(\mathbb{R}^{\mathbb{R}})} = \{\alpha: \mathfrak{R} \rightarrow (\mathbb{R}^{\mathbb{R}}) \mid \text{uncurry}(\alpha) \in \mathbf{QBS}(\mathfrak{R} \times \mathbb{R} \rightarrow \mathbb{R})\}$. For any QBS X , $M_X = \mathbf{QBS}(\mathfrak{R}, X)$.

4.3.3 A Monad of Measures. The following development is novel.

Definition 4.7. A *measure* μ on a quasi-Borel space is a triple (Ω, α, μ) where Ω is a standard Borel space, $\alpha \in \mathbf{QBS}(\Omega, X)$, and μ is a σ -finite measure on Ω .

For example, Ω might be $\overline{\mathbb{R}}^n$ and μ might be the Lebesgue measure. A measure determines an integration operator: if $f \in \mathbf{QBS}(X, \overline{\mathbb{R}}_+)$ then define

$$\int f \, d(\Omega, \alpha, \mu) := \int_{\Omega} f(\alpha(x)) \mu(dx)$$

using Lebesgue integration according to μ . We say that two measures are equal, denoted $(\Omega, \alpha, \mu) \approx (\Omega', \alpha', \mu')$, if they determine the same integration operator. We write $[\Omega, \alpha, \mu]$ for an equivalence class of measures.

As an aside, we note that not every integration operator on $\overline{\mathbb{R}}$ in the classical sense is a measure in the sense of Def. 4.7, because we restrict to σ -finite μ . Technically, the only integration operators that arise in this way are those corresponding to σ -finite measures. This is a class of measures that includes the probability measures, and which works well with iterated integration and probabilistic programming [Staton 2017].

The measures up-to \approx form a monad, as follows. First, the set of all measures M_X forms a QBS by setting $M_{M_X} = \{\lambda r. [D_r, \alpha(r, -), \mu|_{D_r}] \mid \mu \text{ } \sigma\text{-finite on } \Omega, D \subseteq \mathfrak{R} \times \Omega \text{ measurable, } \alpha \in \mathbf{QBS}(D, X)\}$,

where $D_r = \{\omega \mid (r, \omega) \in D\}$. In consequence, when Ω' is a standard Borel space, for every morphism $f : \Omega' \rightarrow MX$, there exist $\Omega, \mu, D \subseteq \Omega' \times \Omega$ and $\alpha \in \mathbf{QBS}(D, X)$ such that $f(\omega') = [D_{\omega'}, \alpha(\omega', -), \mu|_{D_{\omega'}}]$. One intuition is that α is a partial function $\Omega' \times \Omega \rightarrow X$, with domain D .

The unit of the monad, $\text{return} : X \rightarrow MX$, is $\text{return}(x) := [\mathbb{1}, \lambda_{\cdot}. x, \delta_0]$ where δ_0 is the Dirac measure on the one-point space $\mathbb{1}$. We often write $\underline{\delta}_x$ for $\text{return}(x)$. The bind $\gg= : MX \times MY^X \rightarrow MY$ is

$$[\Omega, \alpha, \mu] \gg= \varphi := [D, \beta, (\mu \otimes \mu')|_D]$$

where $\varphi(\alpha(r)) = [D_r, \beta(r, -), \mu']$. Note that $(\varphi \circ \alpha) : \Omega \rightarrow MX$ must be of this form because it is a morphism from a standard Borel space. The measure $\mu \otimes \mu'$ is the product measure, which exists because μ and μ' are σ -finite.

This structure satisfies the monad laws, it is commutative by the Fubini-Tonelli theorem, and it satisfies the biproduct axioms, and so it is a model of synthetic measure theory. Every measure on $\mathbb{1}$ is equivalent to one of the form $([0, r], !, \mu)$ where $r \in \overline{\mathbb{R}}_+$, $! : [0, r] \rightarrow \mathbb{1}$ is the unique such random element, and μ is the Lebesgue measure. Thus $M\mathbb{1} \cong \mathbb{R}_+$.

As another aside, we note that when Ω, Ω' are standard Borel spaces, the Kleisli morphisms $\Omega \rightarrow M\Omega'$ correspond to s -finite kernels, which were shown in [Staton 2017] to provide a fully complete model of first-order probabilistic programming.

5 CONTINUOUS INFERENCE

We now develop the continuous counterpart to Sec. 3. The semantic structure of the category of quasi-Borel spaces allows us to transport many of the definitions with little change. For example, a monadic interface \underline{T} consists of analogous data, but the assignments are indexed by quasi-Borel spaces, T assigns quasi-Borel spaces, and return^T and $\gg=^T$ assign quasi-Borel space morphisms.

Definition 5.1. A *continuous representation* \underline{T} is a tuple $(T, \text{return}^T, \gg=^T, m^T)$ consisting of:

- a monadic interface $(T, \text{return}^T, \gg=^T)$;
- an assignment of a *meaning* morphism $m^T_X : TX \rightarrow MX$ for every space X

such that m^T preserves return^T and $\gg=^T$.

A *sampling* representation is a tuple $(T, \text{return}^T, \gg=^T, m^T, \text{sample}^T)$ such that its first four components form a continuous representation, it has an additional \mathbf{Qbs} -morphism $\text{sample}^T : \mathbb{1} \rightarrow T\mathbb{1}$, and m^T maps $\text{sample}^T()$ to the uniform \mathbf{Qbs} -measure $\mathbf{U} = [\mathbb{1}, \text{id}, \text{Uniform}]$ on the unit interval $\mathbb{1}$, where Uniform is the usual uniform distribution on $\mathbb{1}$.

A *conditioning* representation \underline{T} is similarly a tuple $(T, \text{return}^T, \gg=^T, \text{score}^T, m^T)$, with a \mathbf{Qbs} -morphism $\text{score}^T : \mathbb{R}_+ \rightarrow T\mathbb{1}$ such that for each r , m^T maps $\text{score}^T(r)$ to the r -rescaled unit \mathbf{Qbs} -measure $r \odot \underline{\delta}_0 = [\mathbb{1}, \lambda_{\cdot}. (), r \cdot \delta_0]$.

An *inference* representation \underline{T} is a tuple $(T, \text{return}^T, \gg=^T, \text{sample}^T, \text{score}^T, m^T)$ with the appropriate components forming both a sampling representation and a conditioning representation.

This definition refines Def. 3.2 with sampling and scoring representations, allowing us to talk about inference transformers that augment a representation of one kind into another.

Example 5.2 (Continuous sampler). By analogy with Ex. 3.3, we define in Fig. 8a a sampling representation using the type $\text{Sam } \alpha := \{\text{Return } \alpha \mid \text{Sample } (\mathbb{1} \rightarrow \text{Sam } \alpha)\}$. Validating the preservation of sample and the monadic interface is straightforward. It also follows from more general principles: $\underline{\text{Sam}}$ is the initial monad with an operation $\text{sample} : T\mathbb{1}$.

We define inference transformations between any two representations as in Def. 3.5. We have four kinds of representations, and when defining transformers we can augment a representation with additional capabilities:

instance *Sampling Monad* (Sam) **where**

```

return  $x = \text{Return } x$ 
 $a \gg= f = \text{match } a \text{ with } \{$ 
     $\text{Return } x \rightarrow f(x)$ 
     $\text{Sample } k \rightarrow$ 
     $\text{Sample } (\lambda r. k(r) \gg= f)\}$ 
 $\text{sample} = \text{Sample } \lambda r. (\text{Return } r)$ 
 $m a = \text{match } a \text{ with } \{$ 
     $\text{Return } x \rightarrow \underline{\delta}_x$ 
     $\text{Sample } k \rightarrow \int_{\mathbb{I}} k(x) \mathbf{U}(dx)\}$ 

```

(a) Continuous sampler representation

instance *Cond Trans* (W) **where**

```

return  $w \underline{T} x = \text{return}^T(1, x)$ 
 $a \gg=_{w \underline{T}} f = \underline{T}.do \{(r, x) \leftarrow a;$ 
     $(s, y) \leftarrow f(x);$ 
     $\text{return}(r \cdot s, y)\}$ 
 $(\text{tmap } \underline{t})_X = \underline{t}_{\mathbb{R}_+ * X}$ 
 $\text{lift}_{\underline{T}} a = \underline{T}.do \{x \leftarrow a; \text{return}(1, x)\}$ 
 $m_{w \underline{T}} a = \lambda x. \int_{\mathbb{R}_+ * X} r \odot \underline{\delta}_x m^T(a)(dr, dx)$ 
 $\text{score}_{w \underline{T}} r = \text{return}^T(r, ())$ 

```

(b) Continuous weighting inference transformer

Fig. 8. Continuous representations and transformers

Definition 5.3. Let k_1, k_2 be a pair of kinds of representation. A k_1 to k_2 transformer \underline{F} is a tuple $(F, \text{tmap}^{\underline{F}}, \text{lift}^{\underline{F}})$ consisting of an assignments of:

- a k_2 representation $F \underline{T}$ to every k_1 representation \underline{T} ;
- an inference transformation $\text{tmap}^{\underline{F}} \underline{t} : \underline{F} \underline{T} \rightarrow \underline{F} \underline{S}$ to every transformation $\underline{t} : \underline{T} \rightarrow \underline{S}$; and
- an inference transformation $\text{lift}^{\underline{F}} : \underline{T} \rightarrow \underline{F} \underline{T}$ to every k_1 representation \underline{T} .

When the two kinds k_1, k_2 differ, we say that that the transformer is *augmenting*.

When defining a k_1 to k_2 transformer, we adopt a Haskell-like type-class constraint notation $k_1 \implies k_2$ used for example in Fig. 11a.

Example 5.4. By analogy with Ex. 3.9, Fig. 8b presents the *continuous weighting* transformer structure on $W \underline{T} X := T(\mathbb{R}_+ * X)$. It augments any representation transformer with conditioning capabilities. Each conditioning operation is deferred to the return value, and so we can view this transformer as freely adding a conditioning operation that commutes with all other operations. When the starting representation had conditioning capabilities, we have an inference transformation $wagr : W \underline{T} \rightarrow \underline{T}$, given by $wagr a := \underline{T}.do \{(r, x) \leftarrow a; \text{score}^T r; \text{return } x\}$ which conditions based on the aggregated weight.

Its validity follows from a straightforward calculation using the meaning preservation of \underline{T} .

In the continuous case, the output of the final inference transformation will always be $W \text{Sam } X$ or a similar $\text{Pop Sam } X$ described in the next section. From this representation, we obtain the Monte Carlo approximation to the posterior by using a random number generator to supply the values required by Sam. Interpreting the program directly in $W \text{Sam } X$ and sampling from that would correspond to simple importance sampling from the prior, which usually needs a very large number of samples to give a good approximation to the posterior. Our goal in approximate Bayesian inference is therefore to find another representation for the program and a sequence of inference transformations that map it to $W \text{Sam } X$. While, in principle, this output represents the same posterior distribution, hopefully it uses a representation that requires fewer samples to obtain a good approximation than a direct interpretation in $W \text{Sam } X$. We emphasise that approximation is only done in this final sampling step, while all the inference transformations that happen before it are exact.

6 SEQUENTIAL MONTE CARLO

Sequential Monte Carlo (SMC) is a family of inference algorithms for approximating sequences of distributions. In SMC, every distribution is approximated by a collection of weighted samples called a *population*, with each population obtained from the previous one in the sequence by a suitable transformation. In a sub-class of SMC known as *particle filters*, each distribution in the sequence is generated by a known random process from the previous distribution and we can apply this process to samples in the previous population to obtain the current population. In particle filters the samples in the population are called *particles*.

A common problem with particle filters is that, after multiple steps, a few particles have much larger weights than the remaining ones, effectively reducing the sample size in the population well below the actual number of particles, a phenomenon known as *sample impoverishment*. To ameliorate this problem, particle filters introduce additional *resample* operations after each step in the sequence, which constructs a new population by sampling with replacement from the old one. The new population has uniform weights across its particles. In the setting of probabilistic programming, we use suspended computations as particles, and their associated weight is their currently accumulated likelihood.

We show how to decompose a particle filter into a stack of two transformers: a representation to conditioning transformer for representing a population of particles, and a conditioning to conditioning transformer that allows us to run a particle until its next conditioning operation. We define each step of the SMC algorithm as an inference transformation on this stack. We can then apply this stack of transformers to a sampling representation to obtain a correct by construction variation of SMC. The algorithm we obtain is known as the particle filter with multinomial resampling [Doucet and Johansen 2011] that uses the prior as the proposal distribution, but throughout this paper we refer to it simply as SMC.

6.1 The Population Transformer

Given a representation \underline{T} , we define a representation structure over $\text{Pop } \underline{T} X := T(\text{List}(\mathbb{R}_+ * X))$. We further deconstruct this representation transformer as the composition of two transformers: the continuous weighting transformer W from Ex. 5.4, and Haskell’s notorious ListT transformer.

The negative reputation associated to the transformer $\text{ListT } \underline{T} X := T(\text{List } X)$ stems from its failure to validate the monad laws when \underline{T} is not commutative.² However, it is a perfectly valid representation transformer, described in Fig. 9a, since we do not require that representations satisfy monad laws.

To prove the meaning function preserves return, simply calculate. For $\gg=$ preservation, show:

$$a_s : \text{List}(TX) \vdash m^{\text{ListT } \underline{T}}(\text{sequence } a_s) = \sum_{a \in a_s} m^{\underline{T}}(a)$$

and proceed via straightforward calculation using the linearity of the Kock integral and the commutative (σ -)monoid structure on measures.

By composing the two representation transformers, we obtain the representation to conditioning transformer Pop , given explicitly in Fig. 9b.

Fig. 9c presents a \mathbb{N}_+ -indexed family of inference transformations. Fix any $n \in \mathbb{N}$. The spark function generates a population of particles with the unit value, and the same weight $\frac{1}{n}$. Thus, $\text{spawn}(n, a)$ takes a distribution a over particle populations, sparks n equally weighted particles, and for each of them, samples a population based on a . A straightforward calculation confirms that the meaning of spark is 1, and so $\text{spawn}(n, -) : \text{Pop } \underline{T} \rightarrow \text{Pop } \underline{T}$ is an inference transformation. In

²For a list transformer “done right”, see Jaskelioff’s thesis [2009], and its generalisations [Fiore and Saville 2017; Piróg 2016].

Auxiliary functions:

sequence : List(TX) \rightarrow T(List X)

sequence := foldr (return [])

$$(\lambda(a, r). \underline{T}. \mathbf{do} \{x \leftarrow a; \\ x_s \leftarrow r; \\ \mathbf{return}(x :: x_s)\})$$

concat : List(List X) \rightarrow List X

concat := foldr [] ++

$$\sum_{x \in x_s} f(x) := \text{foldr } 0 \ (\lambda(x, s). f(x) + s) \ x_s$$

instance Rep Trans (ListT) where

return_{ListT \underline{T}} x = return ^{\underline{T}} [x]

a \gg _{ListT \underline{T}} f = $\underline{T}. \mathbf{do} \{x_s \leftarrow a; \\ \mathbf{let} \ b_s = \text{map } f \ x_s \ \mathbf{in} \\ y_{ss} \leftarrow \text{sequence } b_s; \\ \mathbf{return}(\text{concat } y_{ss})\}$

$m_{\text{ListT } \underline{T}} a$ = $\int_{\text{List } X} m^{\underline{T}}(a)(dx_s) \sum_{x \in x_s} \delta_x$

lift_{ListT \underline{T}} a = $\underline{T}. \mathbf{do} \{x \leftarrow a; \mathbf{return} [x]\}$

(**tmap** \underline{t})_X = $\underline{t}_{\text{List } X}$

(a) The list transformer

instance Cond Trans (Pop) where

return_{Pop \underline{T}} = **return**_{(W \circ ListT) \underline{T}}

\gg _{Pop \underline{T}} = \gg _{(W \circ ListT) \underline{T}}

lift_{Pop \underline{T}} = **lift**_{W(ListT) \underline{T}} \circ **lift**_{ListT \underline{T}}

tmap_{Pop \underline{T}} = **tmap**_{W(ListT) \underline{T}} \circ **tmap**_{ListT \underline{T}}

$m_{\text{Pop } \underline{T}}$ = $m_{(W \circ \text{ListT}) \underline{T}}$
 = $\lambda a. \int m^{\underline{T}}(a)(dx_s) \sum_{(r, x) \in x_s} r \odot \delta_x$
 List($\mathbb{R}_+ \times X$)

score_{Pop \underline{T}} = **score**_{(W \circ ListT) \underline{T}}

(b) The population transformer

replicate : $\mathbb{N} * X \rightarrow$ List X

replicate(n, x) := $\mathbb{N}. \mathbf{fold} \lambda \{ \text{Zero} \rightarrow [] \\ \text{Succ } x_s \rightarrow x :: x_s \} n$

spark : $\mathbb{N}_+ \rightarrow$ Pop $\underline{T} \ 1$

spark := return ^{\underline{T}} (replicate(n, ($\frac{1}{n}$, ())))

spawn : $\mathbb{N}_+ * \text{Pop } \underline{T} \ X \rightarrow$ Pop $\underline{T} \ X$

spawn(n, a) := Pop $\underline{T}. \mathbf{do} \{\text{spark } n; a\}$

(c) Spawning new particles

Fig. 9. Representing populations

the version of SMC we consider below, we will only pass to spawn a distribution a over uniformly-weighted single-particle populations.

We use spawn to *resample* a new population. Thinking operationally, we have a population of weighted particles and we obtain a new population by sampling with replacement from the current one, where the probability of selecting a given particle is proportional to its weight. Doing so is equivalent to simulating a discrete weighted sample using a uniform one.

LEMMA 6.1. *There is a Qbs-morphism $\text{dwrand} : \text{List}(\mathbb{R}_+ * X) * \mathbb{I} \rightarrow \{\text{Take } X \mid \text{Fail}\}$ such that:*

- For all x_s for which $\sum_{(r, _) \in x_s} r = 0$, we have $\text{dwrand}(x_s, -) * \mathbb{U} = \delta_{\text{Fail}}$.
- For all x_s for which $w := \sum_{(r, _) \in x_s} r > 0$, we have $\text{dwrand}(x_s, -) * \mathbb{U} = \sum_{(r, x) \in x_s} \frac{r_i}{w} \odot \delta_{\text{Take } x}$.

Fig. 10a presents one such morphism, though its precise implementation does not matter to our development. As a consequence, for every sampling representation \underline{T} for which we have an element fail : TX such that $m^{\underline{T}}(\text{fail}) = 0$, we can define a discrete weighted sampler $\text{dwsample}^{\underline{T}}(x_s) : \text{List}(\mathbb{R}_+ * X) \rightarrow TX$ in Fig. 10b which will then satisfy $m^{\underline{T}}(\text{dwsample}^{\underline{T}}(x_s)) = \sum_{(r, x) \in x_s} r \odot \delta_x$.

The resampling step in Fig. 10c operationally takes the current population, creates a computation/thunk that samples a single particle from this population, and then spawns n new particles that are initialised with this thunk. The morphism $\text{resample}(n, -) : \text{Pop } \underline{T} \rightarrow \text{Pop } \underline{T}$ is an inference transformation because, as we know, $\text{spawn}(n, -)$ is one and $\text{dwsample}^{\text{Pop } \underline{T}} : \text{Pop } \underline{T} \rightarrow \text{Pop } \underline{T}$ samples a population consisting of just a single unit weight particle with a probability proportional to its renormalised weight in the original population.

<pre> dwrand(x_s, r) := let $w = \sum_{(r, _) \in x_s} r$ in if $w = 0$ then Fail else foldr ($w \cdot r$, Fail) ($\lambda((s, x), (fuel, result)).$ if $0 \leq fuel < s$ then $(-1, \text{Take } x)$ else $-- \text{potential underflow}$ ($fuel - s, result$)) x_s </pre> <p>(a) A discrete weighted randomiser</p>	<pre> dwsample^T(x_s) := $\underline{T}.\text{do}$ {score($\sum_{(r, _) \in x_s} r$); $r \leftarrow \text{sample}$; match dwrand(x_s, r) with { Fail $\rightarrow \text{fail}$ Take $x \rightarrow \text{return } x$}} </pre> <p>(b) A discrete weighted sampler</p> <pre> resample : $\mathbb{N}_+ * \text{Pop } \underline{T} X \rightarrow \text{Pop } \underline{T} X$ resample(n, a) := $\underline{T}.\text{do}$ {$x_s \leftarrow a$; spawn($n, \text{dwsample}^{\text{Pop } \underline{T} X} x_s$)} </pre> <p>(c) Resampling</p>
--	---

Fig. 10. The resampling transformation

<pre> instance Cond \implies Cond Trans (Sus) where return_{Sus \underline{T}} x = return_{\underline{T}} (Return x) $a \gg_{\text{Sus } \underline{T}} f$ = fold ($\lambda b. \underline{T}.\text{do}$ { $t \leftarrow b$; match t with { Return $x \rightarrow f(x)$ Yield $c \rightarrow \text{Yield } c$}}) a lift_{Sus \underline{T}} a = $\underline{T}.\text{do}$ {$x \leftarrow a$; return_{Sus \underline{T}} x} (tmap_{Sus \underline{T}} t)_{X} = Sus $\underline{T} X.\text{fold}$ ($\lambda b. m_S(b)$) m_{Sus \underline{T}} a = m_{\underline{T}} (finish_{Sus \underline{T}} (a)) score r = return_{\underline{T}} (Yield lift_{Sus \underline{T}} (score r)) </pre> <p>(a) The suspension transformer</p>	<pre> advance_{\underline{T}} : Sus $\underline{T} X \rightarrow \text{Sus } \underline{T} X$ advance_{\underline{T}} a = $\underline{T}.\text{do}$ { $t \leftarrow a$; match t with { Return $x \rightarrow \text{return } x$ Yield $t \rightarrow t$}} finish_{\underline{T}} : Sus $\underline{T} X \rightarrow \underline{T} X$ finish_{\underline{T}} a = fold $\lambda b. \underline{T}.\text{do}$ { $t \leftarrow b$; match t with { Return $x \rightarrow \text{return } x$ Yield $b \rightarrow b$}} </pre> <p>(b) Suspension operations</p>
--	---

Fig. 11. The suspension transformation

6.2 The Suspension Transformer

The second transformer in the SMC algorithm allows us to suspend computation after each conditioning. The suspension transformer equips the standard *resumption* monad transformer $\text{Sus } \underline{T} X := T\{\text{Return } X \mid \text{Yield } (\text{Sus } \underline{T} X)\}$, presented in Fig. 11a, with inference transformations.

The two transformations on suspended computations in Fig. 11b take one step, and complete the computation, accordingly. As the meaning function for the transformed representation returns the meaning the computation would have if it was allowed to run to completion, these two operations do not change the meaning and so form inference transformations.

We can now put all the components together:

THEOREM 6.2. *Let \underline{T} be a sampling representation. For every pair of natural numbers n, k , the following composite forms an inference transformation:*

$$\text{smc}_{n,k}^{\underline{T}} := (\text{Sus} \circ \text{Pop})\underline{T} \xrightarrow{\text{tmap}_{\text{Sus}} \text{spawn}(n,-)} (\text{Sus} \circ \text{Pop})\underline{T} \xrightarrow{(\text{advance} \circ \text{tmap}_{\text{Sus}} \text{resample}(n,-))^{\circ k}} (\text{Sus} \circ \text{Pop})\underline{T} \xrightarrow{\text{finish}} \text{Pop} \underline{T}$$

In the above $(-)^{\circ -} : X^X \times \mathbb{N} \rightarrow X^X$ denotes n -fold composition. The transformation $\text{smc}_{n,k}^{\underline{T}}$ amounts to running the SMC algorithm with n particles for k steps. If the representation \underline{T} is operational in nature, such as the continuous sampler Sam, we get a sequence of weighted values over the return type when we run the resulting representation. By construction, the distribution on the results, rescaled according to their final weights, would be identical to the desired posterior distribution.

When the representation \underline{T} is not a commutative monad, like the continuous sampler Sam, the resulting representation $\text{Pop} \underline{T}$ is not a monad: the monad laws do not hold. Therefore, to encompass representations of $\text{Pop} \underline{T}$ one must generalise beyond monads.

7 TRACE MARKOV CHAIN MONTE CARLO

Markov Chain Monte Carlo (MCMC) algorithms operate by repeatedly using a transition kernel to generate a new sample from a current one. Thus they can be thought of as performing a random walk around the space they are exploring. If the transition kernel is well-behaved, they are guaranteed to preserve the distribution. A popular MCMC algorithm used for Bayesian inference is Metropolis-Hastings (MH), where the transition kernel consists of a proposal kernel followed by a decision to either accept the proposed sample or keep the old one. The accept or reject step is used to correct for bias introduced by the proposal kernel, thus producing a valid MCMC algorithm for a rich family of proposal kernels.

MH is a general inference method, but it requires specialised knowledge about the space on which they operate on. In the context of a probabilistic programming language, the *Trace MH* algorithm replaces the unknown target space with the space of program *traces*, which are shared by all probabilistic programs. Thus, Trace MH allows probabilistic programming language designers to devise general-purpose kernels to effectively explore traces.

We analyse the the Trace MH as follows. First, we prove a quasi-Borel space counterpart of the *Metropolis-Hastings-Green (MHG) Theorem*, that forms the theoretical foundation for the correctness of MH. We then present the *tracing* representation and show its validity. We present the Trace MH algorithm, parameterised by a proposal kernel for traces, and give sufficient conditions on this kernel for the resulting transformation to be valid. We then give a concrete proposal kernel and show that it satisfies these general conditions.

7.1 Abstract Metropolis-Hastings-Green

In the abstract, the key ingredient in MH is the *Metropolis-Hastings-Green (MHG)* morphism η presented in Fig. 12a, formulated in terms of an arbitrary inference representation \underline{T} . This transformation is usually known as the *update step* of the MH algorithm. It is parameterised by a (representation of a) *proposal* kernel $\psi : X \rightarrow T X$, and by a chosen (representation of a) Radon-Nikodym derivative $\rho : X \times X \rightarrow \mathbb{R}_+$.

To use η in an inference transformation, we need to provide well-behaved parameters ψ, ρ , and their behaviour may depend on the representation of the input distribution a . In particular, the parameter ρ should represent a well-behaved appropriate Radon-Nikodym derivative. To simplify our proofs, we also require that the proposal kernel ψ is Markov, which suffices for our application.

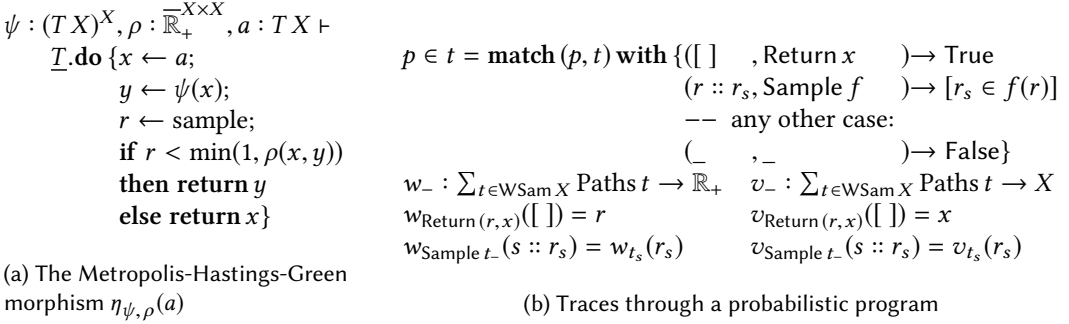


Fig. 12. Basic notions in Trace MH

THEOREM 7.1 (METROPOLIS-HASTINGS-GREEN). *Let X be a qbs, $a \in TX$ a distribution, $\psi : X \rightarrow TX$ a kernel, and $\rho : X \times X \rightarrow \overline{\mathbb{R}}_+$ a Qbs-morphism. Set $k := m^{\perp} \circ \psi$ and $\underline{\mu} := [\rho \neq 0] \odot (m^{\perp}(a) \boxtimes k)$.*

Assume that: (1) k is Markov; (2) $[1 = (\rho \circ \text{swap}) \cdot \rho]$ holds $\underline{\mu}$ -a.e.; (3) ρ is a Radon-Nikodym derivative of $\text{swap}_ \underline{\mu}$ with respect to $\underline{\mu}$; and (4) $\rho(x, y) = 0 \iff \rho(y, x) = 0$ for all $x, y \in X$.*

Then $(m^{\perp} \circ \eta_{\psi, \rho})(a) = m^{\perp}(a)$.

Using Kock’s synthetic measure theory, we were able to follow closely standard measure-theoretic proofs of MHG [Geyer 2011]. The synthetic setting highlights the different roles each of the three abstractions: a.e.-equality, a.e.-properties, and Radon-Nykodim derivatives play in the proof that our formulation exposes (cf. § 4.2.3).

7.2 Tracing Representation

A sampling *trace* is a sequence of samples that occur during the execution of a probabilistic program. We represent such programs as elements of the continuous weighted sampler WSam from (cf. Fig. 8). Consequently, the collection of traces through a program $t \in \text{WSam } X$ is a subset of $\text{List } \mathbb{I}$. Fig. 12b defines a measurable predicate $[\in] : \text{WSam } X \times \text{List } \mathbb{I} \rightarrow \text{bool}$ that tests whether a given sequence p of probabilistic choice forms a complete trace in the program t . Consequently, we can define the set of *paths* through a given program t by $\text{Paths } t := \{p \in \text{List } \mathbb{I} \mid p \in t\} \subseteq \text{List } \mathbb{I}$, and equip it with the subspace structure it inherits from $\text{List } \mathbb{I}$. We can therefore define the set:

$$\sum_{t \in \text{WSam } X} \text{Paths } t := \{(t, p) \in \text{WSam } X \times \text{List } \mathbb{I} \mid p \in t\} \subseteq \text{WSam } X \times \text{List } \mathbb{I}$$

which we can also equip with a subspace structure. We can now define the *weight* w_- and *valuation* v_- morphisms in Fig. 12b that retrieve the likelihood and value at the end of a trace.

We can now define the tracing inference representation. It is parameterised by an inference representation \underline{T} and given for X as the following subspace of $\text{WSam } X \times T(\text{List } \mathbb{I})$:

$$\text{Tr } \underline{T} X := \left\{ (t, a) \in \text{WSam } X \times T(\text{List } \mathbb{I}) \left| \begin{array}{l} [\in t] m_{\underline{T}}(a)\text{-a.e., and} \\ m_{\text{WSam}}(t) = \int_{\text{List } \mathbb{I}} \delta_{v_t(p)} m_{\underline{T}}(a)(dp) \end{array} \right. \right\}.$$

Thus, a representation consists of a program representation t , together with a distribution a on all lists, but maintaining two invariants. First, the lists are $m_{\underline{T}}(a)$ -almost-everywhere paths through t , and so we can indeed think of a as a representation of a distribution over traces. Second, if we calculate the posterior of the paths through t according to $m_{\underline{T}}(a)$, it should have the same meaning as the original program.

<p>instance $\text{Inf} \implies \text{Inf Monad}(\text{Tr } \underline{T})$ where</p> <p>return $x = (\text{return}_{\text{WSam}} x, \text{return}_{\underline{T}} [\])$</p> <p>$(t, a) \gg= (f, g) = (t \gg=\text{WSam } f, \underline{T}. \text{do} \{ p \leftarrow a;$ $q \leftarrow g \circ v_t(p);$ $\text{return}(p \# q) \})$</p> <p>$m(t, a) = m_{\text{WSam}}(t) = \int_{\text{List } \mathbb{I}} \delta_{v_t(p)} m_{\underline{T}}(a)(dp)$</p> <p>tmap $\underline{t} = \text{id} \times \underline{t}_{\text{List } \mathbb{I}}$</p> <p>sample $= (\text{sample}_{\text{WSam}},$ $\underline{T}. \text{do} \{ r \leftarrow \text{sample}; \text{return}[r] \})$</p> <p>score $r = (\text{score}_{\text{WSam}},$ $\underline{T}. \text{do} \{ \text{score } r; \text{return}[\] \})$</p> <p>(a) The tracing inference</p>	<p>$\eta_{\psi, \rho}^{\text{Tr } T} : \text{Tr } T X \rightarrow \text{Tr } T X$</p> <p>$\eta_{\psi, \rho}^{\text{Tr } T}(t, a) := (t, \eta_{\psi_t, \rho_t}(a))$</p> <p>(b) Trace MH update-step</p> <p>$\text{pri}_{\underline{T}} : \text{WSam } X \rightarrow T(\text{List } \mathbb{I})$</p> <p>$\text{pri}_{\underline{T}}(t) := \text{fold}$ $\lambda \{ \text{Return}(r, x) \rightarrow \text{return}_{\underline{T}} [\]$ $\mid \text{Sample } k \rightarrow \underline{T}. \text{do} \{$ $r \leftarrow \text{sample}_{\underline{T}};$ $k(r) \}$</p> <p>(c) Prior representation</p>
---	---

Fig. 13. Building blocks of Trace MH

We stress that an implementation need *not* compute the meaning of the program. But this representation guarantees that the meaning will be preserved by the inference operations.

Note that the integrand in the definition of $(t, a) \in \text{Tr } T X$ is only partially defined. This partiality is not an issue because the first condition guarantees it is $m^T(a)$ -a.e. defined. We can then choose the constantly 0 distribution when $p \notin t$.

Fig. 13a presents the inference representation structure of $\text{Tr } \underline{T}$. Most of the proof revolves around preseving the invariant, i.e., that these definitions define set-theoretic functions.

The inference transformation $\text{marginal}_{\underline{T}} : \text{Tr } T X \rightarrow T X$ marginalises the trace transformer once it is no longer useful. It first samples a path and then uses it to run the program discarding the weight: $\text{marginal}(t, a) = \text{do} \{ x \leftarrow a; \text{return } v_t(x) \}$. Its correctness is precisely the invariant.

7.3 Inference with MHG

The transition from \underline{T} to $\text{Tr } \underline{T}$ still requires a proposal kernel and a representation of the appropriate derivative, but these can now be given in terms of concrete traces.

Given an inference representation \underline{T} , a *trace proposal kernel* is a transformation representing a kernel $\psi : (\sum_{t \in \text{WSam } X} \text{Paths } t) \rightarrow T(\text{List } \mathbb{I})$. A *trace derivative* is a transformation representing the derivative $\rho : (\sum_{t \in \text{WSam } X} \text{Paths } t \times \text{Paths } t) \rightarrow \overline{\mathbb{R}}_+$. Given a trace proposal kernel ψ and a trace derivative ρ , Fig. 13b presents the trace MHG update transformation using the corresponding MHG update on $T(\text{List } \mathbb{I})$.

The Trace MH update step requires some assumptions to form an inference transformation:

THEOREM 7.2 (TRACE METROPOLIS-HASTINGS-GREEN). *Let T be an inference representation, ψ a trace proposal kernel, and ρ a trace derivative. Assume that, for every $(t, a) \in \text{Tr } T X$, letting $k := m^T \circ \psi_t$ and $\underline{\mu} := [\rho_t \neq 0] \odot (m^T(a) \boxtimes k)$: (1) k is Markov; (2) $[1 = \rho_t \cdot (\rho_t \circ \text{swap})]$ holds $\underline{\mu}$ -a.e.; (3) ρ_t is a Radon-Nikodym derivative of $\text{swap}_* \underline{\mu}$ with respect to $\underline{\mu}$; and (4) $\rho_t(p, q) = 0 \iff \rho_t(q, p) = 0$ for all $p, q \in \text{List } \mathbb{I}$. Then $\eta_{\psi, \rho}^{\text{Tr } T} : \text{Tr } T \rightarrow \text{Tr } T$ is a valid inference transformation.*

We will now demonstrate such a simple and generic trace proposal kernel and trace derivative that implement a MHG update step of a popular lightweight Metropolis-Hastings algorithm in several probabilistic programming language systems [Goodman et al. 2008; Goodman and Stuhlmüller 2014; Hur et al. 2015; Wood et al. 2014].

For any inference representation \underline{T} , Fig. 13c defines the morphism $\text{pri}_{\underline{T}}$ that maps a representation $t \in \text{WSam } X$ to its prior distribution on paths over t . Let $\text{U}_D(n) \in M(\mathbb{N})$ be the measure for the uniform discrete distribution with support $\{0, 1, \dots, n\}$. Intuitively, it assigns a probability $\frac{1}{n+1}$ to every element in the support. It be easily defined from sample_M , which denotes the uniform distribution on \mathbb{I} , as in Lemma 6.1.

We now define our concrete proposal ψ_t and derivative, a.k.a. ratio, ρ_t :

$$\begin{aligned} \psi_t &: \text{List}(\mathbb{I}) \rightarrow T(\text{List}(\mathbb{I})) \\ \psi_t(p) &:= \underline{T}.\text{do} \{i \leftarrow \text{U}_D^{\underline{T}}(|p|) \\ &\quad q \leftarrow \text{pri}^T(\text{sub}(t, \text{take}(i, p))) \\ &\quad \text{return}(\text{take}(i, p) + q)\} \end{aligned} \quad \begin{aligned} \rho_t &: \text{List}(\mathbb{I}) \times \text{List}(\mathbb{I}) \rightarrow \overline{\mathbb{R}}_+ \\ \rho_t(p, q) &:= \frac{w_t(q) \cdot (|p|+1)}{w_t(p) \cdot (|q|+1)} \end{aligned}$$

where $\text{sub}(t, x)$ selects a subterm of a given term by following the list x and $\text{take}(i, p)$ retrieves the i -th prefix of p . This proposal and derivative/ratio satisfy the condition in the Trace MH.

Our approach lets us combine MH updates with other inference building block. For example, recall the SMC algorithm from Section 6.2. Each time it performs resampling, multiple particles are given the same values, which results in inadequate coverage of the space, a phenomenon known as *degeneracy*. One way to ameliorate this problem is to apply multiple MH transitions to each particle after resampling in order to spread them across the space, resulting in an algorithm known as resample-move SMC [Doucet and Johansen 2011].

The implementation of resample-move SMC is very similar to that of SMC from Section 6.2, except we introduce an additional layer Tr between Sus and Pop :

THEOREM 7.3. *Let \underline{T} be a sampling representation. For every pair of natural numbers n, k, ℓ the following composite forms an inference transformation:*

$$\begin{aligned} \text{rmsmc}_{n,k,\ell}^{\underline{T}} &:= (\text{Sus} \circ \text{Tr} \circ \text{Pop})\underline{T} \xrightarrow{\text{tmap}_{\text{Sus}} \text{tmap}_{\text{Tr}} \text{spawn}(n, -)} (\text{Sus} \circ \text{Tr} \circ \text{Pop})\underline{T} \\ &\xrightarrow{(\text{advance} \circ \text{tmap}_{\text{Sus}} \eta^{\circ \ell} \circ \text{tmap}_{\text{Sus}} \text{tmap}_{\text{Tr}} \text{resample}(n, -))^{\circ k}} (\text{Sus} \circ \text{Tr} \circ \text{Pop})\underline{T} \xrightarrow{\text{marginal} \circ \text{finish}} \text{Pop } \underline{T} \end{aligned}$$

In the above we apply ℓ MH transitions after each resampling. Our compositional correctness criterion corresponds to a known result that resample-move SMC is an unbiased importance sampler.

8 RELATED WORK AND CONCLUDING REMARKS

The idea of developing a programming language for machine learning and statistics is old, and was explored at least in the early 2000s [Park et al. 2005; Ramsey and Pfeffer 2002; Thomas et al. 1992] as an interesting yet niche research topic. In the past five years, however, designing such a language and building its runtime system has become an active research area, and lead to practical programming languages and libraries [Carpenter et al. 2017; Goodman et al. 2008; Goodman and Stuhlmüller 2014; Gordon et al. 2014; Mansinghka et al. 2014; Minka et al. 2014; Murray 2013; Narayanan et al. 2016; Tran et al. 2017; Wood et al. 2014]. Most of these research efforts have focussed on developing efficient inference algorithms and implementations [Kucukelbir et al. 2015; Le et al. 2017; Tran et al. 2017; Wingate and Weber 2013]. Only a smaller amount of work has been dedicated to justifying the algorithms or other runtime systems of those languages [Borgström et al. 2015; Hur et al. 2015]. Our work contributes to this less-explored line of research by providing novel denotational techniques and tools for specifying and verifying key components of inference algorithms, in particular, those for expressive higher-order probabilistic programming languages. Such specifications can then be combined to construct the correctness argument of a complex inference algorithm, as we have shown in the paper.

The idea of constructing inference algorithms by composing transformations of an intermediate representations is, to the best of our knowledge, relatively recent. In previous work with Gordon [Ścibior et al. 2015], we manipulated a free monad representation to obtain an implementation of SMC. However, we did not implement MH, did not break down SMC further into resampling and suspension, and our semantics was not compositional. Zinkov and Shan [2017] directly manipulate syntax trees of a small language Hakaru. Their semantics is only first-order and they focus on local program transformations corresponding to solving integrals analytically, which is orthogonal to our global transformations relating to sampling algorithms.

Our approach does not yet deal with two important aspects of inference. In practice, one wants *convergence* guarantees for the inference algorithm, estimating the results within an error margin after a given number of inference steps. As any purely-measure theoretic approach, ours does not express such properties. Additionally, we can not express algorithms that rely on derivatives of the density function for the program traces, such as Hamiltonian Monte Carlo or *variational* inference. Developing a theory of differentiation over quasi-Borel spaces might enable us to express such algorithms.

Another interesting direction for future work is to develop a denotational account of some probabilistic programming languages that allow users to select or compose parts of inference algorithms [Mansinghka et al. 2014; Tran et al. 2017]. The exposure of an inference algorithm in such languages breaks the usual abstraction of probabilistic programs as distributions, and causes difficulties of applying existing semantic techniques to such programs. Our more intensional semantics may be able to overcome these difficulties. Finally, it would be interesting to consider indexed or effect-annotated versions of inference representations, transformations and transformers, where indices or annotations ensure that inference components are applied to a certain type of programs. Such a refined version of our results may lead to a way of selectively applying Hamiltonian Monte Carlo or other algorithms that assume the presence of differentiable densities and a fixed length of all paths through the program.

ACKNOWLEDGMENTS

Supported by a Royal Society University Research Fellowship, Institute for Information & Communications Technology Promotion (IITP) grant funded by the Korea government (MSIP) No. R0190-16-2011 ‘Development of Vulnerability Discovery Technologies for IoT Software Security’, Engineering and Physical Sciences Research Council (EPSRC) Early Career Fellowship EP/L002388/1 ‘Combining viewpoints in quantum theory’, an EPSRC studentship and grants EP/N007387/1 ‘Quantum computation as a programming language’ and EP/M023974/1 ‘Compositional higher-order model checking: logics, models, and algorithms’, a Balliol College Oxford Career Development Fellowship, and a University College Oxford Junior Research Fellowship. We would like to thank Samson Abramsky, Thorsten Altenkirch, Bob Coecke, Mathieu Huot, Radha Jagadeesan, Dexter Kozen, Paul B. Levy, and the anonymous reviewers for fruitful discussions and suggestions.

REFERENCES

- R. J. Aumann. 1961. Borel structures for function spaces. *Illinois Journal of Mathematics* 5 (1961), 614–630.
- Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szymczak. 2015. A Lambda-Calculus Foundation for Universal Probabilistic Programming. *CoRR* abs/1512.08990 (2015). <http://arxiv.org/abs/1512.08990>
- Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szymczak. 2016. A lambda-calculus foundation for universal probabilistic programming. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming, ICFP 2016, Nara, Japan, September 18-22, 2016*. 33–46.
- Bob Carpenter, Andrew Gelman, Matthew Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2017. Stan: A Probabilistic Programming Language. *Journal of Statistical Software, Articles* 76, 1 (2017), 1–32. <https://doi.org/10.18637/jss.v076.i01>

- Arnaud Doucet and Adam M. Johansen. 2011. A Tutorial on Particle Filtering and Smoothing: Fifteen years later. In *The Oxford Handbook of Nonlinear Filtering*, Dan Crisan and Boris Rozovskii (Eds.). Oxford University Press, 656–704.
- Matthias Felleisen. 1991. On the Expressive Power of Programming Languages. *Sci. Comput. Program.* 17, 1-3 (1991), 35–75.
- Marcelo Fiore and Philip Saville. 2017. List Objects with Algebraic Structure. In *2st International Conference on Formal Structures for Computation and Deduction, FSCD 2017*.
- Herman Geuvers and Erik Poll. 2007. Iteration and primitive recursion in categorical terms. In *Reflections on Type Theory, Lambda Calculus, and the Mind, Essays Dedicated to Henk Barendregt on the Occasion of his 60th Birthday*. Radboud Universiteit, Nijmegen, 101–114.
- Charles J. Geyer. 2011. Introduction to Markov Chain Monte Carlo. In *Handbook of Markov Chain Monte Carlo*, Steve Brooks, Andrew Gelman, Galin L. Jones, and Xiao-Li Meng (Eds.). Chapman and Hall/CRC, Chapter 1, 3–48.
- Noah Goodman, Vikash Mansinghka, Daniel M Roy, Keith Bonawitz, and Joshua B Tenenbaum. 2008. Church: a language for generative models. In *UAI*.
- Noah Goodman and Andreas Stuhlmüller. 2014. Design and Implementation of Probabilistic Programming Languages. <http://dippl.org>. (2014).
- Andrew D. Gordon, Thore Graepel, Nicolas Rolland, Claudio Russo, Johannes Borgstrom, and John Guiver. 2014. Tabular: A Schema-driven Probabilistic Programming Language. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '14)*. ACM, New York, NY, USA, 321–334. <https://doi.org/10.1145/2535838.2535850>
- Esfandiar Haghverdi and Philip Scott. 2006. A categorical model for the geometry of interaction. *Theoretical Computer Science* 350, 2 (2006), 252 – 274. <https://doi.org/10.1016/j.tcs.2005.10.028> Automata, Languages and Programming: Logic and Semantics (ICALP-B 2004).
- Chris Heunen, Ohad Kammar, Sam Staton, and Hongseok Yang. 2017. A Convenient Category for Higher-Order Probability Theory. In *Proceedings of the 32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS '17, Reykjavik, Iceland, June 20-23, 2017*.
- Chung-Kil Hur, Aditya V. Nori, Sriram K. Rajamani, and Selva Samuel. 2015. A Provably Correct Sampler for Probabilistic Programs. In *35th LARCS Annual Conference on Foundation of Software Technology and Theoretical Computer Science, FSTTCS 2015, December 16-18, 2015, Bangalore, India*. 475–488.
- Graham Hutton. 1999. A Tutorial on the Universality and Expressiveness of Fold. *J. Funct. Program.* 9, 4 (July 1999), 355–372. <https://doi.org/10.1017/S0956796899003500>
- Bart Jacobs. 2017. From Probability Monads to Commutative Effectuses. *Journ. of Logical and Algebraic Methods in Programming* (2017). To appear.
- Mauro Jaskielioff. 2009. *Lifting of Operations in Modular Monadic Semantics*. Ph.D. Dissertation. University of Nottingham.
- G M Kelly. 1980. A unified treatment of transfinite constructions for free algebras, free monoids, colimits, associated sheaves and so on. *Bull. Austral. Math. Soc.* 22 (1980), 1–83.
- Anders Kock. 1972. Strong functors and monoidal monads. *Archiv der Mathematik* 23, 1 (1972), 113–120.
- Anders Kock. 2012. Commutative monads as a theory of distributions. *Theory and Applications of Categories* 26, 4 (2012), 97–131.
- Alp Kucukelbir, Rajesh Ranganath, Andrew Gelman, and David Blei. 2015. Automatic Variational Inference in Stan. In *NIPS*. <https://papers.nips.cc/paper/5758-automatic-variational-inference-in-stan>
- Tuan Anh Le, Atilim Gunes Baydin, and Frank Wood. 2017. Inference Compilation and Universal Probabilistic Programming. In *AISTATS*. <http://www.tuananhle.co.uk/assets/pdf/le2016inference.pdf>
- Vikash K. Mansinghka, Daniel Selsam, and Yura N. Perov. 2014. Venture: a higher-order probabilistic programming platform with programmable inference. *arXiv:1404.0099* (2014). <http://arxiv.org/abs/1404.0099>
- Francisco Marmolejo and Richard J. Wood. 2010. Monads as extension systems — no iteration is necessary. *Theory and Applications of Categories* 24, 4 (2010), 84–113.
- 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). Microsoft Research Cambridge. <http://research.microsoft.com/infernet>.
- Eugenio Moggi. 1989. Computational Lambda-Calculus and Monads. In *LICS*. IEEE Computer Society, USA, 14–23.
- Lawrence M. Murray. 2013. Bayesian State-Space Modelling on High-Performance Hardware Using LibBi. *arXiv:1306.3277*. (2013).
- Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. 2016. Probabilistic inference by program transformation in Hakaru (system description). In *International Symposium on Functional and Logic Programming - 13th International Symposium, FLOPS 2016, Kochi, Japan, March 4-6, 2016, Proceedings*. Springer, 62–79. https://doi.org/10.1007/978-3-319-29604-3_5
- Sungwoo Park, Frank Pfenning, and Sebastian Thrun. 2005. A probabilistic language based upon sampling functions. In *Proceedings of the 32nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2005, Long Beach, California, USA, January 12-14, 2005*. 171–182.

- Maciej Piróg. 2016. Eilenberg-Moore Monoids and Backtracking Monad Transformers. In *Proceedings 6th Workshop on Mathematically Structured Functional Programming, MSFP@ETAPS 2016, Eindhoven, Netherlands, 8th April 2016. (EPTCS)*, Robert Atkey and Neelakantan R. Krishnaswami (Eds.), Vol. 207. 23–56. <https://doi.org/10.4204/EPTCS.207.2>
- Norman Ramsey and Avi Pfeffer. 2002. Stochastic lambda calculus and monads of probability distributions. In *Conference Record of POPL 2002: The 29th SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Portland, OR, USA, January 16-18, 2002*. 154–165.
- Adam Ścibior, Zoubin Ghahramani, and Andrew Gordon. 2015. Practical Probabilistic Programming with Monads. In *Haskell*. <http://dl.acm.org/citation.cfm?id=2804317>
- Sam Staton. 2017. Commutative semantics for probabilistic programming. In *Proc. ESOP 2017*.
- Andrew Thomas, David J. Spiegelhalter, and W. R. Gilks. 1992. BUGS: A program to perform Bayesian inference using Gibbs sampling. *Bayesian statistics 4* (1992), 837–842. Issue 9.
- Dustin Tran, Matthew D. Hoffman, Rif A. Saurous, Eugene Brevdo, Kevin Murphy, and David M. Blei. 2017. Deep Probabilistic Programming. In *ICLR*.
- David Wingate, Andreas Stuhlmüller, and Noah Goodman. 2011. Lightweight Implementations of Probabilistic Programming Languages Via Transformational Compilation. In *AISTATS*. <https://web.stanford.edu/~ngoodman/papers/lightweight-mcmc-aistats2011.pdf> The published version contains a serious bug in the definition of alpha that was fixed in revision 3 available at the given URL.
- David Wingate and Theophane Weber. 2013. Automated Variational Inference in Probabilistic Programming. arXiv:1301.1299. (2013).
- 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.
- Robert Zinkov and Chung-chieh Shan. 2017. Composing inference algorithms as program transformations. In *Proceedings of Uncertainty in Artificial Intelligence*.