



Compiling Probabilistic Programs for Variable Elimination with Information Flow

JIANLIN LI, University of Waterloo, Canada

ERIC WANG, University of Waterloo, Canada

YIZHOU ZHANG, University of Waterloo, Canada

A key promise of probabilistic programming is the ability to specify rich models using an expressive programming language. However, the expressive power that makes probabilistic programming languages enticing also poses challenges to inference, so much so that specialized approaches to inference ban language features such as recursion. We present an approach to variable elimination and marginal inference for probabilistic programs featuring bounded recursion, discrete distributions, and sometimes continuous distributions. A compiler eliminates probabilistic side effects, using a novel information-flow type system to factorize probabilistic computations and hoist independent subcomputations out of sums or integrals. For a broad class of recursive programs with dynamically recurring substructure, the compiler effectively decomposes a global marginal-inference problem, which may otherwise be intractable, into tractable subproblems. We prove the compilation correct by showing that it preserves denotational semantics. Experiments show that the compiled programs subsume widely used PTIME algorithms for recursive models and that the compilation time scales with the size of the inference problems. As a separate contribution, we develop a denotational, logical-relations model of information-flow types in the novel measure-theoretic setting of probabilistic programming; we use it to prove noninterference and consequently the correctness of variable elimination.

CCS Concepts: • **Theory of computation** → *Probabilistic computation; Program semantics; Program reasoning; Type theory*; • **Mathematics of computing** → *Bayesian computation; Statistical software*; • **Computing methodologies** → *Machine learning*; • **Software and its engineering** → *Compilers; Functional languages; Language features; Formal language definitions*.

Additional Key Words and Phrases: Probabilistic programming, probabilistic inference, information flow.

ACM Reference Format:

Jianlin Li, Eric Wang, and Yizhou Zhang. 2024. Compiling Probabilistic Programs for Variable Elimination with Information Flow. *Proc. ACM Program. Lang.* 8, PLDI, Article 218 (June 2024), 26 pages. <https://doi.org/10.1145/3656448>

1 INTRODUCTION

A probabilistic model describes a joint distribution $p(x, z)$ over latent variables z and observations x . Bayesian inference is concerned with computing $p(z|x) = p(x, z)/p(x)$, the posterior distribution of z conditioned on x . Typically, the hard work is in computing the *marginal likelihood* $p(x)$, also known as the *model evidence*. Computing the marginal may be intractable, as it generally requires integration over all possible values of the latent variables: $p(x) = \int p(x, z) dz$.

Probabilistic programming languages (PPLs) are powerful means to specify probabilistic models and solve inference problems. A PPL allows for harnessing the expressivity of a high-level programming language to specify rich Bayesian models, as opposed to using more limiting formalisms such

Authors' addresses: Jianlin Li, jianlin.li@uwaterloo.ca, University of Waterloo, Canada; Eric Wang, e224wang@uwaterloo.ca, University of Waterloo, Canada; Yizhou Zhang, yizhou@uwaterloo.ca, University of Waterloo, Canada.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2024 Copyright held by the owner/author(s).

ACM 2475-1421/2024/6-ART218

<https://doi.org/10.1145/3656448>

as Bayesian networks. However, the expressive power that makes PPLs enticing makes inference even harder. Recent advances in PPL inference often specialize to a particular class of models and impose restrictions on expressible models, prohibiting useful features such as recursion.

Variable elimination (VE) is an effective approach to inference for probabilistic models with discrete random variables (r.v.s) [68]. It works by marginalizing out (i.e., eliminating) discrete r.v.s from a joint distribution, thus producing the marginal likelihood or reducing the inference problem to ones involving only continuous r.v.s.

VE has been generalized to PPLs, but the support for VE does not meet the desired level generality and scalability. For example, Factorie [42] is an interpreted PPL that supports VE for factor graphs, but its VE algorithms make specific, rigid assumptions on the factor-graph structure. SlicStan [29] and PERPL [15] are more recent, compiled approaches to VE and are designed to support broader classes of models. Unfortunately, SlicStan lacks support for recursion—recursion is a natural means to specify models in domains such as language modeling and computational biology. Another issue is that the time SlicStan takes to compile a program does not scale well with the size of the inference problem. PERPL, while supporting (unbounded) recursion, is designed to work in the absence of continuous r.v.s. For exact inference on certain recursive models involving only discrete r.v.s, PERPL does not empirically scale as well as the best algorithms known for the same models.

This paper presents a novel approach to VE for an expressive PPL. While acknowledging that it is an elusive, likely impossible goal for any single inference method to excel at all expressible programs, we aim to achieve good efficiency and scalability across a wide range of programs featuring bounded recursion, discrete distributions, and occasionally continuous distributions, all the while with provable correctness guarantees. We embody this approach in a PPL called MAPPL.

VE as compilation. In MAPPL, a probabilistic computation is compiled into a pure computation of the marginal likelihood. A discrete r.v. is eliminated by summing (over the variable's finite support) the product of all factors dependent on that variable. Control flow, namely branching and function calls, are compiled in continuation-passing style: the compiled branch or function takes as input a continuation representing the product of all factors dependent on the return value.

Decomposition, memoization, and amortization. We observe that many recursive probabilistic programs of interest enjoy the property that the exponentially many possible executions share substructure. For these programs, the VE compilation effectively decomposes, in a recursive manner, a global marginal inference problem into subproblems amenable to dynamic programming [9]. The same subproblem instance may be queried multiple times during inference, so the solution to the inference subproblem can be *memoized* and thus the cost of solving the subproblem *amortized*.

The subproblems are likely easier to solve than the global problem, because they have reduced dimensionality and are free of language constructs such as recursion that are difficult for inference. Some of these subproblems may be solved easily if they happen to contain no continuous r.v.s, some may be solved by existing approximate inference methods that specialize in straight-line programs with continuous r.v.s, and some may sometimes even be solved analytically by capitalizing on advances in symbolic integration. The upshot is that the VE compilation may render an otherwise intractable inference problem solvable in polynomial time.

Factorization by information-flow typing. It is well understood that the effectiveness of VE critically depends on exploiting independence to factorize joint distributions. With recursion, there is more reason for a VE compiler to exploit independence, as decomposition and memoization just would not be as effective if subproblem definitions were too coarse-grained.

To reason about independence, we design an information-flow type system for MAPPL. To eliminate a variable x from a computation, the MAPPL compiler consults information-flow typing to

factorize the computation into two parts, the probabilistic side effects of which are respectively dependent and independent of x . The idea of using information-flow typing [21] to reason about independence is similar to that in SlicStan [29], but MAPPL’s type-system design and formal development differ substantially from SlicStan’s. SlicStan is an imperative while-language where variables must be global and programs must have deterministic support, whereas MAPPL is an expressive functional language allowing recursion and stochastic support. SlicStan is defined with an operational semantics, whereas we adopt a *compositional*, denotational treatment suited for reasoning about independence—and thus for factorizing computations—for open terms under binders.

Generality and scalability. MAPPL generalizes VE compilation and information-flow typing to recursive probabilistic programs—for example, those expressing hidden Markov models (HMMs) and probabilistic context-free grammars (PCFGs).¹ Experiments show that MAPPL’s VE compiler can generate code that recovers widely used polynomial-time inference algorithms: the forward algorithm for HMMs [55] and the inside algorithm for PCFGs [5].

We consider it important to achieve good scalability of not only inference but also compilation. Notably, compared with SlicStan, the increased generality of MAPPL to support recursion has implications for the scalability of compilation. In SlicStan, HMMs have to be expressed by unrolling recursion into a fixed number of iterations. For such models, SlicStan’s compilation time does not scale well with the problem size (e.g., the length of the observed sequence). In MAPPL, by contrast, compilation time stays constant as the problem size increases for such models, because MAPPL can express them as probabilistic recursive functions and compile them to pure recursive functions, without unrolling. In addition, SlicStan uses a semilattice—as opposed to a lattice—of information-flow labels for factorization, which is considered to impede efficient label inference. By contrast, MAPPL uses a simpler, better-behaved two-level lattice.

Correctness guarantees. We want to show that MAPPL’s VE compilation is correct by proving that the compiled program computes the marginal likelihood as defined by the denotational semantics. Since compilation uses information-flow typing to factorize computations, we need to show that our information-flow type system is sound with respect to the denotational semantics. To that end, we contribute a logical-relations model of information-flow types for proving noninterference in the novel, measure-theoretic setting of probabilistic programming.

2 KEY FEATURES, MAIN IDEAS, AND EXAMPLES

We use a simple hidden Markov model as a starting point to illustrate the key features and the main ideas of our approach. Figure 1a models a sequence of observations as being generated by a sequence of hidden states. The recursive function `hmm` takes as input the initial hidden state z_0 and a `data` sequence, which is assumed to be gathered by prepending newer observations to the front of the sequence. The return value of `hmm` is the next hidden state. The probability of transitioning from a state to the next state is given by a pure function `step` : $\mathbb{B} \rightarrow \text{dist}(\mathbb{B})$. The probability of observing a data point in a state is given by an emission function `emit` : $\mathbb{B} \rightarrow \text{dist}(\mathbb{R})$. The HMM is conditioned on observing the `data` sequence.

The inference problem is, given any given `data` sequence, to compute the marginal likelihood of observing it. The MAPPL compiler translates the recursive, probabilistic `hmm` in Figure 1a into the recursive, pure `hmm` in Figure 1b. When the compiled `hmm` is called with the top-level continuation $\lambda_.0$ for the parameter `k`, it computes the desired marginal likelihood. This procedure for exact inference runs in time linear in the length of `data`, recovering the forward algorithm for HMMs.

¹As a caveat, this paper does not address VE for almost surely terminating programs, which use unbounded recursion; we consider the restriction to deterministically bounded recursion to be a reasonable trade-off (Section 8).

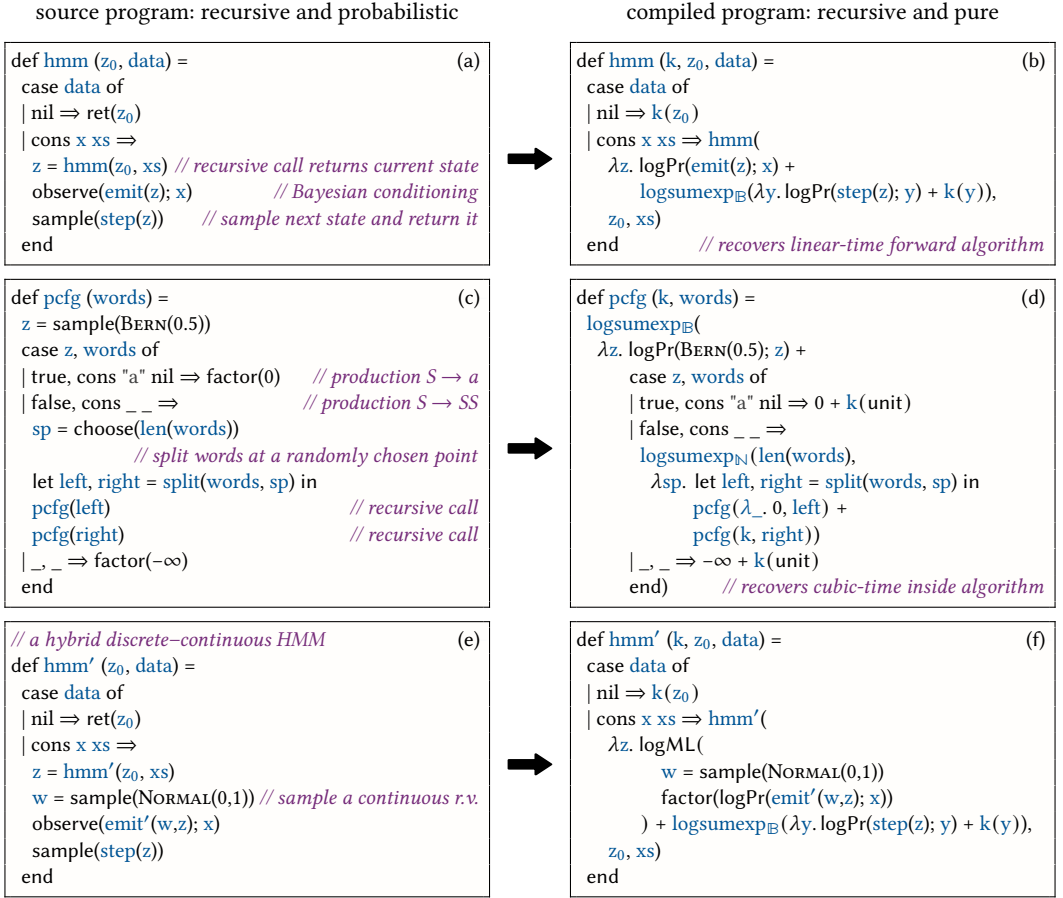


Figure 1. Examples of VE-compiling probabilistic programs in MAPPL. The return value of a multi-line block of terms is that of the last term. The construct $\text{ret}(e)$ is the monadic return that lifts a pure expression to a probabilistic term. The primitive $\text{logPr}(d; \cdot)$ is the log-probability density or mass function of the distribution d . $\text{logsumexp}_{\mathbb{R}} : (\mathbb{B} \rightarrow \mathbb{R}) \rightarrow \mathbb{R}$ and $\text{logsumexp}_{\mathbb{N}} : \mathbb{N} \rightarrow (\mathbb{N} \rightarrow \mathbb{R}) \rightarrow \mathbb{R}$ are the usual log-sum-exp functions for log-domain sums. The choose primitive randomly selects a natural number in a given range, but unlike the sample primitive, choose does not otherwise incur any probabilistic side effects.

Handling expressive language features. The `hmm` example uses recursion, which violates the assumptions of many existing approaches to PPL inference. Instead of using recursion to define the HMM, one could perform exact inference by unfolding the model into a fixed number of iterations and then applying existing inference methods good at nonrecursive programs. But this approach is awkward when the number of iterations is not known statically—namely, when `data` is dynamically sized. A key design goal of MAPPL is to support a broad class of models definable with bounded recursion where the bound may not be known statically.

Unrolling recursion is even more awkward for models that are not iterative but properly recursive and for models where control flow is stochastic. A prime example is PCFGs. Figure 1c shows a PCFG model in MAPPL that samples parse trees for the simple grammar $S \rightarrow a(0.5) \mid SS(0.5)$. The program is conditioned on it generating a parse tree for a sequence of words. The recursion pattern of `pcfg` is more complex than the HMMs. First, it is tree-structured rather than linear. Second, control flow is stochastic—which branch of case is taken depends on the sampled variable `z`

in each recursive call. Yet, the MAPPL compiler can still compile the program into a pure one that computes the marginal likelihood of observing `words`. The compilation applies to *mutually* recursive functions, too, which are useful for expressing more complex PCFGs in practice. The pure `pcfg` in Figure 1d recovers the cubic-time inside algorithm.

An information-flow type system. VE for Bayesian networks has a worst-case exponential running time, but tractable inference is possible for certain models by exploiting independence in the model structure: factors independent of a variable can be factored out of the summation that marginalizes out the variable.

Similarly, VE for recursive programs requires analyzing dependence and independence, for which we use a *static* information-flow analysis. In particular, we design an information-flow type system. Information-flow typing is a compositional, automatable means to reason about dependence, with applications in many different contexts [1], among which the most well-known is language-based security [60]. We repurpose the idea to our expressive PPL.

While information-flow typing for the pure fragment of MAPPL is mostly standard, it is less obvious how to design typing rules for the probabilistic fragment. A principle is that the design of the type system should be guided by the denotational semantics. The denotational semantics of a probabilistic computation in MAPPL is a measure over the space of its possible outcomes, which can be roughly thought of as a function that, given a set of values, produces the unnormalized probability that the computation returns a value in that set. Accordingly, our type system assigns to a probabilistic computation a labeled type A^ℓ : A types the return value of the computation, and the label ℓ classifies the level of information contained in the measure denoting the computation.

For example, consider typing variable bindings in the probabilistic fragment. As expected, the rule requires the composed computation $x = t; m$ to have a label no lower than the labels of the computations t and m being composed. However, it does not explicitly constrain

$$\frac{\Delta; \Psi; \Xi; C; \Gamma \vdash t : A^{\ell_1} \quad \Delta; \Psi; \Xi; C; \Gamma, x : A^{\ell_2} \vdash m : B^{\ell_3}}{\Delta; \Psi; \Xi; C; \Gamma \vdash x = t; m : B^{\ell_1 \sqcup \ell_3}}$$

the label ℓ_2 of the variable x being bound. In particular, it is *not* required that the label of x be at least as high as the label of t . This is in keeping with the denotation of $x = t; m$, which is defined by composing the measures that denote t and m and *marginalizing* over the entire support of x .

Factorizing computations using information-flow typing. To eliminate a variable from a probabilistic computation, MAPPL's VE compiler infers labels for the subcomputations, constraining that the variable being eliminated be labeled H (high). As many subcomputations are inferred to be labeled L (low) as possible. The larger computation can thus be factorized into a H partition and a L partition, and the L partition need not be involved in the elimination of the variable.

For example, in Figure 1a, to eliminate `sample(step(z))` from the `cons` branch, the sampled variable is labeled H, and the information-flow analysis deduces that the probabilistic side effects of the recursive call `hmm(z0, xs)` and the conditioning `observe(emit(z); x)` can be labeled L, while the side effects of `sample(step(z))` and any computations in the caller that depend on the return value z must be labeled H. This factorization indicates that in the compiled program, only the H partition needs to be nested under the `logsumexpB` that marginalizes out `sample(step(z))`.

Compiling with continuations. Continuation-passing style (CPS) transformations [19] are an effective way to eliminate various forms of effects away from a program. It has found applications in the implementation of PPLs [28] and in the cost analysis of randomized algorithms [4, 35]. The MAPPL compiler uses CPS to capture dependence in the presence of functions calls and branching.

For example, in the compiled `hmm`, a continuation of type $\mathbb{B} \rightarrow \mathbb{R}$ represents the dependencies of the recursive call's return value z : the continuation takes z as input and returns a log-likelihood that is the transformation of those caller terms whose probabilistic side effects depend on z . In the

compiled `pcfg`, the continuations passed to recursive calls are less involved, as the information-flow analysis deduces that the return value has no nontrivial dependencies.

Decomposition into subproblems with memoization. The compiled `hmm` and `pcfg` run in polynomial time, despite the exponentially many possible executions of the input programs. This algorithmic efficiency is because the recursive programs have dynamically recurring substructure, on which the VE compilation capitalizes to generate recurring subproblems and memoize their solutions. For example, the compiled `hmm` corresponds to the following recursive equations for computing, in log domain, the marginal likelihood $L(\mathbf{k}, z_0, \text{data})$ of observing `data`:

$$L(\mathbf{k}, z_0, \text{nil}) = k(z_0) \quad \underbrace{\quad}_{\phi_1} \quad \underbrace{\quad}_{\phi_2: \text{ nested inference subproblem (summing out } y)} \\ L(\mathbf{k}, z_0, \text{cons } x \text{ xs}) = L\left(\lambda z. \log \Pr(\text{emit}(z); x) + \log \sum_y \exp\left(\log \Pr(\text{step}(z); y) + k(y)\right), z_0, \text{xs}\right)$$

The generated inference subproblem ϕ_2 eliminates the discrete r.v. `sample(step(z))`, denoted by `y`, by summing over its finite support. Although the subproblem is nested inside a continuation, inference time is linear in the length of `data`: the subproblem needs to be solved at most once for each of the two possible values of `z` and for each continuation `k` created—there can only be as many as the length of `data`. The solution to a subproblem instance, once computed, can be memoized and reused whenever the same subproblem instance is encountered again. This decomposition and memoization is what recovers the dynamic-programming algorithms for HMMs and PCFGs [55, 5]. Contrast this with solving the global problem directly (say, with an enumeration-based approach to exact inference) without first compiling the recursive program to decompose it into subproblems:

$$L(z_0, x_1, \dots, x_n) = \log \sum_{z_1} \sum_{z_2} \dots \sum_{z_n} \prod_{i=0}^{n-1} \Pr(\text{emit}(z_i); x_{i+1}) \Pr(\text{step}(z_i); z_{i+1})$$

The sums over the Boolean variables `z1`, `z2`, ..., `zn` enumerate all possible execution traces of the program. Hence, the inner product has to be computed $O(2^n)$ times, where n is the length of `data`.

Continuous parameters. Consider `hmm'` in Figure 1e, a hybrid discrete–continuous HMM. It is largely the same as `hmm` except that the emission function takes as input an additional, freshly sampled Gaussian variable `w`. Directly solving the marginal-inference problem using a general-purpose inference method such as importance sampling would be intractable. Instead, Figure 1f shows that `hmm'` is compiled similarly to `hmm`, with the factor ϕ_1 replaced by a nested inference problem marginalizing out `w`. Marginal inference with the compiled `hmm'` is efficient, even when the inference subproblem $\log \text{ML}(\dots)$ is solved using general-purpose Monte Carlo methods. The key is that with memoization, this subproblem generated by MAPPL's compiler needs to be solved only $O(n)$ times, once for each of the two values of `z` and for each of the at most n values of `x`:

$$L(\mathbf{k}, z_0, \text{cons } x \text{ xs}) = L\left(\lambda z. \log \int_w \exp\left(\log \Pr(\text{NORMAL}(0, 1); w) + \log \Pr(\text{emit}'(w, z); x)\right) dw + \log \sum_y \exp(\dots), z_0, \text{xs}\right) \\ \underbrace{\quad}_{\phi'_1: \text{ nested inference subproblem (integrating out } w)} \quad \underbrace{\quad}_{\phi_2}$$

A semantic model of information-flow types. Factoring out independent factors is akin to *loop-invariant code motion*, a compiler optimization that moves code outside a loop if it is independent of the loop index. The dependence analysis VE entails is more sophisticated, though, due to the measure-theoretic nature of the *semantics* of probabilistic programs. Fortunately, information-flow typing provides a *syntactic*, principled means to reason about independence.

How can we argue that the syntactic approach of information-flow typing to reasoning about independence is semantically sound in this probabilistic setting? For an information-flow type

system, the usual notion of soundness is noninterference [27]. In our novel setting, we take noninterference to mean that the measure denoting a probabilistic computation of a type labeled low behaves irrespective of substitutions for its high-labeled free variables. To prove noninterference, we adapt the semantic, logical-relations proof technique [61, 1], interpreting labeled types as partial equivalence relations on measures indistinguishable to an observer. We believe that our semantic model and its metatheories are the first to introduce observer-sensitive equational reasoning to a measure-theoretic setting and thus are of independent interest.

3 SYNTAX, TYPE SYSTEM, AND DENOTATIONAL SEMANTICS

Syntax. Figure 2 defines the syntax of MAPPL programs. Local variables and global variables are notated in blue. An overline denotes a sequence of zero or more elements.

MAPPL has a pure, deterministic fragment and a monadic, probabilistic fragment, similar to some prior PPL formalisms [67, 37, 36]. Pure computations take the form of *expressions*. The pure fragment is a simply typed λ -calculus equipped with real numbers, pairs, sums, iso-recursive types, n -ary operations, and two primitive distributions (Bernoulli and normal, representative of discrete and continuous distributions). A special binary operation $\log\text{Pr}(d; e)$ gives the log-probability density or mass of a distribution d at a point e .

Recursive types, sum types, and product types together enable the expression of algebraic data types, including Booleans and lists: $\mathbb{B} \stackrel{\text{def}}{=} \mathbb{U} + \mathbb{U}$, $\text{list}_\tau \stackrel{\text{def}}{=} \mu\alpha. \mathbb{U} + (\tau \times \alpha)$. Distributions have type $\text{dist}(\sigma)$, where $\sigma = \mathbb{B}$ for a Bernoulli distribution and $\sigma = \mathbb{R}$ for a normal distribution.

Probabilistic computations are in the forms of *terms* and *commands*. A command sequences terms. The return value of a command is that of its last term. We will write t for t \$ when it is clear from context that t is being used as the last term of a command. Terms have the following forms: $\text{ret}(e)$ returns the value of a pure expression e , $\text{sample}(d)$ samples from a distribution d , $\text{factor}(e)$ conditions the program using a log-domain expression e , $\text{case}(e; x.m_1; x.m_2)$ branches on a sum-typed expression e , and $f(\bar{e})$ invokes a global function f with arguments \bar{e} . The factor form supports soft constraints, which subsume conditioning on continuous observations—that is, $\text{observe}(d; e)$ can be encoded as $\text{factor}(\log\text{Pr}(d; e))$. For brevity of presentation, we omit hard constraints (i.e., $\text{factor}(-\infty)$), but they are straightforward to incorporate in both the syntax and the semantics. We sometimes write $\text{sample}_\sigma(d)$ and $\log\text{Pr}_\sigma(d; e)$, where σ is \mathbb{B} or \mathbb{R} , to indicate the type of the support of the distribution d .

The pure fragment supports nested marginal inference via the form $\log\text{ML}(m)$. While the command m is probabilistic, $\log\text{ML}(m)$ is a pure expression, since inference handles the probabilistic effects of m . It returns the log-marginal likelihood of the probabilistic computation m .

A program in MAPPL is composed of a set of global definitions and a *main command*. The global definitions can be either pure (\mathcal{G}) or probabilistic (\mathcal{F}). Importantly, MAPPL allows mutual recursion among pure globals and among probabilistic globals.

Recursion leaves open the possibility of nontermination, but for VE in this paper, we do not concern ourselves with programs that have possibility of not terminating.

Type system. The base type system for MAPPL is standard. Figure 2 shows selected rules. Expression typing judgments have the form $\Delta; \Gamma \vdash e : \tau$, where Δ is a context mapping names of pure global definitions to their types, and Γ is a context mapping local variables to their types. Term and command typing judgments have the forms $\Delta; \Psi; \Gamma \vdash t : \tau$ and $\Delta; \Psi; \Gamma \vdash m : \tau$. Computations in the probabilistic fragment can use probabilistic globals, whose types are provided by the context Ψ .

Denotational semantics. We use ω -quasi Borel spaces (ωqbses) [64] as the semantic domain. ωqbses are as a drop-in replacement for measurable spaces. They enable carrying out measure theory in the presence of higher-order types and recursive types, by providing well-behaved

expressions	$e, d ::= \mathbf{x} \mid \text{unit} \mid r \mid \lambda \mathbf{x}. e \mid e_1 e_2 \mid \langle e_1, e_2 \rangle \mid \pi_1(e) \mid \pi_2(e)$ $\mid \text{inl } e \mid \text{inr } e \mid \text{case}(e; \mathbf{x}.e_1; \mathbf{x}.e_2) \mid \text{let}(e_1; \mathbf{x}.e_2)$ $\mid \text{roll}(e) \mid \text{unroll}(e_1; \mathbf{x}.e_2) \mid \text{op}(e_1, \dots, e_n)$ $\mid D(e_1, \dots, e_n) \mid \text{logPr}(d; e) \mid \text{logML}(m)$	types
distributions	$D ::= \text{BERN} \mid \text{NORMAL}$	$\tau, \sigma ::= \mathbb{U} \mid \mathbb{R} \mid \alpha \mid \text{dist}(\tau) \mid \mu\alpha. \tau$ $\mid \tau_1 \rightarrow \tau_2 \mid \tau_1 + \tau_2 \mid \tau_1 \times \tau_2$
terms	$t ::= \text{ret}(e) \mid \text{sample}(d) \mid \text{factor}(e) \mid \mathbf{f}(e_1, \dots, e_n)$ $\mid \text{case}(e; \mathbf{x}.m_1; \mathbf{x}.m_2)$	contexts of locals
commands	$m ::= t \$ \mid \mathbf{x} = t; m$	$\Gamma ::= \emptyset \mid \Gamma, \mathbf{x} : \tau$
pure globals	$\mathcal{G} ::= \text{def } \mathbf{x} = e$	contexts of pure globals
prob. globals	$\mathcal{F} ::= \text{def } \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) = m$	$\Delta ::= \emptyset \mid \Delta, \mathbf{x} : \tau$
programs	$\mathcal{P} ::= \overline{\mathcal{G}}; \overline{\mathcal{F}}; m$	types of probabilistic globals
		$F ::= (\tau_1, \dots, \tau_n) \rightarrow \tau$
		contexts of probabilistic globals
		$\Psi ::= \emptyset \mid \Psi, \mathbf{f} : F$
	$\frac{\Gamma(\mathbf{x}) = \tau}{\Delta; \Gamma \vdash \mathbf{x} : \tau} \quad \frac{\mathbf{x} \notin \text{dom}(\Gamma) \quad \Delta(\mathbf{x}) = \tau}{\Delta; \Gamma \vdash \mathbf{x} : \tau} \quad \frac{\Delta; \Gamma \vdash e : \mathbb{R}}{\Delta; \Gamma \vdash \text{BERN}(e) : \text{dist}(\mathbb{B})} \quad \frac{\Delta; \Gamma \vdash e_1 : \mathbb{R} \quad \Delta; \Gamma \vdash e_2 : \mathbb{R}}{\Delta; \Gamma \vdash \text{NORMAL}(e_1, e_2) : \text{dist}(\mathbb{R})}$	
	$\frac{\Delta; \Gamma \vdash d : \text{dist}(\tau) \quad \Delta; \Gamma \vdash e : \tau}{\Delta; \Gamma \vdash \text{logPr}(d; e) : \mathbb{R}} \quad \frac{\Delta; \emptyset; \Gamma \vdash m : \mathbb{U}}{\Delta; \Gamma \vdash \text{logML}(m) : \mathbb{R}} \quad \frac{\Delta; \Gamma \vdash e : \tau}{\Delta; \Psi; \Gamma \vdash \text{ret}(e) : \tau} \quad \frac{\Delta; \Gamma \vdash e : \mathbb{R}}{\Delta; \Psi; \Gamma \vdash \text{factor}(e) : \mathbb{U}}$	
	$\frac{\Delta; \Gamma \vdash e : \tau_1 + \tau_2 \quad \Delta; \Psi; \Gamma, \mathbf{x} : \tau_i \vdash m_i : \tau \text{ for } i \in \{1, 2\}}{\Delta; \Psi; \Gamma \vdash \text{case}(e; \mathbf{x}.m_1; \mathbf{x}.m_2) : \tau} \quad \frac{\Psi(\mathbf{f}) = (\tau_1, \dots, \tau_n) \rightarrow \tau \quad \Delta; \Gamma \vdash e_i : \tau_i \text{ for } i \in \{1, \dots, n\}}{\Delta; \Psi; \Gamma \vdash \mathbf{f}(e_1, \dots, e_n) : \tau}$	
	$\frac{\Delta; \Gamma \vdash d : \text{dist}(\sigma)}{\Delta; \Psi; \Gamma \vdash \text{sample}(d) : \sigma} \quad \frac{\Delta; \Psi; \Gamma \vdash t : \tau}{\Delta; \Psi; \Gamma \vdash t \$: \tau} \quad \frac{\Delta; \Psi; \Gamma \vdash t : \tau \quad \Delta; \Psi; \Gamma, \mathbf{x} : \tau \vdash m : \tau'}{\Delta; \Psi; \Gamma \vdash \mathbf{x} = t; m : \tau'}$	
	$\frac{\Psi(\mathbf{f}) = (\tau_1, \dots, \tau_n) \rightarrow \tau \quad \Delta; \Psi; \mathbf{x}_1 : \tau_1, \dots, \mathbf{x}_n : \tau_n \vdash m : \tau}{\Delta; \Psi \vdash \text{def } \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) = m : \Psi(\mathbf{f})} \quad \frac{\text{dom}(\Psi) = \{\mathbf{f}_i \mid \mathcal{F}_i \in \overline{\mathcal{F}} \wedge \text{name}(\mathcal{F}_i) = \mathbf{f}_i\} \quad \Delta \vdash \overline{\mathcal{G}} : \Delta \quad \Delta; \Psi \vdash \overline{\mathcal{F}} : \Psi}{\Delta; \Psi \vdash \overline{\mathcal{F}} : \Psi} \quad \frac{\Delta; \Psi; \emptyset \vdash m : \tau}{\vdash \overline{\mathcal{G}}; \overline{\mathcal{F}}; m}$	

Figure 2. Syntax of MAPPL and selected typing rules.

function spaces and ωcpos . Thus, in this paper, we use standard measure-theory notations with the understanding that we are working with ωqbs s. For a measurable function $f : X \rightarrow \mathbb{R}^+$, the Lebesgue integral of f with respect to a measure μ on X is denoted $\int f d\mu$ or $\int f(x)\mu(dx)$.

The semantic interpretation $\llbracket \tau \rrbracket$ of each type τ is an ωqbs . Types in MAPPL are similar to those in the SFPC calculus of Vákár et al. [64], which has function types and iso-recursive types; we refer the reader to their paper for detailed constructions of the ωqbs s. For an ωqbs X , we write X_\perp for the lifting of $\llbracket \tau \rrbracket$ to another ωqbs with an extra element \perp signifying partiality. There is a commutative strong monad of measures on ωqbs s [64]; for an ωqbs X , we write $\text{Meas } X$ for the ωqbs of measures on X .

Figure 3 shows selected interpretations of expressions, terms, and commands. The definitions use the operator \gg to handle partiality. In addition to $\gg : \llbracket \tau_1 \rrbracket_\perp \rightarrow (\llbracket \tau_1 \rrbracket \rightarrow \llbracket \tau_2 \rrbracket_\perp) \rightarrow \llbracket \tau_2 \rrbracket_\perp$, we overload \gg on $\llbracket \tau_1 \rrbracket_\perp \rightarrow (\llbracket \tau_1 \rrbracket \rightarrow \text{Meas } \llbracket \tau_2 \rrbracket_\perp) \rightarrow \text{Meas } \llbracket \tau_2 \rrbracket_\perp$ such that $v \gg f \stackrel{\text{def}}{=} v$ if $v = \perp$ then $\lambda E. \mathbb{1}_E(\perp)$ else $f(v)$. Here, $\mathbb{1}_E(v)$ is the indicator function that is 1 if $v \in E$ and 0 otherwise.

The contract for interpreting expressions is that an expression of type τ is interpreted as an element of $\llbracket \tau \rrbracket_\perp$. The denotation $\llbracket e \rrbracket_{\delta; \gamma}$ of an expression e typed in contexts Δ and Γ is interpreted under substitutions $\delta \in \llbracket \Delta \rrbracket$ and $\gamma \in \llbracket \Gamma \rrbracket$ for the bindings in Δ and Γ . That is, δ and γ provide semantic interpretations for the bound global and local variables. The denotation of a primitive distribution is its probability density or mass function.

The contract for interpreting terms and commands is that a term or command of type τ is interpreted as a measure on $\llbracket \tau \rrbracket_\perp$ —i.e., an element of $\text{Meas } \llbracket \tau \rrbracket_\perp$. The interpretations $\llbracket t \rrbracket_{\delta; \gamma}$ and

$$\begin{array}{l}
\llbracket \mathbf{x} \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \gamma(\mathbf{x}) \text{ if } \mathbf{x} \in \text{dom}(\gamma) \\
\llbracket \mathbf{x} \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \delta(\mathbf{x}) \text{ if } \mathbf{x} \notin \text{dom}(\gamma) \wedge \mathbf{x} \in \text{dom}(\delta) \\
\llbracket \text{inl } e \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \text{inl } \llbracket e \rrbracket_{\delta; \gamma} \\
\llbracket \text{BERN}(e) \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \llbracket e \rrbracket_{\delta; \gamma} \ggg \lambda p. \text{ if } 0 \leq p \leq 1 \\
\quad \text{then } (\lambda v. \text{ if } v \text{ then } p \text{ else } 1 - p) \text{ else } \perp \\
\llbracket \log \text{Pr}(d; e) \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \llbracket d \rrbracket_{\delta; \gamma} \ggg \lambda f. \llbracket e \rrbracket_{\delta; \gamma} \ggg \lambda v. \log f(v) \\
\llbracket \log \text{ML}(m) \rrbracket_{\delta; \gamma} \stackrel{\text{def}}{=} \text{let } r = \llbracket m \rrbracket_{\delta; \varnothing; \gamma}(\llbracket \cup \rrbracket) \text{ in} \\
\quad \text{if } 0 < r < \infty \text{ then } \log r \text{ else } \perp \\
\llbracket \text{ret}(e) \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \lambda E. \mathbb{1}_E(\llbracket e \rrbracket_{\delta; \gamma}) \\
\llbracket \text{case}(e; \mathbf{x}. m_1; \mathbf{x}. m_2) \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \llbracket e \rrbracket_{\delta; \gamma} \ggg \lambda v. \text{ case } v \text{ of} \\
\quad \text{inl } u \Rightarrow \llbracket m_1 \rrbracket_{\delta; \psi; \gamma[x \mapsto u]} \mid \text{inr } u \Rightarrow \llbracket m_2 \rrbracket_{\delta; \psi; \gamma[x \mapsto u]} \text{ end} \\
\llbracket \text{factor}(e) \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \llbracket e \rrbracket_{\delta; \gamma} \ggg \lambda r. \lambda E. \exp(r) \cdot \mathbb{1}_E(\text{unit}) \\
\llbracket \text{sample}_\sigma(d) \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \llbracket d \rrbracket_{\delta; \gamma} \ggg \lambda f. \lambda E. \int_E f d\nu_\sigma \\
\quad \text{where } \nu_\sigma \text{ is the reference measure over } \llbracket \sigma \rrbracket \\
\llbracket \mathbf{f}(e_1, \dots, e_n) \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \llbracket e_1 \rrbracket_{\delta; \gamma} \ggg \lambda v_1. \dots \\
\quad \llbracket e_n \rrbracket_{\delta; \gamma} \ggg \lambda v_n. \psi(\mathbf{f})(v_1, \dots, v_n) \\
\llbracket \mathbf{t} \$ \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \llbracket \mathbf{t} \rrbracket_{\delta; \psi; \gamma} \\
\llbracket \mathbf{x} = \mathbf{t}; \mathbf{m} \rrbracket_{\delta; \psi; \gamma} \stackrel{\text{def}}{=} \lambda E. \\
\quad \int_u \left(u \ggg \lambda v. \llbracket \mathbf{m} \rrbracket_{\delta; \psi; \gamma[x \mapsto v]} \right) (E) \llbracket \mathbf{t} \rrbracket_{\delta; \psi; \gamma} (du) \\
\llbracket \text{def } \mathbf{x} = e \rrbracket_{\delta} \stackrel{\text{def}}{=} \llbracket e \rrbracket_{\delta; \varnothing} \\
\llbracket \text{def } \mathbf{f}(\bar{\mathbf{x}}) = \mathbf{m} \rrbracket_{\delta; \psi} \stackrel{\text{def}}{=} \lambda \bar{v}. \llbracket \mathbf{m} \rrbracket_{\delta; \psi; \{\bar{\mathbf{x}} \mapsto \bar{v}\}} \\
\llbracket \mathcal{G}_1 \dots \mathcal{G}_n \rrbracket_{\delta} \stackrel{\text{def}}{=} \{ \mathbf{g}_1 \mapsto \llbracket \mathcal{G}_1 \rrbracket_{\delta}, \dots, \mathbf{g}_n \mapsto \llbracket \mathcal{G}_n \rrbracket_{\delta} \} \\
\quad \text{where } \text{name}(\mathcal{G}_i) = \mathbf{g}_i \text{ for } i \in \{1, \dots, n\} \\
\llbracket \mathcal{F}_1 \dots \mathcal{F}_n \rrbracket_{\delta; \psi} \stackrel{\text{def}}{=} \{ \mathbf{f}_1 \mapsto \llbracket \mathcal{F}_1 \rrbracket_{\delta; \psi}, \dots, \mathbf{f}_n \mapsto \llbracket \mathcal{F}_n \rrbracket_{\delta; \psi} \} \\
\quad \text{where } \text{name}(\mathcal{F}_i) = \mathbf{f}_i \text{ for } i \in \{1, \dots, n\}
\end{array}$$

Figure 3. Selected definitions of the denotational semantics for MAPPL.

$\llbracket m \rrbracket_{\delta; \psi; \gamma}$ are additionally indexed by a semantic substitution $\psi \in \llbracket \Psi \rrbracket$ for the global variables bound in Ψ . The denotation of $\text{sample}_\sigma(d)$ is a measure over $\llbracket \sigma \rrbracket$, obtained by integrating the density function denoting d with respect to either the Lebesgue measure or the counting measure, depending on whether σ is \mathbb{R} or \mathbb{B} . The denotation of $\mathbf{x} = \mathbf{t}; \mathbf{m}$ composes the denotations of \mathbf{t} and \mathbf{m} , integrating the measure denoting \mathbf{m} with respect to the measure denoting \mathbf{t} . The pure expression $\log \text{ML}(m)$ is denoted by the logarithm of the measure denoting the \cup -typed m on $\llbracket \cup \rrbracket$.

Global definitions are mutually recursive and thus are interpreted under semantic substitutions too. For a program $\bar{\mathcal{G}}; \bar{\mathcal{F}}; m$, the index-free denotations are given by the fixpoints $\llbracket \bar{\mathcal{G}} \rrbracket = \delta_* \stackrel{\text{def}}{=} \text{fix } \delta. \llbracket \bar{\mathcal{G}} \rrbracket_{\delta}$ and $\llbracket \bar{\mathcal{F}} \rrbracket = \psi_* \stackrel{\text{def}}{=} \text{fix } \psi. \llbracket \bar{\mathcal{F}} \rrbracket_{\delta_*; \psi_*}$. We shall write $\llbracket m \rrbracket_{\gamma}$ for $\llbracket m \rrbracket_{\delta_*; \psi_*; \gamma}$; similarly for $\llbracket \mathbf{t} \rrbracket_{\gamma}$ and $\llbracket e \rrbracket_{\gamma}$.

4 INFORMATION-FLOW TYPE SYSTEM

Syntax. Figure 4 shows the syntax of labels, types, and contexts of the information-flow type system. The type system is parameterized by a join semilattice \mathcal{L} of labels, although in this work, we will need only the two-level lattice $\{\mathbb{L}, \mathbb{H}\}$ with $\mathbb{L} \sqsubseteq \mathbb{H}$.

Types are of the form A^ℓ , where A is an unlabeled type that is further constructed from labeled types. Metavariables τ range over (labeled) types. We differentiate them from types τ in the base type system by typesetting labeled types in upright font and in purple. Similarly, we typeset metavariables Δ, Ψ , and Γ differently than Δ, Ψ , and Γ , as Δ, Ψ , and Γ now contain labeled types.

The type system supports label polymorphism, as well as ordering constraints on labels, for functions—that is, the type of a function can be parameterized by label variables η and ordering constraints of the form $\ell_1 \sqsubseteq \ell_2$. Label polymorphism allows more reusable code [45, 44].

Typing the pure fragment. Typing judgments for expressions have the form $\Delta; \Xi; \mathbb{C}; \Gamma \vdash e : \tau$, where Ξ records the label variables in scope and \mathbb{C} the label ordering constraints. These rules are largely standard. In particular, introduction rules (e.g., those for unit, lambdas, pairs, and distributions) do not constrain the label of the expression—the label can be arbitrarily low. A subsumption rule exists to allow weakening (i.e., increasing) the label of an expression. Subtyping rules are standard and thus omitted. For a labeled type A^ℓ , subtyping \leq is covariant in both A and ℓ .

The type system uses fine-grained labeling, in that every type, top-level or nested, is labeled [56]. For example, the type of a pair takes the form $(A_1^{\ell_1} \times A_2^{\ell_2})^{\ell_3}$. The nested labels ℓ_1 and ℓ_2 classify the

labels	$\ell, O ::= \eta \mid l \in \mathcal{L} \mid \ell_1 \sqcup \ell_2$	contexts of label variables	$\Xi ::= \emptyset \mid \Xi, \eta$
unlabeled types	$A, B ::= \mathbb{U} \mid \mathbb{R} \mid \alpha \mid \text{dist}(\tau) \mid \mu\alpha. \tau$ $\mid \forall[\Xi C].\tau_1 \rightarrow \tau_2 \mid \tau_1 + \tau_2 \mid \tau_1 \times \tau_2$	label constraints	$C ::= \emptyset \mid C, \ell_1 \sqsubseteq \ell_2$
labeled types	$\tau ::= A^\ell$	contexts of locals	$\Gamma ::= \emptyset \mid \Gamma, x : \tau$
types of prob. globals	$F ::= \forall[\Xi C].(\tau_1, \dots, \tau_n) \rightarrow \tau_0$	contexts of pure globals	$\Delta ::= \emptyset \mid \Delta, x : \tau$
		contexts of prob. globals	$\Psi ::= \emptyset \mid \Psi, f : F$

$\frac{\Delta; \Xi; C; \Gamma \vdash e : \tau \quad \Xi; C; \mathcal{L} \vdash \tau \leq \tau'}{\Delta; \Xi; C; \Gamma \vdash e : \tau'}$	$\frac{\Delta; \Psi; \Xi; C; \Gamma \vdash t : \tau \quad \Xi; C; \mathcal{L} \vdash \tau \leq \tau'}{\Delta; \Psi; \Xi; C; \Gamma \vdash t : \tau'}$	$\frac{\Delta; \Psi; \Xi; C; \Gamma \vdash m : \tau \quad \Xi; C; \mathcal{L} \vdash \tau \leq \tau'}{\Delta; \Psi; \Xi; C; \Gamma \vdash m : \tau'}$	
$\frac{}{\Delta; \Xi; C; \Gamma \vdash \text{unit} : \mathbb{U}^\ell}$	$\frac{\Gamma(x) = \tau}{\Delta; \Xi; C; \Gamma \vdash x : \tau}$	$\frac{x \notin \text{dom}(\Gamma) \quad \Delta(x) : \tau}{\Delta; \Xi; C; \Gamma \vdash x : \tau}$	$\frac{\Delta; \Xi; C; \Gamma \vdash d : \text{dist}(A^\ell)^\ell \quad \Delta; \Xi; C; \Gamma \vdash e : A^\ell}{\Delta; \Xi; C; \Gamma \vdash \log\text{Pr}(d; e) : \mathbb{R}^\ell}$
$\frac{\Delta; \Xi; C; \Gamma \vdash e : \mathbb{R}^{\ell_1}}{\Delta; \Xi; C; \Gamma \vdash \text{BERN}(e) : \text{dist}(\mathbb{B}^{\ell_1})^{\ell_2}}$	$\frac{\Delta; \Xi; C; \Gamma \vdash e_1 : \mathbb{R}^{\ell_1} \quad \Delta; \Xi; C; \Gamma \vdash e_2 : \mathbb{R}^{\ell_1}}{\Delta; \Xi; C; \Gamma \vdash \text{NORMAL}(e_1, e_2) : \text{dist}(\mathbb{R}^{\ell_1})^{\ell_2}}$	$\frac{\Delta; \emptyset; \Xi; C; \Gamma \vdash m : \mathbb{U}^\ell}{\Delta; \Xi; C; \Gamma \vdash \log\text{ML}(m) : \mathbb{R}^\ell}$	
$\frac{\Delta; \Xi; C; \Gamma \vdash e : \tau}{\Delta; \Psi; \Xi; C; \Gamma \vdash \text{ret}(e) : \tau}$	$\frac{\Delta; \Xi; C; \Gamma \vdash e : \mathbb{R}^\ell}{\Delta; \Psi; \Xi; C; \Gamma \vdash \text{factor}(e) : \mathbb{U}^\ell}$	$\frac{\Delta; \Xi; C; \Gamma \vdash d : \text{dist}(A^\ell)^\ell}{\Delta; \Psi; \Xi; C; \Gamma \vdash \text{sample}(d) : A^\ell}$	
$\frac{\Delta; \Xi; C; \Gamma \vdash e : (\tau_1 + \tau_2)^\ell \quad \Delta; \Psi; \Xi; C; \Gamma, x : \tau_i \vdash m_i : A^\ell \text{ for } i \in \{1, 2\}}{\Delta; \Psi; \Xi; C; \Gamma \vdash \text{case}(e; x.m_1; x.m_2) : A^\ell}$	$\frac{\Psi(f) = \forall[\bar{\eta} C_f].(\tau_1, \dots, \tau_n) \rightarrow \tau_0 \quad \Xi; C; \mathcal{L} \vdash C_f\{\bar{\ell}/\bar{\eta}\} \quad \Delta; \Xi; C; \Gamma \vdash e_i : \tau_i\{\bar{\ell}/\bar{\eta}\} \text{ for } i \in \{1, \dots, n\}}{\Delta; \Psi; \Xi; C; \Gamma \vdash f(e_1, \dots, e_n) : \tau_0\{\bar{\ell}/\bar{\eta}\}}$		
$\frac{\Delta; \Psi; \Xi; C; \Gamma \vdash t : \tau}{\Delta; \Psi; \Xi; C; \Gamma \vdash t \$: \tau}$	$\frac{\Delta; \Psi; \Xi; C; \Gamma \vdash t : A^\ell}{\Delta; \Psi; \Xi; C; \Gamma, x : A^{\ell'} \vdash m : B^\ell}$	$\frac{\Psi(f) = \forall[\Xi C].(\tau_1, \dots, \tau_n) \rightarrow \tau \quad \Delta; \Psi; \Xi; C; x_1 : \tau_1, \dots, x_n : \tau_n \vdash m : \tau}{\Delta; \Psi \vdash \text{def } f(x_1, \dots, x_n) = m : \Psi(f)}$	

Figure 4. Syntax of information-flow types and selected rules of the information-flow type system.

contents of the pair, while the top-level label ℓ_3 classifies the reference to the pair. The distinction enables fine-grained control over the flow of information.

The introduction rules for primitive distributions use the type $\text{dist}(A^{\ell_1})^{\ell_2}$, where A is either \mathbb{B} or \mathbb{R} , ℓ_1 classifies the contents of the distribution (i.e., how the r.v. is distributed), and ℓ_2 classifies the reference to the distribution. For instance, given $x : \mathbb{B}^H$, the expression $\text{BERN}(\text{case}(x; _ .0.7; _ .0.1))$ can be typed at $\text{dist}(\mathbb{B}^H)^\ell$. This fine-grained labeling allows the distribution to be stored in a data structure that can only hold L references, while controlling that when the distribution is eventually retrieved and sampled, the probabilistic effects are classified at H. Whereas coarser-grained type systems trade fine-grained control for reduced label annotation burden, label annotation is not a concern in our setting, because programmers do not specify security policies through labels as they would in a security-typed language. Instead, labels are automatically inferred.

Typing the probabilistic fragment. The design of the type system is guided by the denotational semantics: while the label of an expression e is designed to classify the information the semantic value $\llbracket e \rrbracket_{\delta; \gamma}$ contains, the label of a term t or a command m should classify the information the measure $\llbracket t \rrbracket_{\delta; \psi; \gamma}$ or $\llbracket m \rrbracket_{\delta; \psi; \gamma}$ contains.

Consider typing $\text{sample}(d)$, where d has type $\text{dist}(A^{\ell_1})^{\ell_2}$. Since the contents of the distribution d , as well as the identity of it, determine the measure denoting $\text{sample}(d)$, the term should be typed at a level no lower than $\ell_1 \sqcup \ell_2$. In Figure 4, the typing rule for $\text{sample}(d)$ handles distributions that can be typed at $\text{dist}(A^\ell)^\ell$. This rule suffices, as the type $\text{dist}(A^{\ell_1})^{\ell_2}$ is covariant in both ℓ_1 and ℓ_2 : both labels can be weakened to a label $\ell \sqsupseteq \ell_1 \sqcup \ell_2$ by subsumption. Consider typing $\text{case}(e; x.m_1; x.m_2)$, where e has type $(\tau_1 + \tau_2)^\ell$. Information flows, via a control structure, from e to the measure over the possible outcomes of evaluating the term. So the term should be classified at a level no lower than ℓ .

Typing a call to a probabilistic global function checks, with $\Xi; C; \mathcal{L} \vdash C_f\{\bar{\ell}/\bar{\eta}\}$, that the constraints specified in the function's type (after substitution) are satisfied under the current context.

Now consider typing $x = t; m$. Since the denotation is defined by composing the two measures $\llbracket t \rrbracket$ and $\llbracket m \rrbracket$, the label of $x = t; m$ is required to be no lower than the labels of t and m . It is perhaps surprising that in the typing rule, the label ℓ' of x is not required to be at least as high as the label ℓ of t . An explanation is that denotationally, t merely defines a measure on the possible values x can take; it does not determine the value of x in any given run.

This wrinkle has implications for the precision of the information-flow analysis. Consider the example below. The unnormalized joint density of x and y consists of three factors: $\phi_1(x)\phi_2(x, y)\phi_3(y)$.

$x = \text{sample}(\text{BERN}(.5))$

$y = \text{sample}(\text{BERN}(\text{case}(x; _-.0.7; _-.0.1)))$

$\text{factor}(\text{case}(y; _-.8; _-.2))$

Suppose that we want to marginalize out x . Labeling x at H, we hope to type the third term at L, to justify that it need not be involved in the marginalization of x : $\sum_x \phi_1(x)\phi_2(x, y)\phi_3(y) = \phi_3(y) \sum_x \phi_1(x)\phi_2(x, y)$. The typing rule for $x = t; m$ allows y

to be labeled L, despite that its right-hand side term needs to be typed at H. In contrast, if the typing rule required y to be labeled H, then the third term would also have to be typed at H, which would disallow the third term to be factored out of the sum.

Finally, the expression $\log\text{ML}(m)$ is typed at a label no lower than the command m 's label, since the measure denoting m determines the model evidence of the probabilistic computation m .

Remarks. In the presence of side effects such as mutable state, information-flow type systems often use a program-counter label [22] to lower-bound information leaked through side effects. Indeed, in the imperative while-language of SlicStan [29], typing judgments of commands do use a label to lower-bound the write effects of commands. In MAPPL, however, the label of a command is an upper bound that directly classifies the information that can flow into the command's measure denotation, just as an expression's label upper-bounds the information that can flow into the expression's denotation. In SlicStan, by contrast, factorization must first produce an upper bound by joining the labels of its subexpressions.

5 NONINTERFERENCE VIA A LOGICAL-RELATIONS MODEL

Semantic types. We now establish the soundness of the information-flow type system, by constructing a semantic model of the types. Figure 5 defines our logical-relations model.

The definition uses a function $[\cdot]$ that strips a type of all labels occurring in it; $[\cdot]$ sends a type in the information-flow type system (Section 4) to a type in the base type system (Section 3). The function is overloaded on labeled types τ , unlabeled types A , and contexts Δ, Ψ, Γ .

The main idea behind our model is to interpret each type τ as two binary relations: a value relation $\mathcal{V}\llbracket \tau \rrbracket_{\xi}^{\mathcal{O}}$ for semantic typing of pure expressions, and a measure relation $\mathcal{M}\llbracket \tau \rrbracket_{\xi}^{\mathcal{O}}$ for semantic typing of probabilistic computations. Both relations are parameterized by a label \mathcal{O} that stands for the "security clearance" of an observer, and by a substitution ξ for the label variables occurring free in τ . The model is constructed by first defining, by induction on types, the interpretations $\mathcal{V}\llbracket A \rrbracket_{\xi, \theta}^{\mathcal{O}}$, $\mathcal{V}\llbracket \tau \rrbracket_{\xi, \theta}^{\mathcal{O}}$, and $\mathcal{M}\llbracket \tau \rrbracket_{\xi, \theta}^{\mathcal{O}}$ parameterized by a semantic substitution θ for free type variables, with fixpoints taken in the case of recursive types. The relations $\mathcal{V}\llbracket A \rrbracket_{\xi}^{\mathcal{O}}$, $\mathcal{V}\llbracket \tau \rrbracket_{\xi}^{\mathcal{O}}$, and $\mathcal{M}\llbracket \tau \rrbracket_{\xi}^{\mathcal{O}}$ are then defined on those types without free type variables.

The value relation $\mathcal{V}\llbracket A \rrbracket_{\xi}^{\mathcal{O}}$ relates two semantic values in the semantic domain $\llbracket [A] \rrbracket$ if they are indistinguishable to the observer. What is considered indistinguishable is determined by the observer's label \mathcal{O} and the type A . At a ground type \mathbb{U} or \mathbb{R} , only identical values are related. At a function type, two functions are related if they send related inputs to related outputs. The relations at product, sum, and recursive types are standard as well. The relation at a distribution type $\text{dist}(A^{\ell})$ relates two density functions.

$$\begin{aligned}
\mathcal{V}\llbracket \cup \rrbracket_{\xi;\theta}^O(u_1, u_2) &\stackrel{\text{def}}{=} u_1 = u_2 \in \llbracket \cup \rrbracket \\
\mathcal{V}\llbracket \mathbb{R} \rrbracket_{\xi;\theta}^O(r_1, r_2) &\stackrel{\text{def}}{=} r_1 = r_2 \in \llbracket \mathbb{R} \rrbracket \\
\mathcal{V}\llbracket \text{dist}(A^\ell) \rrbracket_{\xi;\theta}^O(f_1, f_2) &\stackrel{\text{def}}{=} \forall v \in \llbracket [A] \rrbracket. \\
&\quad \mathcal{V}\llbracket \mathbb{R} \rrbracket_{\xi;\theta}^O(f_1(v), f_2(v)) \quad \text{where } A = \mathbb{R} \text{ or } \mathbb{B} \\
\mathcal{V}\llbracket \forall[\Xi|C].\tau \rightarrow \tau' \rrbracket_{\xi;\theta}^O(f_1, f_2) &\stackrel{\text{def}}{=} \forall \xi' \in \llbracket \Xi | \xi C \rrbracket. \\
&\quad \forall (v_1, v_2) \in \mathcal{V}\llbracket \tau \rrbracket_{\xi;\xi';\theta}^O. \mathcal{V}\llbracket \tau' \rrbracket_{\xi;\xi';\theta}^O(f_1(v_1), f_2(v_2)) \\
\mathcal{V}\llbracket \tau + \tau' \rrbracket_{\xi;\theta}^O(v_1, v_2) &\stackrel{\text{def}}{=} \\
&\quad (v_1 = \text{inl } u_1 \wedge v_2 = \text{inl } u_2 \wedge \mathcal{V}\llbracket \tau \rrbracket_{\xi;\theta}^O(u_1, u_2)) \vee \\
&\quad (v_1 = \text{inr } u_1 \wedge v_2 = \text{inr } u_2 \wedge \mathcal{V}\llbracket \tau' \rrbracket_{\xi;\theta}^O(u_1, u_2)) \\
\mathcal{V}\llbracket \tau \times \tau' \rrbracket_{\xi;\theta}^O(v_1, v_2) &\stackrel{\text{def}}{=} v_1 = (u_1, u_1') \wedge v_2 = (u_2, u_2') \wedge \\
&\quad \mathcal{V}\llbracket \tau \rrbracket_{\xi;\theta}^O(u_1, u_1') \wedge \mathcal{V}\llbracket \tau' \rrbracket_{\xi;\theta}^O(u_2, u_2') \\
\mathcal{V}\llbracket \alpha \rrbracket_{\xi;\theta}^O(v_1, v_2) &\stackrel{\text{def}}{=} \theta(\alpha)(v_1, v_2) \\
\mathcal{V}\llbracket \mu\alpha.\tau \rrbracket_{\xi;\theta}^O(v_1, v_2) &\stackrel{\text{def}}{=} \mu R. \mathcal{V}\llbracket \tau \rrbracket_{\xi;\theta[\alpha \rightarrow R]}^O(v_1, v_2) \\
\mathcal{V}\llbracket A^\ell \rrbracket_{\xi;\theta}^O(v_1, v_2) &\stackrel{\text{def}}{=} \begin{cases} \mathcal{V}\llbracket A \rrbracket_{\xi;\theta}^O(v_1, v_2) & \text{if } \xi \ell \sqsubseteq O, \\ v_1, v_2 \in \llbracket [A] \rrbracket & \text{otherwise} \end{cases} \\
\mathcal{E}\llbracket A \rrbracket_{\xi;\theta}^O(E_1, E_2) &\stackrel{\text{def}}{=} E_1, E_2 \in \Sigma\llbracket [A] \rrbracket \wedge \\
&\quad \forall (v_1, v_2) \in \mathcal{V}\llbracket A \rrbracket_{\xi;\theta}^O. v_1 \in E_1 \Leftrightarrow v_2 \in E_2 \\
\mathcal{M}\llbracket A^\ell \rrbracket_{\xi;\theta}^O(\mu_1, \mu_2) &\stackrel{\text{def}}{=} \\
&\quad \begin{cases} \mathcal{E}\llbracket A \rrbracket_{\xi;\theta}^O. \mu_1(E_1) = \mu_2(E_2) & \text{if } \xi \ell \sqsubseteq O, \\ \mu_1, \mu_2 \in \text{Meas } \llbracket [A] \rrbracket & \text{otherwise} \end{cases} \\
\mathcal{V}\llbracket A \rrbracket_{\xi}^O &\stackrel{\text{def}}{=} \mathcal{V}\llbracket A \rrbracket_{\xi;\emptyset}^O & \mathcal{E}\llbracket A \rrbracket_{\xi}^O &\stackrel{\text{def}}{=} \mathcal{E}\llbracket A \rrbracket_{\xi;\emptyset}^O \\
\mathcal{V}\llbracket \tau \rrbracket_{\xi}^O &\stackrel{\text{def}}{=} \mathcal{V}\llbracket \tau \rrbracket_{\xi;\emptyset}^O & \mathcal{M}\llbracket \tau \rrbracket_{\xi}^O &\stackrel{\text{def}}{=} \mathcal{M}\llbracket \tau \rrbracket_{\xi;\emptyset}^O
\end{aligned}$$

$$\begin{aligned}
\mathcal{V}\llbracket A \rrbracket_{\xi}^O(v_1, v_2) &\stackrel{\text{def}}{=} (v_1 = v_2 = \perp) \vee (v_1, v_2) \in \mathcal{V}\llbracket A \rrbracket_{\xi}^O \\
\mathcal{V}\llbracket A^\ell \rrbracket_{\xi}^O(v_1, v_2) &\stackrel{\text{def}}{=} \begin{cases} \mathcal{V}\llbracket A \rrbracket_{\xi}^O(v_1, v_2) & \text{if } \xi \ell \sqsubseteq O, \\ v_1, v_2 \in \llbracket [A] \rrbracket_{\perp} & \text{otherwise} \end{cases} \\
\mathcal{E}\llbracket A \rrbracket_{\xi}^O(E_1, E_2) &\stackrel{\text{def}}{=} E_1, E_2 \in \Sigma\llbracket [A] \rrbracket_{\perp} \wedge \\
&\quad \forall (v_1, v_2) \in \mathcal{V}\llbracket A \rrbracket_{\xi}^O. v_1 \in E_1 \Leftrightarrow v_2 \in E_2 \\
\mathcal{M}\llbracket A^\ell \rrbracket_{\xi}^O(\mu_1, \mu_2) &\stackrel{\text{def}}{=} \\
&\quad \begin{cases} \mathcal{E}\llbracket A \rrbracket_{\xi}^O. \mu_1(E_1) = \mu_2(E_2) & \text{if } \xi \ell \sqsubseteq O, \\ \mu_1, \mu_2 \in \text{Meas } \llbracket [A] \rrbracket_{\perp} & \text{otherwise} \end{cases} \\
\llbracket \Xi | C \rrbracket(\xi) &\stackrel{\text{def}}{=} \text{dom}(\Xi) \subseteq \text{dom}(\xi) \wedge \emptyset; \emptyset; \mathcal{L} \vdash \xi C \\
\llbracket \forall[\Xi|C].\tau_1, \dots, \tau_n \rightarrow \tau \rrbracket(\xi, f_2) &\stackrel{\text{def}}{=} \forall \xi \in \llbracket \Xi | C \rrbracket. \\
&\quad \forall (v_{11}, v_{12}) \in \mathcal{V}\llbracket \tau_1 \rrbracket_{\xi}^O. \dots \forall (v_{n1}, v_{n2}) \in \mathcal{V}\llbracket \tau_n \rrbracket_{\xi}^O. \\
&\quad \mathcal{M}\llbracket \tau \rrbracket_{\xi}^O(f_1(v_{11}, \dots, v_{n1}), f_2(v_{12}, \dots, v_{n2})) \\
\llbracket \Delta \rrbracket(\delta_1, \delta_2) &\stackrel{\text{def}}{=} \forall x \in \text{dom}(\Delta). \mathcal{V}\llbracket \Delta(x) \rrbracket_{\xi}^O(\delta_1(x), \delta_2(x)) \\
\llbracket \Psi \rrbracket(\psi_1, \psi_2) &\stackrel{\text{def}}{=} \forall f \in \text{dom}(\Psi). \llbracket \Psi(f) \rrbracket_{\xi}^O(\psi_1(f), \psi_2(f)) \\
\llbracket \Gamma \rrbracket(\gamma_1, \gamma_2) &\stackrel{\text{def}}{=} \forall x \in \text{dom}(\Gamma). \mathcal{V}\llbracket \Gamma(x) \rrbracket_{\xi}^O(\gamma_1(x), \gamma_2(x)) \\
\Delta; \Xi; C; \Gamma \vDash e_1 \approx e_2 : \tau &\stackrel{\text{def}}{=} \forall O \in \mathcal{L}. \\
&\quad \forall (\delta_1, \delta_2) \in \llbracket \Delta \rrbracket^O. \forall \xi \in \llbracket \Xi | C \rrbracket. \\
&\quad \forall (\gamma_1, \gamma_2) \in \llbracket \Gamma \rrbracket_{\xi}^O. \mathcal{V}\llbracket \tau \rrbracket_{\xi}^O(\llbracket e_1 \rrbracket_{\delta_1; \gamma_1}, \llbracket e_2 \rrbracket_{\delta_2; \gamma_2}) \\
\Delta; \Psi; \Xi; C; \Gamma \vDash t_1 \approx t_2 : \tau &\stackrel{\text{def}}{=} \forall O \in \mathcal{L}. \\
&\quad \forall (\delta_1, \delta_2) \in \llbracket \Delta \rrbracket^O. \forall (\psi_1, \psi_2) \in \llbracket \Psi \rrbracket^O. \forall \xi \in \llbracket \Xi | C \rrbracket. \\
&\quad \forall (\gamma_1, \gamma_2) \in \llbracket \Gamma \rrbracket_{\xi}^O. \mathcal{M}\llbracket \tau \rrbracket_{\xi}^O(\llbracket t_1 \rrbracket_{\delta_1; \psi_1; \gamma_1}, \llbracket t_2 \rrbracket_{\delta_2; \psi_2; \gamma_2}) \\
\Delta; \Psi; \Xi; C; \Gamma \vDash m_1 \approx m_2 : \tau &\stackrel{\text{def}}{=} \forall O \in \mathcal{L}. \\
&\quad \forall (\delta_1, \delta_2) \in \llbracket \Delta \rrbracket^O. \forall (\psi_1, \psi_2) \in \llbracket \Psi \rrbracket^O. \forall \xi \in \llbracket \Xi | C \rrbracket. \\
&\quad \forall (\gamma_1, \gamma_2) \in \llbracket \Gamma \rrbracket_{\xi}^O. \mathcal{M}\llbracket \tau \rrbracket_{\xi}^O(\llbracket m_1 \rrbracket_{\delta_1; \psi_1; \gamma_1}, \llbracket m_2 \rrbracket_{\delta_2; \psi_2; \gamma_2})
\end{aligned}$$

Figure 5. Semantic information-flow types $\llbracket \cdot \rrbracket$ and semantic information-flow typing \vDash .

Having defined the value relation at unlabeled types A , we can define the relation $\mathcal{V}\llbracket A^\ell \rrbracket_{\xi}^O$ at a labeled type A^ℓ , which depends on the labels ℓ and O . If $\xi \ell \sqsubseteq O$ ($\xi \ell$ means applying the substitution function ξ to the label ℓ), the observer is cleared to see values at label $\xi \ell$, so $\mathcal{V}\llbracket A^\ell \rrbracket_{\xi}^O$ contains exactly those values related by $\mathcal{V}\llbracket A \rrbracket_{\xi;\theta}^O$. Otherwise, $\xi \ell \not\sqsubseteq O$ and thus the observer does not have the clearance to see values at $\xi \ell$, so $\mathcal{V}\llbracket A^\ell \rrbracket_{\xi}^O$ is the full relation $\llbracket [A] \rrbracket \times \llbracket [A] \rrbracket$.

The measure relation $\mathcal{M}\llbracket A^\ell \rrbracket_{\xi}^O$ relates two measures in the semantic domain $\text{Meas } \llbracket [A] \rrbracket$ if they are indistinguishable to the observer. Recall that for a probabilistic computation, its label ℓ represents the information contained in the measure denoting it. Hence the following two cases:

- If $\xi \ell \not\sqsubseteq O$, the observer is not classified at a high enough level to differentiate between two measures at $\xi \ell$. So $\mathcal{M}\llbracket A^\ell \rrbracket_{\xi}^O$ is the full relation $\text{Meas } \llbracket [A] \rrbracket \times \text{Meas } \llbracket [A] \rrbracket$.
- Indistinguishability is subtler to define for the case $\xi \ell \sqsubseteq O$. Here, we consider two measures indistinguishable if they agree on related measurable sets. The relation $\mathcal{E}\llbracket A \rrbracket_{\xi}^O$ defines the notion of relatedness for measurable sets. Two measurable sets E_1 and E_2 are related when they are closed to one another under the value relation $\mathcal{V}\llbracket A \rrbracket_{\xi}^O$ —that is, $v_1 \in E_1 \Leftrightarrow v_2 \in E_2$ for all $(v_1, v_2) \in \mathcal{V}\llbracket A \rrbracket_{\xi}^O$.

The relations $\mathcal{V}\llbracket A \rrbracket_{\xi}^O$, $\mathcal{V}\llbracket \tau \rrbracket_{\xi}^O$, $\mathcal{E}\llbracket A \rrbracket_{\xi}^O$, and $\mathcal{M}\llbracket \tau \rrbracket_{\xi}^O$ then lift $\mathcal{V}\llbracket A \rrbracket_{\xi}^O$, $\mathcal{V}\llbracket \tau \rrbracket_{\xi}^O$, $\mathcal{E}\llbracket A \rrbracket_{\xi}^O$, and $\mathcal{M}\llbracket \tau \rrbracket_{\xi}^O$ to account for partiality.

Semantic typing. To define semantic typing, we first define the semantic interpretation of contexts. In particular, two substitutions γ_1 and γ_2 are related at context Γ when for every variable in the domain of Γ , γ_1 and γ_2 map them to related values. And similarly for the interpretation of Δ and Ψ .

Semantic typing judgments have the forms $\Delta; \Xi; C; \Gamma \vdash e_1 \approx e_2 : \tau, \Delta; \Psi; \Xi; C; \Gamma \vdash t_1 \approx t_2 : \tau$, and $\Delta; \Psi; \Xi; C; \Gamma \vdash m_1 \approx m_2 : \tau$. Intuitively, they are defined to hold when the denotations are related under any observer and any related substitutions.

We prove the fundamental property of logical relations, which states that syntactically well-typed expressions, terms, and commands are semantically well-typed.

THEOREM 5.1 (FUNDAMENTAL PROPERTY).

- (1) $\Delta; \Xi; C; \Gamma \vdash e : \tau$ implies $\Delta; \Xi; C; \Gamma \vDash e \approx e : \tau$.
- (2) $\Delta; \Psi; \Xi; C; \Gamma \vdash t : \tau$ implies $\Delta; \Psi; \Xi; C; \Gamma \vDash t \approx t : \tau$.
- (3) $\Delta; \Psi; \Xi; C; \Gamma \vdash m : \tau$ implies $\Delta; \Psi; \Xi; C; \Gamma \vDash m \approx m : \tau$.

The proof is by induction on the (syntactic) typing derivations, with each case proving that semantic typing is compatible with some syntactic typing rule. As with a typical logical-relations proof, the challenge is in setting up the logical-relations model (think of them as induction hypotheses) and the proof is routine. Of special note about the proof is that it involves showing that two integrals are equal $\int f_1(x) \mu_1(dx) = \int f_2(x) \mu_2(dx)$ when the integrands and the measures are not equal point-wise. Lemma 5.2 is a convenient result [18] that enables proving equivalence using a *coarser* structure than point-wise equality.

LEMMA 5.2 (COARSENING). *Let (X, Σ_X) be a measurable space, μ_1, μ_2 be measures on X , and $f_1, f_2 : X \rightarrow \mathbb{R}^+$ be measurable functions. Let $R \subseteq \Sigma_X \times \Sigma_X$ be a binary relation on measurable sets. If (1) μ_1 and μ_2 agree on R -related sets, i.e., $\mu_1(E_1) = \mu_2(E_2)$ for all $(E_1, E_2) \in R$, and (2) if f_1 and f_2 have R -related preimages, i.e., $(f_1^{-1}(S), f_2^{-1}(S)) \in R$ for all $S \in \Sigma_{\mathbb{R}}$, then $\int f_1(x) \mu_1(dx) = \int f_2(x) \mu_2(dx)$.*

The lemma allows proving integrals equal by picking a suitable relation R on measurable sets, for which the relations $\mathcal{E}_{\perp} \llbracket A \rrbracket_{\xi}^O$ on measurable sets will fit the bill.

Noninterference. Noninterference follows from the fundamental property. It guarantees that the measure denoting a L-typed term behaves irrespective of the H-labeled variables in the context.

THEOREM 5.3 (NONINTERFERENCE). *Let Γ_H be a context that binds only H-labeled variables—that is, for all $\mathbf{x} \in \text{dom}(\Gamma_H)$, there is some A such that $\Gamma_H(\mathbf{x}) = A^H$. Let $f : \llbracket \tau \rrbracket_{\perp} \rightarrow \mathbb{R}^+$ be a measurable function. If $\vdash \tau : L, \vdash \Gamma_L : L$, and $\Delta; \Psi; \emptyset; \emptyset; \Gamma_H; \Gamma_L \vdash t : \tau$, then for all $\gamma_1, \gamma_2 \in \llbracket \llbracket \Gamma_H \rrbracket \rrbracket$ and $\gamma_L \in \llbracket \llbracket \Gamma_L \rrbracket \rrbracket$,*

$$\int f(x) \llbracket t \rrbracket_{\gamma_1, \gamma_L}(dx) = \int f(x) \llbracket t \rrbracket_{\gamma_2, \gamma_L}(dx).$$

Here, $\vdash \tau : L$ is defined to mean that all labels occurring in τ , including those nested labels, are L. And $\vdash \Gamma : L$ means that $\vdash \Gamma(z) : L$ for all $z \in \text{dom}(\Gamma)$. It is not a sufficient condition that the *outermost* label of τ and those in Γ are L. For example, the typing $\Psi; \Delta; \emptyset; \emptyset; \mathbf{x} : \mathbb{R}^H \vdash \text{ret}(\langle \mathbf{x}, \mathbf{x} \rangle) : (\mathbb{R}^H \times \mathbb{R}^H)^L$ is valid, but it would be absurd if it implied that $\langle \mathbf{x}, \mathbf{x} \rangle$ behaved irrespective of \mathbf{x} . Similarly, $\mathbf{x} : \mathbb{R}^H, \mathbf{y} : (\mathbb{R}^H \rightarrow \mathbb{R}^L)^L \vdash \text{ret}(\mathbf{y} \mathbf{x}) : \mathbb{R}^L$ is valid typing, but it would be absurd if it implied that the application of \mathbf{y} to \mathbf{x} behaved irrespective of \mathbf{x} for a γ_L such that $\gamma_L(\mathbf{y}) = \lambda x. x \in \llbracket \llbracket (\mathbb{R}^H \rightarrow \mathbb{R}^L)^L \rrbracket \rrbracket$. Theorem 5.3 is a corollary of a more general version of noninterference (given in an appendix [39]) that relaxes the condition $\vdash \Gamma_L : L$ and allows the integrands on the two sides of the equation to be different.

6 VARIABLE-ELIMINATION TRANSFORMATION

Main idea. To generate a pure program, the transformation must compile away all probabilistic-fragment constructs. In particular, (1) it eliminates random variables (r.v.s) by summation or

<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{T}[\mathcal{P}] = \mathcal{P}'$</div> $\frac{\mathbb{K}[m](\lambda x. 0) = e \quad \mathbb{T}[\mathcal{F}] = \mathcal{G}' \text{ for } \mathcal{F} \in \overline{\mathcal{F}} \quad \text{def } \text{logsumexp} = \lambda x. \text{log}(\text{exp}(x \text{ true}) + \text{exp}(x \text{ false})) \in \overline{\mathcal{G}}}{\mathbb{T}[\overline{\mathcal{G}}; \overline{\mathcal{F}}; m] = \overline{\mathcal{G}}, \overline{\mathcal{G}}'; \emptyset; \text{ret}(e)}$	<p style="text-align: center;"><i>CPS-translate a probabilistic program to a pure program</i></p>
<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{T}[\mathcal{F}] = \mathcal{G}$</div> <p style="text-align: center;"><i>CPS-translate a probabilistic global function to a pure one</i></p> $\frac{\mathbb{K}[m] k = e}{\mathbb{T}[\text{def } f(x_1, \dots, x_n) = m] = \text{def } f = \lambda k. \lambda x_1. \dots. \lambda x_n. e}$	<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{C}[m] = e$</div> <p style="text-align: center;"><i>Eliminate probabilistic-level control flow (calls & branching) with CPS</i></p> $\frac{\Delta; \Psi; \gamma : A^H; \Gamma \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L \quad \vdash \Gamma : L \quad \mathbb{C}[m_H] = e_H}{\mathbb{C}[\overline{x} = \bar{t}; m_L; \text{factor}(f(\lambda y. e_H) \bar{e})] = e'}$ $\frac{\mathbb{C}[\overline{x} = \bar{t}; \gamma = f(\bar{e}); m] = e'}{\Delta; \Psi; \gamma : A^H; \Gamma \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L \quad \vdash \Gamma : L \quad \mathbb{C}[m_H] = e_H}$ $\frac{\mathbb{K}[m_i](\lambda y. e_H) = e_i \text{ for } i \in \{1, 2\}}{\mathbb{C}[\overline{x} = \bar{t}; m_L; \text{factor}(\text{case}(e; z.e_1; z.e_2))] = e'}$ $\frac{\mathbb{C}[\overline{x} = \bar{t}; \gamma = \text{case}(e; z.m_1; z.m_2); m] = e'}{\forall i \in \{1, \dots, n\}. t_i \neq f(\dots) \wedge t_i \neq \text{case}(\dots)}$ $\frac{\mathbb{R}[x_1 = t_1; \dots; x_{n-1} = t_{n-1}; t_n] = e}{\mathbb{C}[x_1 = t_1; \dots; x_{n-1} = t_{n-1}; t_n] = e}$
<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{K}[m] e_k = e$</div> <p style="text-align: center;"><i>CPS-translate a command to a pure expression</i></p> $\frac{\mathbb{D}^+[\emptyset \vdash x_1 = t_1; \dots; x_n = t_n; \text{factor}(e_k x_n)] = e}{\mathbb{K}[x_1 = t_1; \dots; x_{n-1} = t_{n-1}; t_n] e_k = e}$	<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{R}[m] = e$</div> <p style="text-align: center;"><i>Transform away any remaining probabilistic-level terms</i></p> $\frac{\mathbb{R}[m] = e'}{\mathbb{R}[\text{ret}(e)] = \text{let}(e; _0) \quad \mathbb{R}[x = \text{ret}(e); m] = \text{let}(e; x.e')}$ $\frac{\mathbb{R}[m\{\text{unit}/x\}] = e'}{\mathbb{R}[\text{factor}(e)] = e \quad \mathbb{R}[x = \text{factor}(e); m] = e + e'}$ $\frac{\Delta; \Psi; \gamma : \mathbb{R}^H; \Gamma \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L \quad \vdash \Gamma : L}{\mathbb{R}[m_L; \text{factor}(\text{logML}(\gamma = \text{sample}_{\mathbb{R}}(d); m_H))] = e}$ $\frac{\mathbb{R}[\gamma = \text{sample}_{\mathbb{R}}(d); m] = e}{\mathbb{R}[m] = e}$
<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{D}^+[\overline{z} \vdash m] = e$</div> <p style="text-align: center;"><i>Accumulate discrete r.v.s into worklist \overline{z}</i></p> $\frac{\mathbb{D}^+[\overline{z}, \gamma \vdash \overline{x} = \bar{t}; \text{factor}(\text{logPr}_{\mathbb{B}}(d; \gamma)); m] = e}{\mathbb{D}^+[\overline{z} \vdash \overline{x} = \bar{t}; \gamma = \text{sample}_{\mathbb{B}}(d); m] = e}$ $\frac{\forall i \in \{1, \dots, n\}. t_i \neq \text{sample}_{\mathbb{B}}(d)}{\mathbb{D}^-[\overline{z} \vdash x_1 = t_1; \dots; x_n = t_n; t_{n+1}] = e}$ $\frac{\mathbb{D}^-[\overline{z} \vdash x_1 = t_1; \dots; x_n = t_n; t_{n+1}] = e}{\mathbb{D}^+[\overline{z} \vdash x_1 = t_1; \dots; x_n = t_n; t_{n+1}] = e}$	<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\mathbb{R}[m] = e$</div> <p style="text-align: center;"><i>Eliminate discrete r.v.s in worklist \overline{z}</i></p> $\frac{\Delta; \Psi; \gamma : \mathbb{B}^H; \Gamma \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L \quad \vdash \Gamma : L \quad \mathbb{C}[m_H] = e_H}{\mathbb{D}^-[\overline{z} \vdash m_L; \text{factor}(\text{logsumexp}(\lambda y. e_H))] = e}$ $\frac{\mathbb{D}^-[\overline{z}, \gamma \vdash m] = e}{\mathbb{C}[m] = e}$ $\frac{\mathbb{C}[m] = e}{\mathbb{D}^-[\emptyset \vdash m] = e}$
<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;">$\Delta; \Psi; \Gamma_H; \Gamma_L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L$</div> <p style="text-align: center;"><i>Factorize m into m_H and m_L with information-flow typing</i></p> <p style="text-align: center;"><i>Invariants: $\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L, \Gamma' \vdash m_H : \mathbb{U}^H$ and $\Delta; \Psi; \emptyset; \emptyset; \Gamma_L \vdash m_L : \mathbb{U}^L$</i></p> $\frac{\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L \vdash t : \mathbb{U}^L \quad \text{Canonicalize}(t, \Gamma_H) = t'}{\Delta; \Psi; \Gamma_H; \Gamma_L \vdash t \rightsquigarrow_{\emptyset} \text{ret}(\text{unit}) * t'}$ $\frac{\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L \vdash t : A^L \quad \vdash A : L \quad \text{Canonicalize}(t, \Gamma_H) = t'}{\Delta; \Psi; \Gamma_H; \Gamma_L, x : A^L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L}$ $\frac{\Delta; \Psi; \Gamma_H; \Gamma_L, x : A^L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L \quad t = \text{sample}_{\mathbb{R}}(d) \Rightarrow x \notin \text{FV}(m_H)}{\Delta; \Psi; \Gamma_H; \Gamma_L \vdash x = t; m \rightsquigarrow_{\Gamma', x: A^L} m_H * (x = t'; m_L)}$	$\frac{\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L \vdash t : \mathbb{U}^H}{\Delta; \Psi; \Gamma_H; \Gamma_L \vdash t \rightsquigarrow_{\emptyset} t * \text{ret}(\text{unit})}$ $\frac{\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L \vdash t : A^L}{\Delta; \Psi; \Gamma_H, x : A^H; \Gamma_L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L}$ $\frac{\Delta; \Psi; \Gamma_H; \Gamma_L \vdash x = t; m \rightsquigarrow_{\Gamma'} (x = t; m_H) * m_L}{\Delta; \Psi; \Gamma_H; \Gamma_L \vdash x = t; m \rightsquigarrow_{\Gamma', x: A^L} m_H * (x = t'; m_L)}$

Figure 6. Variable-elimination transformation. Typing contexts are omitted for brevity (except in factorization judgments). A version of the transformation with complete context information is given in an appendix.

$$\begin{aligned}
\mathbb{T}[\text{def } \text{hmm}(z_0, \text{data}) = m] &\stackrel{\textcircled{1}}{=} \text{def } \text{hmm} = \lambda k. \lambda z_0. \lambda \text{data}. \mathbb{K}[m] k \\
\mathbb{K}[m] k &= \mathbb{K}[\text{case } \text{data} \dots] k \stackrel{\textcircled{2}}{=} \mathbb{D}^+ \left[\emptyset \vdash \frac{y = \text{case } \text{data} \dots}{\text{factor}(k(y))} \right] \stackrel{\textcircled{3}}{=} \mathbb{D}^- \left[\emptyset \vdash \frac{y = \text{case } \text{data} \dots}{\text{factor}(k(y))} \right] \stackrel{\textcircled{4}}{=} \mathbb{C} \left[\frac{y = \text{case } \text{data} \dots}{\text{factor}(k(y))} \right] \\
&\stackrel{\textcircled{5}}{=} \mathbb{C} \left[\frac{\text{factor}(\text{case } \text{data} \text{ of} \right. \\
&\quad \left. \begin{array}{l} \text{nil} \Rightarrow \mathbb{K}[m_{\text{nil}}] (\lambda y. \mathbb{C}[\text{factor}(k(y))]) \\ \text{cons } x \text{ } xs \Rightarrow \mathbb{K}[m_{\text{cons}}] (\lambda y. \mathbb{C}[\text{factor}(k(y))]) \\ \text{end}) \end{array} \right] \stackrel{\textcircled{6}}{=} \mathbb{C} \left[\frac{\text{case } \text{data} \text{ of} \right. \\
&\quad \left. \begin{array}{l} \text{nil} \Rightarrow \mathbb{K}[m_{\text{nil}}] k \\ \text{cons } x \text{ } xs \Rightarrow \mathbb{K}[m_{\text{cons}}] k \\ \text{end} \end{array} \right] \\
\mathbb{K}[m_{\text{nil}}] k &= \mathbb{K}[\text{ret}(z_0)] k \stackrel{\textcircled{7}}{=} \mathbb{D}^+ [\emptyset \vdash \text{factor}(k(z_0))] \stackrel{\textcircled{8}}{=} \mathbb{D}^- [\emptyset \vdash \text{factor}(k(z_0))] \stackrel{\textcircled{9}}{=} \mathbb{C}[\text{factor}(k(z_0))] \stackrel{\textcircled{10}}{=} k(z_0) \\
\mathbb{K}[m_{\text{cons}}] k &= \mathbb{K} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] k \stackrel{\textcircled{11}}{=} \mathbb{D}^+ \left[\emptyset \vdash \frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \stackrel{\textcircled{12}}{=} \mathbb{D}^+ \left[y \vdash \frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \\
&\quad \left[\emptyset \vdash \frac{y = \text{sample}(\text{step}(z))}{\text{factor}(k(y))} \right] \stackrel{\textcircled{13}}{=} \mathbb{D}^- \left[y \vdash \frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \stackrel{\textcircled{14}}{=} \mathbb{D}^- \left[\emptyset \vdash \frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \stackrel{\textcircled{15}}{=} \mathbb{C} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \\
&\quad \left[y \vdash \frac{\text{factor}(\log\text{Pr}(\text{step}(z); y))}{\text{factor}(k(y))} \right] \stackrel{\textcircled{16}}{=} \mathbb{C} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \stackrel{\textcircled{17}}{=} \mathbb{C} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \\
&\quad \left[\frac{\text{factor}(\log\text{sumexp}(\lambda y. \text{factor}(\log\text{Pr}(\text{step}(z); y)))}{\text{factor}(k(y))} \right] \stackrel{\textcircled{18}}{=} \mathbb{C} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \\
&\quad \left[\frac{\text{factor}(\log\text{Pr}(\text{step}(z); y) + k(y))}{\text{factor}(k(y))} \right] \stackrel{\textcircled{19}}{=} \mathbb{C} \left[\frac{z = \text{hmm}(z_0, xs)}{\text{observe}(\text{emit}(z); x)} \right] \\
&\quad \left[\frac{\text{factor}(\text{hmm}(\lambda z. \mathbb{C}[\text{observe}(\text{emit}(z); x)] \text{factor}(\log\text{sumexp}(\lambda y. \log\text{Pr}(\text{step}(z); y) + k(y))))}{\text{factor}(\log\text{sumexp}(\lambda y. \log\text{Pr}(\text{step}(z); y) + k(y)))} \right] \stackrel{\textcircled{20}}{=} \mathbb{C} \left[\frac{\text{hmm}(\lambda z. \log\text{Pr}(\text{emit}(z); x) + \log\text{sumexp}(\lambda y. \log\text{Pr}(\text{step}(z); y) + k(y)))}{\text{factor}(\log\text{sumexp}(\lambda y. \log\text{Pr}(\text{step}(z); y) + k(y)))} \right]
\end{aligned}$$

Figure 7. Compiling the `hmm` function in Figure 1a to that in Figure 1b. The calculation largely follows the rules in Figure 6. To simplify presentation, we use the equality $\mathbb{C}[\text{factor}(e_1); \dots; \text{factor}(e_n)] = e_1 + \dots + e_n$ and standard λ -calculus conversions without detailing the intermediate steps.

integration, creating inference subproblems as a result, and (2) it compiles stochastic control flow to deterministic control flow in continuation-passing style. In both cases, the transformation uses information-flow typing to factorize a command into a H partition and a L partition. The measure denotation of the L partition is guaranteed to be independent of the H-labeled variable—be it bound to a sample, case, or a call term—being eliminated. So the L partition can be factored out of the summation, integration, or continuation indexed by the H-labeled variable.

Program transformation. Figure 6 formalizes the VE transformation as a set of mutually recursive translation functions. The translation is defined for well typed programs (with respect to the base type system in Section 3), so it additionally takes typing contexts as input, but for brevity, we omit them in Figure 6; a version with complete context information can be found in an appendix. As a running example, the step-by-step translation of `hmm` is shown in Figure 7. We now describe the translation rules in Figure 6, referring to steps in Figure 7 as concrete instantiations of the rules.

$\mathbb{T}[\mathcal{F}]$ translates a probabilistic function \mathcal{F} to a pure function additionally parameterized by a continuation k , with the body m of \mathcal{F} translated as $\mathbb{K}[m] k$ (1).

$\mathbb{K}[m] e_k$ translates a command m given a continuation e_k , which is the log-factor dependent on the return value of m . At the top level, the main command m is translated with the top-level continuation $\lambda x. 0$. $\mathbb{K}[m] e_k$ works by applying e_k to m 's return value, appending the resulting factor to m , and then uses \mathbb{D}^+ to further translate the resulting \mathbb{U} -typed command (2)(7)(1).

While \mathbb{K} can be applied to any well typed command, going forward, the other translations (\mathbb{D}^+ , \mathbb{D}^- , \mathbb{C} , and \mathbb{R}) are only defined on commands of the unit type \mathbb{U} , as the commands will already have been CPS-translated.

D^+ does preparatory work for eliminating discrete r.v.s, of which the real work is done by D^- . D^+ accumulates bindings of discrete r.v.s into a worklist, turning all $\text{sample}_{\mathbb{B}}$ terms into factor terms (⑫), so that D^- can eliminate the discrete r.v.s in the worklist in one go without having to worry about factorization possibly creating unbound references to the variables.

$D^- \llbracket \bar{z} \vdash m \rrbracket$ eliminates from m the discrete r.v.s stored in the worklist \bar{z} , one at a time, until the worklist is empty (⑭). Elimination of a variable y from a command m involves factorizing m into m_H and m_L , using the \rightsquigarrow judgment to be defined shortly. D^- eliminates y by summing over y the factor $\lambda y. e_H$ contributed by m_H . Importantly, m_L can be left out of this sum, because it is guaranteed that the measure denotation of m_L is independent of the H-labeled y . The formalized translation does not make memoization explicit, but solutions to this sum may be memoized, with the memo table indexed by the free (discrete) variables in the sum $\text{logsumexp}(\lambda y. e_H)$. These free variables can be thought of as the Markov blanket [49] of y , conditioned on which all other variables are uncorrelated with y . The order in which the variables are eliminated is left unspecified; it is a well studied, orthogonal problem. It is NP-hard to find the optimal ordering that has an elimination width equal to the tree width [20]. In practice, heuristics (e.g., eliminating variables with fewer neighbors first) are effective in giving good orderings with low elimination widths.

C eliminates probabilistic-level control flow, namely call and branching terms. To translate a function call $y = f(\bar{e})$, the variable y is labeled H, and the continuation m to the function call is factorized into m_H and m_L . The factor $\lambda y. e_H$, representing the probabilistic effects of m_H , is passed to the translated pure function f as the continuation argument. The pure function call $f(\lambda y. e_H) \bar{e}$ returns a \mathbb{R} -valued factor, which is then appended to the command to be translated further (⑯). Notice that the measure denotation of m_L is independent of the return value y . So m_L can be left out of the continuation passed to the CPS-translated function call to f .

C translates a branching term $y = \text{case}(e; x.m_1; x.m_2)$ in a similar manner. It factorizes the probabilistic continuation into two partitions, constructs a pure continuation using the H partition only, and CPS-translates the branches by passing this pure continuation (⑤).

When there is no more control-flow terms to eliminate, R takes over to eliminate any remaining probabilistic-level terms, namely $\text{sample}_{\mathbb{R}}$ and factor . A $\text{sample}_{\mathbb{R}}$ term is eliminated by integration (i.e., applying logML), which can be any marginal-inference method of choice. Solutions to this integration may be memoized, with the memo table indexed by the free (discrete) variables in the integral. Unlike $\text{sample}_{\mathbb{B}}$ terms, $\text{sample}_{\mathbb{R}}$ terms are not first converted to factor terms before being eliminated, as integration in general requires *sampling* from prior distributions of continuous r.v.s.

Like in D^- and C , the probabilistic continuation is factorized in R to allow irrelevant terms to be left out of the integration. Unlike in D^- or C , continuous r.v.s are eliminated not always one at a time, but likely simultaneously, to avoid creating unnecessary nested integrals. For instance, the command $m \stackrel{\text{def}}{=} x = \text{sample}(\text{NORMAL}(0, 1)); y = \text{sample}(\text{NORMAL}(3, 1)); \text{factor}(\text{logPr}(\text{NORMAL}(x^2 + y^2, 1); 6))$ is translated as $R[m] = R[\text{factor}(\text{logML}(m))] = \text{logML}(m)$. Factorization makes sure that all three terms belong to the same H partition, despite that information-flow typing does not demand so, so that x and y can be marginalized out simultaneously.

Command factorization via information-flow typing. The judgment $\Delta; \Psi; \Gamma_H; \Gamma_L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L$ factorizes a command m into m_H and m_L . Here, Γ' is the variable bindings declared in m_L and available for use in m_H , such that $\Delta; \Psi; \emptyset; \emptyset; \Gamma_L \vdash m_L : \mathbb{U}^L$ and $\Delta; \Psi; \emptyset; \emptyset; \Gamma_H, \Gamma_L, \Gamma' \vdash m_H : \mathbb{U}^H$. In broad strokes, factorization works by partitioning all H-labeled terms in the command to m_H and all L-labeled terms to m_L (⑤⑭⑯). Label inference is implemented by solving unification constraints.

Step ⑭ performs factorization to eliminate y . With y labeled H, the two factor terms must be typed at H, while the other two terms can be typed at L and left out of the sum over y . Step ⑯ performs factorization to eliminate the recursive call $z = \text{hmm}(z_0, xs)$. With z labeled H, both terms

in the continuation must be typed at H, so their translation is passed to the pure `hmm` as the continuation argument. Step ⑤ performs factorization to eliminate $y = \text{case data } \dots$, which is similar.

Because noninterference guarantees that a L-labeled term behaves irrespective of the value of the H-labeled variables in Γ_H , factorization additionally *canonicalizes* a L-labeled term: it replaces the H-labeled variables with default values. This substitution ensures that m_L does not refer to variables bound in Γ_H and thus is well typed under the context Γ_L . Finally, factorization specially treats `sampleR` terms: it partitions $x = \text{sample}_R(d)$ into m_L only if x is not needed in the m_H , thus scoping the integral over x to one partition only and avoiding unnecessary nested integrals.

LEMMA 6.1. *Let Γ_H be a context that binds only H-labeled variables and Γ_L be a context such that $\vdash \Gamma_L : L$. If $\Delta; \Psi; \Gamma_H; \Gamma_L \vdash m \rightsquigarrow_{\Gamma'} m_H * m_L$, then for any $\gamma_H \in \llbracket \Gamma_H \rrbracket$ and $\gamma_L \in \llbracket \Gamma_L \rrbracket$, it holds that*

$$\llbracket m \rrbracket_{\gamma_H, \gamma_L} (\llbracket \cup \rrbracket) = \oint_{\llbracket \Gamma' \rrbracket} \llbracket m_H \rrbracket_{\gamma_H, \gamma_L, \gamma'} (\llbracket \cup \rrbracket) \llbracket m_L \rrbracket_{\gamma_L} (d\gamma').$$

Lemma 6.1 assures that factorizing a command m produces two partitions that together preserve the semantics of m . The notation $\oint_{\llbracket \Gamma' \rrbracket} f(\gamma') \llbracket m_L \rrbracket_{\gamma_L} (d\gamma')$, defined in an appendix with the proof, is a shorthand for the multiple integral with respect to the measures each denoting an $x_i = t_i$ in m_L . The lemma is a consequence of noninterference (Theorem 5.3).

Correctness. We prove that the variable-elimination transformation is correct. The theorem states that the transformed pure expression, when it terminates, computes the log model evidence of the original probabilistic program. As expected, the proof depends on Lemma 6.1.

THEOREM 6.2. *Let $\overline{\mathcal{G}}; \overline{\mathcal{F}}; m$ be a well typed program where the main command m has type τ . If $\mathbb{T} \llbracket \overline{\mathcal{G}}; \overline{\mathcal{F}}; m \rrbracket = \overline{\mathcal{G}}; \overline{\mathcal{G}}'; \emptyset; \text{ret}(e)$ and $\llbracket e \rrbracket_{\emptyset} \neq \perp$, then $\log \llbracket m \rrbracket_{\emptyset} (\llbracket \tau \rrbracket) = \llbracket e \rrbracket_{\emptyset}$.*

7 EXPERIMENTAL EVALUATION

Scalability of VE compilation. We compare MAPPL and SlicStan [29]. Stan [13], while a popular PPL, does not support discrete parameters. In response, SlicStan features a state-of-the-art VE compiler that performs information-flow analysis and emits variable-eliminated Stan code. As a benchmark, we consider a simple HMM, for which both compilers can generate code whose running time scales linearly with the length of the observed sequence. But compilation time differs.

Figure 8 shows how compilation time scales as the size of the inference problem increases. (All experiments in Section 7 were run on a server with a 3.6GHz CPU and 12GB of RAM.) In SlicStan, models such as HMMs are expressed by unrolling recursion into a fixed number of iterations, so it is expected that compilation time increases as the size of the inference problem increases. Figure 8 confirms this behavior and further shows that SlicStan struggles with large problem sizes: compiling the model conditioned on generating a sequence of length 60 takes over 30 minutes. In contrast, because MAPPL can express the HMM as a recursive program, the compilation time is constant with respect to the problem size.

We also report the time MAPPL takes to compile a version of the HMM with recursion unrolled. Figure 8 (MAPPL*) suggests that the MAPPL compiler exhibits better scalability than SlicStan on the same unrolled model. A probable reason for this speedup is that MAPPL uses a simple two-level lattice in the information-flow analysis, whereas SlicStan uses a meet semilattice, which, as discussed by Gorinova et al. [29], hinders efficient constraint solving.

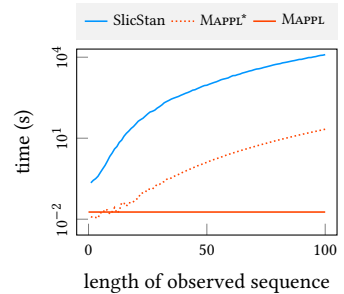


Figure 8. Scaling of compilation time.

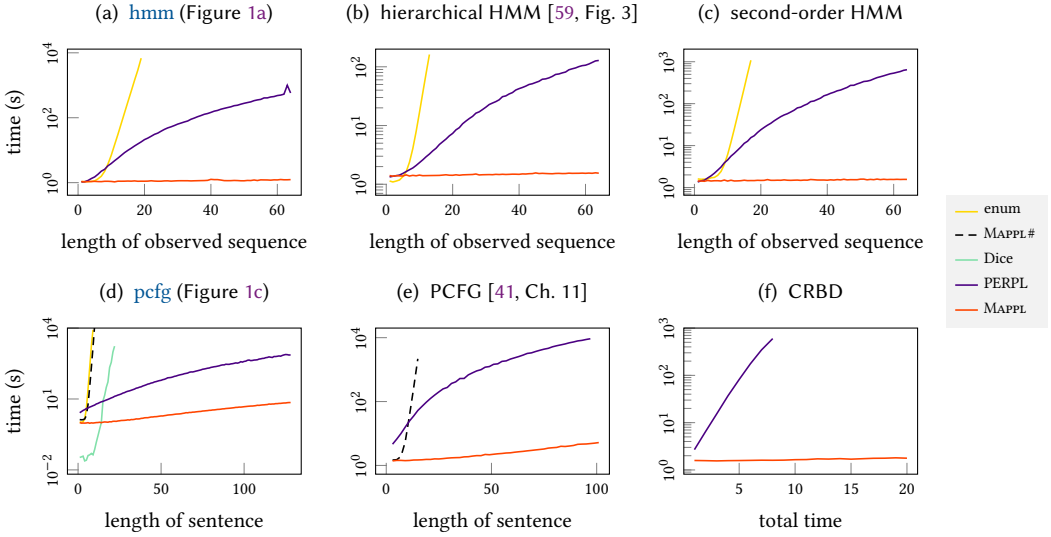


Figure 9. Scaling plots comparing exact-inference methods on recursive programs.

Scalability of exact inference: HMMs, PCFGs, and CRBD. We compare MAPPL and PERPL [15] on recursive programs. PERPL represents a state-of-the-art approach to exact inference for recursive programs, compiling them to factor graph grammars [16] and then to systems of equations.

The benchmarks are the HMM in Figure 1a, a hierarchical HMM, a second-order HMM, the PCFG in Figure 1c, a PCFG with 6 nonterminals and 12 productions, and a discrete-time phylogenetic model. The phylogenetic model generates phylogenetic trees under the constant-rate birth–death (CRBD) assumption. This CRBD model is similar to the PCFGs in that it uses recursion (as opposed to iteration) and exhibits stochastic control flow [58].

Figure 9 shows how the inference running time scales as the size of the inference problem increases. Compilation time is not measured, as it does not vary with the problem size for either MAPPL or PERPL. As PERPL uses a Python back end, to allow a fair comparison, compiled MAPPL programs are further compiled to Python. We use enumeration-based exact inference implemented in Pyro [10] as an additional baseline on some benchmarks; it leads to exponentially increasing running time and runs out of memory quickly on all benchmarks.

On the two HMMs, PERPL leads to running time superlinear in the problem size, whereas MAPPL recovers the linear-time forward algorithm for HMMs. For an observed sequence of length 30, PERPL inference takes over 1 minute, while MAPPL inference takes 1.5 seconds. On the two PCFGs and the CRBD model, PERPL also scales less favorably than MAPPL. We note that PERPL supports unbounded recursion and thus allows the PCFG and CRBD models to be specified in a more declarative way. For example, the CRBD model in PERPL, though complicated by PERPL’s linearity restriction, uses an almost surely terminating function to generate the waiting time until the next speciation or extinction event, whereas the MAPPL version uses a geometric distribution truncated at the remaining time steps to ensure termination.

We also assess, with the PCFGs, the performance implications of the information-flow analysis. Specifically, we evaluate the performance of a version of MAPPL with command factorization disabled (MAPPL#). Disabling factorization means that the VE compilation has to assume correlation between the subparses of a nonterminal, thereby hindering the discovery of recurring substructure

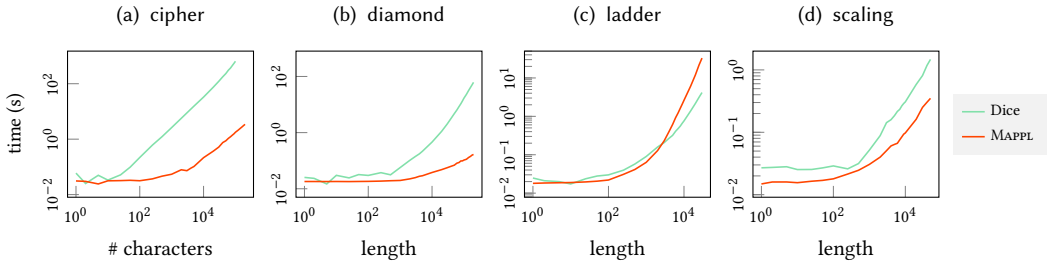


Figure 10. Scaling plots comparing Dice and MAPPL on benchmarks from Holtzen et al. [31, Fig. 10].

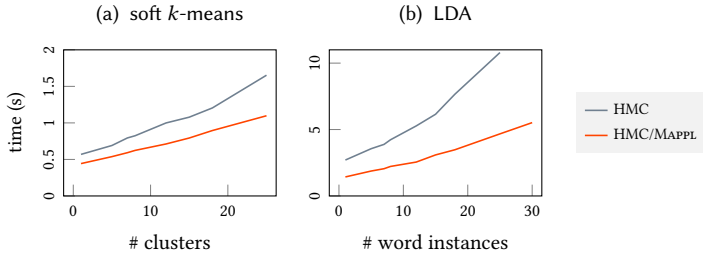


Figure 11. Ahead-of-time VE compilation using MAPPL speeds up HMC in NumPyro. Time per chain is reported.

amenable to dynamic programming. Figures 9d and 9e confirm that without factorization, VE-based inference is intractable.

Scalability of exact inference: Dice benchmarks. We compare MAPPL and Dice [31]. Dice is a state-of-the-art approach to exact inference for discrete, nonrecursive programs. It casts inference to weighted model counting (WMC) on binary decision diagrams (BDDs), exploiting independence structure in programs to create compact BDDs for factorized inference. We use benchmarks [31, Fig. 10] on which Dice has been shown to demonstrate superior scalability over other PPLs that support exact inference. As Dice uses a C library for WMC on BDDs, to allow a fair comparison, compiled MAPPL programs are further compiled to Rust. This Rust back end of MAPPL is not yet full-featured but is sufficient for these Dice benchmarks.

Figure 10 shows how the running time scales as the problem size increases. Given that the Dice running time reported by Holtzen et al. [31] includes the time required for BDD generation, the MAPPL running time reported here includes that for VE compilation. MAPPL is competitive with Dice on these scaling benchmarks, in fact outperforming Dice in three out of four cases. A possible explanation is that since MAPPL can express these benchmarks as recursive programs, compilation time does not increase with the problem size.

We also run Dice on a PCFG. Dice does not support recursion, so we follow the recipe of Chiang et al. [15, App. E] in expressing a PCFG in Dice by manually unfolding a loop that generates a parse from subparses. Figure 9d suggests that PCFGs in this encoding are intractable for Dice.

These benchmarks all contain conditional independence structure as a result of function abstractions. While MAPPL may do better on such programs, we note that Dice performs better on large Bayesian networks (BN). For example, for a BN with ~ 40 nodes, Dice solves the inference problem within 30 ms, while MAPPL takes over 1 s. We conjecture that this is due to the known result that WMC can significantly outperform VE when models contain substantial *local structure* [14].

Approximate inference: Hamiltonian Monte Carlo (HMC). HMC [23] is a powerful sampling method for differentiable models. Discrete latent variables introduce nondifferentiability, posing

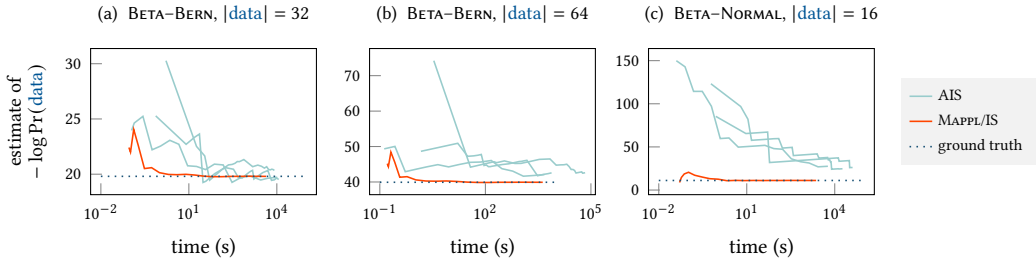


Figure 12. Performance of ML estimation on a family of hybrid discrete–continuous HMM models (Figure 1e), measured by how the negative log ML estimate changes as allowable inference time increases.

challenges to applying HMC to hybrid discrete–continuous models. We consider two such models: a soft k -means clustering model and a latent Dirichlet allocation (LDA) model. One way to handle them is by marginalizing out the discrete variables using VE. For example, Pyro’s support for HMC can handle these models by performing VE on plated factor graphs [48].

Pyro performs VE at run time. So we examine whether the performance of Pyro’s HMC can be improved by ahead-of-time VE compilation through MAPPL. Specifically, we use NumPyro [54], which supports fast HMC inference on top of JAX [24]. We run NumPyro’s HMC on the original model and on the MAPPL-compiled model with necessary syntax adjustments applied (including replacing the top-level logML with an invocation of NumPyro’s HMC).

Figure 11 displays the running time of sampling a single chain consisting of 10,000 samples and 2,500 burn-in samples using the No-U-Turn sampler [32], while varying the number of discrete latent variables in the models. Time is measured after JIT is warmed up. As expected, ahead-of-time VE compilation leads to improved run-time performance.

Approximate inference: marginal-likelihood estimation. We examine the performance implications of MAPPL’s VE compilation for marginal likelihood (ML) estimation, a key task in Bayesian learning and inference. We use a family of hybrid discrete–continuous HMMs (hmm' in Figure 1e) as benchmarks. For the VE-compiled programs, we use importance sampling (IS) with Pyro to solve the inference subproblems (i.e., nested integrals). As a baseline, we use annealed importance sampling (AIS) [47] to solve the global inference problems directly. AIS, generalizing IS, is a widely used sampling method for ML estimation. We assess the convergence rate of the ML estimate as allowable running time increases. Running time roughly translates to the number of importance samples. We experiment with multiple hyperparameter settings for AIS.

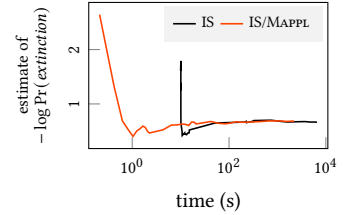
Figures 12a and 12b show the results for an HMM with input data of length 32 and 64, respectively. MAPPL/IS converges within tens of seconds in both cases. In contrast, AIS either takes thousands of seconds to converge or does not show signs of convergence even after thousands of seconds, for each hyperparameter setting tested. The quick convergence with MAPPL/IS is a consequence of the VE compilation eliminating discrete r.v.s and generating single-dimensional subproblems easily solvable by a Monte-Carlo method.

In fact, the compiled program reveals that the inference problems of this experiment have exact solutions: we are able to use Mathematica to obtain closed-form solutions to the generated subproblem integrals. The dotted lines in Figures 12a and 12b represent the exact solutions to the global problems as computed from those to the subproblems. Our approach enables potentially harnessing the power of symbolic-integration engines by translating a recursive program into subproblems that have readily available closed-form solutions.

We also consider a similar HMM where the subproblems in the VE-compiled program are not known to have closed-form solutions—though they can be approximated with arbitrary precision

using numerical methods. Figure 12c shows that MAPPL/IS quickly converges to this approximate solution (the dotted line), whereas AIS fails to converge for the relatively small problem size of 16.

VE does not always lead to improved performance, however. As another benchmark, we consider the aforementioned CRBD model extended with two continuous latent variables for the birth and death rates. The plot to the right shows that IS with the VE-compiled program does not converge faster than IS with the original program. The VE-compiled program consists in a top-level integral that nests an exact-inference problem of the kind addressed in Figure 9f, so each importance sample requires solving a new exact-inference problem and is thus more expensive than an importance sample of the original program.



8 RELATED WORK AND DISCUSSION

VE for probabilistic programs. VE is supported by many PPLs. Here we focus the discussion on those on the more expressive end of the spectrum. Early versions of IBAL support VE for programs with unbounded recursion and use lazy evaluation [51]. This approach allows models such as PCFGs to be specified more declaratively, but it seems to have been abandoned in favor of bounded recursion for correctness and efficiency concerns [50] in later versions of IBAL [52] and Figaro [53].

PERPL [15] supports exact inference for programs with unbounded recursion, by compiling them to monotone systems of polynomial equations. Infinite data types, such as integers and strings, pose challenges to equation solving, as they would lead to infinite systems of equations. In response, PERPL uses whole-program transformations (de- and re-functionalization) to eliminate infinite data types. These transformations further necessitate a linear type system to ensure correctness.

MAPPL shares the restriction of bounded recursion with a few other PPLs that support VE on PCFG-like models. Bounded recursion, while unable to express PCFGs as almost surely terminating programs, is expressive enough for Bayesian-inference queries on these models (see Figure 1c and another encoding given in an appendix), as the observed data is finite. Koller et al. [34] call such queries *evidence-finite*. We consider our choice to restrict attention to bounded recursion a sweet spot in the design space: it aligns well with the evidence-finite nature of many Bayesian inference problems, does not require the programmer to reason about linearity, leads to provably correct VE-compiled code with performance matching the best known PTIME algorithms, and still allows reasonably concise, readable programs.

SlicStan [29] supports VE for an imperative PPL where programs have deterministic support and variables are global. It supports loops but not recursion, and HMMs expressed via loops do not seem to type-check in SlicStan’s information-flow type system. MAPPL, in contrast, is a functional PPL with a wider range of features. The denotational treatment is compositional by nature: it enables a noninterference result on open terms, crucial for eliminating variables in subterms under binders. As Section 7 shows, MAPPL’s support for recursion, as well as its use of a two-level lattice rather than a meet semilattice, avoids the limitation of the SlicStan compiler in scaling to large models.

Solving probabilistic inference problems analytically. Exact, analytical solutions to inference problems are welcome whenever computationally efficient. Some PPLs support exact inference for nonrecursive programs with no or very restricted form of continuous variables, by compiling them into finite graph representations for efficient inference [12, 17, 31, 59]. FSPN [63] and PERPL [15] support exact inference for recursive programs, though they are known to work only for programs with discrete variables. Hakaru [46, 65] and Psi [25, 26] enable exact inference for programs with continuous variables (though still omitting recursion) using computer-algebra solvers.

Delayed sampling in Birch [43] and ProbZelus [8, 3] allow partial analytical solutions to sub-programs by exploiting conjugacy. The similarity to our approach is that both are forms of automatic Rao–Blackwellization [57, §4.2] that analytically reduce an inference problem to a better-behaved one. The distinctions are that delayed sampling derives closed-form posteriors for conjugate priors whereas our approach compiles away discrete variables, that delayed sampling is an *inference-time* approach based on dynamic dependence graphs whereas ours is a *compile-time* transformation, and that delayed sampling is not known to work with recursion.

In practice, no single inference technology is likely to excel at all problems; our approach and existing inference methods are complementary. Identifying independence is generally useful for compile-time Rao–Blackwellization; for example, our information-flow analysis can potentially be used in gradient-based methods to reduce variance for gradient estimators. On the other hand, our approach can potentially capitalize on advances in symbolic integration to solve generated subproblems analytically.

Reasoning about independence in probabilistic programs. Verifying randomized algorithms may require reasoning about independence, for which program logics have been developed [7, 6, 40]. While these program logics enable calculational, largely manual proofs of functional correctness, an information-flow type analysis is more amenable to automation through type inference. Hur et al. [33] and Amtoft and Banerjee [2] study program slicing for probabilistic while-languages. Their reasoning is concerned with determining if two variables are correlated conditioned on the observe statements in a program. Conditional independence can sometimes be determined syntactically for Bayesian networks through the notion of active trails [49]. The idea has been adapted to probabilistic programs [33, 37], though a full soundness result is lacking.

Semantics of probabilistic programs. As a Cartesian-closed alternative to measurable spaces, QBSes are introduced to handle higher-order types [30, 62]. It is further shown that QBSes can be equipped with compatible ω -cpo structures to handle term-level recursion and recursive types [64]. MAPPL’s denotational semantics uses these constructions. Another way to give semantics to PPLs is by first defining a deterministic operational semantics indexed by a randomness source and then integrating over randomness to obtain a measure semantics [11]. Prior work constructs logical relations for program equivalence in this operational setting [18, 66, 69], while we construct logical relations for noninterference in a denotational setting.

9 CONCLUSION

Our approach to variable elimination and marginal inference, presented in the context of MAPPL, represents a generalization and synthesis of several important ideas.

- A compiler eliminates probabilistic effects, generalizing variable elimination from graphical models to a richly expressive PPL with recursion.
- It decomposes a global inference problem into subproblems, recovering and generalizing widely used dynamic-programming algorithms for recursive models.
- It factorizes computations into independent partitions, repurposing information-flow typing to probabilistic programs.
- Its correctness result relies on a logical-relations argument, adapting semantic models for noninterference to a measure-theoretic setting.

The payoff is that MAPPL allows useful recursive models to be expressed in a functional, recursive style, while enabling sound, scalable inference for a broad class of these programs. Future work could explore ways to enable programmer control over the decomposition into inference subproblems and to exploit local structure in certain models to further speed up inference.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their valuable feedback. We thank Oghenevwoaga Ebresafe, Edward Lee, and Cong Ma for discussions and help. This work was supported by the Natural Sciences and Engineering Research Council of Canada. The views and opinions expressed are those of the authors and do not necessarily reflect the position of any funding agency.

DATA AVAILABILITY STATEMENT

The MAPPL compiler implementation is available [38].

REFERENCES

- [1] Martín Abadi, Anindya Banerjee, Nevin Heintze, and Jon G. Riecke. 1999. A core calculus of dependency. In *ACM SIGPLAN Symp. on Principles of Programming Languages (POPL)*. <https://doi.org/10.1145/292540.292555>
- [2] Torben Amtoft and Anindya Banerjee. 2020. A theory of slicing for imperative probabilistic programs. *ACM Tran. on Programming Languages and Systems (TOPLAS)* 42, 2 (April 2020). <https://doi.org/10.1145/3372895>
- [3] Eric Atkinson, Charles Yuan, Guillaume Baudart, Louis Mandel, and Michael Carbin. 2022. Semi-symbolic inference for efficient streaming probabilistic programming. *Proc. of the ACM on Programming Languages (PACMPL)* 6, OOPSLA2 (Oct. 2022). <https://doi.org/10.1145/3563347>
- [4] Martin Avanzini, Gilles Barthe, and Ugo Dal Lago. 2021. On continuation-passing transformations and expected cost analysis. *Proc. of the ACM on Programming Languages (PACMPL)* 5, ICFP (Aug. 2021). <https://doi.org/10.1145/3473592>
- [5] J. K. Baker. 1979. Trainable grammars for speech recognition. *The Journal of the Acoustical Society of America* 65, S1 (1979). <https://doi.org/10.1121/1.2017061>
- [6] Jialu Bao, Simon Docherty, Justin Hsu, and Alexandra Silva. 2021. A bunched logic for conditional independence. In *ACM/IEEE Symp. on Logic In Computer Science (LICS)*. <https://doi.org/10.1109/LICS52264.2021.9470712>
- [7] Gilles Barthe, Justin Hsu, and Kevin Liao. 2019. A probabilistic separation logic. *Proc. of the ACM on Programming Languages (PACMPL)* 4, POPL (Dec. 2019). <https://doi.org/10.1145/3371123>
- [8] Guillaume Baudart, Louis Mandel, Eric Atkinson, Benjamin Sherman, Marc Pouzet, and Michael Carbin. 2020. Reactive probabilistic programming. In *ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/3385412.3386009>
- [9] Richard E. Bellman. 1957. *Dynamic Programming*. Princeton University Press.
- [10] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. 2019. Pyro: Deep universal probabilistic programming. *Journal of Machine Learning Research (JMLR)* 20, 1 (2019). arXiv:1810.09538
- [11] Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szymczak. 2016. A lambda-calculus foundation for universal probabilistic programming. In *ACM SIGPLAN Conf. on Functional Programming (ICFP)*. <https://doi.org/10.1145/2951913.2951942>
- [12] Johannes Borgström, Andrew D. Gordon, Michael Greenberg, James Margetson, and Jurgen Van Gael. 2011. Measure transformer semantics for Bayesian machine learning. In *European Symp. on Programming (ESOP)*. https://doi.org/10.1007/978-3-642-19718-5_5
- [13] Bob Carpenter, Andrew Gelman, Matthew D. 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* 76, 1 (2017). <https://doi.org/10.18637/jss.v076.i01>
- [14] Mark Chavira and Adnan Darwiche. 2008. On probabilistic inference by weighted model counting. *Artificial Intelligence* 172, 6 (2008). <https://doi.org/10.1016/j.artint.2007.11.002>
- [15] David Chiang, Colin McDonald, and Chung-chieh Shan. 2023. Exact recursive probabilistic programming. *Proc. of the ACM on Programming Languages (PACMPL)* 7, OOPSLA1 (April 2023). <https://doi.org/10.1145/3586050> arXiv:2210.01206
- [16] David Chiang and Darcey Riley. 2020. Factor graph grammars. In *Conf. on Neural Information Processing Systems (NeurIPS)*. arXiv:2010.12048
- [17] Guillaume Claret, Sriram K. Rajamani, Aditya V. Nori, Andrew D. Gordon, and Johannes Borgström. 2013. Bayesian inference using data flow analysis. In *Proc. of the 9th Joint Meeting of the European Software Engineering Conf. and the*

- ACM SIGSOFT Symp. on the Foundations of Software Engineering (ESEC/FSE)*. <https://doi.org/10.1145/2491411.2491423>
- [18] Ryan Culpepper and Andrew Cobb. 2017. Contextual equivalence for probabilistic programs with continuous random variables and scoring. In *European Symp. on Programming (ESOP)*. https://doi.org/10.1007/978-3-662-54434-1_14
- [19] Oliver Danvy and Andrzej Filinski. 1992. Representing control: a study of the CPS transformation. *Mathematical Structures in Computer Science* 2, 4 (1992). <https://doi.org/10.1017/S0960129500001535>
- [20] Adnan Darwiche. 2009. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511811357>
- [21] Dorothy E. Denning and Peter J. Denning. 1977. Certification of programs for secure information flow. *Comm. of the ACM (CACM)* 20, 7 (July 1977). <https://doi.org/10.1145/359636.359712>
- [22] Dorothy Elizabeth Robling Denning. 1982. *Cryptography and Data Security*. Addison-Wesley Reading.
- [23] Simon Duane, A.D. Kennedy, Brian J. Pendleton, and Duncan Roweth. 1987. Hybrid Monte Carlo. *Physics Letters B* 195, 2 (1987). [https://doi.org/10.1016/0370-2693\(87\)91197-X](https://doi.org/10.1016/0370-2693(87)91197-X)
- [24] Roy Frostig, Matthew Johnson, and Chris Leary. 2018. Compiling machine learning programs via highlevel tracing. <https://mlsys.org/Conferences/doc/2018/146.pdf>
- [25] Timon Gehr, Sasa Misailovic, and Martin Vechev. 2016. PSI: Exact symbolic inference for probabilistic programs. In *Int'l Conf. on Computer Aided Verification (CAV)*. https://doi.org/10.1007/978-3-319-41528-4_4
- [26] Timon Gehr, Samuel Steffen, and Martin Vechev. 2020. λ PSI: Exact inference for higher-order probabilistic programs. In *ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/3385412.3386006>
- [27] J. A. Goguen and J. Meseguer. 1982. Security policies and security models. In *IEEE Symp. on Security and Privacy*. <https://doi.org/10.1109/SP.1982.10014>
- [28] Noah D. Goodman and Andreas Stuhlmüller. 2014. The design and implementation of probabilistic programming languages. <http://dipl.org>.
- [29] Maria I. Gorinova, Andrew D. Gordon, Charles Sutton, and Matthijs Vákár. 2021. Conditional independence by typing. *ACM Tran. on Programming Languages and Systems (TOPLAS)* 44, 1 (Dec. 2021). <https://doi.org/10.1145/3490421> arXiv:2010.11887
- [30] Chris Heunen, Ohad Kammar, Sam Staton, and Hongseok Yang. 2017. A convenient category for higher-order probability theory. In *ACM/IEEE Symp. on Logic In Computer Science (LICS)*. <https://doi.org/10.5555/3329995.3330072> arXiv:1701.02547
- [31] Steven Holtzen, Guy Van den Broeck, and Todd Millstein. 2020. Scaling exact inference for discrete probabilistic programs. *Proc. of the ACM on Programming Languages (PACMPL)* 4, OOPSLA (Nov. 2020). <https://doi.org/10.1145/3428208> arXiv:2005.09089
- [32] Matthew D. Homan and Andrew Gelman. 2014. The no-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research (JMLR)* 15, 1 (Jan. 2014). <https://doi.org/10.5555/2627435.2638586>
- [33] Chung-Kil Hur, Aditya V. Nori, Sriram K. Rajamani, and Selva Samuel. 2014. Slicing probabilistic programs. In *ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/2594291.2594303>
- [34] Daphne Koller, David McAllester, and Avi Pfeffer. 1997. Effective Bayesian inference for stochastic programs. In *Proc. of the 14th National Conf. on Artificial Intelligence and the 9th Conf. on Innovative Applications of Artificial Intelligence (AAAI/IAAI)*.
- [35] Satoshi Kura. 2023. Higher-order weakest precondition transformers via a CPS transformation. (2023). arXiv:2301.09997
- [36] Alexander K. Lew, Matin Ghavamizadeh, Martin C. Rinard, and Vikash K. Mansinghka. 2023. Probabilistic programming with stochastic probabilities. *Proc. of the ACM on Programming Languages (PACMPL)* 7, PLDI (June 2023). <https://doi.org/10.1145/3591290>
- [37] Jianlin Li, Leni Ven, Pengyuan Shi, and Yizhou Zhang. 2023. Type-preserving, dependence-aware guide generation for sound, effective amortized probabilistic inference. *Proc. of the ACM on Programming Languages (PACMPL)* POPL. <https://doi.org/10.1145/3571243>
- [38] Jianlin Li, Eric Wang, and Yizhou Zhang. 2024. Compiling probabilistic programs for variable elimination with information flow (artifact). <https://doi.org/10.5281/zenodo.10951893>

- [39] Jianlin Li, Eric Wang, and Yizhou Zhang. 2024. *Compiling Probabilistic Programs for Variable Elimination with Information Flow (Extended Version)*. Technical Report CS-2024-03. School of Computer Science, University of Waterloo. <https://cs.uwaterloo.ca/sites/ca.computer-science/files/uploads/files/cs-2024-03.pdf>
- [40] John M. Li, Amal Ahmed, and Steven Holtzen. 2023. Lilac: A modal separation logic for conditional probability. *Proc. of the ACM on Programming Languages (PACMPL)* 7, PLDI (June 2023). <https://doi.org/10.1145/3591226> arXiv:2304.01339
- [41] Christopher Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press.
- [42] Andrew McCallum, Karl Schultz, and Sameer Singh. 2009. FACTORIE: Probabilistic programming via imperatively defined factor graphs. In *Conf. on Neural Information Processing Systems (NIPS)*.
- [43] Lawrence Murray, Daniel Lundén, Jan Kudlicka, David Broman, and Thomas Schön. 2018. Delayed sampling and automatic Rao-Blackwellization of probabilistic programs. In *Int'l Conf. on Artificial Intelligence and Statistics (AISTATS)*. arXiv:1708.07787
- [44] Andrew C. Myers. 1999. JFlow: Practical mostly-static information flow control. In *ACM SIGPLAN Symp. on Principles of Programming Languages (POPL)*. <https://doi.org/10.1145/292540.292561>
- [45] Andrew C. Myers and Barbara Liskov. 1997. A decentralized model for information flow control. In *Proc. of the 16th ACM Symp. on Operating Systems Principles (SOSP)*. <https://doi.org/10.1145/268998.266669>
- [46] Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. 2016. Probabilistic inference by program transformation in Hakaru (system description). In *Functional and Logic Programming*. https://doi.org/10.1007/978-3-319-29604-3_5
- [47] Radford M. Neal. 2001. Annealed importance sampling. *Statistics and Computing* 11 (2001). <https://doi.org/10.1023/A:1008923215028>
- [48] Fritz Obermeyer, Eli Bingham, Martin Jankowiak, Neeraj Pradhan, Justin Chiu, Alexander Rush, and Noah Goodman. 2019. Tensor variable elimination for plated factor graphs. In *Int'l Conf. on Machine Learning (ICML)*. arXiv:1902.03210
- [49] Judea Pearl. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann. <https://doi.org/10.1016/B978-0-08-051489-5.50001-1>
- [50] Avi Pfeffer. [n.d.]. IBAL: An expressive, functional probabilistic modeling language. ([n.d.]). <https://www.cs.tufts.edu/~nr/cs257/archive/avi-pfeffer/ibal-journal.pdf>
- [51] Avi Pfeffer. 2001. IBAL: A probabilistic rational programming language. In *Int'l Joint Conf. on Artificial Intelligence (IJCAI)*.
- [52] Avi Pfeffer. 2007. The design and implementation of IBAL: A general-purpose probabilistic language. In *Introduction to Statistical Relational Learning*. The MIT Press. <https://doi.org/10.7551/mitpress/7432.003.0016>
- [53] Avi Pfeffer. 2016. *Practical Probabilistic Programming*. Manning Publications.
- [54] Du Phan, Neeraj Pradhan, and Martin Jankowiak. 2019. Composable effects for flexible and accelerated probabilistic programming in NumPyro. (2019). arXiv:1912.11554
- [55] Lawrence R. Rabiner and Biing-Hwang Juang. 1986. An introduction to hidden Markov models. *IEEE ASSP Magazine* 3, 1 (1986). <https://doi.org/10.1109/MASSP.1986.1165342>
- [56] Vineet Rajani and Deepak Garg. 2018. Types for information flow control: Labeling granularity and semantic models. In *IEEE Computer Security Foundations Symp. (CSF)*. <https://doi.org/10.1109/CSF.2018.00024>
- [57] Christian P. Robert and George Casella. 1999. *Monte Carlo Statistical Methods*. Springer. <https://doi.org/10.1007/978-1-4757-4145-2>
- [58] Fredrik Ronquist, Jan Kudlicka, Viktor Senderov, Johannes Borgström, Nicolas Lartillot, Daniel Lundén, Lawrence Murray, Thomas B. Schön, and David Broman. 2021. Universal probabilistic programming offers a powerful approach to statistical phylogenetics. *Communications Biology* 4, 1 (2021). <https://doi.org/10.1038/s42003-021-01753-7>
- [59] Feras A. Saad, Martin C. Rinard, and Vikash K. Mansinghka. 2021. SPPL: Probabilistic programming with fast exact symbolic inference. In *ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/3453483.3454078>
- [60] Andrei Sabelfeld and Andrew C. Myers. 2003. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications* 21, 1 (2003). <https://doi.org/10.1109/JSAC.2002.806121>
- [61] Andrei Sabelfeld and David Sands. 2001. A PER model of secure information flow in sequential programs. *Higher-Order and Symbolic Computation* 14 (2001). <https://doi.org/10.1023/A:1011553200337>

- [62] Adam Ścibior, Ohad Kammar, Matthijs Vákár, Sam Staton, Hongseok Yang, Yufei Cai, Klaus Ostermann, Sean K. Moss, Chris Heunen, and Zoubin Ghahramani. 2017. Denotational validation of higher-order Bayesian inference. *Proc. of the ACM on Programming Languages (PACMPL)* 2, POPL (Dec. 2017). <https://doi.org/10.1145/3158148>
- [63] Andreas Stuhlmüller and Noah D. Goodman. 2012. A dynamic programming algorithm for inference in recursive probabilistic programs. (2012). arXiv:1206.3555
- [64] Matthijs Vákár, Ohad Kammar, and Sam Staton. 2019. A domain theory for statistical probabilistic programming. *Proc. of the ACM on Programming Languages (PACMPL)* 3, POPL (2019). <https://doi.org/10.1145/3290349>
- [65] Rajan Walia, Praveen Narayanan, Jacques Carette, Sam Tobin-Hochstadt, and Chung-chieh Shan. 2019. From high-level inference algorithms to efficient code. *Proc. of the ACM on Programming Languages (PACMPL)* 3, ICFP (July 2019). <https://doi.org/10.1145/3341702>
- [66] Mitchell Wand, Ryan Culpepper, Theophilos Giannakopoulos, and Andrew Cobb. 2018. Contextual equivalence for a probabilistic language with continuous random variables and recursion. *Proc. of the ACM on Programming Languages (PACMPL)* 2, ICFP (2018). <https://doi.org/10.1145/3236782>
- [67] Di Wang, Jan Hoffmann, and Thomas Reps. 2021. Sound probabilistic inference via guide types. In *ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*. arXiv:2104.03598
- [68] Nevin Lianwen Zhang and David Poole. 1994. A simple approach to Bayesian network computations. In *Proc. of the 10th Canadian Conference on Artificial Intelligence*.
- [69] Yizhou Zhang and Nada Amin. 2022. Reasoning about “reasoning about reasoning”: Semantics and contextual equivalence for probabilistic programs with nested queries and recursion. *Proc. of the ACM on Programming Languages (PACMPL)* 6, POPL (2022). <https://doi.org/10.1145/3498677>

Received 2023-11-16; accepted 2024-03-31