

Parallelizable Feynman-Kac models for universal probabilistic programming

Michele Boreale

Luisa Collodi

Università degli Studi di Firenze, Dipartimento di Statistica, Informatica, Applicazioni “G. Parenti”,

{michele.boreale,luisa.collodi}@unifi.it

We study provably correct and efficient instantiations of the Sequential Monte Carlo (SMC) inference scheme in the context of formal operational semantics of Probabilistic Programs (PPs). We focus on *universal* PPs featuring sampling from arbitrary measures and conditioning/reweighting in unbounded loops. We first equip Probabilistic Program Graphs (PPGs), an automata-theoretic description format of PPs, with an expectation-based semantics over infinite execution traces, which also incorporates trace weights. We then prove a finite approximation theorem that provides bounds to this semantics based on expectations taken over finite, fixed-length traces. This enables us to frame our semantics within a Feynman-Kac (FK) model, and ensures the consistency of the Particle Filtering (PF) algorithm, an instance of SMC, with respect to our semantics. Building on these results, we introduce VPF, a vectorized version of the PF algorithm tailored to PPGs and our semantics. Experiments conducted with a proof-of-concept implementation of VPF show very promising results compared to state-of-the-art PP inference tools.

Keywords: Probabilistic programming, operational semantics, SIMD parallelism, SMC.

1 Introduction

Probabilistic Programming Languages (PPLs) [25, 7] offer a systematic approach to define arbitrarily complicated probabilistic models. One is typically interested in performing *inference* on these models, given observed data; for example, finding the posterior distribution of the program’s output conditioned on the observed data. Here, in the context of formal operational semantics of Probabilistic Programs, we study provably correct and parallelizable instantiations of the Sequential Monte Carlo (SMC) inference scheme.

In terms of formal semantics of PPLs, the denotational approach introduced by Kozen [30] offers a solid mathematical foundation. However, when it comes to practical algorithms for PPL-based inference, the landscape appears somewhat fragmented. On one hand, *symbolic* and *static analysis* techniques, see e.g. [43, 22, 38, 8, 10, 42], yield results with correctness guarantees firmly grounded in the semantics of PPLs but often struggle with scalability. On the other hand, practical languages and inference algorithms predominantly leverage Monte Carlo (MC) *sampling* techniques (MCMC, SMC), which are more scalable but often lack a clear connection to formal semantics [24, 17, 12]. Notable exceptions to this situation include works such as [40, 51, 33, 13, 32], which are discussed in the related work section.

Establishing the consistency of an inference algorithm with respect to a PPL’s formal semantics is not merely a theoretical pursuit. In the context of *universal* PPLs [23], integration of unbounded loops and conditioning with MC sampling, which requires truncating computations at a finite time, presents significant challenges [10]. Additionally, the interplay between continuous and discrete distributions in these PPLs can lead to complications, potentially causing existing sampling-based algorithms to yield incorrect results [51]. In the present work, we establish a precise connection between *Probabilistic Program Graphs* (PPGs), a general automata-theoretic description format of PPs, and *Feynman-Kac*

(FK) models, a formalism for state-based probabilistic processes and observations defined over a finite time horizon [19, Ch.5]. This connection enables us to prove the consistency for PPGs of the *Particle Filtering (PF)* inference algorithm, one of the incarnations of Sequential Monte Carlo approach [19, Ch.10]. In establishing this connection, we adopt a decisively operational perspective, as explained below.

In a PPG (Section 3), computation (essentially, sampling) progresses in successive stages specified by the direct edges of a graph (transitions), with nodes serving as *checkpoints* between stages for *conditioning* on observed data or more generally updating computation weights. The operational semantics of PPGs is formalized in terms of Markov kernels and score functions. Building on this, we introduce a measure-theoretic, infinite-trace semantics (Section 4, with the necessary measure theory reviewed in Section 2). A finite approximation theorem then allows us to relate this trace semantics precisely to a finite-time horizon FK model (Section 5). PF is known to be *consistent* for FK models asymptotically: as the number N of simulated instances (*particles*) tends to infinity, the distribution of these particles converges to the measure defined by the FK model [19, Ch.11]. Therefore, consistency of PF for PPGs will automatically follow.

Our approach yields additional insights. First, the finite approximation theorem holds for a class of *prefix-closed functions* defined on infinite traces: these are the functions where the output only depends on a finite initial segment of the input argument. The finite approximation theorem implies that the expectation of a prefix-closed function, defined on the probability space of infinite traces, can be approximated by the expectation of functions defined over truncated traces, with respect to a measure defined on a suitable FK model. As expectation in a FK model can be effectively estimated, via PF or other algorithms, our finite approximation result lays a sound basis for the statistical model checking of PPs. Second, the automata-theoretic operational semantics of PPGs translates into a *vectorized* implementation of PF, leveraging the fine-grained, SIMD parallelism existing at the level of particles. Specifically, the PPG’s transition function and the score functions are applied simultaneously to the entire vector of N simulated particles at each step. This is practically significant, as modern CPUs and programming languages offer extensive support for vectorization, that may lead to dramatic speedups. We demonstrate this aspect with a prototype vectorized implementation of a PPG-based PF algorithm using TensorFlow [1], called VPF. Experiments comparing VPF with state-of-the-art PPLs on challenging examples from the literature show very promising results (Section 6). Concluding remarks are provided in the final section (Section 7). Omitted proofs and additional technical material have been reported in an extended version available online [15].

In summary, our main contributions are as follows: (1) A clean semantics for PPGs based on expectation taken over infinite-trace, which incorporates conditioning/reweighting; (2) a finite approximation theorem linking this semantics to finite traces and FK models, thereby establishing the consistency of PF for PPGs; (3) a vectorized version of the PF algorithm based on PPGs.

Related work With few notable exceptions, most work on the semantics of PPL follows the denotational approach initiated by Kozen [30]; see [7, 11, 16, 25, 26, 27, 44, 45, 46, 48] and references therein for representative works in this area. In this context, a general goal orthogonal to ours, is devising methods to combine and reason on densities. Note that we do not require that a PP induces a density on the probability space of infinite traces.

Relevant to our approach is a series of works by Lunden et al. on SMC inference applied to PPLs. In [33], for a lambda-calculus enriched with an explicit resample primitive, consistency of PF is shown to hold, under certain restrictions, independently of the placements of the resamples in the code. Operationally, their functional approach is very different from our automata-theoretic one. In particular,

they handle suspension and resumption of particles in correspondence of resampling via an implicit use of *continuations*, in the style of webPPL [24] and other PPLs. The combination of functional style and continuations does not naturally lend itself to vectorization. For instance, ensuring that all particles are *aligned*, that is are at a resample point of their execution, is an issue that can impact negatively on performance or accuracy. On the contrary, in our automata-theoretic model, placement of resamples and alignment are not issues: resampling always happens after each (vectorized) transition step, so all particles are automatically aligned. Note that in PPGs a transition can group together complicated, conditioning free computations; in any case, consistency of PF is guaranteed. In a subsequent work [35, 34], Lunden et al. study concrete implementation issues of SMC. In [35], they consider *PPL Control-Flow Graphs* (PCFGs), a structure intended as a target for the compilation of high-level PPLs, such as their CorePPL. The PCFG model is very similar in spirit to PPGs, however, it lacks a formal semantics. Lunden et al. also offer an implementation of this framework, designed to take advantage of the potential parallelism existing at the level of particles. We compare our implementation with theirs in Section 6.

Aditya et al. prove consistency of Markov Chain Monte Carlo (MCMC) for their PPL R2 [40], which is based on a big-step sampling semantics that considers finite execution paths. No approximation results bridging finite and infinite traces, and hence unbounded loops, is provided. It is also unclear if a big-step semantics would effectively translate into a SIMD-parallel algorithm. Wu et al. [51] provide the PPL Blog with a rigorous measure-theoretic semantics, formulated in terms of Bayesian Networks, and a very efficient implementation of the PF algorithm tailored to such networks. Again, they do not offer results for unbounded loops. In our previous work [13], we have considered a measure theoretic semantics for a PPL with unbounded loops, and provided a finite approximation result and a SIMD-parallel implementation, with guarantees, of what is in effect a *rejection sampling* algorithm. Rejection may be effective for limited forms of conditioning; but it rapidly becomes wasteful and ineffective as conditioning becomes more demanding, so to speak: e.g. when it is repeated in a loop, or the observed data have a low likelihood in the model. Finally, SMCP3 [32] provides a rich measure-theoretic framework for extending the practical Gen language [20] with expressive proposal distributions.

A rich area in the field of PPL focuses on symbolic, exact techniques [43, 22, 38, 8, 10, 42, 28] aiming to obtain termination certificates, or certified bounds on termination probability of PPs, or even exact representations of the posterior distribution; see also [6, 9, 47, 3, 31, 49] for some recent works in this direction. Our goal and methodology, as already stressed, are rather different, as we focus on scalable inference via sampling and the ensuing consistency issues.

2 Preliminaries on measure theory

We review a few basic concepts from measure theory following closely the presentation in the first two chapters of [2], which is a reference for whatever is not explicitly described below. Given a nonempty set Ω , a *sigma-field* \mathcal{F} on Ω is a collection of subsets of Ω that contains Ω , and is closed under complement and under countable disjoint union. The pair (Ω, \mathcal{F}) is called a *measurable space*. A (total) function $f : \Omega_1 \rightarrow \Omega_2$ is *measurable* w.r.t. the sigma-fields $(\Omega_1, \mathcal{F}_1)$ and $(\Omega_2, \mathcal{F}_2)$ if whenever $A \in \mathcal{F}_2$ then $f^{-1}(A) \in \mathcal{F}_1$. We let $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$ be the set of extended reals, assuming the standard arithmetic for $\pm\infty$ (cf. [2, Sect.1.5.2]), and $\overline{\mathbb{R}}^+$ the set of nonnegative reals including $+\infty$. The *Borel sigma-field* \mathcal{F} on $\Omega = \overline{\mathbb{R}}^m$ is the minimal sigma-field that contains all rectangles of the form $[a_1, b_1] \times \cdots \times [a_n, b_n]$, with $a_i, b_i \in \overline{\mathbb{R}}$. An important case of measurable spaces (Ω, \mathcal{F}) is when $\Omega = \overline{\mathbb{R}}^m$ for some $m \geq 1$ and \mathcal{F} is the Borel sigma-field over Ω . Throughout the paper, “*measurable*” means “*Borel measurable*”, both for sets and for functions. On functions, Borel measurability is preserved by composition and other elementary

operations on functions; continuous real functions are Borel measurable. We will let \mathcal{F}_k denote the Borel sigma-field over $\overline{\mathbb{R}}^k$ ($k \geq 1$) when we want to be specific about the dimension of the space.

A *measure* over a measurable space (Ω, \mathcal{F}) is a function $\mu : \mathcal{F} \rightarrow \overline{\mathbb{R}}^+$ that is countably additive, that is $\mu(\cup_{j \geq 1} A_j) = \sum_{j \geq 1} \mu(A_j)$ whenever A_j 's are pairwise disjoint sets in \mathcal{F} . The *Lebesgue integral* of a Borel measurable function f w.r.t. a measure μ [2, Ch.1.5], both defined over a measure space (Ω, \mathcal{F}) , is denoted by $\int_{\Omega} \mu(d\omega) f(\omega)$, with the subscript Ω omitted when clear from the context. When μ is the standard Lebesgue measure, we may omit μ and write the integral as $\int_{\Omega} d\omega f(\omega)$. For $A \in \mathcal{F}$, $\int_A \mu(d\omega) f(\omega)$ denotes $\int_{\Omega} \mu(d\omega) f(\omega) 1_A(\omega)$, where $1_A(\cdot)$ is the indicator function of the set A . We let δ_v denote Dirac's measure concentrated on v : for each set A in an appropriate sigma-field, $\delta_v(A) = 1$ if $v \in A$, $\delta_v(A) = 0$ otherwise. Otherwise said, $\delta_v(A) = 1_A(v)$. Another measure that arises (in connection with discrete distributions) is the counting measure, $\mu_C(A) := |A|$. In particular, for a nonnegative f , we have the equality $\int_A \mu_C(d\omega) f(\omega) = \sum_{\omega \in A} f(\omega)$. A *probability measure* is a measure μ defined on \mathcal{F} such that $\int \mu(d\omega) = 1$. For a given nonnegative measurable function f defined over Ω , its *expectation* w.r.t. a probability measure ν is just its integral: $E_{\nu}[f] = \int \nu(d\omega) f(\omega)$. The following definition is central.

Definition 1 (Markov kernel) Let $(\Omega_1, \mathcal{F}_1)$ and $(\Omega_2, \mathcal{F}_2)$ be measurable spaces. A function $K : \Omega_1 \times \mathcal{F}_2 \rightarrow \overline{\mathbb{R}}^+$ is a Markov kernel from Ω_1 to Ω_2 if it satisfies the following properties:

1. for each $\omega \in \Omega_1$, the function $K(\omega, \cdot) : \mathcal{F}_2 \rightarrow \overline{\mathbb{R}}^+$ is a probability measure on $(\Omega_2, \mathcal{F}_2)$;
2. for each $A \in \mathcal{F}_2$, the function $K(\cdot, A) : \Omega_1 \rightarrow \overline{\mathbb{R}}^+$ is measurable.

Notationally, we will most often write $K(\omega, A)$ as $K(\omega)(A)$. The following is a standard result about the construction of finite product of measures over a product space¹ $\Omega^t = \Omega \times \cdots \times \Omega$ (t times) for $t \geq 1$ an integer. It is customary to denote the measure μ^t defined by the theorem also as $\mu^1 \otimes K_2 \otimes \cdots \otimes K_t$.

Theorem 1 (product of measures, [2], Th.2.6.7) Let $t \geq 1$ be an integer. Let μ^1 be a probability measure on Ω and K_2, \dots, K_t be $t - 1$ (not necessarily distinct) Markov kernels from Ω to Ω . Then there is a unique probability measure μ^t defined on $(\Omega^t, \mathcal{F}^t)$ such that for every $A_1 \times \cdots \times A_t \in \mathcal{F}^t$ we have: $\mu^t(A_1 \times \cdots \times A_t) = \int_{A_1} \mu^1(d\omega_1) \int_{A_2} K_2(\omega_1)(d\omega_2) \cdots \int_{A_t} K_t(\omega_{t-1})(d\omega_t)$.

3 Probabilistic programs

We first introduce a general formalism for specifying programs, in the form of certain graphs that can be regarded as symbolic finite automata. For this formalism, we introduce an operational semantics in terms of Markov kernels.

Probabilistic Program Graphs In defining probabilistic programs, we will rely on a repertoire of basic distributions: continuous, discrete and mixed distributions will be allowed. A crucial point for expressiveness is that a measure may depend on *parameters*, whose value at runtime is determined by the state of the program. To ensure that the resulting programs define measurable functions (on a suitable space), it is important that the dependence between the parameters and the measure be in turn of measurable type. We will formalize this in terms of Markov kernels. Additionally, we will consider score functions, a generalization of 0/1-valued predicates. Formally, we will consider the two families of functions defined below. In the definitions, we will let $m \geq 1$ denote a fixed integer, representing the

¹We shall freely identify language-theoretic words with *tuples*, hence use the notations $A_1 \cdot A_2 \cdot \cdots \cdot A_k$ and $A_1 \times A_2 \times \cdots \times A_k$ interchangeably. This convention will also apply to infinite words (cf. Section 4).

number of *variables* in the program, conventionally referred to as x_1, \dots, x_m . We will let v range over $\overline{\mathbb{R}}^m$, the content of the program variables in a given state, or *store*.

- *Parametric measures*: Markov kernels $\zeta : \overline{\mathbb{R}}^m \times \mathcal{F}_m \rightarrow [0, 1]$.
- *Score functions*: measurable functions $\gamma : \overline{\mathbb{R}}^m \rightarrow [0, 1]$. A *predicate* is a special case of a score function $\varphi : \overline{\mathbb{R}}^m \rightarrow \{0, 1\}$. An Iverson bracket style notation will be often employed, e.g.: $[x_1 \geq 1]$ is the predicate that on input v yields 1 if $v_1 \geq 1$, 0 otherwise.

For a parametric measure ζ and a store $v \in \overline{\mathbb{R}}^m$, $\zeta(v)$ is a distribution, that can be used to sample a new store $v' \in \overline{\mathbb{R}}^m$ depending on the current program store v . Analytically, ζ may be expressed by, for instance, chaining together sampling of individual components of the store. This can be done by relying on *parametric densities*: measurable functions $\rho : \overline{\mathbb{R}}^m \times \overline{\mathbb{R}} \rightarrow \overline{\mathbb{R}}^+$ such that, for a designated measure μ_ρ , the function $(v, A) \mapsto \int_A \mu_\rho(dr) \rho(v, r)$ ($A \in \mathcal{F}_m$) is a Markov kernel from $\overline{\mathbb{R}}^m$ to $\overline{\mathbb{R}}$. This is explained via the following example.

Example 1 Fix $m = 2$. Consider the Markov kernel defined as follows, for each $x_1, x_2 \in \overline{\mathbb{R}}$ and $A \in \mathcal{F}_2$

$$\zeta(x_1, x_2)(A) := \int \mu_1(dr_1) \left(\rho_1(x_1, x_2, r_1) \cdot \int \mu_2(dr_2) \rho_2(r_1, x_2, r_2) 1_A(r_1, r_2) \right) \quad (1)$$

where: $\mu_1 = \mu_C$ is the counting measure; $\rho_1(x_1, x_2, r) = \frac{1}{2} 1_{\{x_1\}}(r) + \frac{1}{2} 1_{\{x_2\}}(r)$ is the density of a discrete distribution on $\{x_1, x_2\}$; $\mu_2 = \mu_L$ is the ordinary Lebesgue measure; $\rho_2(x_1, x_2, r) = N(x_1, x_2, r) := \frac{1}{|x_2| \sqrt{2\pi}} \exp(-\frac{1}{2}(\frac{r-x_1}{|x_2|})^2)$ is the density of the Normal distribution of mean x_1 and standard deviation² $|x_2|$. The function ζ is a parametric measure: concretely, it corresponds to first sampling uniformly r_1 from the set $\{x_1, x_2\}$, then sampling r_2 from the Normal distribution of mean r_1 and s.d. $|x_2|$ (if $|x_2|$ is positive and finite, otherwise from a default distribution). Rather than via (1), we will describe ζ via the following more handy notation: $r_1 \sim \rho_1(x_1, x_2); r_2 \sim \rho_2(r_1, x_2)$ (or listed top-down). Note that the sampling order from left to right is relevant here.

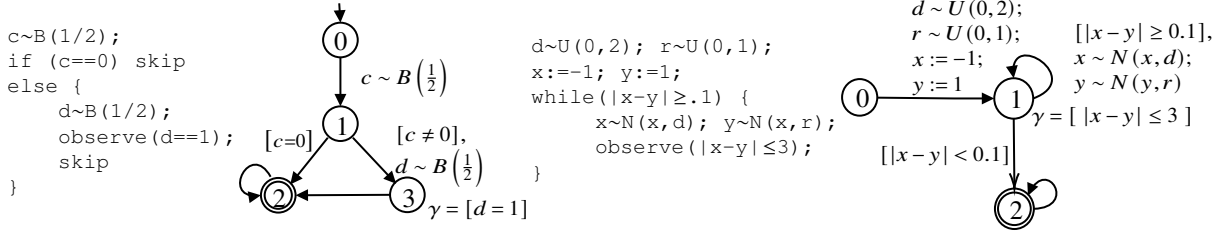
In fact, as far as the formal framework of PPGs introduced below is concerned, how the parametric measures ζ 's are analytically described is irrelevant. From the practical point of view, it is important we know how to (efficiently) sample from the measure $\zeta(v)$, for any v , in order for the inference algorithms to be actually implemented (see Section 5). In concrete terms, $\zeta(v)$ might represent the (possibly unknown) distribution of the outputs in $\overline{\mathbb{R}}^m$ returned by a piece of code, when invoked with input v . Another important special case of parametric measure is the following. For any $v = (v_1, \dots, v_m) \in \overline{\mathbb{R}}^m$, $r \in \overline{\mathbb{R}}$ and $1 \leq i \leq m$, let $v[r@i] := (v_1, \dots, r, \dots, v_m)$ denote the tuple where v_i has been replaced by r . Consider the parametric measure $\zeta(v) = \delta_{v[g(v)@i]}$, where $g : \overline{\mathbb{R}}^m \rightarrow \overline{\mathbb{R}}$ is a measurable function. In programming terms, this corresponds to the deterministic *assignment* of the value $g(v)$ to the variable x_i . We will describe this ζ as: $x_i := g(x_1, \dots, x_m)$.

In the definition of PPG below, one may think of the computation (sampling) taking place in successive stages on the edges (transitions) of the graph, with nodes serving as *checkpoints* (a term we have borrowed from [33]) between stages for conditioning on observed data — or, more generally, re-weighting the score assigned to a computation. The edges also account for the control flow among the different stages via predicates computed on the store of the source nodes.

Definition 2 (PPG) Fix $m \geq 1$. A Probabilistic Program Graph (PPG) on $\overline{\mathbb{R}}^m$ is a 4-tuple $\mathbf{G} = (\mathcal{P}, E, \text{nil}, \text{sc})$ satisfying the following.

- $\mathcal{P} = \{S_1, \dots, S_k\}$ is a finite, nonempty set of program checkpoints (programs, for short).

²With the proviso that, when $x_2 = 0$ or $|x_1|, |x_2| = +\infty$, $N(x_1, x_2, r)$ denotes an arbitrarily fixed, default probability density.



Some additional notational shorthand is in order. First, we identify \mathcal{P} with the finite set of naturals $\{0, \dots, |\mathcal{P}| - 1\}$. With this convention, we have that $\overline{\mathbb{R}}^m \times \mathcal{P} \subseteq \overline{\mathbb{R}}^{m+1}$. Henceforth, we define our state space and sigma-field as follows:

$$\Omega := \overline{\mathbb{R}}^{m+1} \quad \mathcal{F} := \text{Borel sigma-field over } \overline{\mathbb{R}}^{m+1}.$$

We keep the symbol \mathcal{F}_k for the Borel sigma-field over $\overline{\mathbb{R}}^k$, for any $k \geq 1$. For any $S \in \mathcal{P}$ and $A \in \mathcal{F}$, we let $A_S := \{v \in \overline{\mathbb{R}}^m : (v, S) \in A\}$ be the *section* of A at S . Note that $A_S \in \mathcal{F}_m$, as sections of measurable sets are measurable, see [2, Th.2.6.2, proof(1)].

Definition 3 (PPG Markov kernel) *The function $\kappa : \Omega \times \mathcal{F} \rightarrow \mathbb{R}^+$ is defined as follows, for each $\omega \in \Omega$ and $A \in \mathcal{F}$:*

$$\kappa(\omega)(A) := \begin{cases} \delta_\omega(A) & \text{if } \omega \notin \overline{\mathbb{R}}^m \times \mathcal{P} \\ \sum_{(S, \varphi, \zeta, S') \in E_S} \varphi(v) \cdot \zeta(v)(A_{S'}) & \text{if } \omega = (v, S) \in \overline{\mathbb{R}}^m \times \mathcal{P}. \end{cases} \quad (2)$$

Lemma 1 *The function κ is a Markov kernel from Ω to Ω .*

4 Trace semantics and finite approximation for PPGs

Trace semantics In what follows, we fix an arbitrary PPG, $\mathbf{G} = (\mathcal{P}, E, \text{nil}, \text{sc})$ and let κ denote the induced Markov kernel, as per Definition 3. For any $t \geq 1$, we call Ω^t the set of *paths of length t* . Consider now the set of paths of infinite length, Ω^∞ , that is the set of infinite sequences $\tilde{\omega} = (\omega_1, \omega_2, \dots)$ with $\omega_i \in \Omega$. For any $\omega^t \in \Omega^t$ and $\tilde{\omega} \in \Omega^\infty$, we identify the pair $(\omega^t, \tilde{\omega})$ with the element of Ω^∞ in which the prefix ω^t is followed by $\tilde{\omega}$. For $t \geq 1$ and a measurable $B_t \subseteq \Omega^t$, we let $c(B_t) := B_t \cdot \Omega^\infty \subseteq \Omega^\infty$ be the *measurable cylinder* generated by B_t . We let \mathcal{C} be the minimal sigma-field over Ω^∞ generated by all measurable cylinders. Under the same assumptions of Theorem 1 on the measure μ^1 and on the kernels K_2, K_3, \dots there exists a unique measure μ^∞ on \mathcal{C} such that for each $t \geq 1$ and each measurable cylinder $c(B_t)$, it holds that $\mu^\infty(c(B_t)) = \mu^t(B_t)$: see [2, Th.2.7.2], also known as the *Ionescu-Tulcea theorem*. In the definition below, we let $0 = (0, \dots, 0)$ (m times) and consider $\delta_{(0, S)}$, the Dirac's measure on Ω that concentrates all the probability mass in $(0, S)$.

Definition 4 (probability measure induced by S) *Let $S \in \mathcal{P}$. For each integer $t \geq 1$, we let μ_S^t be the probability measure over Ω^t uniquely defined by Theorem 1(a) by letting $\mu^1 = \delta_{(0, S)}$ and $K_2 = \dots = K_t = \kappa$. We let μ_S^∞ be the unique probability measure on \mathcal{C} induced by μ_1 and $K_2 = \dots = K_t = \dots = \kappa$, as determined by the Ionescu-Tulcea theorem.*

In other words, $\mu_S^t = \delta_{(0, S)} \otimes \kappa \otimes \dots \otimes \kappa$ ($t - 1$ times κ). By convention, if $t = 1$, $\mu_S^t = \delta_{(0, S)}$. The measure μ_S^∞ can be informally interpreted as the limit of the measures μ_S^t and represents the semantics of S .

Recall that the *support* of an (extended) real valued function f is the set $\text{supp}(f) := \{z : f(z) \neq 0\}$. In what follows, we shall concentrate on nonnegative measurable functions f to avoid unnecessary complications with the existence of integrals. General functions can be dealt with by the usual trick of decomposing f as $f = f^+ - f^-$, where $f^+ = \max(0, f)$ and $f^- = -\min(0, f)$, and then dealing separately with f^+ and f^- . Let us introduce a *combined score function* $\text{sc} : \Omega \rightarrow [0, 1]$ as follows, for each $\omega = (v, S)$:

$$\text{sc}(\omega) := \begin{cases} \text{sc}(S)(v) & \text{if } \omega = (v, S) \in \overline{\mathbb{R}}^m \times \mathcal{P} \\ 1 & \text{otherwise.} \end{cases} \quad (3)$$

The function $\text{sc}(\cdot)$ is extended to a *weight function* on infinite traces, $w : \Omega^\infty \rightarrow [0, 1]$ by letting³, for any $\tilde{\omega} = (\omega_1, \omega_2, \dots) \in \Omega^\infty$:

$$w(\tilde{\omega}) := \prod_{j \geq 1} \text{sc}(\omega_j). \quad (4)$$

For each $t \geq 1$, we define the weight function truncated at time t , $w_t : \Omega^t \rightarrow [0, 1]$, by $w_t(\omega^t) := \prod_{j=1}^t \text{sc}(\omega_j)$. Both w and w_t ($t \geq 1$) are measurable functions on the respective domains. We arrive at the definition of the semantics of programs. We consider the ratio of the unnormalized semantics ($[S]f$) to the weight of all traces, terminated or not ($[S]w$). In the special case when the score functions represent conditioning, this choice corresponds to quotienting over the probability of *non failed* traces. In PPL, quotienting over non failed states is somewhat standard: see e.g. the discussion in [29, Section 8.3.2].

Definition 5 (trace semantics) *Let f be a nonnegative measurable function defined on Ω^∞ . We let the unnormalized semantics of S and f be $[S]f := E_{\mu_S^\infty}[f] (= \int \mu_S^\infty(d\tilde{\omega})f(\tilde{\omega}))$. We let*

$$[[S]]f := \frac{[S](f \cdot w)}{[S]w} \quad (5)$$

provided the denominator above is > 0 ; otherwise $[[S]]f$ is undefined.

Finite approximation We are mainly interested in $[[S]]f$ in cases where the value of f is, informally speaking, determined by a finite prefix of its argument: we call these functions *prefix-closed*, and will define them further below. We first have to introduce prefix-closed languages⁴, for which some notation on languages of finite and infinite words is useful. Given two words $w, w' \in \Omega^*$, we write $w < w'$ if w is a prefix of w' , i.e. there exists a word $w'' \in \Omega^*$ such that $ww'' = w'$; otherwise we write $w \not< w'$. For $L, L' \subseteq \Omega^*$, we write $L \not< L'$ if for all $w \in L$ and $w' \in L'$ we have $w \not< w'$. A sequence of languages L_0, L_1, \dots such that for each j , $L_j \subseteq \Omega^j$ (with $\Omega^0 := \{\epsilon\}$, the empty sequence) is said to be *prefix-free* if for each $i \neq j$, $L_i \not< L_j$. Note that if $L_0 \neq \emptyset$ then $L_j = \emptyset$ for $j \geq 1$. For the sake of uniform notation, in what follows we convene that $\omega^0 := \epsilon$ and $\mathbf{c}(\{\epsilon\}) := \Omega^\infty$. We say $A \subseteq \Omega^\infty$ is a *prefix closed set* if there is a prefix-free sequence of languages L_0, L_1, \dots such that $A = \bigcup_{j=0}^\infty \mathbf{c}(L_j)$; we call L_j a *j-branch* of A , and refer to L_0, L_1, \dots collectively as *branches of A* . For any $t \geq 1$, we define the following subsets of Ω^t :

$$L^{\leq t} := \bigcup_{j=0}^t L_j \cdot \Omega^{t-j}, \quad L^{> t} := \{\omega^t : \text{there is } t' > t \text{ and } \omega_{t'} \in L_{t'} \text{ s.t. } \omega^t < \omega^{t'}\}.$$

Informally speaking, $L^{\leq t}$ is the set of paths of length t that will become members of A however we extend them to infinite words. $L^{> t}$ is the set of paths of length t for which some infinite extensions, but not all, are in A — they are so to speak “undecided”. Of special interest is the prefix-free sequence of languages defined below.

Definition 6 (termination) *Let $\mathbb{T} := \overline{\mathbb{R}}^m \times \{\text{nil}\}$ be the set of terminated states and let \mathbb{T}^c denote its complement. We let $T_j \subseteq \Omega^j$ ($j \geq 0$) be the set of finite sequences that terminate at time j , that is: $T_0 := \emptyset$ and $T_j := (\mathbb{T}^c)^{j-1} \cdot \mathbb{T}$, for $j \geq 1$. We let $T_\infty := \bigcup_{t \geq 0} \mathbf{c}(T_t) \subseteq \Omega^\infty$ denote the set of infinite sequences that terminate in finite time.*

Note that $\{T_j : j \geq 0\}$ forms a prefix-free sequence, that $T^{\leq t} \subseteq \Omega^t$ is the set of all paths of length t that terminate within time t , while $\mathbf{c}(T_t) \subseteq \Omega^\infty$ is the set of infinite execution paths with termination at time t . The next definition introduces prefix-closed functions. These are functions f with a prefix-free support,

³Note that $w(\tilde{\omega})$ is well-defined because $0 \leq \text{sc}(\omega_j) \leq 1$ for each $j \geq 0$, hence the series of partial products is nonincreasing.

⁴In the context of model checking, these languages arise as complements of Safety properties; see e.g. [5, Def.3.22].

condition (a), additionally satisfying two extra conditions. Condition (b) just states that the value of f on its support is determined by a finite prefix of the input sequence. Condition (c), T-respectfulness, means that a trace that terminates at time j ($\omega^j \in T_j$) cannot lead to $\text{supp}(f)$ at a later time ($\omega^j \notin L^{>j}$). This is a consistency condition, formalizing that the value of f does not depend on, so to speak, what happens *after* termination.

Definition 7 (prefix-closed function) Let $f : \Omega^\infty \rightarrow \overline{\mathbb{R}}^+$ be a nonnegative measurable function and (L_0, L_1, \dots) be a prefix-closed sequence. We say f is a prefix-closed function with branches L_0, L_1, \dots if the following conditions are satisfied.

- (a) $\text{supp}(f)$ is prefix-free with branches L_j ($j \geq 0$).
- (b) for each $j \geq 0$ and $\omega^j \in L_j$, f is constant on $\mathfrak{c}(\{\omega^j\})$.
- (c) $\text{supp}(f)$ is T-respectful: for each $j \geq 0$, $L^{>j} \cap T_j = \emptyset$.

Note that there may be different prefix-free sequences w.r.t. which f is prefix-closed.

Example 4 The indicator function 1_{T_f} is clearly a prefix closed, measurable function with $\text{supp}(1_{T_f}) = T_f$ and branches $L_j = T_j$. For more interesting examples, consider the PPG in Example 3 and the functions f_1 , that returns the time the process terminates, and f_2 that returns the value of d at termination. Here $\text{supp}(f_1) = T_f$ has branches $L_j = T_j$ ($j \geq 0$), instead $\text{supp}(f_2) = \{\tilde{\omega} \in T_f : \text{the first terminated state } \omega \text{ in } \tilde{\omega}, \text{ if it exists, has } \omega(1) = d \in (0, 2]\}$, and $L_0 = \emptyset$, $L_j = (T^c)^{j-1} \cdot (T \cap ((0, 2] \times \overline{\mathbb{R}}^4))$ ($j \geq 1$).

We will now study how to consistently approximate infinite computations (μ_S^∞ semantics) with finite ones (μ_S^t semantics). This will lead to the main result of this section (Theorem 2). As a first step, let us introduce an appropriate notion of finite approximation for functions f defined on the infinite product space Ω^∞ . Fix an arbitrary element $\star \in \Omega$. For each $f : \Omega^\infty \rightarrow \overline{\mathbb{R}}^+$ and $t \geq 1$, let us define the function $f_t : \Omega^t \rightarrow \overline{\mathbb{R}}^+$ by $f_t(\omega^t) := f(\omega^t, \star^\infty)$. The intuition here is that, for a prefix-closed function f , the function f_t approximates correctly f for all finite paths in the L_j -branches of f , for $j \leq t$. Consider for instance the function $f = f_1$ in Example 4. On $L^{\leq t}$, the approximation f_t gives the correct value w.r.t. f in a precise sense: $f_t(\omega^t) = f(\omega^t, \star^\infty) = f(\omega^t, \tilde{\omega}')$ whatever \star and $\tilde{\omega}'$. On the other hand, for finite paths $\omega^t \in L^{>t}$, f_t may not approximate f correctly: we may have $f_t(\omega^t) = f(\omega^t, \star^\infty) \neq f(\omega^t, \tilde{\omega}')$ depending on the specific \star and $\tilde{\omega}'$. The catch is, as t grows large, the set $L^{>t}$ will become thinner and thinner — at least under reasonable assumptions on the measure μ_S^∞ .

It is not difficult to check that, for any t , f_t is measurable over Ω^t . The next result shows how to approximate $[[S]]f$ with quantities defined *only in terms of* f_t, w_t and μ_S^t , which is the basis for the sampling-based inference algorithm in the next section. Formally, for $t \geq 1$ and a measurable function $h : \Omega^t \rightarrow \overline{\mathbb{R}}^+$, we let

$$[S]^t h := E_{\mu_S^t}[h] = \int \mu_S^t(d\omega^t) h(\omega^t).$$

The intuition of the theorem is as follows. Consider a prefix closed function f with branches L_0, L_1, \dots . For any time t , it is not difficult to see that $\mathfrak{c}(L^{\leq t} \cap T^{\leq t}) \subseteq \text{supp}(f) \subseteq \mathfrak{c}(L^{\leq t} \cap T^{\leq t}) \cup (\mathfrak{c}(T^{\leq t}))^c$ (the last inclusion involves T-respectfulness). Since f_t approximates correctly f on $L^{\leq t}$, one sees that the first inclusion leads to the lower bound $[S]^t f_t \cdot 1_{L^{\leq t} \cap T^{\leq t}} \cdot w_t \leq [S]f w$. As for the upper bound, the intuition is that, over $(\mathfrak{c}(T^{\leq t}))^c$, f is upper-bounded by M .

Theorem 2 (finite approximation) Consider $S \in \mathcal{P}$ and $t \geq 1$ such that $[S]^t 1_{T^{\leq t}} \cdot w_t > 0$. Then for any prefix-closed function f with branches L_0, L_1, \dots we have that $[[S]]f$ is well defined. Moreover, given an upper bound $f \leq M$ ($M \in \overline{\mathbb{R}}^+$), for each t large enough and $\alpha_t := \frac{[S]^t w_t}{[S]^t 1_{T^{\leq t}} \cdot w_t}$ we have:

$$\frac{[S]^t f_t \cdot 1_{L^{\leq t} \cap T^{\leq t}} \cdot w_t}{[S]^t w_t} \leq [[S]]f \leq \frac{[S]^t f_t \cdot 1_{L^{\leq t} \cap T^{\leq t}} \cdot w_t}{[S]^t w_t} \alpha_t + M \cdot (\alpha_t - 1). \quad (6)$$

When f is an indicator function, $f = 1_A$, we can of course take $M = 1$ in the theorem above. We first illustrate the above result with a simple example.

Example 5 Consider the PPG of Example 2 (Fig. 1, left). We ask what is the expected value of c upon termination of this program. Formally, we consider the program checkpoint $S = 0$, and the function f on traces that returns the value of c on the first terminated state, if any, and 0 elsewhere. This f is clearly prefix-closed with branches $L_j \subseteq T_j$. We apply Theorem 2 to $[[S]]f$. Fixing the time $t = 4$, we can calculate easily the quantities involved in the approximation of $[[S]]f$ in (6). In doing so, we must consider the finitely many paths of length t of nonzero probability and weight (there are only two of them), their weights and the value of c on their final state when terminated⁵.

$$\begin{aligned} [S]^t f_t \cdot 1_{T \leq t} \cdot w_t &= 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} & [S]^t w_t &= \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4} \\ [S]^t 1_{T \leq t} \cdot w_t &= \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4} & \alpha_t &= 1. \end{aligned}$$

Then, with $M = 1$, the lower and upper bounds in (6) coincide and yield $[[S]]f = \frac{1}{3}$. If we remove conditioning on node 4, then all the paths of length t have weight 1, and a similar calculation yields $[[S]]f = \frac{1}{2}$.

In more complicated cases, we may not be able to calculate exactly the quantities involved in (6), but only to estimate them via sampling. To this purpose, we will introduce Feynman-Kac models and the Particle Filtering algorithm in the next section. The theorem below confirms that the bounds established above are asymptotically tight, at least under the assumption that the program $S \in \mathcal{P}$ terminates with probability 1. In this case, in fact, the probability mass outside $T^{\leq t}$ tends to 0, which leads the lower and the upper bound in (6) to coincide. Moreover, we get a simpler formula in the special case when termination is guaranteed to happen within a fixed time limit; for instance, in the case of acyclic⁶ PPGs.

Theorem 3 (tightness) Assume the same hypotheses as in Theorem 2. Further assume that $\mu_S^\infty(T_f) = 1$. Then both the lower and the upper bounds in (6) tend to $[[S]]f$ as $t \rightarrow +\infty$. In particular, if for some $t \geq 1$ we have $[S]^t 1_{T \leq t} = 1$, then

$$[[S]]f = \frac{[S]^t f_t \cdot w_t}{[S]^t w_t}. \quad (7)$$

Example 6 For the PPG of Example 2 one has $\mu_S^\infty(T_f) = 1$. As already seen in Example 2, lower and upper bounds coincide for $t \geq 4$.

A practically relevant class of closed prefix functions are those where the result $f(\tilde{\omega})$ only depends on computing a function h , defined on Ω , on the first terminated state, if any, of the sequence $\tilde{\omega}$. This way h is lifted to Ω^∞ . This case covers all the examples seen so far. We formally introduce lifting below. Recall that for $t \geq 1$, $T_t = (T^c)^{t-1} \cdot T$.

Definition 8 (lifting) Let $h : \Omega \rightarrow \overline{\mathbb{R}}^+$ a nonnegative measurable function such that $\text{supp}(h) \subseteq T$. The lifting of h is the measurable function $\check{h} : \Omega^\infty \rightarrow \overline{\mathbb{R}}^+$ defined as follows for each $\tilde{\omega} = (\omega_1, \omega_2, \dots)$: $\check{h}(\tilde{\omega}) := \sum_{t \geq 1} 1_{c(T_t)}(\tilde{\omega}) \cdot h(\omega_t)$.

Clearly, any \check{h} is prefix closed with branches $L_0 = \emptyset$ and $L_j = (T^c)^{j-1} \cdot \text{supp}(h) \subseteq T_j$ for $j \geq 1$. In particular, $\text{supp}(\check{h}) \subseteq T_f$. As an example, the indicator function for the set of paths that eventually terminate, $\check{h} = 1_{T_f}$, is clearly the lifting of $h = 1_T$; the functions f_1, f_2 in Example 4 can also be obtained by lifting (details omitted).

⁵Here, we also use the fact that $f_t \cdot 1_{T \leq t} = f_t \cdot 1_{L \leq t \cap T \leq t}$, a consequence of $L_j \subseteq T_j$ for all j s.

⁶Or, more accurately, PPGs where the only loop is the self-loop on the nil state.

5 Feynman-Kac models

In the field of Sequential Monte Carlo methods, Feynman-Kac (FK) models [19, Ch.9] are characterized by the use of *potential* functions. A potential in a Feynman-Kac model is a function that assigns a weight $G_t(x)$ to a *particle* (instance of a random process) in state x at time t . This weight represents how plausible or fit x is at time t based on some observable or conditioning. In other words, G_t modifies the *importance* of particles as the system evolves. For instance, in a model for tracking an object, the potential function could depend on the distance between the predicted particle position and the actual observed position. Particles closer to the observed position get higher weights.

FK models and probabilistic program semantics We first introduce FK models in a general context. Our formulation follows closely [19, Ch.9]. Throughout this and the next section, we let $t \geq 1$ be an arbitrary fixed integer.

Definition 9 (Feynman-Kac models) A Feynman-Kac (FK) model is a tuple $\text{FK} = (X, t, \mu^1, \{K_i\}_{i=2}^t, \{G_i\}_{i=1}^t)$, where $X = \overline{\mathbb{R}}^\ell$ for some $\ell \geq 1$, μ^1 is a probability measure on X and, for $i = 2, \dots, t$: K_i is a Markov kernel from X to X , and $G_i : X \rightarrow \overline{\mathbb{R}}^+$ is a measurable function.

Let μ^t denote the unique product measure on X^t induced by μ^1, K_2, \dots, K_t as per Theorem 1. Let $G := \prod_{i=1}^t G_i$. Provided $0 < E_{\mu^t}[G] < +\infty$, the Feynman-Kac measure induced by FK is defined by the following, for every measurable $A \subseteq X^t$:

$$\phi_{\text{FK}}(A) := \frac{E_{\mu^t}[1_A \cdot G]}{E_{\mu^t}[G]}. \quad (8)$$

We will refer to G in the above definition as the *global potential*. Equality (8) easily generalizes to expectations taken according to ϕ_{FK} . That is, for any measurable nonnegative function g on X^t , we can easily show that:

$$E_{\phi_{\text{FK}}}[g] = \frac{E_{\mu^t}[g \cdot G]}{E_{\mu^t}[G]}. \quad (9)$$

In what follows, we will suppress the subscript FK from ϕ_{FK} in the notation, when no confusion arises. Comparing (9) against the definition (5) suggests that the global potential G should play in FK models a role analogous to the weight function w in probabilistic programs. Note however that there is a major technical difference between the two, because FK models are only defined for a finite time horizon model given by t . A reconciliation between the two is possible thanks to the finite approximation theorem seen in the last section; this will be elaborated further below (see Theorem 4).

We will be particularly interested in the t -th *marginal* of ϕ , that is the probability measure on X defined as ($A \subseteq X$ measurable):

$$\phi_t(A) := \phi(X^{t-1} \times A) = E_\phi[1_{X^{t-1} \times A}]. \quad (10)$$

The measure ϕ_t is called *filtering* distribution (at time t), and can be effectively computed via the Particle Filtering algorithm described in the next subsection. Now let $\mathbf{G} = (\mathcal{P}, E, \text{nil}, \text{sc})$ be an arbitrary fixed PPG. Comparing (9) against e.g. the lower bound in (6) suggests considering the following FK model associated with \mathbf{G} and a checkpoint S .

Definition 10 (FK_S model) Let $t \geq 1$ be an integer and S a program checkpoint of \mathbf{G} . We define FK_S as the FK model where: $X = \Omega$, $\mu^1 = \delta_{(0, S)}$, $K_i = \kappa$ ($i = 2, \dots, t$) and $G_i = \text{sc}$ ($i = 1, \dots, t$). We let ϕ_S denote the measure on Ω^t induced by FK_S .

We now restrict our attention to functions f that are the lifting of a nonnegative h defined on Ω . Let $\phi_{S,t}$ denote the filtering distribution of ϕ_S at time t obtained by (10). In the following theorem we express the bounds in (6) in terms of the measure $\phi_{S,t}$. The whole point and interest of this result is that the bounds are expressed directly as expectations; these are moreover taken w.r.t. a *1-dimensional* filtering distribution ($\phi_{S,t}$), rather than a t -dimensional one (μ_S^t). Importantly, there are well-known algorithms to estimate expectations under a filtering distribution, as we will see in the next subsection.

Theorem 4 (filtering distributions and lifted functions) *Under the same assumptions of Theorem 2, further assume that f is the lifting of h . Then $\alpha_t = \mathbb{E}_{\phi_{S,t}}[1_{\top}]^{-1}$ and*

$$\beta_L := \mathbb{E}_{\phi_{S,t}}[h] \leq [[S]]f \leq \mathbb{E}_{\phi_{S,t}}[h] \cdot \alpha_t + M \cdot (\alpha_t - 1) =: \beta_U. \quad (11)$$

Example 7 Consider again the PPG of Example 2. We can re-compute $[[S]]f$ relying on Theorem 4. Fix $t = 4$. We first compute the filtering distribution ϕ_t on $\mathcal{X} = \overline{\mathbb{R}}^3$ relying on its definition (10). Similarly to what we did in Example 5, we consider the nonzero-weight, nonzero-probability traces of length four. Then we project onto the final (fourth) state, and compute the weights of the resulting triples (c, d, S) , then normalize. There are only two triples (c, d, S) of nonzero probability: $\phi_t(0, 0, 2) = \frac{2}{3}$, $\phi_t(1, 1, 2) = \frac{1}{3}$. The function f considered in Example 5 is the lifting of the function $h(c, d, S) = c \cdot [S = 2]$ defined on $\mathcal{X} = \overline{\mathbb{R}}^3$. We apply Theorem 4 and get $\beta_L = \mathbb{E}_{\phi_t}[h] = \frac{1}{3} \leq [[S]]f$. Moreover $\mathbb{E}_{\phi_t}[1_{\top}] = 1$, hence $\alpha_t = 1$ according to Theorem 4. Hence $\beta_L = \beta_U = [[S]]f = \frac{1}{3}$. This coincides with what found in examples 5 and 6.

We can apply the above theorem to the functions described in Example 4 and to other computationally challenging cases: we will do so in Section 6, after introducing in the next section the Particle Filtering algorithm.

The Particle Filtering algorithm From a computational point of view, our interest in FK models lies in the fact that they allow for a simple, unified presentation of a class of efficient inference algorithms, known as *Particle Filtering (PF)* [19, 37, 50]. For the sake of presentation, we only introduce here the basic version, *Bootstrap PF*, following closely⁷ [19, Ch.11]. Fix a generic FK model, $\text{FK} = (\mathcal{X}, t, \mu^1, \{K_i\}_{i=2}^t, \{G_i\}_{i=1}^t)$. Fix $N \geq 1$, the number of *particles*, that is instances of the random process represented by the K_i 's, we want to simulate. Let $W = W^{1:N} = (W^{(1)}, \dots, W^{(N)})$ be a tuple of N real nonnegative random variables, the *weights*. Denote by \widehat{W} the normalized version of W , that is $\widehat{W}^{(i)} = W^{(i)} / (\sum_{j=1}^N W^{(j)})$. A *resampling scheme* for (N, W) is a N -tuple of random variables $R = (R_1, \dots, R_N)$ taking values on $1..N$ and depending on W , such that, for each $1 \leq i \leq N$, one has: $\mathbb{E}[\sum_{j=1}^N 1_{R^{(j)}=i} | W] = N \cdot \widehat{W}^{(i)}$. In other words, each index $i \in 1..N$ on average is selected in R a number of times proportional to its weight in W . We shall write $R(W)$ to indicate that R depends on a given weight vector W . Various resampling schemes have been proposed in the literature, among which the simplest is perhaps *multinomial resampling*; see e.g. [19, Ch.9] and references therein. Algorithm 1 presents a generic PF algorithm. Resampling here takes place at step 4: its purpose is to give more importance to particles with higher weight, when extracting the next generation of N particles, while discarding particles with lower weight.

The justification and usefulness of this algorithm is that, under mild assumptions, for any measurable function h defined on \mathcal{X} , expectation under ϕ_t , the filtering distribution on \mathcal{X} at time t , in the limit can be expressed a weighted sum with weights $\widehat{W}_t^{(j)}$:

$$\sum_{j=1}^N \widehat{W}_t^{(j)} \cdot h(X_t^{(j)}) \longrightarrow \mathbb{E}_{\phi_t}[h] \quad \text{a.s. as } N \longrightarrow +\infty. \quad (12)$$

⁷Additional details in [15].

Algorithm 1 A generic PF algorithm

Input: $\text{FK} = (\mathcal{X}, t, \mu^1, \{K_k\}_{k=2}^t, \{G_k\}_{k=1}^t)$, a FK model; $N \geq 1$, number of particles.
Output: $X_t^{1:N} \in \mathcal{X}^N$, $W_t^{1:N} \in \mathbb{R}^{+N}$.

```

1:  $X_1^{(j)} \sim \mu^1$  ▷ state initialisation
2:  $W_1^{(j)} := G_1(X_1^{(j)})$  ▷ weight initialisation
3: for  $k = 2, \dots, t$  do
4:    $r_{1:N} \sim R(W_{k-1}^{1:N})$  ▷ resampling
5:    $X_k^{(j)} \sim K_k(X_{k-1}^{(r_j)})$  ▷ state update
6:    $W_k^{(j)} := G_k(X_k^{(j)})$  ▷ weight update
7: end for
8: return  $(X_t, W_t)$ 

```

The practical implication here is that we can estimate quite effectively the expectations involved in (11), for $\phi_t = \phi_{S,t}$, as weighted sums like in (12). Note that in the above consistency statement t is held fixed — it is one of the parameter of the FK model — while the number of particles N tends to $+\infty$.

6 Implementation and experimental validation

Implementation The PPG model is naturally amenable to a vectorized implementation of PF that leverages the fine-grained, SIMD parallelism existing at the level of particles. At every iteration, the state of the N particles, $\omega^N = (\omega_1, \dots, \omega_N)$ with $\omega_i = (v_i, z_i) \in \overline{\mathbb{R}}^{m+1}$, will be stored using a pair of arrays (V, Z) of shape $N \times m$ and $N \times 1$, respectively. The weight vector is stored using another array W of shape $N \times 1$. We rely on vectorization of operations: for a function $f : \overline{\mathbb{R}}^k \rightarrow \overline{\mathbb{R}}$ and a $N \times k$ array U , $f(U)$ will denote the $N \times 1$ array obtained by applying f to each row of U . In particular, we denote by $(Z = s)$ (for any $s \in \mathbb{N}$) the $N \times 1$ array obtained applying element-wise the indicator function $1_{\{s\}}$ to Z element-wise, and by $\varphi(V)$ the $N \times 1$ array obtained by applying the predicate φ to V . For U a $N \times k$ array and W a $N \times 1$ array, $U * W$ denotes the $N \times k$ array obtained by multiplying the j th row of U by the j th element of W , for $j = 1, \dots, N$: when W is a 0/1 vector, this is an instance of *boolean masking*. Abstracting the vectorization primitives of modern CPUs and programming languages, we model the assignments of a vector to an array variable as a single instruction, written $U := Z$. The usual rules for broadcasting scalars to vectors apply, so e.g. $V := S$ for $S \in \overline{\mathbb{R}}$ means filling V with S . Likewise, for ζ a parametric distribution, $U \sim \zeta(V)$ means sampling N times independently from $\zeta(v_1), \dots, \zeta(v_N)$, and assigning the resulting matrix to U : this too counts as a single instruction.

Based on the above idealized model of vectorized computation, we present VPF, a vectorized version of the PF algorithm for PPGs, as Algorithm 2. Here it is assumed that $\mathcal{P} \subseteq \mathbb{N}$, while $\text{sc}(s) = \gamma_s$. On line 4, $\text{Resampling}(\cdot)$ denotes the result of applying a generic resampling algorithm based on weights W to the current particles' state, represented by the pair of vectors (V, Z) . With respect to the generic PF Algorithm 1, here in the returned output, (V, Z) corresponds to X_t and W to W_t .

Experimental validation We illustrate some experimental results obtained with a proof-of-concept TensorFlow-based [1] implementation of Algorithm 2 (VPF). We have considered a number of challenging probabilistic programs that feature conditioning inside loops. For all these programs, we will estimate $[[S]]f$, for given functions f , relying on the bounds provided by Theorem 4 in terms of expectations w.r.t. filtering distributions. Such expectations will be estimated via VPF. We also compare VPF with two state-of-the-art PPLs, webPPL [24] and CorePPL [35]. In [35], a comparison of CorePPL

Algorithm 2 VPF, a Vectorized PF algorithm for PPGs.

Input: $G = (\mathcal{P}, E, \text{nil}, \text{sc})$, a PPG; $S \in \mathcal{P}$, initial program checkpoint; $t \geq 1$, time horizon; $N \geq 1$, number of particles.

Output: $V \in \mathbb{R}^{m \times N}$, $Z, W \in \mathbb{R}^{1 \times N}$.

```

1:  $V := S$ ;  $Z := S$                                 ▶ state initialisation
2:  $W := \gamma_S(Z)$                                 ▶ weight initialisation
3: for  $t - 1$  times do
4:    $(V, Z) := \text{Resampling}((V, Z), W)$                                 ▶ resampling
5:   for  $(s, \varphi, \zeta, s') \in E$  do
6:      $M_{s, \varphi} := \varphi(V) * (Z = s)$                                 ▶ mask computation
7:   end for
8:    $V \sim \sum_{(s, \varphi, \zeta, s') \in E} \zeta(V) * M_{s, \varphi}$ ;  $Z := \sum_{(s, \varphi, \zeta, s') \in E} s' \cdot M_{s, \varphi}$                                 ▶ state update
9:    $W := \sum_{s \in \mathcal{P}} \gamma_s(V) * (Z = s)$                                 ▶ weight update
10: end for
11: return  $(V, Z, W)$ 

```

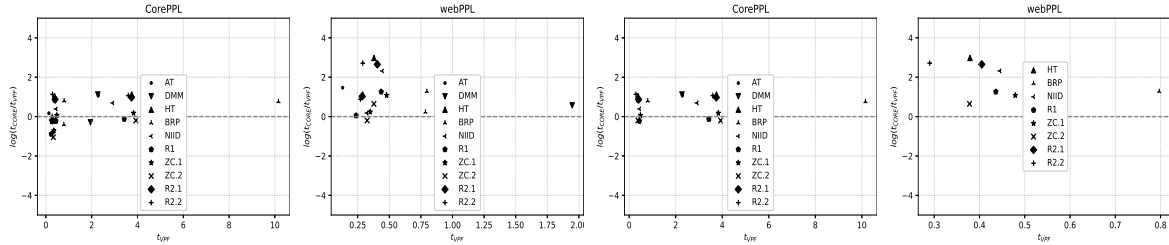


Figure 2: For $N = 10^5, 10^6$, scatterplots of the log-ratios of execution times, $\log_{10}(\text{time}_{\text{tool}}/\text{time}_{\text{VPF}})$, based on the data points of Table 1. From left to right: $N = 10^5$, $\text{tool} = \text{CorePPL}$; $N = 10^5$, $\text{tool} = \text{WebPPL-smc}$; $N = 10^6$, $\text{tool} = \text{CorePPL}$; $N = 10^6$, $\text{tool} = \text{WebPPL-smc}$. For $N = 10^6$, the vast majority of the data points lie above the x-axis, indicating superior performance of VPF across different examples.

with webPPL, Pyro [12] and other PPLs in terms of performance shows the superiority of CorePPL SMC-based inference across a number of benchmarks.

At least for $N \geq 10^5$, the tools tend generally to return similar estimates of the expected value, which we take as an empirical evidence of accuracy. Additional insight into accuracy is obtained by directly comparing the results of VPF with those of webPPL-rejection (when available), which is an exact inference algorithm. The expected values estimated by webPPL-rejection are consistently in line to those of VPF. In terms of performance, a graphical representation of our results is provided in Figure 2, with scatterplots showing the ratio of execution times ($\text{time}_{\text{other-tool}}/\text{time}_{\text{VPF}}$) on a log scale. In the case of WebPPL, nearly all data points lie above the x-axis, indicating superiority of VPF. In the case of CorePPL, for $N = 10^5$ the data points are quite uniformly distributed across the x-axis, indicating basically a tie. For $N = 10^6$, we have a majority of points above the x-axis, indicating again superiority of VPF by *orders of magnitude*; additional details can be found in Appendix A.

7 Conclusion

We study correct and efficient implementations of Sequential Monte Carlo inference algorithms for universal probabilistic programs. We offer a clean trace-based operational semantics for PPGs, a finite approximation theorem and consistency of the PF algorithm via a connection to FK models. Experiments conducted with VPF, a vectorized version of PF tailored to PPGs, show very promising results.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Josh Levenberg, Manjunath Kudlur, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu & Xiaoqiang Zheng (2016): *TensorFlow: A System for Large-Scale Machine Learning*. *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16)*. arXiv:1605.08695.
- [2] Robert B. Ash (1972): *Real Analysis and Probability*. Academic Press, Inc., New York, NY, USA, doi:10.1016/C2013-0-07155-1.
- [3] Martin Avanzini, Georg Moser & Michael Schaper (2023): *Automated Expected Value Analysis of Recursive Programs*. *roc. ACM Program. Lang.*, 7, PLDI, doi:10.1145/3591263.
- [4] Alexander Bagnall, Gordon Stewart & Anindya Banerjee (2023): *Formally Verified Samplers from Probabilistic Programs with Loops and Conditioning*. *Proc. ACM Program. Lang.*, 7, 1–24, doi:10.1145/3591220.
- [5] Christel Baier & Joost-Pieter Katoen (2008): *Principles of model checking*. MIT press.
- [6] Jialu Bao, Nitesh Trivedi, Drashti Pathak, Justin Hsu & Subhajit Roy (2022): *Data-driven invariant learning for probabilistic programs*. In: *CAV. Lecture Notes in Computer Science*, vol. 13371, pp. 33–54. Springer, doi:0.1007/s10703-024-00466-x.
- [7] Gilles Barthe, Joost-Pieter Katoen & Alexandra Silva (2020): *Foundations of Probabilistic Programming*. Cambridge University Press, doi:10.1017/9781108770750.
- [8] Ezio Bartocci, Laura Kovács & Miroslav Stankovic (2020): *Mora-automatic generation of moment-based invariants*. *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 492–498, doi:10.1007/978-3-030-45190-5_28.
- [9] Kevin Batz, Mingshuai Chen, Sebastian Junges, Benjamin Lucien Kaminski, Joost-Pieter Katoen & Christoph Matheja (2023): *Probabilistic Program Verification via Inductive Synthesis of Inductive Invariants*. In Sankaranarayanan, S., Sharygina, N. (eds) *Tools and Algorithms for the Construction and Analysis of Systems. TACAS 2023. Lecture Notes in Computer Science*, vol 13994. Springer, doi:10.1007/978-3-031-30820-8_25.
- [10] Raven Beutner, Luke Ong & Fabian Zaiser (2022): *Guaranteed bounds for posterior inference in universal probabilistic programming*. In *Proceedings of the 43rd ACM SIGPLAN international conference on programming language design and implementation, (PLDI 22)*, doi:10.1145/3519939.3523721.
- [11] Sooraj Bhat, Johannes Borgström, Andrew D. Gordon & Claudio Russo (2013): *Deriving probability density functions from probabilistic functional programs*. In *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2013)*, doi:10.23638/LMCS-13(2:16)2017.
- [12] Eli Bingham, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall & Noah D Goodman (2019): *Pyro: Deep universal probabilistic programming*. *The Journal of Machine Learning Research* 20, 1, 973–978.
- [13] Michele Boreale & Luisa Collodi (2023): *Guaranteed inference for probabilistic programs: a parallelisable, small-step operational approach*. *Lecture Notes in Computer Science*, Vol. 14500, Springer, doi:0.1007/978-3-031-50521-8_7.
- [14] Michele Boreale & Luisa Collodi (2025): *Code for the experiments described in the present paper*. <https://github.com/Luisa-unifi/Vectorized-Particle-Filtering>.
- [15] Michele Boreale & Luisa Collodi (2025): *Full version of the present paper*. Available from <https://github.com/Luisa-unifi/Vectorized-Particle-Filtering>.
- [16] Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon & Marcin Szymczak (2016): *A lambda-calculus foundation for universal probabilistic programming*. *ACM SIGPLAN Notices* 51.9: 33–46, doi:10.1145/3022670.2951942.

- [17] Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li & Allen Riddell (2017): *Stan: A probabilistic programming language*. *Journal of Statistical Software* 76(1), doi:10.18637/jss.v076.i01.
- [18] Krishnendu Chatterjee, Amir K. Goharshady, Tobias Meggendorfer & Đorđe Žikelić. (2022): *Sound and Complete Certificates for Quantitative Termination Analysis of Probabilistic Programs*. In: Shoham, S., Vitzel, Y. (eds) *Computer Aided Verification. CAV 2022. Lecture Notes in Computer Science*, vol 13371. Springer, doi:10.1007/978-3-031-13185-1_4.
- [19] Nicolas Chopin & Omiros Papaspiliopoulos (2021): *An Introduction to Sequential Monte Carlo*. Springer Nature Switzerland AG, doi:10.1007/978-3-030-47845-2.
- [20] Marco F. Cusumano-Towner, Feras A. Saad, Alexander K. Lew & Vikash K. Mansinghka (2019): *Gen: A General-purpose Probabilistic Programming System with Programmable Inference*. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)* pp. 221–236. New York, NY, USA, doi:10.1145/3314221.3314642.
- [21] Víctor Elvira, Luca Martino & Christian P. Robert. (2018): *Rethinking the Effective Sample Size*. *arXiv: 1809.04129*, doi:10.1111/insr.12500.
- [22] Timon Gehr, Sasa Misailovic & Martin T. Vechev (2016): *Psi: Exact symbolic inference for probabilistic programs*. *Proceedings of the 28th International Conference in Computer Aided Verification (CAV 2016)*, Toronto, page 62–83, doi:10.1007/978-3-319-41528-4_4.
- [23] Noah D. Goodman, Vikash Mansinghka, Daniel Roy, Keith Bonawitz & Joshua Tenenbaum (2008): *Church: a language for generative models*. *Proc. Uncertainty in Artificial Intelligence*.
- [24] Noah D. Goodman & Andreas Stuhlmüller (2016): *The design and implementation of probabilistic programming languages*. *Proceedings of the 28th International Conference in Computer Aided Verification (CAV 2016)*, Toronto, page 62–83.
- [25] Andrew D. Gordon, Thomas A. Henzinger, Aditya V. Nori & Sriram K. Rajamani (2014): *Probabilistic programming*. *Proceedings of Future of Software Engineering Proceedings (FOSE2014)*, pages 167–181, doi:10.1145/2593882.2593900.
- [26] Maria I. Gorinova, Andrew D. Gordon & Charles Sutton (2019): *Probabilistic programming with densities in SlicStan: efficient, flexible, and deterministic*. *Proceedings of the ACM on Programming Languages* 3, POPL.1–30, doi:10.1145/3290348.
- [27] Maria I. Gorinova, Andrew D. Gordon, Charles Sutton & Matthijs Vákár (2021): *Conditional independence by typing*. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 44.1: 1–54, doi:10.1145/3490421.
- [28] Michikazu Hirata, Yasuhiko Minamide & Tetsuya Sato (2023): *Semantic Foundations of Higher-Order Probabilistic Programs in Isabelle/HOL*. *14th International Conference on Interactive Theorem Proving (ITP), 2023. Leibniz International Proceedings in Informatics (LIPIcs)*, Volume 268, pp. 18:1–18:18, Schloss Dagstuhl – Leibniz-Zentrum für Informatik, doi:10.4230/LIPIcs.ITP.2023.18.
- [29] Benjamin L. Kaminski (2019): *Advanced weakest precondition calculi for probabilistic programs*. *Doctoral thesis (Ph.D)*, RWTH Aachen University, doi:10.18154/RWTH-2019-01829.
- [30] Dexter Katoen (1981): *Semantics of probabilistic programs*. *Journal of Computer and System Sciences*, 22(3):328–350, doi:10.1016/0022-0000(81)90036-2.
- [31] Lutz Klinkenberg, Christian Blumenthal, Mingshuai Chen, Darion Haase & Joost-Pieter Katoen (2024): *Exact Bayesian Inference for Loopy Probabilistic Programs using Generating Functions*. In *Proceedings of the ACM on Programming Languages (OPSLA)*, Volume 8: 127,923–953, doi:10.1145/3649844.
- [32] Alexander K. Lew, George Matheos, Tan Zhi-Xuan, Matin Ghavamizadeh, Nishad Gothoskar, Stuart Russell & Vikash K. Mansinghka (2023): *SMCP3: Sequential Monte Carlo with Probabilistic Program Proposals*. *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics (AISTATS)*. PMLR 206:7061–7088.

- [33] Daniel Lundén, Johannes Borgström & David Broman (2021): *Correctness of sequential Monte Carlo inference for probabilistic programming languages*. *Programming Languages and Systems*, pp. 404–431. Springer International Publishing, doi:10.1007/978-3-030-72019-3_15.
- [34] Daniel Lundén, David Broman, Fredrik Ronquist & Lawrence M. Murray (2023): *Automatic alignment of sequential Monte Carlo inference in higher-order probabilistic programs*. In *ESOP*, pages 535–563.
- [35] Daniel Lundén, Joey Ohman, Jan Kudlicka, Viktor Senderov, Fredrik Ronquist & David Broman (2022): *Compiling universal probabilistic programming languages with efficient parallel sequential Monte Carlo inference*. In *ESOP*, pages 29–56, doi:10.1007/978-3-030-99336-8_2.
- [36] Tshepo Mokoena (2014): *Urban Myths: Are you never more than 6ft from a rat in a city?* *The Guardian online*, <https://www.theguardian.com/cities/2014/feb/13/urban-myths-6ft-from-a-rat>.
- [37] Pierre Del Moral, Arnaud Doucet & Ajay Jasra (2006): *Sequential Monte Carlo Samplers*. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, Vol. 68, No. 3, pp. 411–436, doi:0.1111/j.1467-9868.2006.00553.x.
- [38] Praveen Narayanan, Jacques Carette, Chungchieh Shan Wren Romano & Robert Zinkov (2016): *Probabilistic inference by program transformation in Hakaru (system description)*. In *Proceedings of the 13th International Symposium on Functional and Logic Programming (FLOPS 2016)*, pages 62–79, doi:10.1007/978-3-319-29604-3_5.
- [39] Praveen Narayanan & Chung chieh Shan (2020): *Symbolic disintegration with a variety of base measures*. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 42 (2), 1–60, doi:10.1145/3374208.
- [40] Aditya V. Nori, Chung-Kil Hur, Sriram K. Rajamani & Selva Samuel (2014): *R2: An efficient mcmc sampler for probabilistic programs*. *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, page 2476–2482, doi:10.1609/aaai.v28i1.9060.
- [41] Alexey Radul & Boris Alexeev (2021): *The Base Measure Problem and its Solution*. *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, PMLR 130:3583–3591.
- [42] Francesca Randone, Luca Bortolussi, Emilio Incerto & Mirco Tribastone (2024): *Inference of Probabilistic Programs with Moment-Matching Gaussian Mixtures*. *Proceedings of the ACM on Programming Languages (POPL)*:1882–1912, doi:10.1145/3632905.
- [43] Sriram Sankaranarayanan, Aleksandar Chakarov & Sumit Gulwani (2013): *Static analysis for probabilistic programs: Inferring whole program properties from finitely many paths*. *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation.*, doi:10.1145/2499370.2462179.
- [44] Adam Scibior, Ohad Kammar, Matthijs Vákár, Sam Staton, Hongseok Yang, Yufei Cai, Klaus Ostermann, Sean K. Moss, Chris Heunen & Zoubin Ghahramani (2017): *Denotational validation of higher-order Bayesian inference*. *Proceedings of the ACM on Programming Languages* 2. POPL. 1–29.
- [45] Sam Staton (2017): *Commutative semantics for probabilistic programming*. *Proceedings of the 26th European Symposium on Programming (ESOP2017)*, Uppsala, Sweden, doi:10.1007/978-3-662-54434-1_32.
- [46] Sam Staton, Frank Wood, Hongseok Yang, Chris Heunen & Ohad Kammar (2016): *Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints*. *Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, doi:10.1145/2933575.2935313.
- [47] Zachary Susag, Sumit Lahiri, Justin Hsu & Subhajit Roy (2022): *Symbolic execution for randomized programs*. *Proc. ACM Program. Lang.*, 6, OOPSLA, 1583–1612, doi:10.1145/3563344.
- [48] Joseph Tassarotti & Jean-Baptiste Tristan (2023): *Verified density compilation for a probabilistic programming language*. *Proceedings of the ACM on Programming Languages* 7, PLDI. 615–637, doi:10.1145/3591245.
- [49] Peixin Wang, Tengshun Yang, Hongfei Fu & C.-H. Luke Ong Guanyan Li (2024): *Static Posterior Inference of Bayesian Probabilistic Programming via Polynomial Solving*. *Proceedings of the ACM on Programming Languages*, Volume 8, Issue PLDI’24, 202, 1361 – 1386, doi:10.1145/3656432.

- [50] Frank Wood, Jan Willem Meent & Vikash Mansinghka (2014): *A New Approach to Probabilistic Programming Inference*. *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, PMLR 33:1024-1032.
- [51] Yi Wu, Siddharth Srivastava, Nicholas Hay, Simon Du & Stuart Russell (2018): *Discrete-Continuous Mixtures in Probabilistic Programming: Generalized Semantics and Inference Algorithms*. *Proceedings of the 35th International Conference on Machine Learning*, PMLR 80:5343-5352.

A Experimental validation

We illustrate some experimental results obtained with a proof-of-concept TensorFlow-based [1] implementation of Algorithm 2. We still refer to this implementation as VPF. We have considered a number of challenging probabilistic programs that feature conditioning inside loops. For all these programs, we will estimate $[[S]]f$, for given functions f , relying on the bounds provided by Theorem 4 in terms of expectations w.r.t. filtering distributions. Such expectations will be estimated via VPF. We also compare VPF with two state-of-the-art PPLs, webPPL [24] and CorePPL [35]. webPPL is a popular PPL supporting several inference algorithms, including SMC, where resampling is handled via continuation passing. We have chosen to consider CorePPL as it supports a very efficient implementation of PF. In [35], a comparison of CorePPL with webPPL, Pyro [12] and other PPLs in terms of performance shows the superiority of CorePPL SMC-based inference across a number of benchmarks. As discussed in the Introduction, CorePPL’s implementation is based on a compilation into an intermediate format, conceptually similar to our PPGs⁸.

Models For our experiments we have considered the following programs: *Aircraft tracking* (AT, [51]), *Drunk man and mouse* (DMM, Example 3), *Hare and tortoise* (HT, e.g. [4]), *Bounded retransmission protocol* (BRP, [31]), *Non-i.i.d. loops* (NIID, e.g. [31]), the *ZeroConf* protocol (ZC, [9]), and two variations of *Random Walks*, RW1 ([13], Example 2) and RW2 in the following. In particular, AT is a model where a single aircraft is tracked in a 2D space using noisy measurements from six radars. HT simulates a race between a hare and a tortoise on a discrete line. BRP models a scenario where multiple packets are transmitted over a lossy channel. NIID describes a process that keeps tossing two fair coins until both show tails. ZC is an idealized version of the network connection protocol by the same name. RW1, RW2 are random walks with Gaussian steps. The pseudo-code of these models is reported in Appendix B.

These programs feature conditioning/scoring inside loops. In particular, DMM, HT and NIID feature unbounded loops: for these three programs, in the case of VPF we have truncated the execution after $k = 1000, 100, 100$ iterations, respectively, and set the time parameter t of Theorem 4 accordingly, which allows us to deduce bounds on the value of $[[S]]f$.⁹ For the other tools, we just consider the truncated estimate returned at the end of k iterations. AT, BRP, ZC, RW1 and RW2 feature bounded loops, but are nevertheless quite challenging. In particular, AT features multiple conditioning inside a for-loop, sampling from a mix of continuous and discrete distributions, and noisy observations. Below, we discuss the obtained experimental results in terms of accuracy, performance, scalability.

Accuracy We report in Table 1 the execution time, the estimated expected value and the Effective Sample Size (ESS, a measure of diversity of particles, the higher the better; see [15]) for VPF, CorePPL and webPPL, as the number N of particles increases. At least for $N \geq 10^5$, the tools tend generally to return similar estimates of the expected value, which we take as an empirical evidence of accuracy. Additional insight into accuracy is obtained by directly comparing the results of VPF with those of

⁸Direct compilation of CorePPL to GPU via the intermediate-level format RootPPL is also supported. However, the results we have obtained with RootPPL are generally worse in terms of execution time, and not presented here. Our PC configuration is as follows. OS: Windows 10; CPU: 2.8 GHz Intel Core i7; GPU: Nvidia T500, driver v. 522.06; TF: v. 2.10.1; CUDA Toolkit v. 11.8; cuDNN SDK v. 8.6.0.

⁹For the precise definition of f in each case, see Appendix B.

webPPL-rejection (when available), which is an exact inference algorithm. The expected values estimated by webPPL-rejection are consistently in line to those of VPF. In terms of ESS, the difference across the tools is significant. Except for model RW1, VPF yields ESS that are higher or comparable to those of the other tools.

Performance For larger values of N VPF generally outperforms the other considered tools in terms of execution time. The difference is especially noticeable for $N = 10^6$. A graphical representation of the data in Table 1 is provided in Figure 2 (see Figure 4 for an enlarged version of the plots), with scatterplots showing the ratio of execution times ($time_{other-tool}/time_{VPF}$) on a log scale. In the case of WebPPL, nearly all data points lie above the x-axis, indicating superiority of VPF. In the case of CorePPL, for $N = 10^5$ the data points are quite uniformly distributed across the x-axis, indicating basically a tie. For $N = 10^6$, we have a majority of points above the x-axis, indicating again superiority of VPF.

A closer look in the $N = 10^6$ case reveals that the only programs where CorePPL beats VPF are RW1 and ZC. This is most likely due to the low probability of conditioning in these programs; for instance in RW1 just a single final conditioning is performed. As in CorePPL resampling is only performed following a conditioning, this may explain its lower execution times in these cases. To further investigate this issue, we consider RW2, where the probability of conditioning is governed by a parameter $\lambda \in [0, 1]$, and run it for different values of λ . The obtained results are showed in Figure 3. We observe that for both CorePPL and WebPPL execution time tends to increase as the probability λ of conditioning increases; on the contrary, the execution time of VPF appears to be insensitive to λ . This suggests that VPF has a definite advantage over tools with explicit resample, on models with heavy conditioning.

Scalability The plot on the right shows the behaviour the *average unit cost (per particle)* of VPF, CorePPL and WebPPL across all the models we analyzed for $N = 10^3, \dots, 10^6$ on a log-scale. Here, for each N the average unit cost (in seconds) is $(t_1 + t_2 + \dots + t_k)/(N \cdot k)$, with t_i the execution time of the i -th example. We can observe that the cost of VPF decreases as the number of samples increases, whereas the cost of the other tools remains constant or increases (webPPL).

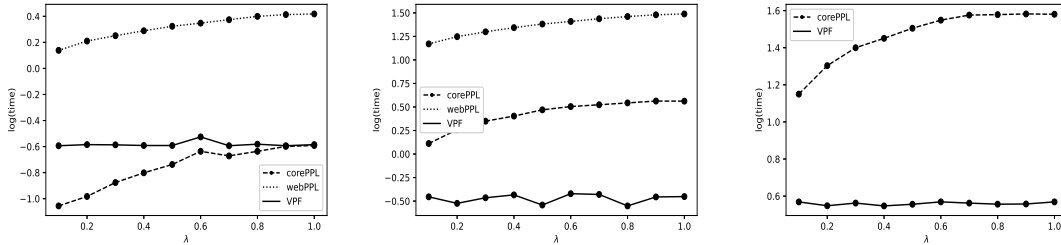
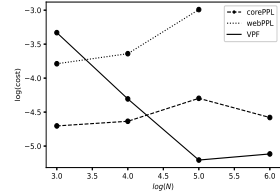


Figure 3: Execution times (in seconds) for the RW2 program, as a function of the probability λ of conditioning on external data for $N = 10^4$ (left), $N = 10^5$ (center) and $N = 10^6$ (right). webPPL missing from the right-most plot due to time-out. Execution times of VPF are basically insensitive to λ .

A.1 Table 1

		AT				DMM				HT				BRP				NIID			
		VPF	CorePPL	webPPL-smc	VPF	CorePPL	webPPL-smc	VPF	CorePPL	webPPL-smc	VPF	CorePPL	webPPL-smc	VPF	CorePPL	webPPL-smc	VPF	CorePPL	webPPL-smc		
$N = 10^3$	time	0.009	0.014	0.190	2.348	0.050	0.474	0.274	0.015	0.152	0.676	0.021	0.155	0.240	0.010	0.061					
	EV	6.805	6.955	6.696	0.478±0.110	0.501	0.427	32.834	33.683	32.368	0.018	0.016	0.023	3.594	2.694	3.473					
	ESS	1000	1000	999	994.6	726.9	973	955.0	758.9	951.1	1000	1000	974.5	1000	846.6	726.9					
$N = 10^4$	time	0.131	0.194	3.842	1.947	0.988	7.190	0.290	0.180	3.839	0.786	0.309	1.328	0.323	0.058	0.490					
	EV	6.817	6.967	6.760	0.539±0.115	0.498	0.481	32.725	33.474	32.702	0.029	0.025	0.024	3.364	2.766	3.417					
	ESS	10 ⁴	10 ⁴	9975	9984.5	7797.4	69417	9445.0	7692.9	9476.2	10 ⁴	10 ⁴	9745.8	10 ⁴	8555.6	7560.5					
$N = 10^5$	time	0.354	2.252	-	2.268	29.936	-	0.379	4.225	361.792	0.797	5.010	15.038	0.445	1.083	92.419					
	EV	6.818	6.970	-	0.537±0.120	0.507	-	33.128	33.545	32.560	0.024	0.025	0.026	3.467	2.772	3.430					
	ESS	10 ⁵	9.9e10 ⁴	-	≈ 10 ⁵	7.6e10 ⁴	-	≈ 10 ⁵	7.7e10 ⁴	≈ 10 ⁵	10 ⁵	10 ⁵	9.7e10 ⁴	10 ⁵	8.5e10 ⁴	7.6e10 ⁴					
$N = 10^6$	time	2.286	26.481	-	38.977	-	-	3.749	49.493	-	10.155	58.448	-	2.916	14.323	-					
	EV	6.834	6.980	-	0.541±0.111	-	-	33.432	33.606	-	0.024	0.025	-	3.413	2.774	-					
	ESS	10 ⁶	9.9e10 ⁵	-	≈ 10 ⁶	-	-	≈ 10 ⁶	7.7e10 ⁵	-	10 ⁶	10 ⁶	-	10 ⁶	8.5e10 ⁵	-					
webPPL-rej		-				0.494				32.683				0.023				3.414			

		RW1				ZC.1				ZC.2				RW2.1				RW2.2			
		VPF	CorePPL	webPPL	VPF	CorePPL	webPPL	VPF	CorePPL	webPPL	VPF	CorePPL	webPPL	VPF	CorePPL	webPPL	VPF	CorePPL	webPPL		
$N = 10^3$	time	0.192	0.009	0.045	0.206	0.018	0.083	0.262	0.016	0.049	0.231	0.024	0.232	0.231	0.021	0.187					
	EV	0.323	0.324	0.343	0.212	0.142	0.250	0.514	0.477	0.483	1.046	0.642	0.912	0.677	0.750	1.092					
	ESS	537.0	1000	46.739	1000	1000	392.2	1000	1000	245.2	992.0	780.0	479.9	999.0	997.0	133.9					
$N = 10^4$	time	0.238	0.031	0.269	0.349	0.068	0.610	0.325	0.029	0.207	0.285	0.186	3.043	0.271	0.275	2.081					
	EV	0.334	0.328	0.336	0.242	0.129	0.232	0.483	0.474	0.478	1.367	0.856	1.066	0.982	0.929	1.083					
	ESS	5163.0	10000	562.795	10000	10000	4263.0	10000	10000	2446.8	9967.0	7529.0	4026.6	9701.0	9350.9	582.5					
$N = 10^5$	time	0.436	0.260	7.956	0.479	0.558	5.778	0.378	0.243	1.673	0.405	3.003	181.802	0.290	3.887	149.831					
	EV	0.337	0.328	0.332	0.174	0.131	0.233	0.493	0.479	0.479	1.009	0.998	0.982	1.040	0.982	1.037					
	ESS	5.1e10 ⁴	10 ⁵	5.7e10 ³	10 ⁵	10 ⁵	4.2e10 ⁴	10 ⁵	10 ⁵	2.4e10 ⁴	≈ 10 ⁵	7.3e10 ⁴	4.7e10 ⁴	≈ 10 ⁵	9.5e10 ⁴	2.0e10 ³					
$N = 10^6$	time	3.422	2.569	-	3.829	5.790	-	3.928	2.485	-	3.742	35.271	-	3.595	42.134	-					
	EV	0.329	0.329	-	0.245	0.130	-	0.479	0.480	-	1.011	1.001	-	1.023	1.002	-					
	ESS	5.2e10 ⁵	10 ⁶	-	10 ⁶	10 ⁶	-	10 ⁶	10 ⁶	-	≈ 10 ⁶	7.3e10 ⁵	-	≈ 10 ⁶	9.4e10 ⁵	-					
webPPL-rej		-				0.235				0.479				1.022				1.061			

Table 1: Execution time (*time*) in seconds, estimated expected value (*EV*) and effective sample size ($ESS := (\sum_{i=1}^N W_i)^2 / (\sum_{i=1}^N W_i^2)$); the higher the better, see e.g. [21]) as the number of particles (N) increases, for VPF, CorePPL and webPPL, when applied on Aircraft tracking (AT), Drunk man and mouse (DMM), Hare and tortoise (HT), Bounded retransmission protocol (BRP), Non-i.i.d. loops (NIID), ZeroConf (ZC.1, ZC.2) and Random Walks (RW1 and RW2.1, RW2.2). For VPF, with reference to Theorem 4: for the bounded loops AT, BRP, RW1, RW2.1, RW2.2, ZC.1 and ZC.2, we have $EV = \beta_L = \beta_U$ (as $\alpha_i = 1$); for HT and NIID, we only provide β_L , as β_U is vacuous ($M = +\infty$). For DMM we give the midpoint of the interval $[\beta_L, \beta_U] \pm$ its half-width. Best results for *time* and *ESS* for each example and value of N are marked in **boldface**. Everywhere, ‘-’ means no result due to out-of-memory or timeout (500s). The results for DDM, especially for smaller values of N , exhibit a significant empirical variance: those reported in the table are obtained by averaging over 10 runs of each algorithm. Generally, there is an agreement across the tools about the estimates *EV*: an exception is NIID, where CorePPL returns values significantly different from the other tools’ and from the exact value $\frac{24}{7} = 3.428\dots$, cf. [31]. Also, for DMM the EV estimates returned by CorePPL and webPPL appear to be consistently lower than the midpoint of the interval returned by VPF.

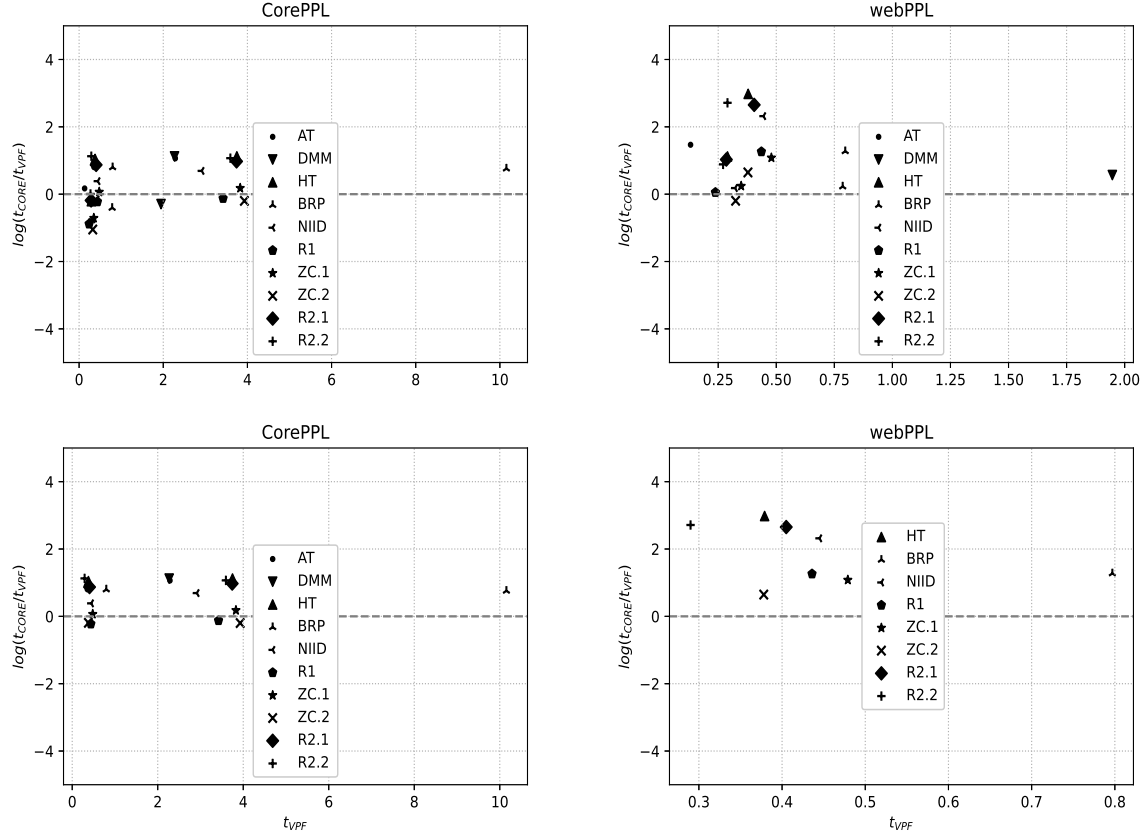


Figure 4: Enlarged version of plots in Fig. 2.

B Probabilistic programs pseudo-code

We consider the following probabilistic models described for example in [25]. For convenience, the programs are described in the language of [25], which is based on sequential composition, but they are easy to translate into our language. The `return` statement at the end of each program describes the function f considered in the estimation of $[[S]]f$; e.g. the DMM example, `return r` means f is the lifting of the function $(d, r, x, y, S) \mapsto r$.

```

1: while(time<=8){
2:   float[] radius;
3:   float obs-dist;
4:   if(time==0){
5:     x = Gaussian(2,1);
6:     y = Gaussian(-1.5,1);
7:   }else{
8:     x = Gaussian(x,2);
9:     y = Gaussian(y,2);
10:    for i in (0,6){
11:      d= compute-distance(i,x,y);
12:      if(d>radius[i]){
13:        flag=Bernoulli(0.999);
14:        if (flag==true){
15:          obs-dist=radius[i];
16:        }else{
17:          obs-dist=radius[i]+0.001*
18:            trunc-gauss(0,radius[i]);
19:        }
20:      }else{
21:        obs-dist=d+0.1*trunc-gauss(0,radius[i]);
22:      }
23:      obs-dist1 = Gaussian(obs-dist,0.01);
24:      observe(obs-dist1==...); //evidence numbers
25:      omitted
26:    }
27:    time=time+1;
28:  }
29: return x

```

(a) Aircraft Tracking (AT).

```

1: d=uniform(0,2);
2: r=uniform(0,1);
3: x=-1;
4: y=1;
5: while(|x-y|<1/10){
6:   x=Gaussian(x,d);
7:   y=Gaussian(y,r);
8:   observe(|x-y|<=3);
9: }
10: return r;
11:
12:
13:
14:
15:
16:
17:
18:
19:
20:
21:
22:
23:
24:
25:
26:
27:
28:
29:

```

(b) Drunk Man and Mouse (DMM).

```

1: initialPos=uniform(0,10);
2: tortoise=initialPos;
3: hare=0;
4: n=0;
5: while(hare<tortoise){
6:   n=n+1;
7:   tortoise=tortoise+1
8:   flag=Bernoulli(2/5);
9:   if (flag==true){
10:    hare=hare+Gaussian(4,2);
11:  }
12:  observe(|hare-tortoise|<=10);
13: }
14: observe((n>=20));
15: return hare;
16:
17:
18:

```

(c) Hare and Tortoise (HT).

```

1: initialPos=uniform(0,10);
2: s=100;
3: f=0;
4: t=0;
5: n=0;
6: while(s>0 && f<=4 && t<=280){
7:   t=t+1;
8:   flag=Bernoulli(0.2);
9:   if (flag==1){
10:    f=f+1;
11:    n=n+1;
12:    observe((s<=80));
13:   }else{
14:    f=0;
15:    s=s-1;
16:   }
17: }
18: return s>0;

```

(d) Bounded Retransmission Protocol (BRP).

```

1: a0=1;
2: b0=1;
3: c0=1;
4: d0=1;
5: n=0;
6: while((a0==1||b0==1)){
7:   a1=Bernoulli(0.5);
8:   b1=Bernoulli(0.5);
9:   observe(c0==a1 || d0==b1);
10:  c1=a1;
11:  d1=b1;
12:  n=n+1;
13: }
14: return n

```

(e) Non-i.i.d. loops (NIID).

```

1: r=uniform(0,1);
2: y=0;
3: n=0;
4: while(|y|<1 && (n<=100)){
5:   y=Gaussian(y,2r);
6:   n=n+1;
7: }
8: observe(n>=3);
9: return r;
10:
11:
12:
13:
14:

```

(f) Random Walk (RW1).

```

1: p=uniform(0,1);
2: start=1;
3: curprobe,established=0;
4: while(curprobe < 100 && established <=0 && start <= 1){
5:   if(start == 1){
6:     flag=Bernoulli(p)
7:     if (flag==false){
8:       established=1;
9:     }
10:    start=0;
11:  }else{
12:    flag=Bernoulli(lambda)
13:    if (flag==true){
14:      curprobe := curprobe + 1;
15:    }else{
16:      observe(curprobe>=20);
17:      start=1;
18:      curprobe=0;
19:    }
20:  }
21: }
22: return p;

```

(g) ZeroConf (ZC.1: $\lambda = 0.99$, ZC.2: $\lambda = 0.5$).

```

1: var=uniform(0,7);
2: y=1;
3: prob=0.5;
4: i=0;
5: while(i <= 100){
6:   oldy=y;
7:   y=Gaussian(oldy,2*var);
8:   flag=Bernoulli(lambda);
9:   if(flag){
10:    observe(|y-oldy|<2);
11:    i=i+1;
12:  }
13:  return y;
14:
15:
16:
17:
18:
19:
20:
21:

```

(h) Random Walk (RW2.1: $\lambda = 0.5$, RW2.2: $\lambda = 0.9999$).