

Scalable Verification of Probabilistic Networks*

Steffen Smolka
Cornell University
Ithaca, NY, USA

Praveen Kumar
Cornell University
Ithaca, NY, USA

David M. Kahn[†]
Carnegie Mellon University
Pittsburgh, PA, USA

Nate Foster
Cornell University
Ithaca, NY, USA

Justin Hsu[†]
University of Wisconsin
Madison, WI, USA

Dexter Kozen
Cornell University
Ithaca, NY, USA

Alexandra Silva
University College London
London, UK

Abstract

This paper presents McNetKAT, a scalable tool for verifying probabilistic network programs. McNetKAT is based on a new semantics for the guarded and history-free fragment of Probabilistic NetKAT in terms of finite-state, absorbing Markov chains. This view allows the semantics of all programs to be computed exactly, enabling construction of an automatic verification tool. Domain-specific optimizations and a parallelizing backend enable McNetKAT to analyze networks with thousands of nodes, automatically reasoning about general properties such as probabilistic program equivalence and refinement, as well as networking properties such as resilience to failures. We evaluate McNetKAT’s scalability using real-world topologies, compare its performance against state-of-the-art tools, and develop an extended case study on a recently proposed data center network design.

CCS Concepts • **Theory of computation** → **Automated reasoning; Program semantics**; Random walks and Markov chains; • **Networks** → *Network properties*; • **Software and its engineering** → *Domain specific languages*.

Keywords Network verification, Probabilistic Programming

ACM Reference Format:

Steffen Smolka, Praveen Kumar, David M. Kahn, Nate Foster, Justin Hsu, Dexter Kozen, and Alexandra Silva. 2019. Scalable Verification of Probabilistic Networks. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI ’19)*, June 22–26, 2019, Phoenix, AZ, USA. ACM, New York, NY, USA, 23 pages. <https://doi.org/10.1145/3314221.3314639>

1 Introduction

Networks are among the most complex and critical computing systems used today. Researchers have long sought

*Extended version with appendix.

[†]Work performed at Cornell University.

PLDI ’19, June 22–26, 2019, Phoenix, AZ, USA

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI ’19)*, June 22–26, 2019, Phoenix, AZ, USA, <https://doi.org/10.1145/3314221.3314639>.

to develop automated techniques for modeling and analyzing network behavior [51], but only over the last decade has programming language methodology been brought to bear on the problem [6, 7, 36], opening up new avenues for reasoning about networks in a rigorous and principled way [4, 14, 25, 27, 33]. Building on these initial advances, researchers have begun to target more sophisticated networks that exhibit richer phenomena. In particular, there is renewed interest in *randomization* as a tool for designing protocols and modeling behaviors that arise in large-scale systems—from uncertainty about the inputs, to expected load, to likelihood of device and link failures.

Although programming languages for describing randomized networks exist [13, 17], support for automated reasoning remains limited. Even basic properties require quantitative reasoning in the probabilistic setting, and seemingly simple programs can generate complex distributions. Whereas state-of-the-art tools can easily handle deterministic networks with hundreds of thousands of nodes, probabilistic tools are currently orders of magnitude behind.

This paper presents McNetKAT, a new tool for reasoning about probabilistic network programs written in the guarded and history-free fragment of Probabilistic NetKAT (ProbNetKAT) [4, 13, 14, 46]. ProbNetKAT is an expressive programming language based on Kleene Algebra with Tests, capable of modeling a variety of probabilistic behaviors and properties including randomized routing [30, 48], uncertainty about demands [40], and failures [19]. The history-free fragment restricts the language semantics to input-output behavior rather than tracking paths, and the guarded fragment provides conditionals and while loops rather than union and iteration operators. Although the fragment we consider is a restriction of the full language, it is still expressive enough to encode a wide range of practical networking models. Existing deterministic tools, such as Anteater [35], HSA [25], and Veriflow [27], also use guarded and history-free models.

To enable automated reasoning, we first reformulate the semantics of ProbNetKAT in terms of finite state Markov chains. We introduce a *big-step* semantics that models programs as Markov chains that transition from input to output in a single step, using an auxiliary *small-step* semantics to compute the closed-form solution for the semantics of

the iteration operator. We prove that the Markov chain semantics coincides with the domain-theoretic semantics for ProbNetKAT developed in previous work [13, 46]. Our new semantics also has a key benefit: the limiting distribution of the resulting Markov chains can be computed exactly in closed form, yielding a concise representation that can be used as the basis for building a practical tool.

We have implemented McNetKAT in an OCaml prototype that takes a ProbNetKAT program as input and produces a stochastic matrix that models its semantics in a finite and explicit form. McNetKAT uses the UMFPACK linear algebra library as a back-end solver to efficiently compute limiting distributions [8], and exploits algebraic properties to automatically parallelize the computation across multiple machines. To facilitate comparisons with other tools, we also developed a back-end based on PRISM [31].

To evaluate the scalability of McNetKAT, we conducted experiments on realistic topologies, routing schemes, and properties. Our results show that McNetKAT scales to networks with thousands of switches, and performs orders of magnitude better than a state-of-the-art tool based on general-purpose symbolic inference [17, 18]. We also used McNetKAT to carry out a case study of the resilience of a fault-tolerant data center design proposed by Liu et al. [34].

Contributions and outline. The central contribution of this paper is the development of a *scalable probabilistic network verification tool*. We develop a new, tractable semantics that is sound with respect to ProbNetKAT’s original denotational model. We present a prototype implementation and evaluate it on a variety of scenarios drawn from real-world networks. In §2, we introduce ProbNetKAT using a running example. In §3, we present a semantics based on *finite stochastic matrices* and show that it fully characterizes the behavior of ProbNetKAT programs (Theorem 3.1). In §4, we show how to compute the matrix associated with iteration in closed form. In §5, we discuss our implementation, including symbolic data structures and optimizations that are needed to handle the large state space efficiently. In §6, we evaluate the scalability of McNetKAT on a common data center design and compare its performance against state-of-the-art probabilistic tools. In §7, we present a case study using McNetKAT to analyze resilience in the presence of link failures. We survey related work in §8 and conclude in §9. We defer proofs to the appendix.

2 Overview

This section introduces a running example that illustrates the main features of the ProbNetKAT language as well as some quantitative network properties that arise in practice.

Background on ProbNetKAT. Consider the network in Figure 1, which connects a source to a destination in a topology with three switches. We will first introduce a program that

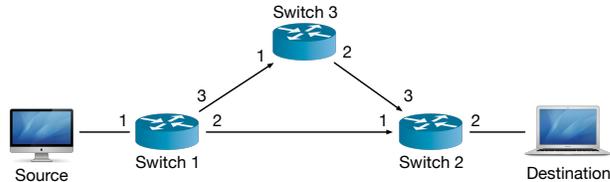


Figure 1. Network topology for running example.

forwards packets from the source to the destination, and then verify that it correctly implements the desired behavior. Next, we will show how to enrich our program to model the possibility of link failures, and develop a fault-tolerant forwarding scheme that automatically routes around failures. Using a quantitative version of program refinement, we will show that the fault-tolerant program is indeed more resilient than the initial program. Finally, we will show how to compute the expected degree of resilience analytically.

To a first approximation, a ProbNetKAT program can be thought of as a randomized function that maps input packets to sets of output packets. Packets are modeled as records, with fields for standard headers—such as the source (*src*) and destination (*dst*) addresses—as well as two fields *switch* (*sw*) and *port* (*pt*) encoding the current location of the packet. ProbNetKAT provides several primitives for manipulating packets: a *modification* $f \leftarrow n$ returns the input packet with field f updated to n , while a *test* $f = n$ returns either the input packet unmodified if the test succeeds, or the empty set if the test fails. The primitives *skip* and *drop* behave like a test that always succeeds and fails, respectively. In the guarded fragment of the language, programs can be composed sequentially ($p ; q$), using conditionals (**if** p **then** q_1 **else** q_2), while loops (**while** p **do** q), or probabilistic choice ($p \oplus q$).

Although ProbNetKAT programs can be freely constructed by composing primitive operations, a typical network model is expressed using two programs: a *forwarding* program (sometimes called a *policy*) and a *link* program (sometimes called a *topology*). The forwarding program describes how packets are transformed locally by the switches at each hop. In our running example, to route packets from the source to the destination, switches 1 and 2 can simply forward all incoming packets out on port 2 by modifying the port field (*pt*). This program can be encoded in ProbNetKAT by performing a case analysis on the location of the input packet, and then setting the port field to 2:

$$p \triangleq \text{if } sw=1 \text{ then } pt \leftarrow 2 \text{ else} \\ \text{if } sw=2 \text{ then } pt \leftarrow 2 \text{ else drop}$$

The final drop at the end of this program encodes the policy for switch 3, which is unreachable.

We can model the topology as a cascade of conditionals that match packets at the end of each link and update their

locations to the link's destination:

$$t \triangleq \text{if } sw=1; pt=2 \text{ then } sw \leftarrow 2; pt \leftarrow 1 \text{ else } \dots$$

To build the overall network model, we first define predicates for the ingress and egress locations,

$$in \triangleq sw=1; pt=1 \quad out \triangleq sw=2; pt=2$$

and then combine the forwarding policy p with the topology t . More specifically, a packet traversing the network starts at an ingress and is repeatedly processed by switches and links until it reaches an egress:

$$M(p, t) \triangleq in; p; \text{while } \neg out \text{ do } (t; p)$$

We can now state and prove properties about the network by reasoning about this model. For instance, the following equivalence states that p forwards all packets to the destination:

$$M(p, t) \equiv in; sw \leftarrow 2; pt \leftarrow 2$$

The program on the right can be regarded as an ideal specification that “teleports” each packet to its destination. Such equations were also used in previous work to reason about properties such as waypointing, reachability, isolation, and loop freedom [4, 14].

Probabilistic reasoning. Real-world networks often exhibit nondeterministic behaviors such as fault tolerant routing schemes to handle unexpected failures [34] and randomized algorithms to balance load across multiple paths [30]. Verifying that networks behave as expected in these more complicated scenarios requires a form of probabilistic reasoning, but most state-of-the-art network verification tools model only deterministic behaviors [14, 25, 27].

To illustrate, suppose we want to extend our example with link failures. Most modern switches execute low-level protocols such as Bidirectional Forwarding Detection (BFD) that compute real-time health information about the link connected to each physical port [5]. We can enrich our model so that each switch has a boolean flag up_i that indicates whether the link connected to the switch at port i is up. Then, we can adjust the forwarding logic to use backup paths when the link is down: for switch 1,

$$\hat{p}_1 \triangleq \text{if } up_2=1 \text{ then } pt \leftarrow 2 \text{ else} \\ \text{if } up_2=0 \text{ then } pt \leftarrow 3 \text{ else drop}$$

and similarly for switches 2 and 3. As before, we can package the forwarding logic for all switches into a single program:

$$\hat{p} \triangleq \text{if } sw=1 \text{ then } \hat{p}_1 \text{ else if } sw=2 \text{ then } \hat{p}_2 \text{ else } \hat{p}_3$$

Next, we update the encoding of our topology to faithfully model link failures. Links can fail for a wide variety of reasons, including human errors, fiber cuts, and hardware faults. A natural way to model such failures is with a *probabilistic model*—i.e., with a distribution that captures how often

certain links fail:

$$f_0 \triangleq up_2 \leftarrow 1; up_3 \leftarrow 1$$

$$f_1 \triangleq \oplus \{ f_0 @ \frac{1}{2}, (up_2 \leftarrow 0; up_3 \leftarrow 1) @ \frac{1}{4}, (up_2 \leftarrow 1; up_3 \leftarrow 0) @ \frac{1}{4} \}$$

$$f_2 \triangleq (up_2 \leftarrow 1 \oplus_{.8} up_2 \leftarrow 0); (up_3 \leftarrow 1 \oplus_{.8} up_3 \leftarrow 0)$$

Intuitively, in f_0 no links fail, in f_1 the links ℓ_{12} and ℓ_{13} fail with probability 25% but at most one link fails, while in f_2 the links fail independently with probability 20%. Using the up flags, we can model a topology with possibly faulty links like so:

$$\hat{t} \triangleq \text{if } sw=1; pt=2; up_2=1 \text{ then } sw \leftarrow 2; pt \leftarrow 1 \text{ else } \dots$$

Combining the policy, topology, and failure model yields a model of the entire network:

$$\hat{M}(p, t, f) \triangleq \text{var } up_2 \leftarrow 1 \text{ in} \\ \text{var } up_3 \leftarrow 1 \text{ in} \\ M((f; p), t)$$

This refined model \hat{M} wraps our previous model M with declarations of the two local fields up_2 and up_3 and executes the failure model (f) at each hop before executing the programs for the switch (p) and topology (t).

Now we can analyze our resilient routing scheme \hat{p} . As a sanity check, we can verify that it delivers packets to their destinations in the absence of failures. Formally, it behaves like the program that teleports packets to their destinations:

$$\hat{M}(\hat{p}, \hat{t}, f_0) \equiv in; sw \leftarrow 2; pt \leftarrow 2$$

More interestingly, \hat{p} is 1-resilient—i.e., it delivers packets provided at most one link fails. Note that this property does *not* hold for the original, naive routing scheme p :

$$\hat{M}(\hat{p}, \hat{t}, f_1) \equiv in; sw \leftarrow 2; pt \leftarrow 2 \not\equiv \hat{M}(p, \hat{t}, f_1)$$

While \hat{p} is not fully resilient under failure model f_2 , which allows two links to fail simultaneously, we can still show that the refined routing scheme \hat{p} performs strictly better than the naive scheme p by checking

$$\hat{M}(p, \hat{t}, f_2) < \hat{M}(\hat{p}, \hat{t}, f_2)$$

where $p < q$ intuitively means that q delivers packets with higher probability than p .

Going a step further, we might want to compute more general quantitative properties of the distributions generated for a given program. For example, we might compute the probability that each routing scheme delivers packets to the destination under f_2 (i.e., 80% for the naive scheme and 96% for the resilient scheme), potentially valuable information to help an Internet Service Provider (ISP) evaluate a network design to check that it meets certain service-level agreements (SLAs). With this motivation in mind, we aim to build a scalable tool that can carry out automated reasoning on probabilistic network programs expressed in ProbNetKAT.

3 ProbNetKAT Syntax and Semantics

This section reviews the syntax of ProbNetKAT and presents a new semantics based on finite state Markov chains.

Preliminaries. A *packet* π is a record mapping a finite set of fields f_1, f_2, \dots, f_k to bounded integers n . As we saw in the previous section, fields can include standard header fields such as source (`src`) and destination (`dst`) addresses, as well as logical fields for modeling the current location of the packet in the network or variables such as up_j . These logical fields are not present in a physical network packet, but they can track auxiliary information for the purposes of verification. We write $\pi.f$ to denote the value of field f of π and $\pi[f:=n]$ for the packet obtained from π by updating field f to hold n . We let Pk denote the (finite) set of all packets.

Syntax. ProbNetKAT terms can be divided into two classes: *predicates* (t, u, \dots) and *programs* (p, q, \dots). Primitive predicates include *tests* ($f=n$) and the Boolean constants *false* (`drop`) and *true* (`skip`). Compound predicates are formed using the usual Boolean connectives: disjunction ($t \& u$), conjunction ($t ; u$), and negation ($\neg t$). Primitive programs include *predicates* (t) and *assignments* ($f \leftarrow n$). The original version of the language also provides a *dup* primitive, which logs the current state of the packet, but the history-free fragment omits this operation. Compound programs can be formed using *parallel composition* ($p \& q$), *sequential composition* ($p ; q$), and *iteration* (p^*). In addition, the *probabilistic choice* operator $p \oplus_r q$ executes p with probability r and q with probability $1 - r$, where r is rational, $0 \leq r \leq 1$. We sometimes use an n -ary version and omit the r 's: $p_1 \oplus \dots \oplus p_n$ executes a p_i chosen uniformly at random. In addition to these core constructs (summarized in Figure 2), many other useful constructs can be derived. For example, mutable local variables (e.g., up_j , used to track link health in §2), can be desugared into the language:

$$\mathbf{var} \ f \leftarrow n \ \mathbf{in} \ p \triangleq f \leftarrow n ; p ; f \leftarrow 0$$

Here f is a field that is local to p . The final assignment $f \leftarrow 0$ sets the value of f to a canonical value, “erasing” it after the field goes out of scope. We often use local variables to record extra information for verification—e.g., recording whether a packet traversed a given switch allows reasoning about simple waypointing and isolation properties, even though the history-free fragment of ProbNetKAT does not model paths directly.

Guarded fragment. Conditionals and while loops can be encoded using union and iteration:

$$\begin{aligned} \mathbf{if} \ t \ \mathbf{then} \ p \ \mathbf{else} \ q &\triangleq t ; p \& \neg t ; q \\ \mathbf{while} \ t \ \mathbf{do} \ p &\triangleq (t ; p)^* ; \neg t \end{aligned}$$

Note that these constructs use the predicate t as a *guard*, resolving the inherent nondeterminism in the union and iteration operators. Our implementation handles programs

Naturals	$n ::= 0 \mid 1 \mid 2 \mid \dots$	
Fields	$f ::= f_1 \mid \dots \mid f_k$	
Packets	$\text{Pk} \ni \pi ::= \{f_1 = n_1, \dots, f_k = n_k\}$	
Probabilities	$r \in [0, 1] \cap \mathbb{Q}$	
Predicates	$t, u ::= \text{drop}$	<i>False</i>
	$\mid \text{skip}$	<i>True</i>
	$\mid f = n$	<i>Test</i>
	$\mid t \& u$	<i>Disjunction</i>
	$\mid t ; u$	<i>Conjunction</i>
	$\mid \neg t$	<i>Negation</i>
Programs	$p, q ::= t$	<i>Filter</i>
	$\mid f \leftarrow n$	<i>Assignment</i>
	$\mid p \& q$	<i>Union</i>
	$\mid p ; q$	<i>Sequence</i>
	$\mid p \oplus_r q$	<i>Choice</i>
	$\mid p^*$	<i>Iteration</i>

Figure 2. ProbNetKAT Syntax.

in the guarded fragment of the language—i.e., with loops and conditionals but without union and iteration—though we will develop the theory in full generality here, to make connections to previous work on ProbNetKAT clearer. We believe this restriction is acceptable from a practical perspective, as the main purpose of union and iteration is to encode forwarding tables and network-wide processing, and the guarded variants can often perform the same task. A notable exception is multicast, which cannot be expressed in the guarded fragment.

Semantics. Previous work on ProbNetKAT [13] modeled history-free programs as maps $2^{\text{Pk}} \rightarrow \mathcal{D}(2^{\text{Pk}})$, where $\mathcal{D}(2^{\text{Pk}})$ denotes the set of probability distributions on 2^{Pk} . This semantics is useful for establishing fundamental properties of the language, but we will need a more explicit representation to build a practical verification tool. Since the set of packets is finite, probability distributions over sets of packets are discrete and can be characterized by a *probability mass function*, $f : 2^{\text{Pk}} \rightarrow [0, 1]$ such that $\sum_{b \subseteq \text{Pk}} f(b) = 1$. It will be convenient to view f as a *stochastic vector* of non-negative entries that sum to 1.

A program, which maps inputs a to distributions over outputs, can then be represented by a square matrix indexed by Pk in which the stochastic vector corresponding to input a appears as the a -th row. Thus, we can interpret a program p as a matrix $\mathcal{B}[[p]] \in [0, 1]^{2^{\text{Pk}} \times 2^{\text{Pk}}}$ indexed by packet sets, where the matrix entry $\mathcal{B}[[p]]_{ab}$ gives the probability that p produces output $b \in 2^{\text{Pk}}$ on input $a \in 2^{\text{Pk}}$. The rows of the matrix $\mathcal{B}[[p]]$ are stochastic vectors, each encoding the distribution produced for an input set a ; such a matrix is called *right-stochastic*, or simply *stochastic*. We write $\mathbb{S}(2^{\text{Pk}})$ for the set of right-stochastic matrices indexed by 2^{Pk} .

$$\boxed{\mathcal{B}[[p]] \in \mathbb{S}(2^{\text{Pk}})}$$

$$\begin{aligned} \mathcal{B}[[\text{drop}]]_{ab} &\triangleq [b = \emptyset] \\ \mathcal{B}[[\text{skip}]]_{ab} &\triangleq [a = b] \\ \mathcal{B}[[f=n]]_{ab} &\triangleq [b = \{\pi \in a \mid \pi.f = n\}] \\ \mathcal{B}[[\neg t]]_{ab} &\triangleq [b \subseteq a] \cdot \mathcal{B}[[t]]_{a,a-b} \\ \mathcal{B}[[f \leftarrow n]]_{ab} &\triangleq [b = \{\pi[f := n] \mid \pi \in a\}] \\ \mathcal{B}[[p \& q]]_{ab} &\triangleq \sum_{c,d} [c \cup d = b] \cdot \mathcal{B}[[p]]_{a,c} \cdot \mathcal{B}[[q]]_{c,d} \\ \mathcal{B}[[p; q]] &\triangleq \mathcal{B}[[p]] \cdot \mathcal{B}[[q]] \\ \mathcal{B}[[p \oplus_r q]] &\triangleq r \cdot \mathcal{B}[[p]] + (1-r) \cdot \mathcal{B}[[q]] \\ \mathcal{B}[[p^*]]_{ab} &\triangleq \lim_{n \rightarrow \infty} \mathcal{B}[[p^{(n)}]]_{ab} \end{aligned}$$

Figure 3. ProbNetKAT Semantics. The notation $\mathcal{B}[[p]]_{ab}$ denotes the probability that p produces b on input a .

Figure 3 defines an interpretation of ProbNetKAT programs as stochastic matrices; the Iverson bracket $[\varphi]$ is 1 if φ is true, and 0 otherwise. Deterministic program primitives are interpreted as $\{0, 1\}$ -matrices—e.g., the program primitive drop is interpreted as the following stochastic matrix:

$$\mathcal{B}[[\text{drop}]] = \begin{array}{c} \emptyset \quad b_2 \quad \dots \quad b_n \\ \emptyset \quad \begin{bmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_n & 1 & 0 & \dots & 0 \end{bmatrix} \end{array} \quad \begin{array}{c} \boxed{a_2} \xrightarrow{1} \boxed{a_1 = \emptyset} \\ \vdots \\ \boxed{a_n} \xrightarrow{1} \boxed{a_1 = \emptyset} \end{array} \quad (1)$$

which assigns all probability mass to the \emptyset -column. Similarly, skip is interpreted as the identity matrix. Sequential composition can be interpreted as matrix product,

$$\mathcal{B}[[p; q]]_{ab} = \sum_c \mathcal{B}[[p]]_{ac} \cdot \mathcal{B}[[q]]_{cb} = (\mathcal{B}[[p]] \cdot \mathcal{B}[[q]])_{ab}$$

which reflects the intuitive semantics of composition: to step from a to b in $\mathcal{B}[[p; q]]$, one must step from a to an intermediate state c in $\mathcal{B}[[p]]$, and then from c to b in $\mathcal{B}[[q]]$.

As the picture in (1) suggests, a stochastic matrix $B \in \mathbb{S}(2^{\text{Pk}})$ can be viewed as a *Markov chain* (MC)—i.e., a probabilistic transition system with state space 2^{Pk} . The B_{ab} entry gives the probability that the system transitions from a to b .

Soundness. The matrix $\mathcal{B}[[p]]$ is equivalent to the denotational semantics $[[p]]$ defined in previous work [13].

Theorem 3.1 (Soundness). *Let $a, b \in 2^{\text{Pk}}$. The matrix $\mathcal{B}[[p]]$ satisfies $\mathcal{B}[[p]]_{ab} = [[p]](a)(\{b\})$.*

Hence, checking program equivalence for p and q reduces to checking equality of the matrices $\mathcal{B}[[p]]$ and $\mathcal{B}[[q]]$.

Corollary 3.2. $[[p]] = [[q]]$ if and only if $\mathcal{B}[[p]] = \mathcal{B}[[q]]$.

In particular, because the Markov chains are all finite state, the transition matrices are finite dimensional with rational

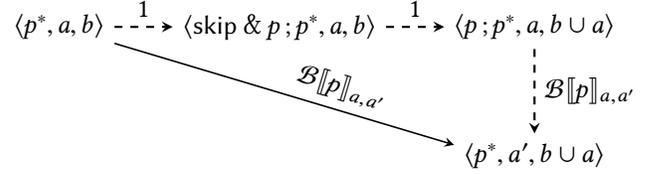


Figure 4. The small-step semantics is given by a Markov chain with states $\langle \text{program}, \text{input set}, \text{output accumulator} \rangle$. The three dashed arrows can be collapsed into the single solid arrow, rendering the program component superfluous.

entries. Accordingly, program equivalence and other quantitative properties can be automatically verified provided we can compute the matrices for given programs. This is relatively straightforward for program constructs besides $\mathcal{B}[[p^*]]$, whose matrix is defined in terms of a limit. The next section presents a closed-form definition of the stochastic matrix for this operator.

4 Computing Stochastic Matrices

The semantics developed in the previous section can be viewed as a “big-step” semantics in which a single step models the execution of a program from input to output. To compute the semantics of p^* , we will introduce a finer, “small-step” chain in which a transition models one iteration of the loop.

To build intuition, consider simulating p^* using a transition system with states given by triples $\langle p, a, b \rangle$ in which p is the program being executed, a is the set of (input) packets, and b is an accumulator that collects the output packets generated so far. To model the execution of p^* on input a , we start from the initial state $\langle p^*, a, \emptyset \rangle$ and unroll p^* one iteration according to the characteristic equation $p^* \equiv \text{skip} \& p; p^*$, yielding the following transition:

$$\langle p^*, a, \emptyset \rangle \xrightarrow{1} \langle \text{skip} \& p; p^*, a, \emptyset \rangle$$

Next, we execute both skip and $p; p^*$ on the input set and take the union of their results. Executing skip yields the input set as output, with probability 1:

$$\langle \text{skip} \& p; p^*, a, \emptyset \rangle \xrightarrow{1} \langle p; p^*, a, a \rangle$$

Executing $p; p^*$, executes p and feeds its output into p^* :

$$\forall a' : \langle p; p^*, a, a \rangle \xrightarrow{\mathcal{B}[[p]]_{a,a'}} \langle p^*, a', a \rangle$$

At this point we are back to executing p^* , albeit with a different input set a' and some accumulated output packets. The resulting Markov chain is shown in Figure 4.

Note that as the first two steps of the chain are deterministic, we can simplify the transition system by collapsing all three steps into one, as illustrated in Figure 4. The program component can then be dropped, as it now remains constant

across transitions. Hence, we work with a Markov chain over the state space $2^{\text{Pk}} \times 2^{\text{Pk}}$, defined formally as follows:

$$\mathcal{S}[[p]] \in \mathbb{S}(2^{\text{Pk}} \times 2^{\text{Pk}})$$

$$\mathcal{S}[[p]]_{(a,b),(a',b')} \triangleq [b' = b \cup a] \cdot \mathcal{B}[[p]]_{a,a'}.$$

We can verify that the matrix $\mathcal{S}[[p]]$ defines a Markov chain.

Lemma 4.1. $\mathcal{S}[[p]]$ is stochastic.

Next, we show that each step in $\mathcal{S}[[p]]$ models an iteration of p^* . Formally, the $(n+1)$ -step of $\mathcal{S}[[p]]$ is equivalent to the big-step behavior of the n -th unrolling of p^* .

Proposition 4.2. $\mathcal{B}[[p^{(n)}]]_{a,b} = \sum_{a'} \mathcal{S}[[p]]_{(a,\emptyset),(a',b)}^{n+1}$

Direct induction on the number of steps $n \geq 0$ fails because the hypothesis is too weak. We generalize from start states with empty accumulator to arbitrary start states.

Lemma 4.3. Let p be program. Then for all $n \in \mathbb{N}$ and $a, b, b' \subseteq \text{Pk}$, we have

$$\sum_{a'} [b' = a' \cup b] \cdot \mathcal{B}[[p^{(n)}]]_{a,a'} = \sum_{a'} \mathcal{S}[[p]]_{(a,b),(a',b')}^{n+1}.$$

Proposition 4.2 then follows from Lemma 4.3 with $b = \emptyset$.

Intuitively, the long-run behavior of $\mathcal{S}[[p]]$ approaches the big-step behavior of p^* : letting (a_n, b_n) denote the random state of the Markov chain $\mathcal{S}[[p]]$ after taking n steps starting from (a, \emptyset) , the distribution of b_n for $n \rightarrow \infty$ is precisely the distribution of outputs generated by p^* on input a (by Proposition 4.2 and the definition of $\mathcal{B}[[p^*]]$).

Closed form. The limiting behavior of finite state Markov chains has been well studied in the literature (e.g., see Kemeny and Snell [26]). For so-called *absorbing* Markov chains, the limit distribution can be computed exactly. A state s of a Markov chain T is *absorbing* if it transitions to itself with probability 1,

$$\begin{array}{c} \textcircled{s} \\ \curvearrowright \end{array} 1 \quad (\text{formally: } T_{s,s'} = [s = s'])$$

and a Markov chain $T \in \mathbb{S}(S)$ is *absorbing* if each state can reach an absorbing state:

$$\forall s \in S. \exists s' \in S, n \geq 0. T_{s,s'}^n > 0 \text{ and } T_{s',s'} = 1$$

The non-absorbing states of an absorbing MC are called *transient*. Assume T is absorbing with n_t transient states and n_a absorbing states. After reordering the states so that absorbing states appear first, T has the form

$$T = \begin{bmatrix} I & 0 \\ R & Q \end{bmatrix}$$

where I is the $n_a \times n_a$ identity matrix, R is an $n_t \times n_a$ matrix giving the probabilities of transient states transitioning to absorbing states, and Q is an $n_t \times n_t$ matrix specifying the probabilities of transitions between transient states. Since absorbing states never transition to transient states by definition, the upper right corner contains a $n_a \times n_t$ zero matrix.

From any start state, a finite state absorbing MC always ends up in an absorbing state eventually, i.e. the limit $T^\infty \triangleq \lim_{n \rightarrow \infty} T^n$ exists and has the form

$$T^\infty = \begin{bmatrix} I & 0 \\ A & 0 \end{bmatrix}$$

where the $n_t \times n_a$ matrix A contains the so-called *absorption probabilities*. This matrix satisfies the following equation:

$$A = (I + Q + Q^2 + \dots)R$$

Intuitively, to transition from a transient state to an absorbing state, the MC can take an arbitrary number of steps between transient states before taking a single—and final—step into an absorbing state. The infinite sum $X \triangleq \sum_{n \geq 0} Q^n$ satisfies $X = I + QX$, and solving for X yields

$$X = (I - Q)^{-1} \quad \text{and} \quad A = (I - Q)^{-1}R. \quad (2)$$

(We refer the reader to Kemeny and Snell [26] for the proof that the inverse exists.)

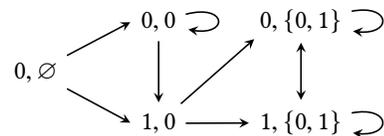
Before we apply this theory to the small-step semantics $\mathcal{S}[[\cdot]]$, it will be useful to introduce some MC-specific notation. Let T be an MC. We write $s \xrightarrow{T} s'$ if s can reach s' in precisely n steps, i.e. if $T_{s,s'}^n > 0$; and we write $s \xrightarrow{T} s'$ if s can reach s' in some number of steps, i.e. if $T_{s,s'}^n > 0$ for some $n \geq 0$. Two states are said to *communicate*, denoted $s \leftrightarrow s'$, if $s \xrightarrow{T} s'$ and $s' \xrightarrow{T} s$. The relation \leftrightarrow is an equivalence relation, and its equivalence classes are called *communication classes*. A communication class is *absorbing* if it cannot reach any states outside the class. Let $\text{Pr}[s \xrightarrow{T} s']$ denote the probability $T_{s,s'}^n$. For the rest of the section, we fix a program p and abbreviate $\mathcal{B}[[p]]$ as B and $\mathcal{S}[[p]]$ as S . We also define *saturated states*, those where the accumulator has stabilized.

Definition 4.4. A state (a, b) of S is called *saturated* if b has reached its final value, i.e. if $(a, b) \xrightarrow{S} (a', b')$ implies $b' = b$.

After reaching a saturated state, the output of p^* is fully determined. The probability of ending up in a saturated state with accumulator b , starting from an initial state (a, \emptyset) , is

$$\lim_{n \rightarrow \infty} \sum_{a'} S_{(a,\emptyset),(a',b)}^n$$

and, indeed, this is the probability that p^* outputs b on input a by Proposition 4.2. Unfortunately, we cannot directly compute this limit since saturated states are not necessarily absorbing. To see this, consider $p^* = (f \leftarrow 0 \oplus_{1/2} f \leftarrow 1)^*$ over a single $\{0, 1\}$ -valued field f . Then S has the form

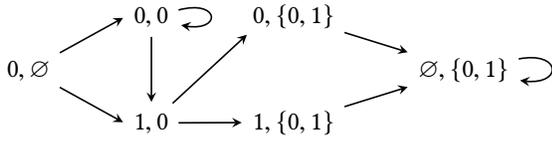


where all edges are implicitly labeled with $\frac{1}{2}$, and 0 and 1 denote the packets with f set to 0 and 1 respectively. We omit states not reachable from $(0, \emptyset)$. The right-most states are saturated, but they communicate and are thus not absorbing.

To align saturated and absorbing states, we can perform a quotient of this Markov chain by collapsing the communicating states. We define an auxiliary matrix,

$$U_{(a,b),(a',b')} \triangleq [b' = b] \cdot \begin{cases} [a' = \emptyset] & \text{if } (a, b) \text{ is saturated} \\ [a' = a] & \text{else} \end{cases}$$

which sends a saturated state (a, b) to a canonical saturated state (\emptyset, b) and acts as the identity on all other states. In our example, the modified chain SU is as follows:



and indeed is absorbing, as desired.

Lemma 4.5. S, U , and SU are monotone in the sense that: $(a, b) \xrightarrow{S} (a', b')$ implies $b \subseteq b'$ (and similarly for U and SU).

Proof. By definition (S and U) and by composition (SU). \square

Next, we show that SU is an absorbing MC:

Proposition 4.6. Let $n \geq 1$.

1. $(SU)^n = S^n U$
2. SU is an absorbing MC with absorbing states $\{(\emptyset, b)\}$.

Arranging the states (a, b) in lexicographically ascending order according to \subseteq and letting $n = |2^{\text{Pk}}|$, it then follows from Proposition 4.6.2 that SU has the form

$$SU = \begin{bmatrix} I_n & 0 \\ R & Q \end{bmatrix}$$

where, for $a \neq \emptyset$, we have

$$(SU)_{(a,b),(a',b')} = [R \ Q]_{(a,b),(a',b')}.$$

Moreover, SU converges and its limit is given by

$$(SU)^\infty \triangleq \begin{bmatrix} I_n & 0 \\ (I - Q)^{-1}R & 0 \end{bmatrix} = \lim_{n \rightarrow \infty} (SU)^n. \quad (3)$$

Putting together the pieces, we can use the modified Markov chain SU to compute the limit of S .

Theorem 4.7 (Closed Form). Let $a, b, b' \subseteq \text{Pk}$. Then

$$\lim_{n \rightarrow \infty} \sum_{a'} S_{(a,b),(a',b')}^n = (SU)_{(a,b),(\emptyset,b')}^\infty.$$

The limit exists and can be computed exactly, in closed-form.

5 Implementation

We have implemented McNetKAT as an embedded DSL in OCaml in roughly 10KLoC. The frontend provides functions for defining and manipulating ProbNetKAT programs and for generating such programs automatically from network topologies encoded using Graphviz. These programs can then be analyzed by one of two backends: the *native backend* (PNK), which compiles programs to (symbolically represented) stochastic matrices; or the *PRISM-based backend* (PPNK), which emits inputs for the state-of-the-art probabilistic model checker PRISM [32].

Pragmatic restrictions. Although our semantics developed in §3 and §4 theoretically supports computations on sets of packets, a direct implementation would be prohibitively expensive—the matrices are indexed by the powerset 2^{Pk} of the universe of all possible packets! To obtain a practical analysis tool, we restrict the state space to single packets. At the level of syntax, we restrict to the guarded fragment of ProbNetKAT, *i.e.* to programs with conditionals and while loops, but without union and iteration. This ensures that no proper packet sets are ever generated, thus allowing us to work over an exponentially smaller state space. While this restriction does rule out some uses of ProbNetKAT—most notably, modeling multicast—we did not find this to be a serious limitation because multicast is relatively uncommon in probabilistic networking. If needed, multicast can often be modeled using multiple unicast programs.

5.1 Native Backend

The native backend compiles a program to a symbolic representation of its big step matrix. The translation, illustrated in Figure 5, proceeds as follows. First, we translate atomic programs to Forwarding Decision Diagrams (FDDs), a symbolic data structure based on Binary Decision Diagrams (BDDs) that encodes sparse matrices compactly [45]. Second, we translate composite programs by first translating each sub-program to an FDD and then merging the results using standard BDD algorithms. Loops require special treatment: we (i) convert the FDD for the body of the loop to a sparse stochastic matrix, (ii) compute the semantics of the loop by using an optimized sparse linear solver [8] to solve the system from §4, and finally (iii) convert the resulting matrix back to an FDD. We use exact rational arithmetic in the frontend and FDD-backend to preempt concerns about numerical precision, but trust the linear algebra solver UMFPACK (based on 64 bit floats) to provide accurate solutions.¹ Our implementation relies on several optimizations; we detail two of the more interesting ones below.

Probabilistic FDDs. Binary Decision Diagrams [1] and variants thereof [15] have long been used in verification and

¹UMFPACK is a mature library powering widely-used scientific computing packages such as MATLAB and SciPy.

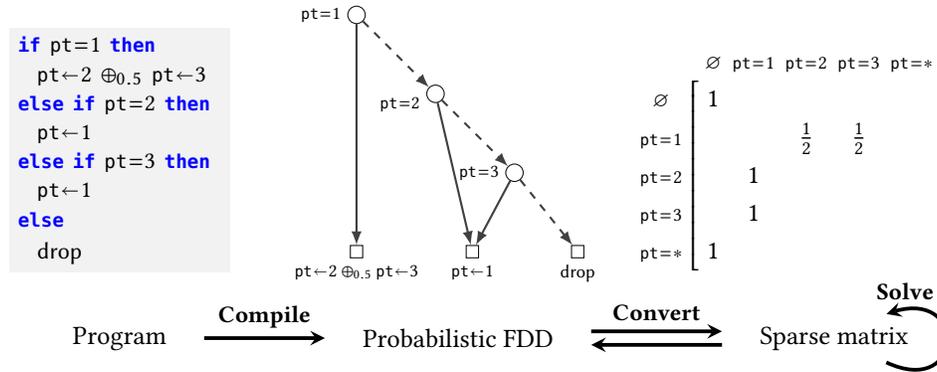


Figure 5. Implementation using FDDs and a sparse linear algebra solver.

model checking to represent large state spaces compactly. A variant called Forwarding Decision Diagrams (FDDs) [45] was previously developed specifically for the networking domain, but only supported deterministic behavior. In this work, we extended FDDs to probabilistic FDDs. A probabilistic FDD is a rooted directed acyclic graph that can be understood as a control-flow graph. Interior nodes test packet fields and have outgoing true- and false- branches, which we visualize by solid lines and dashed lines in Figure 5. Leaf nodes contain distributions over actions, where an action is either a set of modifications or a special action drop. To interpret an FDD, we start at the root node with an initial packet and traverse the graph as dictated by the tests until a leaf node is reached. Then, we apply each action in the leaf node to the packet. Thus, an FDD represents a function of type $P_k \rightarrow \mathcal{D}(P_k + \emptyset)$, or equivalently, a stochastic matrix over the state space $P_k + \emptyset$ where the \emptyset -row puts all mass on \emptyset by convention. Like BDDs, FDDs respect a total order on tests and contain no isomorphic subgraphs or redundant tests, which enables representing sparse matrices compactly.

Dynamic domain reduction. As Figure 5 shows, we do not have to represent the state space $P_k + \emptyset$ explicitly even when converting into sparse matrix form. In the example, the state space is represented by symbolic packets $pt = 1$, $pt = 2$, $pt = 3$, and $pt = *$, each representing an equivalence class of packets. For example, $pt = 1$ can represent all packets π satisfying $\pi.pt = 1$, because the program treats all such packets in the same way. The packet $pt = *$ represents the set $\{\pi \mid \pi.pt \notin \{1, 2, 3\}\}$. The symbol $*$ can be thought of as a wildcard that ranges over all values not explicitly represented by other symbolic packets. The symbolic packets are chosen dynamically when converting an FDD to a matrix by traversing the FDD and determining the set of values appearing in each field, either in a test or a modification. Since FDDs never contain redundant tests or modifications, these sets are typically of manageable size.

5.2 PRISM backend

PRISM is a mature probabilistic model checker that has been actively developed and improved for the last two decades. The tool takes as input a Markov chain model specified symbolically in PRISM’s input language and a property specified using a logic such as Probabilistic CTL, and outputs the probability that the model satisfies the property. PRISM supports various types of models including finite state Markov chains, and can thus be used as a backend for reasoning about ProbNetKAT programs using our results from §3 and §4. Accordingly, we implemented a second backend that translates ProbNetKAT to PRISM programs. While the native backend computes the big step semantics of a program—a costly operation that may involve solving linear systems to compute fixed points—the PRISM backend is a purely syntactic transformation; the heavy lifting is done by PRISM itself.

A PRISM program consists of a set of bounded variables together with a set of transition rules of the form

$$\phi \rightarrow p_1 \cdot u_1 + \dots + p_k \cdot u_k$$

where ϕ is a Boolean predicate over the variables, the p_i are probabilities that must sum up to one, and the u_i are sequences of variable updates. The predicates are required to be mutually exclusive and exhaustive. Such a program encodes a Markov chain whose state space is given by the finite set of variable assignments and whose transitions are dictated by the rules: if ϕ is satisfied under the current assignment σ and σ_i is obtained from σ by performing update u_i , then the probability of a transition from σ to σ_i is p_i .

It is easy to see that any PRISM program can be expressed in ProbNetKAT, but the reverse direction is slightly tricky: it requires the introduction of an additional variable akin to a program counter to emulate ProbNetKAT’s control flow primitives such as loops and sequences. As an additional challenge, we must be economical in our allocation of the program counter, since the performance of model checking is very sensitive to the size of the state space.

We address this challenge in three steps. First, we translate the ProbNetKAT program to a finite state machine using a

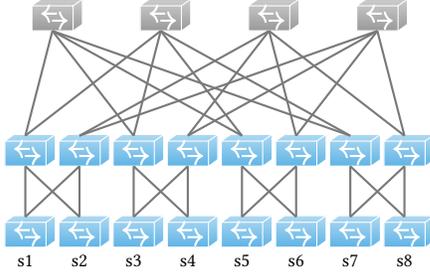


Figure 6. A FatTree topology with $p = 4$.

Thompson-style construction [47]. Each edge is labeled with a predicate ϕ , a probability p_i , and an update u_i , subject to the following well-formedness conditions:

1. For each state, the predicates on its outgoing edges form a partition.
2. For each state and predicate, the probabilities of all outgoing edges guarded by that predicate sum to one.

Intuitively, the state machine encodes the control-flow graph.

This intuition serves as the inspiration for the next translation step, which collapses each basic block of the graph into a single state. This step is crucial for reducing the state space, since the state space of the initial automaton is linear in the size of the program. Finally, we obtain a PRISM program from the automaton as follows: for each state s with adjacent predicate ϕ and ϕ -guarded outgoing edges $s \xrightarrow{\phi/p_i/u_i} t_i$ for $1 \leq i \leq k$, produce a PRISM rule

$$(pc=s \wedge \phi) \rightarrow p_1 \cdot (u_1 ; pc \leftarrow t_1) + \dots + p_k \cdot (u_k ; pc \leftarrow t_k).$$

The well-formedness conditions of the state machine guarantee that the resulting program is a valid PRISM program. With some care, the entire translation can be implemented in linear time. Indeed, McNetKAT translates all programs in our evaluation to PRISM in under a second.

6 Evaluation

To evaluate McNetKAT we conducted experiments on several benchmarks including a family of real-world data center topologies and a synthetic benchmark drawn from the literature [17]. We evaluated McNetKAT’s scalability, characterized the effect of optimizations, and compared performance against other state-of-the-art tools. All McNetKAT running times we report refer to the time needed to compile programs to FDDs; the cost of comparing FDDs for equivalence and ordering, or of computing statistics of the encoded distributions, is negligible. All experiments were performed on machines with 16-core, 2.6 GHz Intel Xeon E5-2650 processors with 64 GB of memory.

Scalability on FatTree topologies. We first measured the scalability of McNetKAT by using it to compute network models for a series of FatTree topologies of increasing size. FatTrees [2] (see also Figure 6) are multi-level, multi-rooted

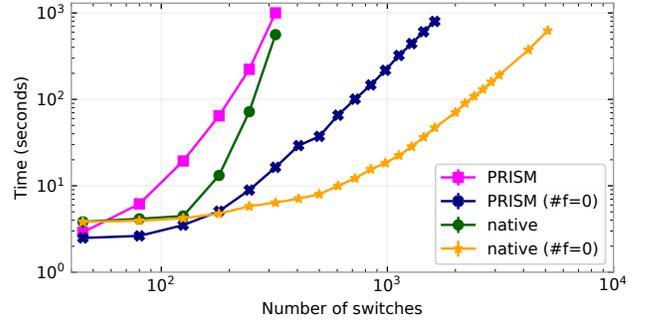


Figure 7. Scalability on a family of data center topologies.

trees that are widely used as topologies in modern data centers. FatTrees can be specified in terms of a parameter p corresponding to the number of ports on each switch. A p -ary FatTree connects $\frac{1}{4}p^3$ servers using $\frac{5}{4}p^2$ switches. To route packets, we used a form of Equal-Cost Multipath Routing (ECMP) that randomly maps traffic flows onto shortest paths. We measured the time needed to construct the stochastic matrix representation of the program on a single machine using two backends (native and PRISM) and under two failure models (no failures and independent failures with probability $1/1000$).

Figure 7 depicts the results, several of which are worth discussing. First, the native backend scales quite well: in the absence of failures ($f = 0$), it scales to a network with 5000 switches in approximately 10 minutes. This result shows that McNetKAT is able to handle networks of realistic size. Second, the native backend consistently outperforms the PRISM backend. We conjecture that the native backend is able to exploit algebraic properties of the ProbNetKAT program to better parallelize the job. Third, performance degrades in the presence of failures. This is to be expected—failures lead to more complex probability distributions which are nontrivial to represent and manipulate.

Parallel speedup. One of the contributors to McNetKAT’s good performance is its ability to parallelize the computation of stochastic matrices across multiple cores in a machine, or even across machines in a cluster. Intuitively, because a network is a large collection of mostly independent devices, it is possible to model its global behavior by first modeling the behavior of each device in isolation, and then combining the results to obtain a network-wide model. In addition to speeding up the computation, this approach can also reduce memory usage, often a bottleneck on large inputs.

To facilitate parallelization, we added an n -ary disjoint branching construct to ProbNetKAT:

```

case  $sw=1$  then  $p_1$  else
case  $sw=2$  then  $p_2$  else
    ...
case  $sw=n$  then  $p_n$ 
    
```

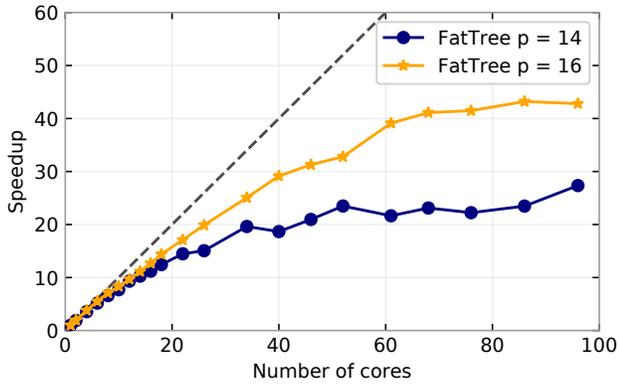


Figure 8. Speedup due to parallelization.

Semantically, this construct is equivalent to a cascade of conditionals; but the native backend compiles it in parallel using a map-reduce-style strategy, using one process per core by default.

To evaluate the impact of parallelization, we compiled two representative FatTree models ($p = 14$ and $p = 16$) using ECMP routing on an increasing number of cores. With m cores, we used one master machine together with $r = \lceil m/16 - 1 \rceil$ remote machines, adding machines one by one as needed to obtain more physical cores. The results are shown in Figure 8. We see near linear speedup on a single machine, cutting execution time by more than an order of magnitude on our 16-core test machine. Beyond a single machine, the speedup depends on the complexity of the submodels for each switch—the longer it takes to generate the matrix for each switch, the higher the speedup. For example, with a $p = 16$ FatTree, we obtained a 30x speedup using 40 cores across 3 machines.

Comparison with other tools. Bayonet [17] is a state-of-the-art tool for analyzing probabilistic networks. Whereas McNetKAT has a native backend tailored to the networking domain and a backend based on a probabilistic model checker, Bayonet programs are translated to a general-purpose probabilistic language which is then analyzed by the symbolic inference engine PSI [18]. Bayonet’s approach is more general, as it can model queues, state, and multi-packet interactions under an asynchronous scheduling model. It also supports Bayesian inference and parameter synthesis. Moreover, Bayonet is fully symbolic whereas McNetKAT uses a numerical linear algebra solver [8] (based on floating point arithmetic) to compute limits.

To evaluate how the performance of these approaches compares, we reproduced an experiment from the Bayonet paper that analyzes the reliability of a simple routing scheme in a family of “chain” topologies indexed by k , as shown in Figure 9.

For $k = 1$, the network consists of four switches organized into a diamond, with a single link that fails with probability

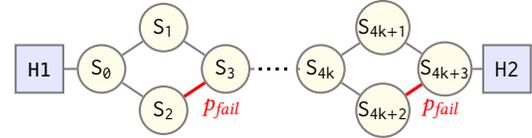


Figure 9. Chain topology

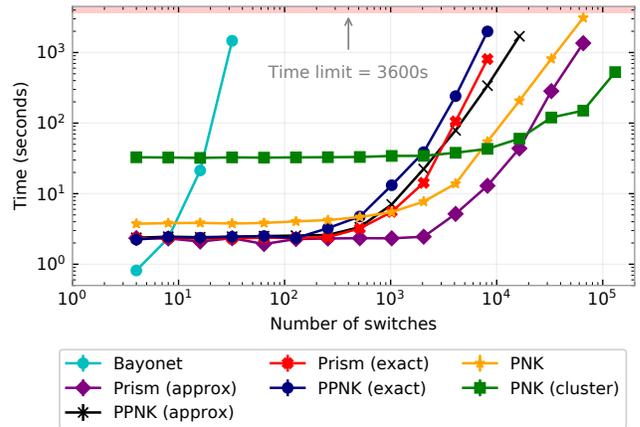


Figure 10. Scalability on chain topology.

$p_{fail} = 1/1000$. For $k > 1$, the network consists of k diamonds linked together into a chain as shown in Figure 9. Within each diamond, switch S_0 forwards packets with equal probability to switches S_1 and S_2 , which in turn forward to switch S_3 . However, S_2 drops the packet if the link to S_3 fails. We analyze the probability that a packet originating at H1 is successfully delivered to H2. Our implementation does not exploit the regularity of these topologies.

Figure 10 gives the running time for several tools on this benchmark: Bayonet, hand-written PRISM, ProbNetKAT with the PRISM backend (PPNK), and ProbNetKAT with the native backend (PNK). Further, we ran the PRISM tools in exact and approximate mode, and we ran the ProbNetKAT backend on a single machine and on the cluster. Note that both axes in the plot are log-scaled.

We see that Bayonet scales to 32 switches in about 25 minutes, before hitting the one hour time limit and 64 GB memory limit at 48 switches. ProbNetKAT answers the same query for 2048 switches in under 10 seconds and scales to over 65000 switches in about 50 minutes on a single core, or just 2.5 minutes using a cluster of 24 machines. PRISM scales similarly to ProbNetKAT, and performs best using the hand-written model in approximate mode.

Overall, this experiment shows that for basic network verification tasks, ProbNetKAT’s domain-specific backend based on specialized data structures and an optimized linear-algebra library [8] can outperform an approach based on a general-purpose solver.

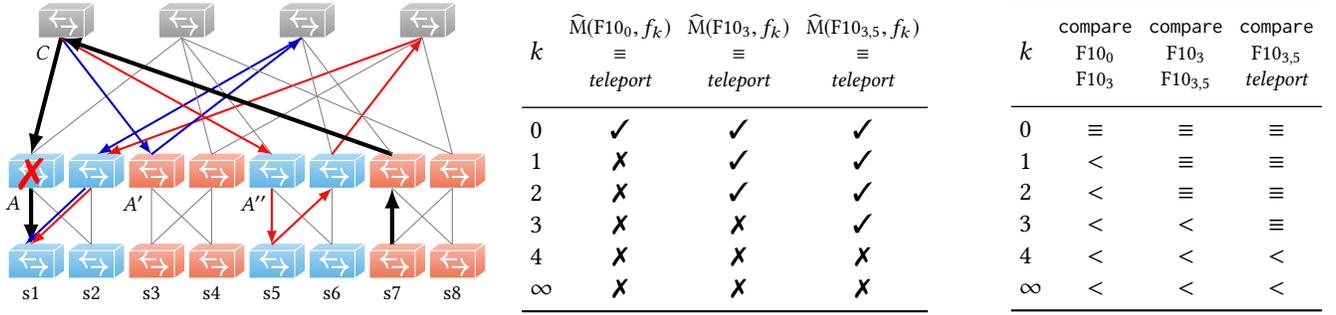


Figure 11. (a) AB FatTree topology with $p = 4$. (b) Evaluating k -resilience. (c) Comparing schemes under k failures.

7 Case Study: Data Center Fault-Tolerance

In this section, we go beyond benchmarks and present a case study that illustrates the utility of McNetKAT for probabilistic reasoning. Specifically, we model the F10 [34] data center design in ProbNetKAT and verify its key properties.

Data center resilience. An influential measurement study by Gill et al. [19] showed that data centers experience frequent failures, which have a major impact on application performance. To address this challenge, a number of data center designs have been proposed that aim to simultaneously achieve high throughput, low latency, and fault tolerance.

F10 topology. F10 uses a novel topology called an *AB FatTree*, see Figure 11(a), that enhances a traditional FatTree [2] with additional backup paths that can be used when failures occur. To illustrate, consider routing from s_7 to s_1 in Figure 11(a) along one of the shortest paths (in thick black). After reaching the core switch C in a standard FatTree (recall Figure 6), if the aggregation switch on the downward path failed, we would need to take a 5-hop detour (shown in red) that goes down to a different edge switch, up to a different core switch, and finally down to s_1 . In contrast, an AB FatTree [34] modifies the wiring of the aggregation later to provide shorter detours—e.g., a 3-hop detour (shown in blue) for the previous scenario.

F10 routing. F10’s routing scheme uses three strategies to re-route packets after a failure occurs. If a link on the current path fails and an equal-cost path exists, the switch simply re-routes along that path. This approach is also known as *equal-cost multi-path routing* (ECMP). If no shortest path exist, it uses a 3-hop detour if one is available, and otherwise falls back to a 5-hop detour if necessary.

We implemented this routing scheme in ProbNetKAT in several steps. The first, $F10_0$, approximates the hashing behavior of ECMP by randomly selecting a port along one of the shortest paths to the destination. The second, $F10_3$, improves the resilience of $F10_0$ by augmenting it with 3-hop re-routing—e.g., consider the blue path in Figure 11(a). We find a port on C that connects to a different aggregation switch A' and forward the packet to A' . If there are multiple

such ports which have not failed, we choose one uniformly at random. The third, $F10_{3,5}$, attempts 5-hop re-routing in cases where $F10_3$ is unable to find a port on C whose adjacent link is up—e.g., consider the red path in Figure 11(a). The 5-hop rerouting strategy requires a flag to distinguish packets taking a detour from regular packets.

F10 network and failure model. We model the network as discussed in §2, focusing on packets destined to switch 1:

$$M(p) \triangleq \text{in}; \text{do } (p; t) \text{ while } (\neg \text{sw}=1)$$

McNetKAT automatically generates the topology program t from a Graphviz description. The ingress predicate in is a disjunction of switch-port tests over all ingress locations. Adding the failure model and some setup code to declare local variables tracking the health of individual links yields the complete network model:

$$\widehat{M}(p, f) \triangleq \text{var } \text{up}_1 \leftarrow 1 \text{ in } \dots \text{var } \text{up}_d \leftarrow 1 \text{ in } M(f; p)$$

Here, d is the maximum degree of a topology node. The entire model measures about 750 lines of ProbNetKAT code.

To evaluate the effect of different kinds of failures, we define a family of failure models f_k indexed by the maximum number of failures $k \in \mathbb{N} \cup \{\infty\}$ that may occur, where links fail otherwise independently with probability pr ; we leave pr implicit. To simplify the analysis, we focus on failures occurring on downward paths (note that $F10_0$ is able to route around failures on the upward path, unless the topology becomes disconnected).

Verifying refinement. Having implemented F10 as a series of three refinements, we would expect the probability of packet delivery to increase in each refinement, but not to achieve perfect delivery in an unbounded failure model f_∞ . Formally, we should have

$$\begin{aligned} \text{drop} &< \widehat{M}(F10_0, f_\infty) < \widehat{M}(F10_3, f_\infty) \\ &< \widehat{M}(F10_{3,5}, f_\infty) < \text{teleport} \end{aligned}$$

where teleport moves the packet directly to its destination, and $p < q$ means the probability assigned to every input-output pair by q is greater than the probability assigned by p . We confirmed that these inequalities hold using McNetKAT.

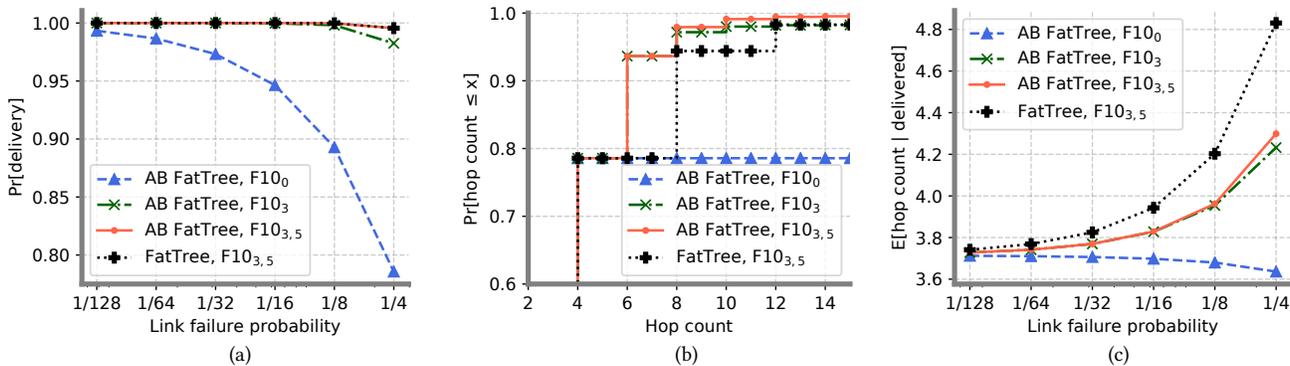


Figure 12. Case study results ($k = \infty$): (a) Probability of delivery vs. link-failure probability; (b) Increased path length due to resilience ($pr = 1/4$); (c) Expected hop-count conditioned on delivery.

Verifying k -resilience. Resilience is the key property satisfied by F10. By using McNetKAT, we were able to automatically verify that F10 is resilient to up to three failures in the AB FatTree Figure 11(a). To establish this property, we increased the parameter k in our failure model f_k while checking equivalence with teleportation (i.e., perfect delivery), as shown in Figure 11(b). The simplest scheme F10₀ drops packets when a failure occurs on the downward path, so it is 0-resilient. The F10₃ scheme routes around failures when a suitable aggregation switch is available, hence it is 2-resilient. Finally, the F10_{3,5} scheme routes around failures as long as any aggregation switch is reachable, hence it is 3-resilient. If the schemes are not equivalent to *teleport*, we can still compare the relative resilience of the schemes using the refinement order, as shown in Figure 11(c). Our implementation also enables precise, quantitative comparisons. For example, Figure 12(a) considers a failure model in which an unbounded number of failures can occur. We find that F10₀'s delivery probability dips significantly as the failure probability increases, while both F10₃ and F10_{3,5} continue to ensure high delivery probability by routing around failures.

Analyzing path stretch. Routing schemes based on detours achieve a higher degree of resilience at the cost of increasing the lengths of forwarding paths. We can quantify this increase by augmenting our model with a counter that is incremented at each hop and analyzing the expected path length. Figure 12(b) shows the cumulative distribution function of latency as the fraction of traffic delivered within a given hop count. On AB FatTree, F10₀ delivers $\approx 80\%$ of the traffic in 4 hops, since the maximum length of a shortest path from any edge switch to s_1 is 4 and F10₀ does not attempt to recover from failures. F10₃ and F10_{3,5} deliver the same amount of traffic when limited to at most 4 hops, but they can deliver significantly more traffic using 2 additional hops by using 3-hop and 5-hop paths to route around failures. F10₃ also delivers more traffic with 8 hops—these are the cases when F10₃ performs 3-hop re-routing twice for a

single packet as it encountered failure twice. We can also show that on a standard FatTree, F10_{3,5} failures have a higher impact on latency. Intuitively, the topology does not support 3-hop re-routing. This finding supports a key claim of F10: the topology and routing scheme should be co-designed to avoid excessive path stretch. Finally, Figure 12(c) shows the expected path length conditioned on delivery. As the failure probability increases, the probability of delivery for packets routed via the core layer decreases for F10₀. Thus, the distribution of delivered packets shifts towards 2-hop paths via an aggregation switch, so the expected hop-count decreases.

8 Related Work

The most closely related system to McNetKAT is Bayonet [17]. In contrast to the domain-specific approach followed in this paper, Bayonet uses a general-purpose probabilistic programming language and inference tool [18]. Such an approach, which reuses existing techniques, is naturally appealing. In addition, Bayonet is more expressive than McNetKAT: it supports asynchronous scheduling, stateful transformations, and probabilistic inference, making it possible to model richer phenomena, such as congestion due to packet-level interactions in queues. Of course, the extra generality does not come for free. Bayonet requires programmers to supply an upper bound on loops as the implementation is not guaranteed to find a fixed point. As discussed in §5, McNetKAT scales better than Bayonet on simple benchmarks. Another issue is that writing a realistic scheduler appears challenging, and one might also need to model host-level congestion control protocols to obtain accurate results. Currently Bayonet programs use deterministic or uniform schedulers and model only a few packets at a time [16].

Prior work on ProbNetKAT [46] gave a measure-theoretic semantics and an implementation that approximated programs using sequences of monotonically improving estimates. While these estimates were proven to converge in the limit, [46] offered no guarantees about the convergence

rate. In fact, there are examples where the approximations do not converge after any finite number of steps, which is obviously undesirable in a tool. The implementation only scaled to 10s of switches. In contrast, this paper presents a straightforward and implementable semantics; the implementation computes limits precisely in closed form, and it scales to real-world networks with thousands of switches. McNetKAT achieves this by restricting to the guarded and history-free fragment of ProbNetKAT, sacrificing the ability to reason about multicast and path-properties directly. In practice this sacrifice seems well worth the payoff: multicast is somewhat uncommon, and we can often reason about path-properties by maintaining extra state in the packets. In particular, McNetKAT can still model the examples studied in previous work by Smolka et al. [46].

Our work is the latest in a long line of techniques using Markov chains as a tool for representing and analyzing probabilistic programs. For an early example, see the seminal paper of Sharir et al. [43]. Markov chains are also used in many probabilistic model checkers, such as PRISM [31].

Beyond networking applications, there are connections to other work on verification of probabilistic programs. Di Pierro, Hankin, and Wiklicky used probabilistic abstract interpretation to statically analyze probabilistic λ -calculus [9]; their work was extended to a language $pWhile$, using a store and program location state space similar to Sharir et al. [43]. However, they do not deal with infinite limiting behavior beyond stepwise iteration, and do not guarantee convergence. Olejnik, Wicklicky, and Cheraghchi provided a probabilistic compiler pwc for a variation of $pWhile$ [38]; their optimizations could potentially be useful for McNetKAT. A recent survey by Gordon et al. [21] shows how to give semantics for probabilistic processes using stationary distributions of Markov chains, and studies convergence. Similar to our approach, they use absorbing strongly connected components to represent termination. Finally, probabilistic abstract interpretation is also an active area of research [49]; it would be interesting to explore applications to ProbNetKAT.

9 Conclusion

This paper presents a scalable tool for verifying probabilistic networks based on a new semantics for the history-free fragment of ProbNetKAT in terms of Markov chains. Natural directions for future work include further optimization of our implementation—e.g., using Bayesian networks to represent joint distributions compactly. We are also interested in applying McNetKAT to other systems that implement algorithms for randomized routing [30, 44], load balancing [11], traffic monitoring [42], anonymity [10], and network neutrality [52], among others.

Acknowledgments

We are grateful to the anonymous reviewers and our shepherd Michael Greenberg for their feedback and help in improving the paper. Thanks also to Jonathan DiLorenzo for suggesting improvements to the paper and for helping us locate a subtle performance bug, and to the Bellairs Research Institute of McGill University for providing a wonderful research environment. This work was supported in part by the National Science Foundation under grants NeTS-1413972 and AiTF-1637532, by the European Research Council under grant 679127, by a Facebook TAV award, by a Royal Society Wolfson fellowship, and a gift from Keysight.

References

- [1] S. B. Akers. 1978. Binary Decision Diagrams. *IEEE Trans. Comput.* 27, 6 (June 1978), 509–516. <https://doi.org/10.1109/TC.1978.1675141>
- [2] Mohammad Al-Fares, Alexander Loukissas, and Amin Vahdat. 2008. A Scalable, Commodity Data Center Network Architecture. In *ACM SIGCOMM Computer Communication Review*, Vol. 38. ACM, 63–74.
- [3] Rajeev Alur, Dana Fisman, and Mukund Raghothaman. 2016. Regular programming for quantitative properties of data streams. In *ESOP 2016*. 15–40.
- [4] Carolyn Jane Anderson, Nate Foster, Arjun Guha, Jean-Baptiste Jeannin, Dexter Kozen, Cole Schlesinger, and David Walker. 2014. NetKAT: Semantic Foundations for Networks. In *POPL*. 113–126.
- [5] Manav Bhatia, Mach Chen, Sami Boutros, Marc Binderberger, and Jeffrey Haas. 2014. Bidirectional Forwarding Detection (BFD) on Link Aggregation Group (LAG) Interfaces. RFC 7130. <https://doi.org/10.17487/RFC7130>
- [6] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, and David Walker. 2014. P4: Programming Protocol-Independent Packet Processors. *SIGCOMM CCR* 44, 3 (July 2014), 87–95.
- [7] Martin Casado, Nate Foster, and Arjun Guha. 2014. Abstractions for Software-Defined Networks. *CACM* 57, 10 (Oct. 2014), 86–95.
- [8] Timothy A. Davis. 2004. Algorithm 832: UMFPACK V4.3—an Unsymmetric-pattern Multifrontal Method. *ACM Trans. Math. Softw.* 30, 2 (June 2004), 196–199. <https://doi.org/10.1145/992200.992206>
- [9] Alessandra Di Pierro, Chris Hankin, and Herbert Wiklicky. 2005. Probabilistic λ -calculus and quantitative program analysis. *Journal of Logic and Computation* 15, 2 (2005), 159–179. <https://doi.org/10.1093/logcom/exi008>
- [10] Roger Dingleline, Nick Mathewson, and Paul Syverson. 2004. Tor: The Second-generation Onion Router. In *USENIX Security Symposium (SSYM)*. 21–21.
- [11] A. Dixit, P. Prakash, Y. C. Hu, and R. R. Kompella. 2013. On the impact of packet spraying in data center networks. In *IEEE INFOCOM*. 2130–2138.
- [12] Manfred Droste, Werner Kuich, and Heiko Vogler. 2009. *Handbook of Weighted Automata*. Springer.
- [13] Nate Foster, Dexter Kozen, Konstantinos Mamouras, Mark Reitblatt, and Alexandra Silva. 2016. Probabilistic NetKAT. In *ESOP*. 282–309. https://doi.org/10.1007/978-3-662-49498-1_12
- [14] Nate Foster, Dexter Kozen, Matthew Milano, Alexandra Silva, and Laure Thompson. 2015. A Coalgebraic Decision Procedure for NetKAT. In *POPL*. ACM, 343–355.
- [15] M. Fujita, P. C. McGeer, and J. C.-Y. Yang. 1997. Multi-Terminal Binary Decision Diagrams: An Efficient DataStructure for Matrix Representation. *Form. Methods Syst. Des.* 10, 2-3 (April 1997), 149–169. <https://doi.org/10.1023/A:1008647823331>

- [16] Timon Gehr, Sasa Misailovic, Petar Tsankov, Laurent Vanbever, Pascal Wiesmann, and Martin T. Vechev. 2018. Bayonet: Probabilistic Computer Network Analysis. Available at <https://github.com/eth-sri/bayonet/>.
- [17] Timon Gehr, Sasa Misailovic, Petar Tsankov, Laurent Vanbever, Pascal Wiesmann, and Martin T. Vechev. 2018. Bayonet: probabilistic inference for networks. In *ACM SIGPLAN PLDI*. 586–602.
- [18] Timon Gehr, Sasa Misailovic, and Martin T. Vechev. 2016. PSI: Exact Symbolic Inference for Probabilistic Programs. 62–83.
- [19] Phillipa Gill, Navendu Jain, and Nachiappan Nagappan. 2011. Understanding Network Failures in Data Centers: Measurement, Analysis, and Implications. In *ACM SIGCOMM*. 350–361.
- [20] Michele Giry. 1982. A categorical approach to probability theory. In *Categorical aspects of topology and analysis*. Springer, 68–85. <https://doi.org/10.1007/BFb0092872>
- [21] Andrew D Gordon, Thomas A Henzinger, Aditya V Nori, and Sriram K Rajamani. 2014. Probabilistic programming. In *Proceedings of the on Future of Software Engineering*. ACM, 167–181. <https://doi.org/10.1145/2593882.2593900>
- [22] Timothy V Griffiths. 1968. The unsolvability of the equivalence problem for Λ -free nondeterministic generalized machines. *Journal of the ACM* 15, 3 (1968), 409–413.
- [23] Tero Harju and Juhani Karhumäki. 1991. The equivalence problem of multitape finite automata. *Theoretical Computer Science* 78, 2 (1991), 347–355.
- [24] David M. Kahn. 2017. Undecidable Problems for Probabilistic Network Programming. In *MFCS 2017*. <http://hdl.handle.net/1813/51765>
- [25] Peyman Kazemian, George Varghese, and Nick McKeown. 2012. Header Space Analysis: Static Checking for Networks. In *USENIX NSDI 2012*. 113–126. <https://www.usenix.org/conference/nsdi12/technical-sessions/presentation/kazemian>
- [26] John G Kemeny and James Laurie Snell. 1960. *Finite markov chains*. Vol. 356. van Nostrand Princeton, NJ.
- [27] Ahmed Khurshid, Wenxuan Zhou, Matthew Caesar, and Brighton Godfrey. 2012. Veriflow: Verifying Network-Wide Invariants in Real Time. In *ACM SIGCOMM*. 467–472.
- [28] Dexter Kozen. 1981. Semantics of probabilistic programs. *J. Comput. Syst. Sci.* 22, 3 (1981), 328–350. [https://doi.org/10.1016/0022-0000\(81\)90036-2](https://doi.org/10.1016/0022-0000(81)90036-2)
- [29] Dexter Kozen. 1997. Kleene algebra with tests. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 19, 3 (May 1997), 427–443. <https://doi.org/10.1145/256167.256195>
- [30] Praveen Kumar, Yang Yuan, Chris Yu, Nate Foster, Robert Kleinberg, Petr Lapukhov, Chiun Lin Lim, and Robert Soulé. 2018. Semi-Oblivious Traffic Engineering: The Road Not Taken. In *USENIX NSDI*.
- [31] M. Kwiatkowska, G. Norman, and D. Parker. 2011. PRISM 4.0: Verification of Probabilistic Real-time Systems. In *Proc. 23rd International Conference on Computer Aided Verification (CAV'11) (LNCS)*, G. Gopalakrishnan and S. Qadeer (Eds.), Vol. 6806. Springer, 585–591. https://doi.org/10.1007/978-3-642-22110-1_47
- [32] Marta Z. Kwiatkowska, Gethin Norman, and David Parker. 2011. PRISM 4.0: Verification of Probabilistic Real-Time Systems. In *CAV*. 585–591.
- [33] Jed Liu, William Hallahan, Cole Schlesinger, Milad Sharif, Jeongkeun Lee, Robert Soulé, Han Wang, Calin Cascaval, Nick McKeown, and Nate Foster. 2018. p4v: Practical Verification for Programmable Data Planes. In *SIGCOMM*. 490–503.
- [34] Vincent Liu, Daniel Halperin, Arvind Krishnamurthy, and Thomas E Anderson. 2013. F10: A Fault-Tolerant Engineered Network. In *USENIX NSDI*. 399–412.
- [35] Hao-hui Mai, Ahmed Khurshid, Rachit Agarwal, Matthew Caesar, P. Brighton Godfrey, and Samuel Talmadge King. 2011. Debugging the Data Plane with Anteater. In *ACM SIGCOMM*. 290–301.
- [36] Nick McKeown, Tom Anderson, Hari Balakrishnan, Guru Parulkar, Larry Peterson, Jennifer Rexford, Scott Shenker, and Jonathan Turner. 2008. OpenFlow: Enabling Innovation in Campus Networks. *SIGCOMM CCR* 38, 2 (2008), 69–74.
- [37] Mehryar Mohri. 2000. Generic ϵ -removal algorithm for weighted automata. In *CIAA 2000*. Springer, 230–242.
- [38] Maciej Olejnik, Herbert Wiklicky, and Mahdi Cheraghchi. 2016. Probabilistic Programming and Discrete Time Markov Chains. <http://www.imperial.ac.uk/media/imperial-college/faculty-of-engineering/computing/public/MaciejOlejnik.pdf>
- [39] Michael O Rabin and Dana Scott. 1959. Finite automata and their decision problems. *IBM Journal of Research and Development* 3, 2 (1959), 114–125.
- [40] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. 2015. Inside the Social Network’s (Datacenter) Network. In *ACM SIGCOMM*. 123–137.
- [41] N. Saheb-Djahromi. 1980. CPOs of measures for nondeterminism. *Theoretical Computer Science* 12 (1980), 19–37. [https://doi.org/10.1016/0304-3975\(80\)90003-1](https://doi.org/10.1016/0304-3975(80)90003-1)
- [42] Vyas Sekar, Michael K. Reiter, Walter Willinger, Hui Zhang, Ramana Rao Kompella, and David G. Andersen. 2008. CSAMP: A System for Network-wide Flow Monitoring. In *USENIX NSDI*. 233–246.
- [43] Micha Sharir, Amir Pnueli, and Sergiu Hart. 1984. Verification of probabilistic programs. *SIAM J. Comput.* 13, 2 (1984), 292–314. <https://doi.org/10.1137/0213021>
- [44] Rachee Singh, Manya Ghobadi, Klaus-Tycho Foerster, Mark Filer, and Phillipa Gill. 2018. RADWAN: Rate Adaptive Wide Area Network. In *ACM SIGCOMM*.
- [45] Steffen Smolka, Spiros Eliopoulos, Nate Foster, and Arjun Guha. 2015. A Fast Compiler for NetKAT. In *ICFP 2015*. <https://doi.org/10.1145/2784731.2784761>
- [46] Steffen Smolka, Praveen Kumar, Nate Foster, Dexter Kozen, and Alexandra Silva. 2017. Cantor Meets Scott: Semantic Foundations for Probabilistic Networks. In *POPL 2017*. <https://doi.org/10.1145/3009837.3009843>
- [47] Ken Thompson. 1968. Regular Expression Search Algorithm. *Commun. ACM* 11, 6 (1968), 419–422. <https://doi.org/10.1145/363347.363387>
- [48] L. Valiant. 1982. A Scheme for Fast Parallel Communication. *SIAM J. Comput.* 11, 2 (1982), 350–361.
- [49] Di Wang, Jan Hoffmann, and Thomas Reps. 2018. PMAF: An Algebraic Framework for Static Analysis of Probabilistic Programs. In *POPL 2018*. <https://www.cs.cmu.edu/~janh/papers/WangHR17.pdf>
- [50] James Worrell. 2013. Revisiting the equivalence problem for finite multitape automata. In *International Colloquium on Automata, Languages, and Programming (ICALP)*. Springer, 422–433.
- [51] Geoffrey G. Xie, Jibin Zhan, David A. Maltz, Hui Zhang, Albert G. Greenberg, Gísli Hjálmtýsson, and Jennifer Rexford. 2005. On static reachability analysis of IP networks. In *INFOCOM*.
- [52] Zhiyong Zhang, Ovidiu Mara, and Katerina Argyraki. 2014. Network Neutrality Inference. In *ACM SIGCOMM*. 63–74.

<p>Semantics $\boxed{\llbracket p \rrbracket \in 2^{\text{Pk}} \rightarrow \mathcal{D}(2^{\text{Pk}})}$</p> <p>$\llbracket \text{drop} \rrbracket(a) \triangleq \delta(\emptyset)$ $\llbracket \text{skip} \rrbracket(a) \triangleq \delta(a)$ $\llbracket f = n \rrbracket(a) \triangleq \delta(\{\pi \in a \mid \pi.f = n\})$ $\llbracket f \leftarrow n \rrbracket(a) \triangleq \delta(\{\pi[f := n] \mid \pi \in a\})$ $\llbracket \neg t \rrbracket(a) \triangleq \mathcal{D}(\lambda b. a - b)(\llbracket t \rrbracket(a))$ $\llbracket p \& q \rrbracket(a) \triangleq \mathcal{D}(\cup)(\llbracket p \rrbracket(a) \times \llbracket q \rrbracket(a))$ $\llbracket p ; q \rrbracket(a) \triangleq \llbracket q \rrbracket^\dagger(\llbracket p \rrbracket(a))$ $\llbracket p \oplus_r q \rrbracket(a) \triangleq r \cdot \llbracket p \rrbracket(a) + (1 - r) \cdot \llbracket q \rrbracket(a)$ $\llbracket p^* \rrbracket(a) \triangleq \bigsqcup_{n \in \mathbb{N}} \llbracket p^{(n)} \rrbracket(a)$ where $p^{(0)} \triangleq \text{skip}$, $p^{(n+1)} \triangleq \text{skip} \& p ; p^{(n)}$</p>	<p>(Discrete) Probability Monad \mathcal{D}</p> <p>Unit $\delta : X \rightarrow \mathcal{D}(X)$ $\delta(x) \triangleq \delta_x$ Bind $-^\dagger : (X \rightarrow \mathcal{D}(Y)) \rightarrow \mathcal{D}(X) \rightarrow \mathcal{D}(Y)$ $f^\dagger(\mu)(A) \triangleq \sum_{x \in X} f(x)(A) \cdot \mu(x)$</p>
--	--

Figure 13. ProbNetKAT semantics.

A ProbNetKAT Denotational Semantics

In the original ProbNetKAT language, programs manipulate sets of *packet histories*—non-empty, finite sequences of packets modeling trajectories through the network [13, 46]. The resulting state space is uncountable and modeling the semantics properly requires full-blown measure theory as some programs generate continuous distributions. In the history-free fragment, programs manipulate sets of packets and the state space is finite, which makes the semantics considerably simpler.

Proposition A.1. *Let $\llbracket - \rrbracket$ denote the semantics defined in Smolka et al. [46]. Then for all dup-free programs p and inputs $a \in 2^{\text{Pk}}$, we have $\llbracket p \rrbracket(a) = \langle p \rangle(a)$, where we identify packets and histories of length one.*

Throughout this paper, we can work in the discrete space 2^{Pk} , i.e., the set of sets of packets. An *outcome* (denoted by lowercase variables a, b, c, \dots) is a set of packets and an *event* (denoted by uppercase variables A, B, C, \dots) is a set of outcomes. Given a discrete probability measure on this space, the probability of an event is the sum of the probabilities of its outcomes.

ProbNetKAT programs are interpreted as *Markov kernels* on the space 2^{Pk} . A Markov kernel is a function $2^{\text{Pk}} \rightarrow \mathcal{D}(2^{\text{Pk}})$ where \mathcal{D} is the probability (or Giry) monad [20, 28]. Thus, a program p maps an input set of packets $a \in 2^{\text{Pk}}$ to a *distribution* $\llbracket p \rrbracket(a) \in \mathcal{D}(2^{\text{Pk}})$ over output sets of packets. The semantics uses the following probabilistic constructions:²

- For a discrete measurable space X , $\mathcal{D}(X)$ denotes the set of probability measures over X ; that is, the set of countably additive functions $\mu : 2^X \rightarrow [0, 1]$ with $\mu(X) = 1$.
- For a measurable function $f : X \rightarrow Y$, $\mathcal{D}(f) : \mathcal{D}(X) \rightarrow \mathcal{D}(Y)$ denotes the *pushforward* along f ; that is, the function that maps a measure μ on X to

$$\mathcal{D}(f)(\mu) \triangleq \mu \circ f^{-1} = \lambda A \in \Sigma_Y. \mu(\{x \in X \mid f(x) \in A\})$$

which is called the *pushforward measure* on Y .

- The *unit* $\delta : X \rightarrow \mathcal{D}(X)$ of the monad maps a point $x \in X$ to the point mass (or *Dirac measure*) $\delta_x \in \mathcal{D}(X)$. The Dirac measure is given by

$$\delta_x(A) \triangleq [x \in A]$$

That is, the Dirac measure is 1 if $x \in A$ and 0 otherwise.

- The *bind* operation of the monad,

$$-^\dagger : (X \rightarrow \mathcal{D}(Y)) \rightarrow \mathcal{D}(X) \rightarrow \mathcal{D}(Y)$$

lifts a function $f : X \rightarrow \mathcal{D}(Y)$ with deterministic inputs to a function $f^\dagger : \mathcal{D}(X) \rightarrow \mathcal{D}(Y)$ that takes random inputs. Intuitively, this is achieved by averaging the output of f when the inputs are randomly distributed according to μ . Formally,

$$f^\dagger(\mu)(A) \triangleq \sum_{x \in X} f(x)(A) \cdot \mu(x).$$

²These can also be defined for uncountable spaces, as would be required to handle the full language.

- Given two measures $\mu \in \mathcal{D}(X)$ and $\nu \in \mathcal{D}(Y)$, $\mu \times \nu \in \mathcal{D}(X \times Y)$ denotes their *product measure*. This is the unique measure satisfying

$$(\mu \times \nu)(A \times B) = \mu(A) \cdot \nu(B)$$

Intuitively, it models distributions over pairs of independent values.

Using these primitives, we can now make our operational intuitions precise (see Figure 13 for formal definitions). A predicate t maps the set of input packets $a \in 2^{\text{Pk}}$ to the subset of packets $b \subseteq a$ satisfying the predicate (with probability 1). Hence, drop drops all packets (i.e., it returns the empty set) while skip keeps all packets (i.e., it returns the input set). The test $f=n$ returns the subset of input packets whose f -field is n . Negation $\neg t$ filters out the packets returned by t .

Parallel composition $p \& q$ executes p and q independently on the input set, then returns the union of their results. Note that packet sets do *not* model nondeterminism, unlike the usual situation in Kleene algebras—rather, they model collections of packets traversing possibly different portions of the network simultaneously. In particular, the union operation is *not* idempotent: $p \& p$ need not have the same semantics as p . Probabilistic choice $p \oplus_r q$ feeds the input to both p and q and returns a convex combination of the output distributions according to r . Sequential composition $p; q$ can be thought of as a two-stage probabilistic process: it first executes p on the input set to obtain a random intermediate result, then feeds that into q to obtain the final distribution over outputs. The outcome of q is averaged over the distribution of intermediate results produced by p .

We say that two programs are *equivalent*, denoted $p \equiv q$, if they denote the same Markov kernel, i.e. if $\llbracket p \rrbracket = \llbracket q \rrbracket$. As usual, we expect Kleene star p^* to satisfy the characteristic fixed point equation $p^* \equiv \text{skip} \& p; p^*$, which allows it to be unrolled ad infinitum. Thus we define it as the supremum of its finite unrollings $p^{(n)}$; see Figure 13. This supremum is taken in a CPO $(\mathcal{D}(2^{\text{Pk}}), \sqsubseteq)$ of distributions that is described in more detail in Appendix A.1. The partial ordering \sqsubseteq on packet set distributions gives rise to a partial ordering on programs: we write $p \leq q$ iff $\llbracket p \rrbracket(a) \sqsubseteq \llbracket q \rrbracket(a)$ for all inputs $a \in 2^{\text{Pk}}$. Intuitively, $p \leq q$ iff p produces any particular output packet π with probability at most that of q for any fixed input— q has a larger probability of delivering more output packets.

A.1 The CPO $(\mathcal{D}(2^{\text{Pk}}), \sqsubseteq)$

The space 2^{Pk} with the subset order forms a CPO $(2^{\text{Pk}}, \subseteq)$. Following Saheb-Djahromi [41], this CPO can be lifted to a CPO $(\mathcal{D}(2^{\text{Pk}}), \sqsubseteq)$ on distributions over 2^{Pk} . Because 2^{Pk} is a finite space, the resulting ordering \sqsubseteq on distributions takes a particularly easy form:

$$\mu \sqsubseteq \nu \iff \mu(\{a\}^\uparrow) \leq \nu(\{a\}^\uparrow) \text{ for all } a \subseteq \text{Pk}$$

where $\{a\}^\uparrow \triangleq \{b \mid a \subseteq b\}$ denotes upward closure. Intuitively, ν produces more outputs than μ . As was shown in Smolka et al. [46], ProbNetKAT satisfies various monotonicity (and continuity) properties with respect to this ordering, including

$$a \subseteq a' \implies \llbracket p \rrbracket(a) \sqsubseteq \llbracket p \rrbracket(a') \quad \text{and} \quad n \leq m \implies \llbracket p^{(n)} \rrbracket(a) \sqsubseteq \llbracket p^{(m)} \rrbracket(a).$$

As a result, the semantics of p^* as the supremum of its finite unrollings $p^{(n)}$ is well-defined.

While the semantics of full ProbNetKAT requires more domain theory to give a satisfactory characterization of Kleene star, a simpler characterization suffices for the history-free fragment.

Lemma A.2 (Pointwise Convergence). *Let $A \subseteq 2^{\text{Pk}}$. Then for all programs p and inputs $a \in 2^{\text{Pk}}$,*

$$\llbracket p^* \rrbracket(a)(A) = \lim_{n \rightarrow \infty} \llbracket p^{(n)} \rrbracket(a)(A).$$

B Omitted Proofs

Lemma B.1. *Let A be a finite boolean combination of basic open sets, i.e. sets of the form $B_a = \{a\}^\uparrow$ for $a \in \wp_\omega(\text{H})$, and let $\llbracket - \rrbracket$ denote the semantics from Smolka et al. [46]. Then for all programs p and inputs $a \in 2^{\text{H}}$,*

$$\llbracket p^* \rrbracket(a)(A) = \lim_{n \rightarrow \infty} \llbracket p^{(n)} \rrbracket(a)(A)$$

Proof. Using topological arguments, the claim follows directly from previous results: A is a Cantor-clopen set by Smolka et al. [46] (i.e., both A and \bar{A} are Cantor-open), so its indicator function $\mathbf{1}_A$ is Cantor-continuous. But $\mu_n \triangleq \llbracket p^{(n)} \rrbracket(a)$ converges weakly to $\mu \triangleq \llbracket p^* \rrbracket(a)$ in the Cantor topology [13, Theorem 4], so

$$\lim_{n \rightarrow \infty} \llbracket p^{(n)} \rrbracket(a)(A) = \lim_{n \rightarrow \infty} \int \mathbf{1}_A d\mu_n = \int \mathbf{1}_A d\mu = \llbracket p^* \rrbracket(a)(A)$$

(To see why A and \bar{A} are open in the Cantor topology, note that they can be written in disjunctive normal form over atoms $B_{\{h\}}$.) \square

Predicates in ProbNetKAT form a Boolean algebra.

Lemma B.2. *Every predicate t satisfies $\llbracket t \rrbracket(a) = \delta_{a \cap b_t}$ for a certain packet set $b_t \subseteq \text{Pk}$, where*

- $b_{\text{drop}} = \emptyset$,
- $b_{\text{skip}} = \text{Pk}$,
- $b_{f=n} = \{\pi \in \text{Pk} \mid \pi.f = n\}$,
- $b_{\neg t} = \text{Pk} - b_t$,
- $b_{t \& u} = b_t \cup b_u$, and
- $b_{t;u} = b_t \cap b_u$.

Proof. For drop, skip, and $f=n$, the claim holds trivially. For $\neg t$, $t \& u$, and $t;u$, the claim follows inductively, using that $\mathcal{D}(f)(\delta_b) = \delta_{f(b)}$, $\delta_b \times \delta_c = \delta_{(b,c)}$, and that $f^\dagger(\delta_b) = f(b)$. The first and last equations hold because $\langle \mathcal{D}, \delta, -^\dagger \rangle$ is a monad. \square

Proof of Proposition A.1. We only need to show that for dup-free programs p and history-free inputs $a \in 2^{\text{Pk}}$, $\llbracket p \rrbracket(a)$ is a distribution on packets (where we identify packets and singleton histories). We proceed by structural induction on p . All cases are straightforward except perhaps the case of p^* . For this case, by the induction hypothesis, all $\llbracket p^{(n)} \rrbracket(a)$ are discrete probability distributions on packet sets, therefore vanish outside 2^{Pk} . By Lemma B.1, this is also true of the limit $\llbracket p^* \rrbracket(a)$, as its value on 2^{Pk} must be 1, therefore it is also a discrete distribution on packet sets. \square

Proof of Lemma A.2. This follows directly from Lemma B.1 and Proposition A.1 by noticing that any set $A \subseteq 2^{\text{Pk}}$ is a finite boolean combination of basic open sets. \square

Proof of Theorem 3.1. It suffices to show the equality $\mathcal{B}\llbracket p \rrbracket_{ab} = \llbracket p \rrbracket(a)(\{b\})$; the remaining claims then follow by well-definedness of $\llbracket - \rrbracket$. The equality is shown using Lemma A.2 and a routine induction on p :

For $p = \text{drop}$, skip, $f=n$, $f \leftarrow n$ we have

$$\llbracket p \rrbracket(a)(\{b\}) = \delta_c(\{b\}) = [b = c] = \mathcal{B}\llbracket p \rrbracket_{ab}$$

for $c = \emptyset$, a , $\{\pi \in a \mid \pi.f = n\}$, $\{\pi[f := n] \mid \pi \in a\}$, respectively.

For $\neg t$ we have,

$$\begin{aligned} \mathcal{B}\llbracket \neg t \rrbracket_{ab} &= [b \subseteq a] \cdot \mathcal{B}\llbracket t \rrbracket_{a, a-b} \\ &= [b \subseteq a] \cdot \llbracket t \rrbracket(a)(\{a-b\}) && \text{(IH)} \\ &= [b \subseteq a] \cdot [a-b = a \cap b_t] && \text{(Lemma B.2)} \\ &= [b \subseteq a] \cdot [a-b = a - (H - b_t)] \\ &= [b = a \cap (H - b_t)] \\ &= \llbracket \neg t \rrbracket(a)(b) && \text{(Lemma B.2)} \end{aligned}$$

For $p \& q$, letting $\mu = \llbracket p \rrbracket(a)$ and $\nu = \llbracket q \rrbracket(a)$ we have

$$\begin{aligned} \llbracket p \& q \rrbracket(a)(\{b\}) &= (\mu \times \nu)(\{(b_1, b_2) \mid b_1 \cup b_2 = b\}) \\ &= \sum_{b_1, b_2} [b_1 \cup b_2 = b] \cdot (\mu \times \nu)(\{(b_1, b_2)\}) \\ &= \sum_{b_1, b_2} [b_1 \cup b_2 = b] \cdot \mu(\{b_1\}) \cdot \nu(\{b_2\}) \\ &= \sum_{b_1, b_2} [b_1 \cup b_2 = b] \cdot \mathcal{B}\llbracket p \rrbracket_{ab_1} \cdot \mathcal{B}\llbracket q \rrbracket_{ab_2} && \text{(IH)} \\ &= \mathcal{B}\llbracket p \& q \rrbracket_{ab} \end{aligned}$$

where we use in the second step that $b \subseteq \text{Pk}$ is finite, thus $\{(b_1, b_2) \mid b_1 \cup b_2 = b\}$ is finite.

For $p; q$, let $\mu = \llbracket p \rrbracket(a)$ and $\nu_c = \llbracket q \rrbracket(c)$ and recall that μ is a discrete distribution on 2^{Pk} . Thus

$$\begin{aligned} \llbracket p; q \rrbracket(a)(\{b\}) &= \sum_{c \in 2^{\text{Pk}}} \nu_c(\{b\}) \cdot \mu(\{c\}) \\ &= \sum_{c \in 2^{\text{Pk}}} \mathcal{B}\llbracket q \rrbracket_{c,b} \cdot \mathcal{B}\llbracket p \rrbracket_{a,c} \\ &= \mathcal{B}\llbracket p; q \rrbracket_{ab}. \end{aligned}$$

For $p \oplus_r q$, the claim follows directly from the induction hypotheses.

Finally, for p^* , we know that $\mathcal{B}[\![p^{(n)}]\!]_{ab} = \llbracket p^{(n)} \rrbracket(a)(\{b\})$ by induction hypothesis. The key to proving the claim is [Lemma A.2](#), which allows us to take the limit on both sides and deduce

$$\mathcal{B}[\![p^*]\!]_{ab} = \lim_{n \rightarrow \infty} \mathcal{B}[\![p^{(n)}]\!]_{ab} = \lim_{n \rightarrow \infty} \llbracket p^{(n)} \rrbracket(a)(\{b\}) = \llbracket p^* \rrbracket(a)(\{b\}). \quad \square$$

Proof of Lemma 4.1. For arbitrary $a, b \subseteq \text{Pk}$, we have

$$\begin{aligned} \sum_{a', b'} \mathcal{S}[\![p]\!]_{(a,b),(a',b')} &= \sum_{a', b'} [b' = a \cup b] \cdot \mathcal{B}[\![p]\!]_{a,a'} \\ &= \sum_{a'} \left(\sum_{b'} [b' = a \cup b] \right) \cdot \mathcal{B}[\![p]\!]_{a,a'} \\ &= \sum_{a'} \mathcal{B}[\![p]\!]_{a,a'} = 1 \end{aligned}$$

where in the last step, we use that $\mathcal{B}[\![p]\!]$ is stochastic ([Theorem 3.1](#)). □

Proof of Lemma 4.3. By induction on $n \geq 0$. For $n = 0$, we have

$$\begin{aligned} \sum_{a'} [b' = a' \cup b] \cdot \mathcal{B}[\![p^{(n)}]\!]_{a,a'} &= \sum_{a'} [b' = a' \cup b] \cdot \mathcal{B}[\![\text{skip}]\!]_{a,a'} \\ &= \sum_{a'} [b' = a' \cup b] \cdot [a = a'] \\ &= [b' = a \cup b] \\ &= [b' = a \cup b] \cdot \sum_{a'} \mathcal{B}[\![p]\!]_{a,a'} \\ &= \sum_{a'} \mathcal{S}[\![p]\!]_{(a,b),(a',b')} \end{aligned}$$

In the induction step ($n > 0$),

$$\begin{aligned} &\sum_{a'} [b' = a' \cup b] \cdot \mathcal{B}[\![p^{(n)}]\!]_{a,a'} \\ &= \sum_{a'} [b' = a' \cup b] \cdot \mathcal{B}[\![\text{skip} \ \& \ p; p^{(n-1)}]\!]_{a,a'} \\ &= \sum_{a'} [b' = a' \cup b] \cdot \sum_c [a' = a \cup c] \cdot \mathcal{B}[\![p; p^{(n-1)}]\!]_{a,c} \\ &= \sum_c \left(\sum_{a'} [b' = a' \cup b] \cdot [a' = a \cup c] \right) \cdot \sum_k \mathcal{B}[\![p]\!]_{a,k} \cdot \mathcal{B}[\![p^{(n-1)}]\!]_{k,c} \\ &= \sum_{c,k} [b' = a \cup c \cup b] \cdot \mathcal{B}[\![p]\!]_{a,k} \cdot \mathcal{B}[\![p^{(n-1)}]\!]_{k,c} \\ &= \sum_k \mathcal{B}[\![p]\!]_{a,k} \cdot \sum_{a'} [b' = a' \cup (a \cup b)] \cdot \mathcal{B}[\![p^{(n-1)}]\!]_{k,a'} \\ &= \sum_k \mathcal{B}[\![p]\!]_{a,k} \cdot \sum_{a'} \mathcal{S}[\![p]\!]_{(k, a \cup b),(a',b')}^n \\ &= \sum_{a'} \sum_{k_1, k_2} [k_2 = a \cup b] \cdot \mathcal{B}[\![p]\!]_{a,k_1} \cdot \mathcal{S}[\![p]\!]_{(k_1, k_2),(a',b')}^n \\ &= \sum_{a'} \sum_{k_1, k_2} \mathcal{S}[\![p]\!]_{(a,b)(k_1, k_2)} \cdot \mathcal{S}[\![p]\!]_{(k_1, k_2),(a',b')}^n \\ &= \sum_{a'} \mathcal{S}[\![p]\!]_{(a,b),(a',b')}^{n+1} \end{aligned} \quad \square$$

Lemma B.3. *The matrix $X = I - Q$ in Equation (2) of §4 is invertible.*

Proof. Let S be a finite set of states, $|S| = n$, M an $S \times S$ substochastic matrix ($M_{st} \geq 0$, $M\mathbf{1} \leq \mathbf{1}$). A state s is *defective* if $(M\mathbf{1})_s < 1$. We say M is *stochastic* if $M\mathbf{1} = \mathbf{1}$, *irreducible* if $(\sum_{i=0}^{n-1} M^i)_{st} > 0$ (that is, the support graph of M is strongly connected), and *aperiodic* if all entries of some power of M are strictly positive.

We show that if M is substochastic such that every state can reach a defective state via a path in the support graph, then the spectral radius of M is strictly less than 1. Intuitively, all weight in the system eventually drains out at the defective states.

Let e_s , $s \in S$, be the standard basis vectors. As a distribution, e_s^T is the unit point mass on s . For $A \subseteq S$, let $e_A = \sum_{s \in A} e_s$. The L_1 -norm of a substochastic vector is its total weight as a distribution. Multiplying on the right by M never increases total weight, but will strictly decrease it if there is nonzero weight on a defective state. Since every state can reach a defective state, this must happen after n steps, thus $\|e_s^T M^n\|_1 < 1$. Let $c = \max_s \|e_s^T M^n\|_1 < 1$. For any $y = \sum_s a_s e_s$,

$$\begin{aligned} \|y^T M^n\|_1 &= \|(\sum_s a_s e_s)^T M^n\|_1 \\ &\leq \sum_s |a_s| \cdot \|e_s^T M^n\|_1 \leq \sum_s |a_s| \cdot c = c \cdot \|y^T\|_1. \end{aligned}$$

Then M^n is contractive in the L_1 norm, so $|\lambda| < 1$ for all eigenvalues λ . Thus $I - M$ is invertible because 1 is not an eigenvalue of M . \square

Proof of Proposition 4.6.

1. It suffices to show that $USU = SU$. Suppose that

$$\Pr[(a, b) \xrightarrow{USU}_1 (a', b')] = p > 0.$$

It suffices to show that this implies

$$\Pr[(a, b) \xrightarrow{SU}_1 (a', b')] = p.$$

If (a, b) is saturated, then we must have $(a', b') = (\emptyset, b)$ and

$$\Pr[(a, b) \xrightarrow{USU}_1 (\emptyset, b)] = 1 = \Pr[(a, b) \xrightarrow{SU}_1 (\emptyset, b)]$$

If (a, b) is not saturated, then $(a, b) \xrightarrow{U}_1 (a, b)$ with probability 1 and therefore

$$\Pr[(a, b) \xrightarrow{USU}_1 (a', b')] = \Pr[(a, b) \xrightarrow{SU}_1 (a', b')]$$

2. Since S and U are stochastic, clearly SU is a MC. Since SU is finite state, any state can reach an absorbing communication class. (To see this, note that the reachability relation \xrightarrow{SU} induces a partial order on the communication classes of SU . Its maximal elements are necessarily absorbing, and they must exist because the state space is finite.) It thus suffices to show that a state set $C \subseteq 2^{\text{Pk}} \times 2^{\text{Pk}}$ in SU is an absorbing communication class iff $C = \{(\emptyset, b)\}$ for some $b \subseteq \text{Pk}$.

“ \Leftarrow ”: Observe that $\emptyset \xrightarrow{B}_1 a'$ iff $a' = \emptyset$. Thus $(\emptyset, b) \xrightarrow{S}_1 (a', b')$ iff $a' = \emptyset$ and $b' = b$, and likewise $(\emptyset, b) \xrightarrow{U}_1 (a', b')$ iff $a' = \emptyset$ and $b' = b$. Thus (\emptyset, b) is an absorbing state in SU as required.

“ \Rightarrow ”: First observe that by monotonicity of SU (Lemma 4.5), we have $b = b'$ whenever $(a, b) \xleftrightarrow{SU} (a', b')$; thus there exists a fixed b_C such that $(a, b) \in C$ implies $b = b_C$.

Now pick an arbitrary state $(a, b_C) \in C$. It suffices to show that $(a, b_C) \xrightarrow{SU} (\emptyset, b_C)$, because that implies $(a, b_C) \xleftrightarrow{SU} (\emptyset, b_C)$, which in turn implies $a = \emptyset$. But the choice of $(a, b_C) \in C$ was arbitrary, so that would mean $C = \{(\emptyset, b_C)\}$ as claimed.

To show that $(a, b_C) \xrightarrow{SU} (\emptyset, b_C)$, pick arbitrary states such that

$$(a, b_C) \xrightarrow{S} (a', b') \xrightarrow{U}_1 (a'', b'')$$

and recall that this implies $(a, b_C) \xrightarrow{SU} (a'', b'')$ by claim (1). Then $(a'', b'') \xrightarrow{SU} (a, b_C)$ because C is absorbing, and thus $b_C = b' = b''$ by monotonicity of S , U , and SU . But (a', b') was chosen as an arbitrary state S -reachable from (a, b_C) , so (a, b) and by transitivity (a', b') must be saturated. Thus $a'' = \emptyset$ by the definition of U . \square

Proof of Theorem 4.7. Using Proposition 4.6.1 in the second step and equation (3) in the last step,

$$\begin{aligned} \lim_{n \rightarrow \infty} \sum_{a'} S_{(a,b),(a',b')}^n &= \lim_{n \rightarrow \infty} \sum_{a'} (S^n U)_{(a,b),(a',b')} \\ &= \lim_{n \rightarrow \infty} \sum_{a'} (SU)_{(a,b),(a',b')}^n \\ &= \sum_{a'} (SU)_{(a,b),(a',b')}^\infty = (SU)_{(a,b),(\emptyset,b')}^\infty \end{aligned}$$

$(SU)^\infty$ is computable because S and U are matrices over \mathbb{Q} and hence so is $(I - Q)^{-1}R$. \square

Corollary B.4. *For programs p and q , it is decidable whether $p \equiv q$.*

Proof of Corollary B.4. Recall from Corollary 3.2 that it suffices to compute the finite rational matrices $\mathcal{B}[p]$ and $\mathcal{B}[q]$ and check them for equality. But Theorem 4.7 together with Proposition 4.2 gives us an effective mechanism to compute $\mathcal{B}[-]$ in the case of Kleene star, and $\mathcal{B}[-]$ is straightforward to compute in all other cases. Summarizing the full chain of equalities, we have:

$$\llbracket p^* \rrbracket(a)(\{b\}) = \mathcal{B}[p^*]_{a,b} = \lim_{n \rightarrow \infty} \mathcal{B}[p^{(n)}]_{a,b} = \lim_{n \rightarrow \infty} \sum_{a'} \mathcal{S}[p]_{(a,\emptyset),(a',b)}^n = (SU)_{(a,\emptyset),(\emptyset,b)}^\infty$$

following from Theorem 3.1, Definition of $\mathcal{B}[-]$, Proposition 4.2, and finally Theorem 4.7. \square

C Handling Full ProbNetKAT: Obstacles and Challenges

History-free ProbNetKAT can describe sophisticated network routing schemes under various failure models, and the program semantics can be computed exactly. Performing quantitative reasoning in full ProbNetKAT appears significantly more challenging. We illustrate some of the difficulties in deciding program equivalence; recall that this is decidable for the history-free fragment (Corollary B.4).

The main difference in the original ProbNetKAT language is an additional primitive `dup`. Intuitively, this command duplicates a packet $\pi \in \text{Pk}$ and outputs the word $\pi\pi \in \text{H}$, where $\text{H} = \text{Pk}^*$ is the set of non-empty, finite sequences of packets. An element of H is called a *packet history*, representing a log of previous packet states. ProbNetKAT policies may only modify the first (*head*) packet of each history; `dup` fixes the current head packet into the log by copying it. In this way, ProbNetKAT policies can compute distributions over the paths used to forward packets, instead of just over the final output packets.

However, with `dup`, the semantics of ProbNetKAT becomes significantly more complex. Policies p now transform sets of packet histories $a \in 2^{\text{H}}$ to distributions $\llbracket p \rrbracket(a) \in \mathcal{D}(2^{\text{H}})$. Since 2^{H} is uncountable, these distributions are no longer guaranteed to be discrete, and formalizing the semantics requires full-blown measure theory (see prior work for details [46]).

Without `dup`, policies operate on sets of packets 2^{Pk} ; crucially, this is a *finite* set and we can represent each set with a single state in a finite Markov chain. With `dup`, policies operate on sets of packet histories 2^{H} . Since this set is not finite—in fact, it is not even countable—encoding each packet history as a state would give a Markov chain with infinitely many states. Procedures for deciding equivalence are not known for such systems in general.

While in principle there could be a more compact representation of general ProbNetKAT policies as finite Markov chains or other models where equivalence is decidable, (e.g., weighted or probabilistic automata [12] or quantitative variants of regular expressions [3]), we suspect that deciding equivalence in the presence of `dup` may be intractable. As circumstantial evidence, ProbNetKAT policies can simulate a probabilistic variant of multitape automaton originally introduced by Rabin and Scott [39]. We specialize the definition here to two tapes, for simplicity, but ProbNetKAT programs can encode any multitape automata with any fixed number of tapes.

Definition C.1. Let A be a finite alphabet. A *probabilistic multitape automaton* is defined by a tuple (S, s_0, ρ, τ) where S is a finite set of states; $s_0 \in S$ is the initial state; $\rho : S \rightarrow (A \cup \{_ \})^2$ maps each state to a pair of letters (u, v) , where either u or v may be a special blank character $_$; and the transition function $\tau : S \rightarrow \mathcal{D}(S)$ gives the probability of transitioning from one state to another.

The semantics of an automaton can be defined as a probability measure on the space $A^\infty \times A^\infty$, where A^∞ is the set of finite and (countably) infinite words over the alphabet A . Roughly, these measures are fully determined by the probabilities of producing any two finite prefixes of words $(w, w') \in A^* \times A^*$.

Presenting the formal semantics would require more concepts from measure theory and take us far afield, but the basic idea is simple to describe. An infinite trace of a probabilistic multitape automaton over states s_0, s_1, s_2, \dots gives a sequence of pairs

of (possibly blank) letters:

$$\rho(s_0), \rho(s_1), \rho(s_2) \dots$$

By concatenating these pairs together and dropping all blank characters, a trace induces two (finite or infinite) words over the alphabet A . For example, the sequence,

$$(a_0, _), (a_1, _), (_, a_2), \dots$$

gives the words $a_0a_1\dots$ and $a_2\dots$. Since the traces are generated by the probabilistic transition function τ , each automaton gives rise to a probability measure over pairs of infinite words.

Probabilistic multitape automata can be encoded as ProbNetKAT policies with `dup`. We sketch the idea here, deferring further details to [Appendix D](#). Suppose we are given an automaton (S, s_0, ρ, τ) . We build a ProbNetKAT policy over packets with two fields, **st** and **id**. The first field **st** ranges over the states S and the alphabet A , while the second field **id** is either 1 or 2; we suppose the input set has exactly two packets labeled with **id** = 1 and **id** = 2. In a set of packet history, the two active packets have the same value for **st** $\in S$ —this represents the current state in the automaton. Past packets in the history have **st** $\in A$, representing the words produced so far; the first and second components of the output are tracked by the histories with **id** = 1 and **id** = 2. We can encode the transition function τ as a probabilistic choice in ProbNetKAT, updating the current state **st** of all packets, and recording non-blank letters produced by ρ in the two components by applying `dup` on packets with the corresponding value of **id**.

Intuitively, a set of packet histories generated by the resulting ProbNetKAT term describes a pair of words generated by the original automaton. With a bit more bookkeeping (see [Appendix D](#)), we can show that two probabilistic multitape automata are equivalent if and only if their encoded ProbNetKAT policies are equivalent. Thus, deciding equivalence for ProbNetKAT with `dup` is harder than deciding equivalence for probabilistic multitape automata; similar reductions have been considered before for showing undecidability of related problems about KAT [29] and probabilistic NetKAT [24].

Deciding equivalence between probabilistic multitape automata is a challenging open problem. In the special case where only one word is generated (say, when the second component produced is always blank), these automata are equivalent to standard automata with ε -transitions (e.g., see Mohri [37]). In this setting, non-productive steps can be eliminated and the automata can be modeled as finite state Markov chains, where equivalence is decidable. In our setting, however, steps producing blank letters in one component may produce non-blank letters in the other. As a result, it is not clear how to eliminate these steps and encode our automata as Markov chains. Removing probabilities, it is known that equivalence between non-deterministic multitape automata is undecidable [22]. Deciding equivalence of deterministic multitape automata remained a challenging open question for many years, until Harju and Karhumäki [23] surprisingly settled the question positively; Worrell [50] later gave an alternative proof. If equivalence of probabilistic multitape automata is undecidable, then equivalence is undecidable for ProbNetKAT programs as well. However if equivalence turns out to be decidable, the proof technique may shed light on how to decide equivalence for the full ProbNetKAT language.

D Encoding 2-Generative Automata in Full ProbNetKAT

To keep notation light, we describe our encoding in the special case where the alphabet $A = \{x, y\}$, there are four states $S = \{s_1, s_2, s_3, s_4\}$, the initial state is s_1 , and the output function ρ is

$$\rho(s_1) = (x, _) \quad \rho(s_2) = (y, _) \quad \rho(s_3) = (_, x) \quad \rho(s_4) = (_, y).$$

Encoding general automata is not much more complicated. Let $\tau : S \rightarrow \mathcal{D}(S)$ be a given transition function; we write $p_{i,j}$ for $\tau(s_i)(s_j)$. We will build a ProbNetKAT policy simulating this automaton. Packets have two fields, **st** and **id**, where **st** ranges over $S \cup A \cup \{\bullet\}$ and **id** ranges over $\{1, 2\}$. Define:

$$p \triangleq \mathbf{st}=s_1 ; \mathbf{loop}^* ; \mathbf{st}\leftarrow\bullet$$

The initialization keeps packets that start in the initial state, while the final command marks histories that have exited the loop by setting **st** to be the special letter \bullet .

The main program **loop** first branches on the current state **st**:

$$\mathbf{loop} \triangleq \text{case} \begin{cases} \mathbf{st}=s_1 : \mathbf{state1} \\ \mathbf{st}=s_2 : \mathbf{state2} \\ \mathbf{st}=s_3 : \mathbf{state3} \\ \mathbf{st}=s_4 : \mathbf{state4} \end{cases}$$

Then, the policy simulates the behavior from each state. For instance:

$$\mathbf{state1} \triangleq \bigoplus \begin{cases} (\mathbf{if\ id=1\ then\ st}\leftarrow x; \mathbf{dup\ else\ skip}); \mathbf{st}\leftarrow s_1 @ p_{1,1}, \\ (\mathbf{if\ id=1\ then\ st}\leftarrow y; \mathbf{dup\ else\ skip}); \mathbf{st}\leftarrow s_2 @ p_{1,2}, \\ (\mathbf{if\ id=2\ then\ st}\leftarrow x; \mathbf{dup\ else\ skip}); \mathbf{st}\leftarrow s_3 @ p_{1,3}, \\ (\mathbf{if\ id=2\ then\ st}\leftarrow y; \mathbf{dup\ else\ skip}); \mathbf{st}\leftarrow s_4 @ p_{1,4} \end{cases}$$

The policies **state2**, **state3**, **state4** are defined similarly.

Now, suppose we are given two probabilistic multitape automata W, W' that differ only in their transition functions. For simplicity, we will further assume that both systems have strictly positive probability of generating a letter in either component in finitely many steps from any state. Suppose they generate distributions μ, μ' respectively over pairs of infinite words $A^\omega \times A^\omega$. Now, consider the encoded ProbNetKAT policies p, p' . We argue that $\llbracket p \rrbracket = \llbracket p' \rrbracket$ if and only if $\mu = \mu'$.³

First, it can be shown that $\llbracket p \rrbracket = \llbracket p' \rrbracket$ if and only if $\llbracket p \rrbracket(e) = \llbracket p' \rrbracket(e)$, where

$$e \triangleq \{\pi\pi \mid \pi \in \text{Pk}\}.$$

Let $v = \llbracket p \rrbracket(e)$ and $v' = \llbracket p' \rrbracket(e)$. The key connection between the automata and the encoded policies is the following equality:

$$\mu(S_{u,v}) = v(T_{u,v}) \quad (4)$$

for every pair of finite prefixes $u, v \in A^*$. In the automata distribution on the left, $S_{u,v} \subseteq A^\omega \times A^\omega$ consists of all pairs of infinite strings where u is a prefix of the first component and v is a prefix of the second component. In the ProbNetKAT distribution on the right, we first encode u and v as packet histories. For $i \in \{1, 2\}$ representing the component and $w \in A^*$ a finite word, define the history

$$h_i(w) \in \text{H} \triangleq (\mathbf{st} = \bullet, \mathbf{id} = i), (\mathbf{st} = w[|w|], \mathbf{id} = i), \dots, (\mathbf{st} = w[1], \mathbf{id} = i), (\mathbf{st} = s_1, \mathbf{id} = i).$$

The letters of the word w are encoded in reverse order because by convention, the head/newest packet is written towards the left-most end of a packet history, while the oldest packet is written towards the right-most end. For instance, the final letter $w[|w|]$ is the most recent (*i.e.*, the latest) letter produced by the policy. Then, $T_{u,v}$ is the set of all history sets including $h_1(u)$ and $h_2(v)$:

$$T_{u,v} \triangleq \{a \in 2^{\text{H}} \mid h_1(u) \in a, h_2(v) \in a\}.$$

Now $\llbracket p \rrbracket = \llbracket p' \rrbracket$ implies $\mu = \mu'$, since (4) gives

$$\mu(S_{u,v}) = \mu'(S_{u,v}).$$

The reverse implication is a bit more delicate. Again by (4), we have

$$v(T_{u,v}) = v'(T_{u,v}).$$

We need to extend this equality to all cones, defined by packet histories h :

$$B_h \triangleq \{a \in 2^{\text{H}} \mid h \in a\}.$$

This follows by expressing B_h as boolean combinations of $T_{u,v}$, and observing that the encoded policy produces only sets of encoded histories, *i.e.*, where the most recent state **st** is set to \bullet and the initial state **st** is set to s_1 .

E Background on Datacenter Topologies

Data center topologies typically organize the network fabric into several levels of switches.

FatTree. A FatTree [2] is perhaps the most common example of a multi-level, multi-rooted tree topology. Figure 6 shows a 3-level FatTree topology with 20 switches. The bottom level, *edge*, consists of top-of-rack (ToR) switches; each ToR switch connects all the hosts within a rack (not shown in the figure). These switches act as ingress and egress for intra-data center traffic. The other two levels, *aggregation* and *core*, redundantly connect the switches from the edge layer.

The redundant structure of a FatTree makes it possible to implement fault-tolerant routing schemes that detect and automatically route around failed links. For instance, consider routing from a source to a destination along shortest paths—*e.g.*, the green links in the figure depict one possible path from (s_7) to (s_1). On the way from the ToR to the core switch, there are multiple paths that could be used to carry the traffic. Hence, if one of the links goes down, the switches can route around the failure by simply choosing a different path. Equal-cost multi-path (ECMP) routing is widely used—it automatically chooses among the available paths while avoiding longer paths that might increase latency.

³We will not present the semantics of ProbNetKAT programs with *dup* here; instead, the reader should consult earlier papers [13, 46] for the full development.

However, after reaching a core switch, there is a *unique* shortest path down to the destination. Hence, ECMP no longer provides any resilience if a switch fails in the aggregation layer (*cf.* the red cross in Figure 6). A more sophisticated scheme could take a longer (5-hop) detour going all the way to another edge switch, as shown by the red lines in the figure. Unfortunately, such detours can lead to increased latency and congestion.

AB FatTree. The long detours on the downward paths in FatTrees are dictated by the symmetric wiring of aggregation and core switches. AB FatTrees [34] alleviate this by using two types of subtrees, differing in their wiring to higher levels. Figure 11(a) shows how to rewire a FatTree to make it an AB FatTree. The two types of subtrees are as follows:

- i) Type A: switches depicted in blue and wired to core using dashed lines.
- ii) Type B: switches depicted in red and wired to core using solid lines.

Type A subtrees are wired in a way similar to FatTree, but Type B subtrees differ in their connections to core switches. In our diagrams, each aggregation switch in a Type A subtree is wired to adjacent core switches, while each aggregation switch in a Type B subtree is wired to core switches in a staggered manner. (See the original paper by Liu et al. [34] for the general construction.)

This slight change in wiring enables much shorter detours around failures in the downward direction. Consider again routing from source ($s7$) to destination ($s1$). As before, we have multiple options going upwards when following shortest paths (*e.g.*, the one depicted in green), as well as a unique downward path. But unlike FatTree, if the aggregation switch on the downward path fails, there is a short detour, as shown in blue. This path exists because the core switch, which needs to re-route traffic, is connected to aggregation switches of both types of subtrees. More generally, aggregation switches of the same type as the failed switch provide a 5-hop detour; but aggregation switches of the opposite type provide an efficient 3-hop detour.