

Flock: Accurate network fault localization at scale

Vipul Harsh, Tong Meng, Kapil Agrawal, P. Brighten Godfrey
University of Illinois at Urbana-Champaign

Abstract

Inferring the root cause of failures among thousands of components in a data center network is challenging, especially for “gray” failures that are not reported directly by switches. Faults can be localized through end-to-end measurements, but past localization schemes are either too slow for large-scale networks or sacrifice accuracy. We describe Flock, a network fault localization algorithm and system that achieves both high accuracy and speed at datacenter scale. Flock uses a probabilistic graphical model (PGM) to achieve high accuracy, coupled with new techniques to dramatically accelerate inference in discrete-valued Bayesian PGMs. Large-scale simulations and experiments in a hardware testbed show Flock speeds up inference by $>10^4$ x compared to past PGM methods, and improves accuracy over the best previous datacenter fault localization approaches, reducing inference error by $1.19 - 11\times$ on the same input telemetry, and by $1.2 - 55\times$ after incorporating passive telemetry. We also prove Flock’s inference is optimal in restricted settings.

1 Introduction

Datacenters often comprise of tens of thousands of network components. Failures in such large networks are common, arising due to software bugs, misconfiguration, and faulty hardware, among other reasons [33]. In many cases, a device will directly report a failure, e.g., a switch may report that one of its line cards is non-responsive or that an interface has a certain packet loss rate. Network operators utilize monitoring software to collect these metrics and raise alerts. However, datacenters also experience significant network downtime and SLO violations from *gray failures* whose root cause is obscure [33, 63]. For example, the reason for poor performance of a distributed service could be a link silently dropping a small fraction of packets without updating switch counters [54], or a driver bug in a virtualized firewall. Diagnosing such performance anomalies is very hard for network operators [11]. With programmable switches, more advanced monitoring is possible [34, 62], but these methods either do not eliminate gray failures [45] or come at a high cost in switch resources [62], and require deployment of programmable switches which is not generally available.

An alternate approach, which is the domain of this paper, is to infer the root cause via end-to-end measurements, which we refer to as *fault localization* [4, 11, 27, 29, 32, 54]. Fault localization has been deployed by large cloud providers [12, 29]. At the heart of fault localization is an *inference model and algorithm* that uses end-to-end observations (e.g., packet

loss rate or latency of TCP connections) to infer a set of faulty components (links or switches). The key challenge is to do this both accurately and quickly.

The most powerful class of inference techniques builds a probabilistic graphical model (PGM) and performs a form of maximum likelihood estimation (MLE): finding the set of components that, if they failed, maximizes the probability of having produced the given end-to-end observations. An early such system was Sherlock [12], whose primary context was inferring faults among dependent services. However, deriving the MLE for a PGM can be computationally intensive. With $\leq k$ failures among n components, the solution space is exponentially large in k ($O(n^k)$). In a datacenter with millions of TCP flows, switches, links, etc. and multiple simultaneous failures, Sherlock’s MLE can require several hours.

Therefore, fault localization schemes intended for datacenter networks move away from PGMs to other techniques – using scores to rank links [11, 32] or solving for drop rates via a system of equations [54]. As we will show, these compromise accuracy and flexibility of PGMs in favor of performance.

In this paper, we present Flock, a fault localization system for datacenter networks that seeks to maximize accuracy, with sufficiently high speed (i.e., seconds). Flock’s core innovation is a novel MLE inference algorithm for PGMs that offers substantial acceleration for the kind of models encountered in fault localization, leveraging two key acceleration approaches: (i) a technique we call *joint likelihood exploration* maintains an array of hypotheses that it can update en masse to find the likelihood of a set of *new but similar* hypotheses, more quickly than computing their likelihoods individually from scratch and (ii) we use a greedy algorithm which builds its solution link by link; this part of the algorithm is simple, of course, but importantly, we prove a sufficient condition for optimality and verify through experiments that it does find the MLE in practice. These two optimizations each individually provide asymptotic speedups, and together allow Flock to use a PGM at scale. The result is that Flock is several orders of magnitude faster than past PGM-based fault localization [12], and is substantially more accurate than past non-PGM-based fault localization [11, 54], on the same input data.

Moreover, the PGM-based approach allows Flock to use different types of input telemetry. Recent datacenter fault localization schemes [11, 54] use observations of *active probes* of the network (which can be constructed to have known paths and uniform distribution) but do not incorporate *passive flow monitoring*, i.e., observations of all ongoing traffic, obtained

via NetFlow, IPFIX, or INT.¹ The large volume of passive data makes it potentially informative. But including passive monitoring would be hard for non-PGM methods because it requires more discerning modeling and inference to handle skew in traffic patterns and path uncertainty.² Although PGMs are generally more flexible in incorporating data of different types, it would be hard for past PGM approaches because of computational difficulty, due to the immense number of flows and because a flow with 10 possible paths is roughly 10× costlier to model than a flow with a known path. Flock’s combination of flexible PGM-based modeling and speed enables it to utilize passive information.

In summary, this paper’s key contributions are as follows:

- **Inference algorithm.** We develop Flock, a new fast PGM MLE inference algorithm for the type of PGM necessary for fault localization, namely discrete-valued Bayesian PGMs (§ 3.3). We analyze a sufficient condition for this algorithm’s accuracy on Flock’s model, providing intuition for why Flock works well in practice (§ 4.2).
- **System implementation.** We implement the Flock system (§ 3), a new end-to-end fault localization system. The Flock algorithm forms the heart of the Flock system, allowing it to employ a PGM and naturally incorporate various kinds of dependence and uncertainty such as unknown paths.
- **Evaluation suite.** We create an open evaluation suite [1] for fault localization, which includes (a) implementations of algorithms from NetBouncer [54], 007 [11], Sherlock [12], and Flock, (b) an implementation of end-host telemetry agents and a collector, (c) telemetry data for six different fault scenarios from a simulated data center and a hardware testbed, and (d) scaling tests. We believe this suite is of independent interest, as it is the first such open data set and expands on the fault scenarios evaluated by past work.
- **Performance evaluation.** For a Clos network with 88K links and 9.5M flows, Flock is empirically $> 10^4\times$ faster than Sherlock’s PGM-based method [12], scanning $\sim 3.5\text{M}$ hypotheses in 17 sec, while achieving the same or better inference results (§ 7.8). In fact, Flock is $\approx 4.5\times$ faster than the non-PGM approach of NetBouncer [54] on the same input. 007 [11] is the fastest of the lot, but its time savings (< 1 sec) is not a good tradeoff with accuracy.
- **Accuracy evaluation.** With the same (active probe) input measurements as past work, Flock reduced inference error by 1.8 – 8× compared to 007 and by 1.19 – 11× compared to NetBouncer.
- **The value of passive information/INT.** Incorporating passive monitoring reduced error even further, by up to 5.3×

¹[11] incorporates only a limited amount of passive monitoring; §2.3.

²Most data centers use non-deterministic ECMP multipath routing, so that only a *set of possible paths* is known when flows are monitored with NetFlow/IPFIX. Paths can be known with INT-based monitoring, but INT is not generally deployed.

compared to Flock with only active monitoring. We also evaluate the value of INT telemetry input.

2 Background and Motivation

2.1 When fault localization is useful

Ideally, network faults are directly reported by the faulty component; e.g., switches typically track interface utilization, packet drops, up/down status, queue length, etc. These metrics are commonly collected from network switches via SNMP [15], polling [17, 52], or streaming telemetry [10, 23].

Fault localization becomes useful when such direct monitoring is not enough. Problems that are not directly reported by the faulty components are known as *gray failures* [33]. For example, silent inter-switch or inter-card drops [11, 50, 54, 62, 63] are extremely challenging to detect, constituting 50% of faults that took > 3 hours to diagnose in [62]. Other examples of gray failures include corruption in TCAM-based forwarding tables causing black holes [40, 63] or loss [18], and a misconfigured switch causing high latency [63] (see [62] for more cases). Gray failures also occur in the numerous software packet processing components present in modern data centers. Software bugs can silently drop or corrupt packets in host virtualization [55, 56], server software [63], and virtualized network functions like software firewalls [47]. All the above faults can be “silent” (the device does not realize the error occurred). Further, the symptoms could manifest at a different location than the problem, e.g. packet data could be corrupted by an intermediate switch but the corruption is discovered only at the receiver. End-to-end observations are well suited to detect such problems (§ 6.4).

Another alternative is to use programmable switches to obtain more information, as in Omnimon [34], FANcY [43] (for ISPs), or NetSeer [62] which runs a packet sequencing protocol across neighboring hops to find silent drops. This can be quite accurate, though it comes at the cost of significant switch resources ($\sim 100\%$ overall PHV usage and 40% ALU usage [62]). But more importantly, although use of programmable switches is growing, they still have very limited deployment (13% estimated market share for 2023 [31]). Both programmable and traditional switching environments are valuable use cases, but *the latter is the target of this paper*; schemes like [34, 43, 62] are out of the scope of our work. As an exception, we will consider the use of data collected with In-band Network Telemetry (INT) [13, 45], which does not directly report gray failures, but does record packets’ paths. INT can be implemented with programmable switches, but similar path data can be obtained other ways; see § 6.2.

Thus, it is very hard to guarantee that every packet processing element detects all faults locally (indeed, the end-to-end principle [51] applies here). Also, even if a fault is reported, the operator may want a way to cross-check. For these reasons, we see fault localization based on end-to-end observations as

an important tool in infrastructure engineers’ toolbox for the foreseeable future.

2.2 Problem setup and goals

The input to a network fault localization algorithm is the network topology, and input telemetry (flow measurements). Each flow measurement includes one or more metrics (TCP loss rate, mean latency, throughput, etc.) and a set of paths through the topology that the flow may have traversed. Depending on the monitoring method, this set may have size one (the exact path is known) or greater than one (typically, a set of possible ECMP paths is known). Given this input, the fault localization algorithm should output a set of links or devices it believes to be faulty, while meeting two goals.

Performance: Network operators often have to resolve a reported problem quickly. For example, a managed service from BT has a 15 minute response time in its SLA [14] and Gartner chose a threshold of 3 minutes in defining “real time” network data analysis [9]. Thus, fault localization within minutes is critical and within a few seconds is ideal.

Accuracy: False positives can bury true problems among several alerts [49]. False negatives could send engineers down the wrong track of investigation. There may be a tradeoff between accuracy and performance. As long as results are available within a few minutes, accuracy (minimizing false positives and false negatives) is of primary importance.

2.3 Existing fault localization approaches

Several past approaches address above goals, with different types of *input telemetry* and different *inference algorithms*.

Input telemetry: Administrators commonly [9] use passive monitoring of flows via NetFlow [20] or IPFIX [21] to understand overall network health. However, this has not been relied on for automated fault localization because vendor-specific ECMP hashing obscures flows’ exact paths.

Recent approaches use active probes with known paths. NetBouncer [54] sends probes uniformly from hosts to core switches in a Clos topology, via special switch support. 007 [11] uses active probes with assistance from passive monitoring: end-host agents monitor production traffic, and when they detect flows with anomalous performance, the agents traceroute the flagged flows’ paths and report the flagged flows’ metrics for analysis. 007 does not incorporate passive monitoring of non-flagged flows. In both systems, the volume of active probes is limited (which assists inference performance, and minimizes host/network overhead, and the exact path of monitored flows is known (which assists accuracy). But active probes do not eliminate all uncertainty: even if the path is known, the culprit link(s) are not; and flows can experience packet loss or latency on non-faulty links (e.g., due to congestion). Hence, good inference is still needed. Further, active probes often take a different data path than

regular traffic and may fail to reproduce problems faced by production flows [11].

INT has, to our knowledge, not been specifically used by past work for end-to-end fault localization. INT is similar to passive NetFlow telemetry in that it can observe a large volume of actual application traffic, if deployed for all flows. However, like active probes, it can trace the traversed path. It also does not eliminate all uncertainty (e.g., a silent packet corruption will not trigger an INT action, and even if the path is known, the faulty link is not).

We observe that *different deployments are likely to have different available information*. 007 requires host support and NetBouncer and INT require switch support. INT is now available in some switches, but is not deployed in most networks, and might be deployed selectively. Furthermore, this technology landscape may evolve. A scheme that is flexible in its input telemetry is thus preferred.

Inference algorithms: Sherlock [12], NetSonar [58] and Shrink [36] use a PGM with a form of maximum likelihood estimation (MLE), but are far too slow. Sherlock targeted a small use case (358 components), and NetSonar, which uses Sherlock’s inference, targeted smaller inter-DC networks. In a data center network with tens of thousands of components, they require several hours (§ 7.8), even with $K \leq 2$ concurrent failures (they scan $O(n^K)$ possibilities, where n = number of components). Furthermore, inevitable concurrent failures [11, 50, 54] may make $K > 3$ essential.

Hence, state-of-the-art data center fault localization schemes move away from PGM-based MLE. 007 [11] uses a scoring function to rank links while NetBouncer [54] optimizes an objective function to solve for drop rates. These are reasonably fast, but as we will see, they fall short on accuracy.

Summary: Our goal is to localize “gray” failures in datacenter networks, using end-to-end information, achieving as high accuracy as possible within roughly a few seconds (including monitoring and inference), making best use of available information, including active probes or INT (where paths are known) and passive monitoring (exact paths are unknown).

3 Flock Design

We present Flock in three components. First (§ 3.1), Flock monitors end-to-end flows at the endpoints (e.g. hosts) of the network. The monitored flow observations, comprised of metrics from active probes and (when available) passive flow monitoring and/or INT, are then sent to a central collector. Next, the collector periodically constructs a probabilistic graphical model (PGM) from the flow observations (§ 3.2). This PGM captures uncertainty in how the observations may have resulted from underlying network components. Finally, Flock periodically performs MLE inference (§ 3.3) on the PGM model, searching through the exponentially-large hypothesis space (i.e., sets of possible faulty components) to find

the hypothesis that maximizes the probability of the observed flow metrics in the model.

The first two components (monitoring and the PGM model) are not the significant contributions of this paper. Our PGM is similar to that of Sherlock and NetSonar [58], with useful adaptations to fit our environment. We describe those pieces for completeness. Our key contribution lies in the MLE inference algorithm, comprised of multiple optimizations to make it scale.

3.1 Flow monitoring

Monitoring of end-to-end flows may occur in an agent process running on hosts, in the hypervisor of a virtualized data center, or potentially at edge/top-of-rack switches (e.g., with INT). Our inference algorithms are agnostic to where exactly these measurements come from. Still, for concreteness, we describe the operation of an endpoint monitoring agent (§5).

The agent periodically actively probes the network and may optionally passively observe performance of ongoing flows. Note that 007 also observes ongoing traffic, to decide what active probes to launch. Metrics from both active and passive monitoring are aggregated by flow, and optionally randomly sampled to reduce volume. Periodically, the agent sends these reports to the collector. For the rest of this section, we assume flow reports include RTT, source, destination, retransmissions, packets sent and the path if known via active probing/INT, else the set of possible paths.

3.2 Inference graph model

After telemetry is collected, Flock builds a probabilistic model of the network using two inputs: (1) telemetry reports as described above, potentially including active probes, INT, or passive telemetry, and (2) the network topology and routes. The latter could be provided by an SDN controller or a topology discovery tool and is used to determine a *set* of paths (typically, ECMP paths) that each flow may have traversed in the specific case of passive data with unknown paths.

Flock employs a probability model based on a Bayesian network that utilizes flow level metrics. A Bayesian network is a probabilistic graphical model that defines the probability distribution of a set of observed variables in terms of a set of unknown (or *hidden*) variables. It consists of a labeled directed acyclic graph (DAG) where each node represents a random variable, unknown or observed, whose distribution can be specified completely as a function of its immediate ancestors in the DAG. The goal of the inference is to estimate the unknown variables (failed/not failed status of each component) given the observed variables (retransmissions, RTT, packets sent etc. associated with each flow).

We use a 3-layer graph model. Fig. 1 shows an example network with two flows and its corresponding model. Formally, the conversion from datacenter network to model is as follows.

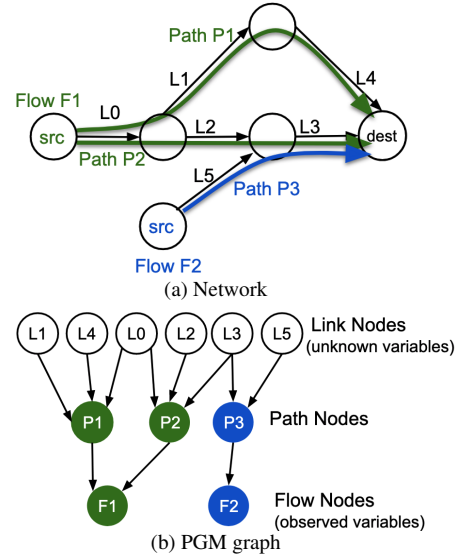


Figure 1: (a) A network with flows F1 and F2 traversing several possible paths; (b) the corresponding PGM model.

The **top layer** nodes in the graphical model represent individual links in the datacenter, referred to as *link-nodes*. Each link-node is a hidden binary variable which is 0 if the link has failed and 1 otherwise. The **intermediate layer** consists of nodes corresponding to paths in the datacenter, referred to as *path-nodes*. For each link ℓ in path p , there is an edge from the link-node of ℓ to the path-node of p in the Bayesian graph. A path-node represents an intermediate binary variable which is 0 if the path consists of a faulty link and 1 otherwise. Finally, the **bottom layer** consists of *flow-nodes* – one for every flow – which are the observed variables. A flow-node has an edge from each path-node in the path-set of that flow. We define the value of the flow-node variable as the number of *bad* packets – packets which experienced a problem. We describe two ways of setting the bad packets variable:

- Per packet analysis: For packet loss and corruption, we set the number of *bad* packets as the number of retransmissions which serves as a proxy for lost packets.
- Per flow analysis: To capture symptoms of high latency, we use a “per-flow” analysis which is in effect a special case of the per-packet model where the number of packets sent is 1, and the number of bad packets is set to 1 if the flow’s RTT is above a threshold and 0 otherwise.

Flock’s model assumes a flow F is routed via ECMP; F takes one of w paths chosen uniformly at random, and packets experience problems independently and uniformly at random. Thus, the probability of observing r bad packets out of the t sent can be given as:

$$P[F = (r, t)] = \frac{1}{w} \sum_{i=1}^w (1 - \gamma_i) p_b^r (1 - p_b)^{t-r} + \gamma_i p_g^r (1 - p_g)^{t-r} \quad (1)$$

where γ_i is the value of the i^{th} possible path of F ($\gamma_i = 0$ if a failed link is on the i^{th} possible path and 1 otherwise). p_g and p_b are model hyperparameters: p_b represents the probability of a packet experiencing problems when taking a bad path (i.e., with at least one faulty link), and p_g represents the probability of a packet experiencing problems on a good path (no faulty links). Intuitively, a packet going through a failed link is much more likely to experience a problem, hence $p_b \gg p_g$. § 5.2 describes how to pick p_g and p_b . Equation 1 can also be adapted to include path weights, like in WCMP [61].

A *hypothesis* is an assignment $H \in \{0, 1\}^n$ for all link-nodes. Equivalently, we can think of H as a *set of links* that are deemed to be failed, with all other link-nodes being not failed. The goal of the inference is to recover the hypothesis that consists of all truly failed links and only those links. Conditioned on a hypothesis H , the probability of the set of flow observations taking on the observed assignment of values (a certain number of bad packets r_i out of t_i packets sent, for each flow i) is simply the product of probabilities of all individual flow probabilities:

$$P[F_1, F_2, \dots, F_m | H] = \prod_{F_i \in \text{flows}} P[F_i = (r_i, t_i) | H] = \prod_{F_i \in \text{flows}} P[F_i | H]$$

where F_i is shorthand for the event that flow i takes on the observed metric values, i.e., $F_i = (r_i, t_i)$.

Incorporating Priors. We assign a prior belief about failures by assuming that, a priori, any link can fail with probability ρ . The priors reduce the false positive rate by effectively assigning a lower prior to hypotheses with more links, thus favouring hypotheses with fewer failed links. If hypothesis H contains $|H|$ candidate failed links and there are n total links, then the likelihood of H after incorporating priors is given as:

$$\text{Prior} * \prod_{F_i \in \text{flows}} P[F_i | H]; \text{Prior} = \rho^{|H|} (1 - \rho)^{n - |H|}$$

Model extensions. The top layer nodes can include other component types besides links; we add device nodes, treating a device as another component in a flow’s path, exactly as links. We found that a device prior that is 5× larger on log-scale, worked well in practice, as it forces Flock to detect a device failure only when there is stronger evidence for it than a link failure. Other components (line cards, racks, pods etc.) can be modeled in a similar way, but is beyond the scope of this paper.

Model Intuition. The model effectively incorporates several kinds of uncertainty. Given an observation of bad packets in a flow, we may not know what path is responsible (modeled via flow nodes having multiple path parents). Even if we do know the path, as we do for active probes and INT, we don’t

know what link is responsible³ (this is modeled via path nodes having multiple link parents). Even if we know what link is responsible, an observed bad packet might or might not mean there is a faulty link (this is modeled via both good and bad links having non-zero probability of bad packets).

Differences from Sherlock’s PGM. Sherlock was intended to model application-level failures, and thus includes elements that we don’t need such as services and load balancers. Sherlock uses three node states – working, failed, and partially working; we omit the last, as Flock models some packet loss even for working links. Starting from Sherlock’s model and making these changes results in a PGM that is very close to Flock’s, except for the probability formulations.

Model Assumptions. Like any other model [12, 36, 58], ours has assumptions: fixed packet fail probabilities (as in [36]), classifying paths as failed on just the absence/presence of at least one failed link (as in [58]), and packets getting affected independently (as in [12, 36, 58]). We tried different assumptions – one with variable fail rates obtaining MLE using Nesterov’s Gradient Descent algorithm, but it was too slow; another variation that treats paths differently depending on their number of failed links, but its accuracy was worse.⁴ We present the model that gave the best results in experiments based on real environments that don’t adhere to any model assumptions. We also give theoretical (§ 4.2) evidence supporting Flock’s effectiveness with these assumptions. Finally, it’s important to keep in mind that the model does not need to match reality perfectly, *it only needs to be “close enough” that the most likely explanation in the model is the right one.* See Fig: 6 in appendix for an illustration of how the PGM-inference can localize more accurately than past schemes.

3.3 Inference algorithm

We describe the inference algorithm next. For ease of exposition, this description only refers to links as components that can fail. Devices are treated exactly in the same way as links and our implementation handles both links and devices.

Recall that a hypothesis H is a candidate set of failed links. From the model, we can compute the probability of the flow variables taking their observed values (number of bad/sent packets for each flow) given H ; this is the likelihood of H . We denote this likelihood as $L(H) = P[F_1 | H] P[F_2 | H] \dots P[F_m | H]$. The maximum likelihood estimator \mathcal{H} is the hypothesis that maximizes $L(H)$, or equivalently log likelihood $LL(H) = \log L(H)$, that is $\mathcal{H} = \arg \max_{H \subseteq \text{links}} LL(H)$.

³While INT can capture per-hop metrics, it may not detect where a gray failure happens. For example, a silent corruption of packet data likely would not be detected until the packet reaches the receiver’s host stack.

⁴Note that Flock still localizes multiple failures on a path since there are flows that transit one failed link, but not the other.

We normalize all likelihoods by the likelihood of the no-failure hypothesis (i.e., $H_0 = \{\}$) to cancel out any flow whose path set does not include any failed links in a hypothesis H .

The goal of MLE inference is to compute \mathcal{H} . Simply computing the likelihood of each possible hypothesis would be impractically slow because there are 2^n hypotheses, where n is the number of links. Sherlock [12] and NetSonar [58] limit max concurrent failures to k , reducing the search space to $O(n^k)$ hypotheses. However, in our setting, this is still far too slow, even for $k = 2$ (§ 7.8). Further, a datacenter can have many concurrent failures making $k > 3$ important.

We introduce two algorithmic techniques to accelerate MLE inference in PGMs. *Greedy search* reduces the number of hypotheses examined. While Greedy MLE is a simple idea, our main contribution is showing that it finds good solutions even after narrowing the hypothesis space: we provide theory (§ 4.2) and experiments (§ 6). *Joint likelihood exploration* (JLE) is a new algorithmic technique, that reduces the time required per examined hypothesis. Both techniques provide significant speedups individually and together speedup the inference by several orders of magnitude (§ 7.8).

Greedy Search: We start from the no-failure hypothesis and extend it one link at a time. Specifically, we maintain a current hypothesis H . Initially, $H = \{\}$. In each iterative step, we scan over each link $l \notin H$ and calculate $LL(H \cup \{l\})$. If one of these log likelihoods improves over $LL(H)$, we set $H := H \cup \{l^*\}$ where l^* is the link offering the biggest improvement, and continue iterating. When no added link failure improves the log likelihood of the current hypothesis H , the search terminates and returns $H_{greedy} = H$.

There is still a performance challenge with Greedy MLE. Each iterative step requires evaluating close to n hypotheses (specifically $n - |H|$) to find l^* . Even with 40 cores, greedy search took over 3 hours for a medium-sized datacenter (§ 7.8). This motivates our second key optimization.

Joint likelihood exploration (JLE): We devise an additional acceleration technique for inference algorithms for PGMs with discrete variables. We use it to speed up each iteration of the greedy algorithm by a $O(n)$ factor, where n is the number of components (links and switches). Note that greedy+JLE produces the exact same solutions as greedy.

Suppose we are given the current best hypothesis H which has the maximum likelihood among hypotheses searched till now. Joint likelihood exploration is a technique to quickly explore all “neighbors” of H – all assignments that are different from H in the inclusion or exclusion of exactly one link. Note that there are n such neighbor hypotheses of H .

DEFINITION 1. Let $H \oplus l$ denote the hypothesis obtained by flipping the status of link l in H , i.e., if $l \in H$ then $H \oplus l = H \setminus \{l\}$ and otherwise $H \oplus l = H \cup \{l\}$. Let $\Delta_H(l) =$

$LL(H \oplus l) - LL(H)$ represent the difference in log likelihoods of hypotheses $(H \oplus l)$ and H .

JLE Intuition: We explain JLE by showing how it accelerates the Greedy algorithm (although it can also accelerate exhaustive search). In each iteration, Greedy computes each $LL(H \cup \{l\})$ for each l , to find the link l^* offering the most improvement. Note that maximizing $LL(H \cup \{l\})$ is equivalent to maximizing the difference in log likelihoods, $\Delta_H(l)$. So, in each iteration of Greedy, we could compute an n -element array Δ_H whose elements are the values $\Delta_H(l)$ for each link l , and then scan the array to find the largest value. This will find the same l^* as in the original greedy algorithm.

But how do we compute the array Δ_H ? If we do it in the obvious way, by iterating over each l and directly computing $LL(H \oplus l) - LL(H)$, then this is nearly identical to the original greedy algorithm which computed $LL(H \cup \{l\})$, with no appreciable change in runtime. What JLE offers is a *faster way to compute* Δ_H , with careful algorithm engineering. This involves two somewhat different types of computation: quickly preparing the array Δ_{H_0} for the first iteration; and quickly iteratively *updating* the array for each subsequent iteration.

To initially create Δ_{H_0} , note each entry $\Delta_{H_0}(l)$ represents the difference in log likelihood due to failing link l , compared to no failures. To compute these differences, we only have to look at the effect on flows whose paths intersect l . Furthermore, there is an important opportunity for memoization in computing the log likelihood-difference formula across different links l : the effect on a flow’s likelihood depends only on the number of failed paths, not the specific failed links.

Now suppose we have an existing Δ_H and the search algorithm is about to move from hypothesis H to hypothesis H' in its next iteration. We need to compute the array $\Delta_{H'}$. To do this, we track the *difference in the difference arrays* (Δ_H vs. $\Delta_{H'}$) rather than directly computing the difference in likelihoods of H and $H \oplus l$ for every l . The key insight here is that each entry of the difference array Δ_H can be written as a sum of contributions for all flows and only some of the terms need to be updated after moving to a new hypothesis H' . The fact that we can do this faster than creating the array from scratch is the key to JLE’s acceleration.

JLE formalization: First the algorithm needs to compute the array Δ_{H_0} for the first iteration of Greedy. We omit these details due to space; see function ComputeInitialDelta in Algorithm 2 in appendix which follows the intuition above.

Next, we describe how the Δ_H array can be updated when the greedy algorithm moves to a new hypothesis H' . We first note that $LL(H)$ is a sum of contributions from all flows and can be written as $LL(H) = \sum_{F \in \text{flows}} LL_F(H)$, where $LL_F(H) =$

$\log P[F|H]$. We have

$$\begin{aligned} \Delta_H(l) &= LL(H \oplus l) - LL(H) \\ &= \sum_{F \in \text{flows}} LL_F(H \oplus l) - \sum_{F \in \text{flows}} LL_F(H) = \sum_{F \in \text{flows}} \Delta_H(l, F) \end{aligned}$$

where $\Delta_H(l, F) = LL_F(H \oplus l) - LL_F(H)$. Next we derive a useful property about the individual flow contributions $\Delta_H(l, F)$.

DEFINITION 2. *A flow F intersects with link l if at least one of the possible paths for F has link l .*

THEOREM 1. *For a link l' and hypothesis H , let $H' = H \oplus l'$. Then for all links l and flows F ,*

- (i) *If F does not intersect with l , then $\Delta_H(l, F) = 0$*
- (ii) *If F does not intersect with l' , then $\Delta_{H'}(l, F) = \Delta_H(l, F)$.*

PROOF. This can be easily seen by expanding $\Delta_{H'}(l, F)$ and $\Delta_H(l, F)$. We note that for any H , the log likelihood of a flow F , given by $LL_F(H) = \log P[F|H]$, does not depend on a link l if F does not intersect with l (that is, $LL_F(H) = LL_F(H \oplus l)$).

For (i), if link l does not intersect with F , then flipping l 's status does not affect F : $LL_F(H) = LL_F(H \oplus l) \Rightarrow \Delta_H(l, F) = 0$.

For (ii), when l' does not intersect with flow F , $LL_F(H') = LL_F(H \oplus l') = LL_F(H)$ and $LL_F(H' \oplus l) = LL_F(H \oplus l' \oplus l) = LL_F(H \oplus l)$. Consequently, we get $\Delta_{H'}(l, F) = \Delta_H(l, F)$. \square

Hence, to obtain $\Delta_{H'}(l)$ from $\Delta_H(l)$, we only need to update the terms $\Delta_H(l, F)$ for flows F that intersect with *both* links l' and l since all other flow contributions to $\Delta_{H'}(l)$ remain unchanged from $\Delta_H(l)$. Let $\text{flows}(l', l)$ denote the set of flows that intersect with both l and l' . After updating the current hypothesis from H to H' , we can compute the new entry $\Delta_{H'}(l)$ for link l using Theorem 1:

$$\begin{aligned} \Delta_{H'}(l) &= \sum_{F \in \text{flows}} \Delta_{H'}(l, F) = \sum_{F \in \text{flows}} \Delta_H(l, F) + \sum_{F \in \text{flows}(l', l)} \Delta_{H'}(l, F) - \Delta_H(l, F) \\ &\Rightarrow \boxed{\Delta_{H'}(l) = \Delta_H(l) + \sum_{F \in \text{flows}(l', l)} \Delta_{H'}(l, F) - \Delta_H(l, F)} \quad (2) \end{aligned}$$

Once we have equation ??, the algorithm to update the Δ array, for all n entries, is simple to state. After moving to $H' = H \oplus l'$, we iterate over all flows that intersect with l' . For each such flow F , let L_F be the set of links that intersect with F . For each $l \in L_F$, we update F 's contribution to $\Delta_{H'}(l)$. With memoization, one can update all entries $\Delta_{H'}(l, F)$ for all $l \in L_F$ in a couple of passes over L_F , similar to how we initially computed Δ_{H_0} . The crux of the greedy+JLE algorithm is outlined in Algorithm 1 in appendix and the full pseudocode in Algorithm 2 is outlined in appendix .

Given $LL(H)$, an alternate approach is to compute $LL(H \oplus l)$ without JLE, individually for each l , as in [12, 58]. This requires updating the contribution of all flows that intersect with l since their likelihoods $LL_F(H \oplus l)$ would have changed after

flipping the status of link l . Thus, the number of flows whose contributions need to be updated for computing just one entry $LL(H \oplus l)$ for a single l is the same as that for computing all n entries of the Δ array jointly with JLE. Thus, JLE results in in a $O(n)$ speedup. The reason for this large improvement is that JLE tracks the change in the Δ 's (i.e., the difference in the differences: $LL_F(H \oplus l) - LL_F(H)$) across iterations which allows reuse of computation from the previous iteration.

Besides greedy search, JLE can apply to any algorithm which explores a hypothesis H 's neighbors: $H \oplus l$ for all l . This includes brute force, Sherlock and NetSonar's inference, and MCMC techniques (e.g. Gibbs sampling). Using JLE, we were able to accelerate (i) Sherlock's inference (Alg. 3 in appendix), and (ii) Gibbs sampling for Flock, both by multiple orders of magnitude. We ended up using Greedy for Flock because (i) Sherlock's inference can not detect $K > 2$ concurrent failures and was still slow with JLE (§ 7.8) and (ii) for Gibbs sampling, it's hard to bound the number of iterations required for convergence. Gibbs sampling [46] without JLE was too slow for our purposes.

4 Flock: Analysis

4.1 Runtime analysis

Let n be the number of links, m be the number of flows, T be an upper bound on the number of links that any flow intersects with, D be an upper bound on the number of flows that any link intersects with and K be the maximum number of concurrent failures (note Flock's inference does not know K). The runtime of Greedy inference with JLE is $O(n+mT+(K-1)DT)$. If we had used just Greedy without JLE (computing likelihood of each hypothesis individually), the runtime would be $O(n+mT+(K-1)nDT)$. In contrast, Sherlock's runtime is $O(n^K DT)$. JLE can improve Sherlock's runtime by a factor of n , to $O(n^{K-1}DT)$. From the analysis above (and our experiments later), it can be seen that Greedy + JLE is dramatically faster. See § C of appendix for derivations of these results.

4.2 Accuracy analysis

We now analyze conditions in which Greedy returns the true MLE hypothesis. To make the problem tractable, we restrict the analysis to cases where path taken is known (true for active probes and INT) and packets crossing a link get dropped independently according to a (unknown) drop probability of that link (inference is NP hard if packets get dropped adversarially, see § 3 in appendix). Theorem 2 gives a sufficient condition on the traffic pattern for Flock's inference to correctly recover the set of failed links, providing intuition for why Flock's model and inference work well in practice.

DEFINITION 3. *For given traffic T , let $T(\{l_1, l_2, \dots, l_k\})$ denote the number of packets that each go through all of the links $\{l_1, l_2, \dots, l_k\}$. If $T(\{l_1, l_2\})/T(\{l_1\}) \leq \epsilon$ for all links l_1 and l_2 , then we say T is ϵ -skewed.*

THEOREM 2. *For any topology with $(1/\alpha)$ -skewed traffic, with high probability, Flock’s inference returns the set of all failed links if the number of failures is $\leq \alpha/2$, the number of packets T_{min} crossing every link is larger than a certain threshold, and the drop probabilities are $< p_g$ on all good links and $> p_b$ on all failed links where (p_g, p_b) satisfy the condition $5p_g < p_b < 0.05$.*

Proof in appendix. Theorem 2 has an intuitive interpretation: Consider two links l_1 and l_2 , where l_1 has failed and l_2 works correctly. Intuitively, at most $\frac{1}{\alpha}$ fraction of the dropped packets on l_1 will transit l_2 . One can think of this as l_2 getting “ $\frac{1}{\alpha}$ fraction of the blame” for the dropped packets on l_1 . If there are $(f\alpha)$ failures and these are arranged adversarially so that all of them add blame to l_2 , then in the worst case l_2 can get $(f\alpha) \cdot \frac{1}{\alpha} = f$ times as much blame as a true failed link. Intuitively, a large enough constant f would confuse the algorithm into classifying l_2 as failed. The theorem proves that for $f < 0.5$, the algorithm outputs the true failed links.

5 Implementation

5.1 Agent and inference engine

We implement an agent that runs on end-hosts and collects flow statistics via a lightweight packet dumping tool based on the PF_RING module [5]. The agent periodically encapsulates the collected flow statistics (52 bytes per flow) into export IPFIX messages, and sends it to the collector. As the specifics of the agent/collector don’t affect our results, we leave more optimized designs for future work. Note that commercial solutions also exist [3, 6], possibly employing alternate approaches such as pulling flow statistics directly from the kernel via eBPF or using multiple collectors at scale.

Flock’s inference engine, written in C++, (i) collects IPFIX flow reports from agents and (ii) periodically runs inference on the collected input. Currently, the network topology is static, but the inference engine can be modified to obtain the topology from a controller. The engine periodically reads flow reports from the queue, every 30 seconds, and runs the inference algorithm of § 3.3.

5.2 Parameter Calibration

A important problem we encountered, with both Flock and the other competing schemes, is how to set their hyperparameters. Flock has 3 hyperparameters (p_g, p_b, ρ) , NetBouncer has 3, and 007 has 1. In real deployed systems, parameters and thresholds are quite common, and are set based on past deployment experience. However, manual tuning puts extra onus on the user and makes it difficult to evaluate systems. The manually set parameters of past schemes 007 [11] and NetBouncer [54] gave suboptimal results in our environments, across different topology scales and failure scenarios.

Thus, we design an automated parameter calibration method that we use for all schemes. We use a training set

of monitoring data to search for the parameter settings that obtain the best precision and recall in the training set. It is tempting to obtain the training set from historical monitoring data, which is generally available in deployed systems. However, this can be tricky, since: (a) historical data may not have faults labeled; and (b) faults, especially rare faults, may have diverse and unpredictable types so that historical data is not entirely representative of the next upcoming incident.

To solve (a), we leverage simulations to obtain the training set. We calibrate parameters once using this training set and use those parameters across our experiments, unless stated otherwise. Note that if this method were used in a deployment, it would increase concerns with (b) since simulations may not match the real world. To address (b) we will experimentally quantify the robustness of each scheme to scenarios where the train and test sets are drawn from *different environments* (different topology, monitoring duration, fault rate, fault type).

Once we have the training set, we use the following calibration method. For each hyperparameter, we choose equally-spaced values in a reasonable range of possible values. We fix a minimum precision P and find the parameters which, in a training set, yielded highest recall and had precision $> P$. Varying P produces a set of parameters that are Pareto-optimal along the precision/recall tradeoff curve. Our evaluation will apply these parameters to a separate test data set.

To choose a single parameter setting (rather than a tradeoff curve), we set $P = 98\%$ and find the setting that maximizes recall (in the training set); if no such point exists or recall is too low ($< 25\%$), then we subtract 5% from P and try again, repeating until a setting is found. This method lays more emphasis on precision, which is usually desirable.

6 Evaluation Methodology

6.1 Systems evaluated

We compare Flock with several state-of-the-art datacenter fault localization schemes. Our codebase consists of 7K LOC including C++ implementations of all inference algorithms.

We implemented the “Ferret” inference algorithm of **Sherlock** (Sec. 3.2 of [12]). For a fair comparison, we run Ferret on the same PGM as Flock (which is anyway similar to Sherlock’s; § 3.2). Sherlock can not detect $K > 2$ failures, but (as expected) resulted in the same accuracy as Flock for $K \leq 2$ failures at small scale. Hence, we only show performance differences.

We implemented **NetBouncer’s** algorithm (Figure 5 in [54]) and **007** (Algorithm 1 in [11]) We verified our 007 implementation by matching scores outputted by the publicly available 007 code. We were also able to reproduce Figure 10 in [11] with our implementation/setup.

For all schemes, we calibrate parameters once using simulations of random packet drops and use those parameters by

default, unless stated otherwise. In some cases, we also show results with parameters calibrated on that environment.

6.2 Input telemetry types

We use four different kinds of input for inference:

- **A1:** Active probes between end-hosts and the core switches with known paths, as designed for NetBouncer [54].
- **A2:** Reports about flows with ≥ 1 retransmission, along with their (actively-probed) paths, as designed for 007 [11].
- **P:** Passive information consisting of reports about regular (application) data flows, whose traffic matrix is thus dictated by the network environment (§ 6.3). A set of possible paths is known (based on ECMP multipath).
- **INT:** We assume INT [2] provides reports, including paths, for both A1 and P (which thus becomes a superset of A2).

Note the last case is intended to test full deployment of INT; in other deployment modes it could trace just a subset of traffic. Similar reports could be obtained from other recent packet marking [50] or mirroring [30, 34, 35, 53, 63] methods. Since Flock can incorporate different kinds of inputs, we compare accuracy across input types: Flock (A1) vs NetBouncer (A1), Flock (INT) vs NetBouncer (INT) and Flock (A2) vs 007 (A2). We also quantify the accuracy boost Flock obtains from additional passive information (A1+P, A2+P, A1+A2+P). NetBouncer and 007 cannot trivially ingest the passive telemetry as they do not model path uncertainty. Finally, the passive flow telemetry can be downsampled in a large datacenter with high link speeds to reduce volume of the monitoring data.

6.3 Network environments

NS3 simulations. We set up a NS3 simulation to output a trace consisting of flow metrics (retransmissions/packets sent). We feed this trace as input to inference. We use a standard 3-tiered Clos topology [7] with 2500 40Gbps links, ECMP routing and 3x oversubscription at ToRs. Like [54], we set drop rates on all non-failed links between 0 – 0.01% chosen independently and uniformly at random to model occasional drops on good links⁵. For all our experiments, half the traces used uniform random traffic and the other half used a skewed traffic pattern where 50% of the traffic is concentrated among 5% of the racks, randomly chosen. Flow sizes were drawn from a Pareto distribution (mean: 200KB, scale:1.05) to mimic irregular flow sizes in a typical datacenter [8].

Large scale simulation. NS3 was too slow for large scale simulations. Hence, we use a flow level simulator (similar to [11]), that drops each packet as per preset drop probabilities on links but does not model queuing or TCP. We use this simulator for scaling experiments (§ 7.8).

Hardware test cluster. We set up a physical testbed with 10 switches and 48 emulated hosts, each with its own dedicated hardware NIC port and one CPU core. One of the 48 hosts⁵TCP can tolerate such low rates, hence it’s reasonable to qualify such links as non-faulty. These rates are not central to Flock

runs Flock’s collector. We use a standard 2-tier Clos topology with 2 spines, 8 leaf racks and 6 hosts per rack. We provision 1 Gbps link speeds to emulate as many hosts as possible. Schemes with A1 are omitted from our testbed results since our switches don’t have the in network IP-in-IP feature for A1 [54].

6.4 Failure scenarios

In simulation:

- **Silent link packet drops:** a link drops a small fraction of packets without updating switch counters. Silent drops are a common problem in the industry [11, 50, 54].
- **Silent device failure:** An error in a device component (e.g., memory, line card) causes silent packet drops. This differs from the prior scenario because it affects many or all links on the device.

In hardware test cluster:

- **Queue misconfiguration:** A WRED queue drops packets with probability p when the queue length is above a configurable threshold w . We misconfigure WRED queues [24] on switches, setting $p = 1\%$ (available choices: 1-100%) and $w = 0$ (so, the link works normally if the queue is empty).
- **Link flap:** We pull out a cable manually and quickly put it back in to emulate link flaps [19]. In our setup, link flaps caused the latency of the flows transiting the link to spike, but did not produce any significant increase in retransmissions (i.e., the link was buffering packets).

Evaluation metrics: We use **Precision** (fraction of predicted failed links/devices that actually failed) and **Recall** (fraction of failed links/devices that were correctly predicted as failed), to quantify false positives and false negatives respectively. We use the standard **Fscore** measure (harmonic mean of precision and recall) when we need a combined measure of accuracy.

7 Evaluation Results

The first goal of our evaluation is to investigate Flock’s accuracy compared to NetBouncer and 007. We compare various input types (INT, A1, A2, P) to quantify benefits obtained from incorporating passive data and INT. As expected, Flock and Sherlock had the same accuracy in small scale experiments (with max $K = 2$ failures), where Sherlock finished in reasonable time. Hence we don’t show accuracy comparisons with Sherlock.

Next, we investigate performance of Flock, including inference algorithm compared to Sherlock, NetBouncer and 007 and the runtime benefits of using JLE to speed up inference.

7.1 Silent packet drops

We generated 63 traces via NS3, each with 1 to 8 failed links with drop rate on each failed link chosen uniformly at random between 0.1% and 1% [54] (drops on good links are set as in § 6.3). For each trace, we send active flows

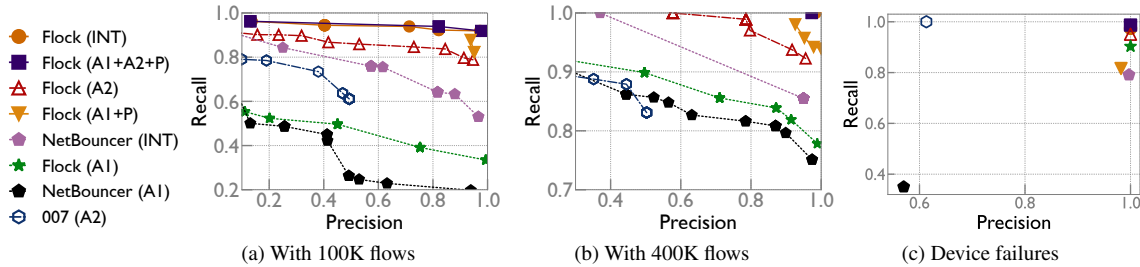


Figure 2: Accuracy for silent packet drops. (a), (b): Tradeoff curves for NetBouncer, 007 and Flock for silent drops, varying hyperparameters for each scheme. Schemes are annotated with the input information they use. (c): Accuracy on device failures.

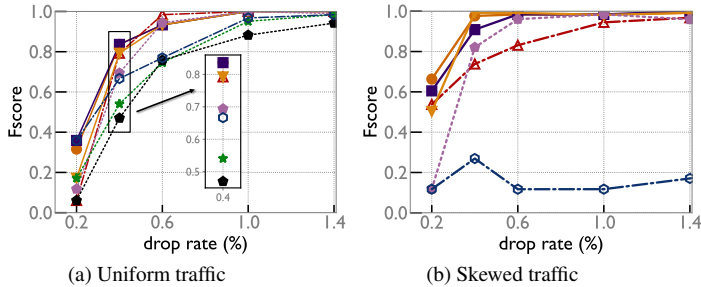


Figure 3: The extent of drop rates that each scheme can detect.

between the hosts and the core switches (A1), sending 40 packets per second and 400K passive flows per second across all hosts. Fig. 2 shows precision-recall tradeoff curves, with different calibrated parameters (§5.2) after 100K and 400K flows. Using the chosen point for each scheme (§ 5.2), we draw the following conclusions for each input type:

- **A1:** Flock reduces error rate over NetBouncer by roughly 45% (fscore: 0.5 vs 0.27). With 4x more probes, Flock still reduces the error rate by >20% compared to NetBouncer (fscore: 0.87 vs 0.84).
- **A2:** Flock (fscore: 0.93) reduces error-rate over 007 (fscore: 0.61) by 5.5x after 400K flows.
- **A1+P and A1+A2+P:** When active probes (A1, A2) are augmented with passive information (P), Flock achieves very high accuracy (fscore with A1+A2+P: 0.98, A1+P: 0.93) after 400K flows (see Fig. 2a), suggesting that these schemes require less data than active-only schemes for localization.
- **INT:** Flock (INT) achieves the best accuracy (fscore 0.99) reducing error over NetBouncer (INT) (fscore 0.88) by 12x. The last two points highlight the benefits of incorporating passive data for accuracy. 007’s performance can be attributed to it being sensitive to traffic skew (§ 7.3).

7.2 Device failures

Using the same setup as § 7.1, we simulate a device failure by failing $f\%$ of a faulty device’s links. We generate 64 traces, each consisting of up to 2 device failures, varying f across traces from 25% to 100%. A subset of links failing on a device is similar to the behaviour of a faulty line card on that device. For all schemes, we used the same parameters as in § 7.1 (we

calibrated NetBouncer’s threshold for the number of problematic flows crossing a device). As shown in Fig. 2c, Flock outperforms NetBouncer and 007 for all types of information. Flock (INT) achieves $\approx 100\%$ recall, compared to 80% recall of NetBouncer (INT). Flock (A2) reduces error-rate compared to 007 by 8x (fscore 0.97 vs 0.76). Flock (A1+P) has poorer precision for device failures than link failures.

7.3 Soft gray failures

We vary the drop rate on a single failed link to test what drop rates Flock can detect. A useful metric is the ratio of drop-rate on a failed link and the maximum drop-rate on a functioning link, which we call Signal to Noise Ratio (SNR), by a slight abuse of terminology. We use the same setup and parameters as § 7.2 (except 007 which had to be calibrated separately for skewed traffic since otherwise it had poor recall). We used 32 traces for each data point. From Figs. 3a and 3b, we conclude that Flock can detect links with > 1% drop rate (or SNR > 100) with high recall, with A2. With uniform traffic, 007’s accuracy is good when SNR > 100 (consistent with the SNR in Fig. 10 of [11]). 007’s recall gets affected significantly with skewed traffic. After adding passive telemetry with either INT or (A1+A2+P), Flock’s accuracy gets boosted and it is able to detect > 0.4% drop rate reliably. NetBouncer’s accuracy is slightly worse than Flock for A1, but its accuracy becomes even worse for multiple concurrent failures with different drop rates (§7.1). Schemes utilizing A1 (active probes) are unaffected by skew in the application traffic and hence omitted from Fig. 3b.

7.4 Misconfigured queue

We move now from simulated faults to our testbed, beginning with misconfigured queues. Flock achieves high accuracy with all information types (see Fig. 4a). Using the same parameters as in § 7.1 for all schemes, Flock (INT) had higher precision and recall than NetBouncer (INT) (16x less error in fscore), whereas Flock (A2) had higher precision than 007 (A2) (the solid markers in Fig. 4a). For comparison, we also show precision recall tradeoff curves with parameters calibrated on the testbed with real examples. In this case, Flock (INT) had 7x less error than NetBouncer (INT) (fscore: 0.87 vs. 0.98) and Flock (A2) had 16x lower error compared to 007

(fscore: 0.97 vs. 0.5) (the hollow markers in Fig. 4a). Flock (A2+P) gets very close to Flock (INT), consistent with §7.1.

7.5 Link flaps

We use a per-flow analysis (§ 3.2), classifying a flow as problematic if its RTT is > 10 ms. Since the analysis is per-flow and not per-packet, we had to recalibrate parameters. Flock does not model acks traversing the reverse path, which is important in this case. Accounting for that with a revised model would be possible, we leave that for future work. Even with a somewhat inaccurate model, Flock (INT) reduces the error rate by 1.66x over NetBouncer (INT) (fscore: 0.81 vs 0.69) and Flock (A2) reduces the error rate over 007 by 1.8x (see Fig. 4b).

7.6 Irregular Clos topologies

Real world datacenters are rarely perfectly symmetric like a Clos topology and typically have asymmetries due to failures, policies, piecemeal upgrades, etc. To see the effect of topology irregularity, we omit links from the fat tree. We recalibrated parameters on the irregular topologies for all schemes since the topology can be known in advance. (we also tested without the recalibration; Flock’s improvement over other schemes was even higher in this case. We omit the results for brevity.)

Fig. 5a, 5b show that Flock’s accuracy is robust to topology irregularity. 007 is sensitive to topology irregularity, probably because its effect is similar as having traffic skew. We omit A1 since its active probing mechanism is designed only for regular Clos topologies.

Flock with passive only input (P): Some networks may only have passive information available. Past fault localization schemes can not be applied to passive only input since they don’t handle path uncertainty. Operators in this situation resort to manual troubleshooting (traceroutes, adjustments to routing or taking links offline, etc.) which can take days. Flock with only passive input (P) can provide partial analysis (Fig. 5a, 5b). Interestingly, Flock (P)’s accuracy actually *improves* with more links removed. This is because in a symmetric Clos topology, there are equivalence classes of links (e.g. all links from a leaf switch to the spine layer) that cannot be differentiated because they participate in the same ECMP paths. As the topology becomes irregular, this breaks symmetry, and Flock’s inference algorithm automatically takes advantage of this.

Fig. 5c shows an even more difficult fully passive scenario where the failed link is one of several symmetric links in a Clos topology, the network has little irregularity for Flock to leverage ($< 5\%$ omitted links) and absence of active probes or path tracing. Flock (P) achieved $>75\%$ recall and $>40\%$ precision. We also show the theoretical maximum precision (calculated from the topology’s link equivalence classes). Note that 40% precision means Flock has narrowed down the faulty link

Parameters calibrated for → (D: different, S: same)	Different topology	Different failure rate		Different monitoring interval		Different failure scenario		Aggregate score (average Fscore)		
		p	r	p	r	p	r			
Flock (A1+A2+P)	D	0.92	0.98	0.97	1	0.98	1	0.95	0.99	0.973
	S	0.96	0.98	0.97	1	0.99	1	0.96	0.99	0.981
Flock (A2)	D	0.90	0.99	0.96	1	0.86	0.97	0.94	0.98	0.948
	S	0.94	0.98	1	1	0.94	0.92	0.94	0.98	0.962
Flock (INT)	D	0.92	0.99	0.96	1	0.98	1	0.95	0.99	0.973
	S	0.96	0.99	0.96	1	0.99	1	0.96	0.99	0.981
007 (A2)	D	0.75	1	0.87	1	0.51	0.82	1	0.33	0.728
	S	0.82	0.74	1	1	0.47	0.87	0.82	0.74	0.792
NetBouncer (INT)	D	0.25	0.33	0.95	1	1	0.66	0.28	0.33	0.589
	S	0.81	0.9	1	1	0.95	0.85	0.81	0.9	0.901

Table 1: Evaluating parameter robustness. For each scheme, we show accuracy when its parameters are calibrated on a different environment than the test data set (D) and when they are calibrated in the same environment (S).

to about 2-3 possibilities. We believe this will be an extremely helpful starting point for operators.

7.7 Parameter calibration robustness

As discussed in § 5.2, all the schemes we consider have parameters that must be set, and we have calibrated them based on simulations with known ground truth. What happens when these systems encounter unexpected situations?

To test robustness, we trained each scheme on one environment and tested in a different environment. (This is effectively a strong form of cross-validation where not only are the train and test sets different, they are drawn from *different distributions*.) Specifically, we created different types of differences between train and test, changing the (a) duration of monitoring, (b) topology, (c) failure rate (training set has failed links with significantly different drop rates), and (d) failure type. In particular, for (b), the schemes were calibrated in our simulator with random packet drops, and tested on misconfigured queues in a 20x smaller topology in our physical testbed. Table 1 shows the accuracy in these cases, both when the train and test sets are drawn from different distributions (“D”) and when they are drawn from the same distribution (“S”). The table’s aggregate score column shows Flock was fairly robust to parameter calibration in a different environment, with under 2% loss in accuracy. 007 was also robust (6% loss) while NetBouncer was more sensitive (31% loss).

We also tested Flock’s parameter sensitivity, i.e., how precision and recall vary with perturbations of its parameters. Accuracy remained high for many choices of parameters. See Fig. 8a in appendix .

7.8 Running time and scalability

Algorithm runtime: Flock’s main algorithmic innovation is its fast PGM inference compared to Sherlock’s PGM inference. Fig. 4c shows Flock’s inference is more than 4 orders of magnitude of faster than Sherlock, whose runtime on a large network was estimated to be 19 days, based on extrapolating a partial run. Recall Flock employs two optimizations: greedy

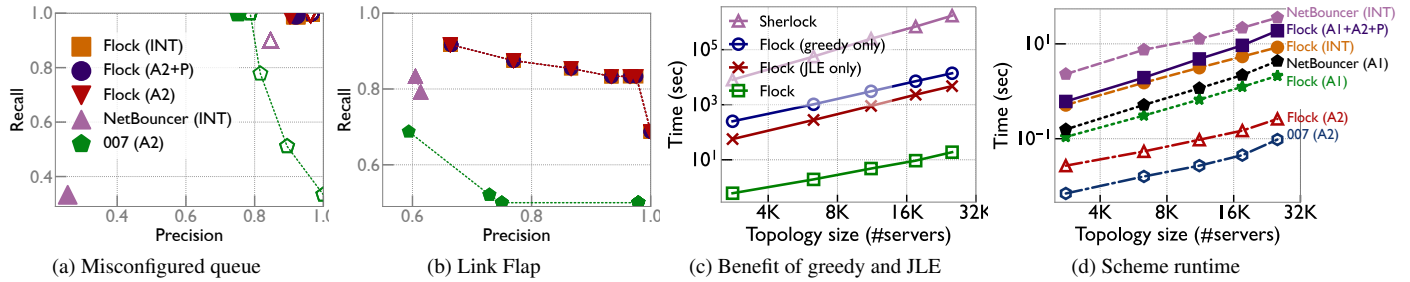


Figure 4: (a), (b) Accuracy on failure scenarios in testbed (solid markers). For (a), for comparison, we also show precision recall tradeoff curves with recalibrated parameters (hollow markers) (c) Running time vs. a past PGM scheme (Sherlock); Flock achieves the same accuracy while being $>10^4$ x faster. Also shown is the effect of Flock’s two optimizations (JLE, greedy) alone. (d) Running time of all schemes on various topology sizes.

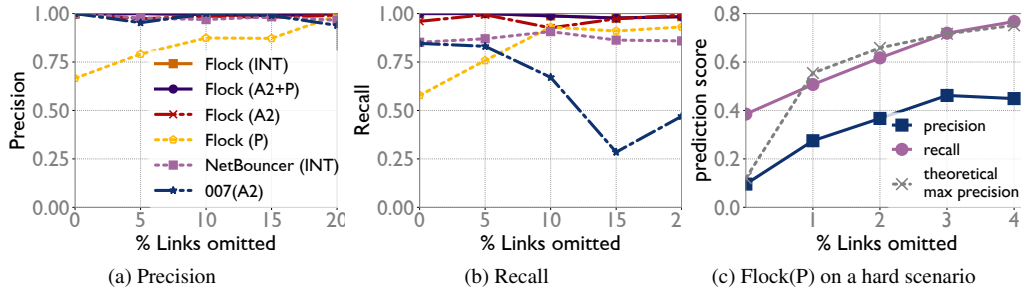


Figure 5: (a),(b) Accuracy on “irregular” Clos networks with a few links omitted. (c) Flock (P) produces useful results in a difficult scenario, where other schemes don’t apply.

and JLE. Fig. 4c shows each of these optimizations alone yields a ≈ 100 x improvement over Sherlock.

Fig. 4d compares Flock to the non-PGM schemes. Flock is faster than NetBouncer on the same input data. 007 is the fastest but its time savings (< 1 sec) does not trade-off well with accuracy.

Agent/Collector : Refer to Appendix A for the scalability of our agent/collector. As other commercial solutions also exist [3, 6], we leave more optimized agent/collector designs for future work.

8 Related work

Many other works have studied fault localization outside of datacenter networks – for troubleshooting reachability, black holes in IP networks [22, 28, 40, 58], virtual disk failures [59], performance problems in distributed services [12, 26, 44] and application performance anomalies [57]. Some of these can benefit from PGM-based inference, accelerated via JLE. We leave this for future work.

Detecting entire flow drops, for e.g., caused by a misconfigured ACL or a forwarding loop, is challenging for end-to-end schemes (as noted in [54]), partly due to (un)available input when paths are fully blocked. Other schemes such as NetSeer [62] and Omnimon [34] or network verification [25, 37, 38, 42] are more suitable for this class of faults.

Several works orchestrate active probes for inference [4, 16, 29, 32, 36, 39, 41, 54, 60]. Flock can handle active probes and outperforms one such approach (NetBouncer). Additionally, it can use passive data for accuracy gains. deTector [48],

MaxCoverage [40], Tomo [22] and Score [39] find a minimal set of components that explain most of the problems (e.g. packet drops). We expect them to run into similar problems as 007, since they don’t account for traffic skew. Simon [27] reconstructs queuing times from active probes. This allows it to diagnose high latency, but not silent packet drops. Packet mirroring [63] can catch packet drops that happen in the switch pipeline (e.g. congestion drops), but can not detect silent interswitch or silent intercard drops. Flock is well suited to detect such problems. Pingmesh [29] and NetNorad [4] use active pings, but do not provide complete localization. [50] identifies anomalies among symmetric links in a Clos network using statistical tests. It is sensitive to topology irregularity and requires path information.

9 Conclusion

Flock is a fault localization system for large datacenter networks based on end-to-end information. Flock’s key innovation is an optimized MLE inference algorithm which allows it to use a PGM at scale, achieving both high accuracy and speed, where past work achieved only one of the two.

Acknowledgements

We thank Radhika Mittal and members of the SysNet group at UIUC for providing feedback on an earlier version of the paper. We also thank CoNext 2023 reviewers for their reviews and feedback.

References

[1] Flock code. <https://github.com/netarch/FaultLocalization>.

- [2] In-band network telemetry (int) dataplane specification. https://github.com/p4lang/p4-applications/blob/master/docs/INT_latest.pdf.
- [3] Manage engine traffic analyzer. <https://www.manageengine.com/products/netflow/>.
- [4] Netnorad: Troubleshooting networks via end-to-end probing. <https://code.fb.com/networking-traffic/netnorad-troubleshooting-networks-via-end-to-end-probing/>.
- [5] Pf_ring by ntop software. github.com/ntop/PF_RING.
- [6] Solarwinds traffic analyzer. <https://www.solarwinds.com/netflow-traffic-analyzer>.
- [7] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. In *Proceedings of the ACM SIGCOMM 2008 Conference on Data Communication, SIGCOMM '08*, pages 63–74, New York, NY, USA, 2008. ACM.
- [8] M. Alizadeh, A. Kabbani, T. Edsall, B. Prabhakar, A. Vahdat, and M. Yasuda. Less is more: Trading a little bandwidth for ultra-low latency in the data center. In *Presented as part of the 9th USENIX Symposium on Networked Systems Design and Implementation (NSDI 12)*, pages 253–266, San Jose, CA, 2012. USENIX.
- [9] Andrew Lerner. Inclusion Criteria for the 2016 NPMD Magic Quadrant. <https://blogs.gartner.com/andrew-lerner/2015/06/29/gotnmpmd/>, 2015.
- [10] Arista. Arista Network Telemetry. <https://www.arista.com/en/solutions/software-defined-network-telemetry>, Accessed 2021-01-27.
- [11] B. Arzani, S. Ciraci, L. Chamon, Y. Zhu, H. H. Liu, J. Padhye, B. T. Loo, and G. Outhred. 007: Democratically finding the cause of packet drops. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 419–435, Renton, WA, 2018. USENIX Association.
- [12] P. Bahl, R. Chandra, A. Greenberg, S. Kandula, D. A. Maltz, and M. Zhang. Towards highly reliable enterprise network services via inference of multi-level dependencies. In *Proceedings of the 2007 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications, SIGCOMM '07*, pages 13–24, New York, NY, USA, 2007. ACM.
- [13] R. Ben Basat, S. Ramanathan, Y. Li, G. Antichi, M. Yu, and M. Mitzenmacher. PINT: probabilistic in-band network telemetry. In *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, pages 662–680, 2020.
- [14] British Telecommunications. Contract for BT Managed WAN Services. <https://business.bt.com/content/dam/terms/it-solutions-support/bt1190.pdf>, 2018.
- [15] J. Case, M. Fedor, M. Schoffstall, and J. Davin. A Simple Network Management Protocol (SNMP). In *RFC 1157*. Internet Engineering Task Force, 1990.
- [16] Y. Chen, D. Bindel, H. Song, and R. H. Katz. An algebraic approach to practical and scalable overlay network monitoring. In *Proceedings of the 2004 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications, SIGCOMM '04*, pages 55–66, New York, NY, USA, 2004. ACM.
- [17] Cisco. Monitoring and Troubleshooting With Cisco Prime LAN Management Solution 4.1. https://www.cisco.com/c/en/us/td/docs/net_mgmt/ciscoworks_lan_management_solution/4-1/user/guide/monitoring_troubleshooting/mnt_ug/SNMPInfo.html, 2018.
- [18] Cisco. Cisco Bug: CSCvn56156 - Silent packet drops may occur on FXOS platforms due to classifier table entry corruption. <https://quickview.cloudapps.cisco.com/quickview/bug/CSCvn56156>, 2020.
- [19] Cisco. Configure Link Flap Prevention on a Cisco Business Switch using CLI. <https://www.cisco.com/c/en/us/support/docs/smb/switches/Cisco-Business-Switching/kmgmt-2249-configure-the-link-flap-prevention-settings-on-a-switch-thro.html>, July 2020.
- [20] B. Claise. Cisco Systems NetFlow services export version 9. *RFC 3954 (Internet Standard)*, Internet Engineering Task Force, 2004.
- [21] B. Claise, B. Trammell, and P. Aitken. Specification of the ip flow information export (ipfix) protocol for the exchange of flow information. *RFC 7011 (Internet Standard)*, Internet Engineering Task Force, 2013.
- [22] A. Dhamdhere, R. Teixeira, C. Dovrolis, and C. Diot. Netdiagnoser: Troubleshooting network unreachabilities using end-to-end probes and routing data. In *Proceedings of the 2007 ACM CoNEXT Conference, CoNEXT '07*, New York, NY, USA, 2007. Association for Computing Machinery.
- [23] Divya Rao. Hot off the press: Introducing OpenConfig Telemetry on NX-OS with gNMI and Telegraf! https://www.cisco.com/c/en/us/td/docs/net_mgmt/ciscoworks_lan_management_solution/4-1/user/guide/monitoring_troubleshooting/mnt_ug/SNMPInfo.html, July 2020.
- [24] S. Floyd and V. Jacobson. Random early detection gateways for congestion avoidance. *IEEE/ACM Trans. Netw.*, 1(4):397–413, Aug. 1993.
- [25] A. Fogel, S. Fung, L. Pedrosa, M. Walraed-Sullivan, R. Govindan, R. Mahajan, and T. Millstein. A general approach to network configuration analysis. In *12th {USENIX} symposium on networked systems design and implementation ({NSDI} 15)*, pages 469–483, 2015.
- [26] Y. Gan, M. Liang, S. Dev, D. Lo, and C. Delimitrou. Sage: Practical and scalable ml-driven performance debugging in microservices. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2021*, page 135–151, New York, NY, USA, 2021. Association for Computing Machinery.
- [27] Y. Geng, S. Liu, Z. Yin, A. Naik, B. Prabhakar, M. Rosenblum, and A. Vahdat. {SIMON}: A simple and scalable method for sensing, inference and measurement in data center networks. In *16th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 19)*, pages 549–564, 2019.
- [28] D. Ghita, H. Nguyen, M. Kurant, K. Argyraki, and P. Thiran. Netscope: Practical network loss tomography. In *Proceedings of the 29th Conference on Information Communications, INFOCOM '10*, pages 1262–1270, Piscataway, NJ, USA, 2010. IEEE Press.
- [29] C. Guo, L. Yuan, D. Xiang, Y. Dang, R. Huang, D. Maltz, Z. Liu, V. Wang, B. Pang, H. Chen, Z.-W. Lin, and V. Kurien. Pingmesh: A large-scale system for data center network latency measurement and analysis. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication, SIGCOMM '15*, pages 139–152, New York, NY, USA, 2015. ACM.
- [30] A. Gupta, R. Harrison, M. Canini, N. Feamster, J. Rexford, and W. Willinger. Sonata: Query-driven streaming network telemetry. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication, SIGCOMM '18*, page 357–371, New York, NY, USA, 2018. Association for Computing Machinery.
- [31] M. Hamblen. Programmable chips for data center switches catch fire with 20% annual growth. <https://www.fierceelectronics.com/electronics/programmable-chips-for-data-center-switches-catch-fire-20-annual-growth>, August 2019.
- [32] H. Herodotou, B. Ding, S. Balakrishnan, G. Outhred, and P. Fitter. Scalable near real-time failure localization of data center networks. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pages 1689–1698, New York, NY, USA, 2014. ACM.
- [33] P. Huang, C. Guo, L. Zhou, J. R. Lorch, Y. Dang, M. Chintalapati, and R. Yao. Gray failure: The achilles' heel of cloud-scale systems. In *Proceedings of the 16th Workshop on Hot Topics in Operating Systems, HotOS '17*, pages 150–155, New York, NY, USA, 2017. ACM.
- [34] Q. Huang, H. Sun, P. P. C. Lee, W. Bai, F. Zhu, and Y. Bao. Omnimon: Re-architecting network telemetry with resource efficiency and full

- accuracy. SIGCOMM '20, page 404–421, New York, NY, USA, 2020. Association for Computing Machinery.
- [35] Q. Huang, H. Sun, P. P. C. Lee, W. Bai, F. Zhu, and Y. Bao. Omnimon: Re-architecting network telemetry with resource efficiency and full accuracy. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM '20, page 404–421, New York, NY, USA, 2020. Association for Computing Machinery.
- [36] S. Kandula, D. Katabi, and J.-P. Vasseur. Shrink: A tool for failure diagnosis in ip networks. In *Proceedings of the 2005 ACM SIGCOMM Workshop on Mining Network Data*, MineNet '05, pages 173–178, New York, NY, USA, 2005. ACM.
- [37] P. Kazemian, G. Varghese, and N. McKeown. Header space analysis: Static checking for networks. In *9th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 12)*, pages 113–126, 2012.
- [38] A. Khurshid, X. Zou, W. Zhou, M. Caesar, and P. B. Godfrey. Veriflow: Verifying network-wide invariants in real time. In *10th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 13)*, pages 15–27, 2013.
- [39] R. R. Kompella, J. Yates, A. Greenberg, and A. C. Snoeren. Ip fault localization via risk modeling. In *Proceedings of the 2Nd Conference on Symposium on Networked Systems Design & Implementation - Volume 2*, NSDI'05, pages 57–70, Berkeley, CA, USA, 2005. USENIX Association.
- [40] R. R. Kompella, J. Yates, A. Greenberg, and A. C. Snoeren. Detection and localization of network black holes. In *Proceedings of the IEEE INFOCOM 2007 - 26th IEEE International Conference on Computer Communications*, pages 2180–2188, Washington, DC, USA, 2007. IEEE Computer Society.
- [41] L. Ma, T. He, A. Swami, D. Towsley, K. K. Leung, and J. Lowe. Node failure localization via network tomography. In *Proceedings of the 2014 Conference on Internet Measurement Conference, IMC '14*, pages 195–208, New York, NY, USA, 2014. ACM.
- [42] H. Mai, A. Khurshid, R. Agarwal, M. Caesar, P. B. Godfrey, and S. T. King. Debugging the data plane with anteater. *ACM SIGCOMM Computer Communication Review*, 41(4):290–301, 2011.
- [43] E. C. Molero, S. Vissicchio, and L. Vanbever. Fast in-network gray failure detection for isps. In *Proceedings of the ACM SIGCOMM 2022 Conference*, SIGCOMM '22, page 677–692, New York, NY, USA, 2022. Association for Computing Machinery.
- [44] R. N. Mysore, R. Mahajan, A. Vahdat, and G. Varghese. Gestalt: Fast, unified fault localization for networked systems. In *2014 USENIX Annual Technical Conference (USENIX ATC 14)*, pages 255–267, Philadelphia, PA, 2014. USENIX Association.
- [45] P4.org Applications Working Group. In-band Network Telemetry (INT) Dataplane Specification Version 2.1. https://github.com/p4lang/p4-applications/blob/master/docs/INT_v2_1.pdf, 2020.
- [46] V. N. Padmanabhan, L. Qiu, and H. J. Wang. Passive network tomography using bayesian inference. In *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet Measurement, IMW '02*, page 93–94, New York, NY, USA, 2002. Association for Computing Machinery.
- [47] Palo Alto Networks. Critical Issues Addressed in PAN-OS Releases. <https://knowledgebase.paloaltonetworks.com/KCSAArticleDetail?id=kA10g000000Cm68CAC>, 2020.
- [48] Y. Peng, J. Yang, C. Wu, C. Guo, C. Hu, and Z. Li. detector: a topology-aware monitoring system for data center networks. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*, pages 55–68, Santa Clara, CA, 2017. USENIX Association.
- [49] A. Roy, D. Bansal, D. Brumley, H. K. Chandrappa, P. Sharma, R. Tewari, B. Arzani, and A. C. Snoeren. Cloud datacenter sdn monitoring: Experiences and challenges. In *Proceedings of the Internet Measurement Conference 2018, IMC '18*, page 464–470, New York, NY, USA, 2018. Association for Computing Machinery.
- [50] A. Roy, H. Zeng, J. Bagga, and A. C. Snoeren. Passive realtime datacenter fault detection and localization. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 595–612, Boston, MA, 2017. USENIX Association.
- [51] J. H. Saltzer, D. P. Reed, and D. D. Clark. End-to-end arguments in system design. *ACM Transactions on Computer Systems (TOCS)*, 2(4):277–288, 1984.
- [52] SolarWinds. Configure polling statistics intervals in the Orion Platform. https://documentation.solarwinds.com/en/Success_Center/orionplatform/content/core-polling-statistics-intervals-sw1829.htm, Accessed 2021-01-24.
- [53] P. Tammana, R. Agarwal, and M. Lee. Simplifying datacenter network debugging with pathdump. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 233–248, Savannah, GA, 2016. USENIX Association.
- [54] C. Tan, Z. Jin, C. Guo, T. Zhang, H. Wu, K. Deng, D. Bi, and D. Xiang. Netbouncer: Active device and link failure localization in data center networks. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 599–614, Boston, MA, 2019. USENIX Association.
- [55] VMware. Possible data corruption after a Windows 2012 virtual machine network transfer. <https://kb.vmware.com/s/article/2058692>, 2017.
- [56] VMware. Network timeouts or packet drops with VMware Tools 11.x with Guest Introspection Driver on ESXi 6.5/6.7. <https://kb.vmware.com/s/article/79185>, 2021.
- [57] M. Yu, A. Greenberg, D. Maltz, J. Rexford, L. Yuan, S. Kandula, and C. Kim. Profiling network performance for multi-tier data center applications. In *Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation, NSDI'11*, page 57–70, USA, 2011. USENIX Association.
- [58] H. Zeng, R. Mahajan, N. McKeown, G. Varghese, L. Yuan, and M. Zhang. Measuring and troubleshooting large operational multipath networks with gray box testing. Technical Report MSR-TR-2015-55, June 2015.
- [59] Q. Zhang, G. Yu, C. Guo, Y. Dang, N. Swanson, X. Yang, R. Yao, M. Chintalapati, A. Krishnamurthy, and T. Anderson. Deepview: Virtual disk failure diagnosis and pattern detection for azure. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 519–532, Renton, WA, Apr. 2018. USENIX Association.
- [60] Y. Zhao, Y. Chen, and D. Bindel. Towards unbiased end-to-end network diagnosis. In *Proceedings of the 2006 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, SIGCOMM '06, pages 219–230, New York, NY, USA, 2006. ACM.
- [61] J. Zhou, M. Tewari, M. Zhu, A. Kabbani, L. Poutievski, A. Singh, and A. Vahdat. Wcmp: Weighted cost multipathing for improved fairness in data centers. In *Proceedings of the Ninth European Conference on Computer Systems, EuroSys '14*, pages 5:1–5:14, New York, NY, USA, 2014. ACM.
- [62] Y. Zhou, C. Sun, H. H. Liu, R. Miao, S. Bai, B. Li, Z. Zheng, L. Zhu, Z. Shen, Y. Xi, P. Zhang, D. Cai, M. Zhang, and M. Xu. Flow event telemetry on programmable data plane. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM '20, page 76–89, New York, NY, USA, 2020. Association for Computing Machinery.
- [63] Y. Zhu, N. Kang, J. Cao, A. Greenberg, G. Lu, R. Mahajan, D. Maltz, L. Yuan, M. Zhang, B. Y. Zhao, and H. Zheng. Packet-level telemetry in large datacenter networks. In *Proceedings of the 2015 ACM Conference*

on Special Interest Group on Data Communication, SIGCOMM '15, pages 479–491, New York, NY, USA, 2015. ACM.

A Proofs

If the topology, flow statistics and path taken for each flow are chosen arbitrarily, then finding the MLE for such adversarial inputs, unsurprisingly, is NP-hard (§ appendix A).

DEFINITION 4. *Define adversarial inference as inference when the input is arbitrary, consisting of topology, flow metrics with arbitrarily chosen source and destinations, paths and arbitrarily chosen flow metrics (packets sent/dropped).*

However, we can make the following assumption, in the context of packet drops, to alleviate intractability for adversarial inputs: for each link, there is a (unknown) ground truth drop probability. Rather than adversarially, each packet crossing that link is dropped independently according to the drop probability of the link. First, we prove that adversarial inference is NP-hard.

THEOREM 3. *Adversarial inference is NP-hard.*

PROOF. We reduce the problem of finding a minimum vertex cover in a graph to the adversarial inference problem. For the proof, we will design an algorithm for min-vertex cover that makes polynomial number of accesses to an oracle O^{AI} for adversarial inference.

Let's say, we're given a graph $G = (V, E)$ where V is the set of vertices and E is the set of edges and our goal is to find a minimum vertex cover in this graph. We create topology \mathcal{T} for O^{AI} as follows-

- We create “vertex”-nodes n_v^1 and n_v^2 in \mathcal{T} for each vertex $v \in V$ and attach them with a directed link $l_v: (n_v^1 \rightarrow n_v^2)$. Note that all vertex-nodes in \mathcal{T} have exactly one link, either incoming or outgoing.
- For each edge $e \in E$, we create an “edge”-flow f_e , where f_e goes through links $(l_{v_1}$ and $l_{v_2})$, where v_1 and v_2 are the endpoints of e . In order for this to be a legitimate path, we connect the endpoints of $(l_{v_1}$ and $l_{v_2})$ to a special node.
- For each link l in \mathcal{T} , we create a “link”-flow f_l . going through just link l .

Finally we need to assign each flow, some number of packets sent and dropped. We first derive an expression for the return value of oracle O^{AI} . The output of O^{AI} is a binary assignment to links in \mathcal{T} where we interpret a value of 0 as that link being failed and 1 as the link being up. Let $E_{\mathcal{T}}$ be the number of links in \mathcal{T} and $\gamma_f^A = \prod_{l \in Path(f)} l$ denote the status of the path taken by flow f with $\gamma_f^A=0$ being path failed and 1 being path not failed. Recall that a path is deemed to be failed if it contains at least one failed link. For flow f , if n and r be the number of packets sent/dropped by the flow respectively, then the likelihood for flow f can be written as:

$$\begin{aligned} & \frac{P[f|H = \{l_1, l_2, \dots, l_{E_{\mathcal{T}}}\}]}{P[f|l_1 = 1, l_2 = 1, \dots, l_{E_{\mathcal{T}}} = 1]} \\ &= \frac{\gamma_f^A p_g^r (1 - p_g)^{n-r} + (1 - \gamma_f^A) p_b^r (1 - p_b)^{n-r}}{p_g^r (1 - p_g)^{n-r}} \\ &= \frac{p_b^r (1 - p_b)^{n-r}}{p_g^r (1 - p_g)^{n-r}} \left(1 + \gamma_f^A \left(\frac{p_g^r (1 - p_g)^{n-r}}{p_b^r (1 - p_b)^{n-r}} - 1 \right) \right) \\ &= \frac{1 + \alpha_f \gamma_f^A}{\alpha_f + 1} = \frac{(1 + \alpha_f \prod_{l \in Path(f)} l)}{\alpha_f + 1} \end{aligned}$$

where, $\alpha_f = \frac{p_g^r (1 - p_g)^{n-r}}{p_b^r (1 - p_b)^{n-r}} - 1 \in (-1, \infty)$ is a constant for flow f irrespective of H . The oracle O^{AI} then returns

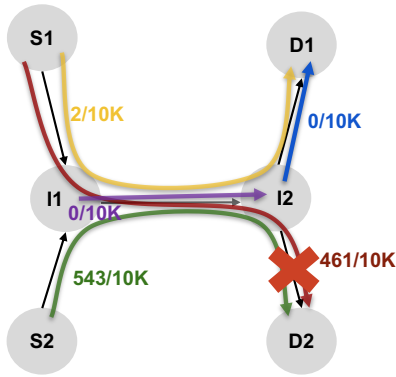
$$\arg \max_{H \in \{0,1\}^{E_{\mathcal{T}}}} \prod_{f: \text{flows}} P[f|H] = \arg \max_{H \in \{0,1\}^{E_{\mathcal{T}}}} \prod_{f: \text{flows}} (1 + \alpha_f \prod_{l \in Path(f)} l^H)$$

Where l^H is the status (0/1) of link l as per hypothesis H . Given the above expression, we set the number of packets sent/dropped for all flows in the following way

- (1) For all edge-flows f_e , we set n, r : the number of packets sent/dropped such that $1 + \alpha_{f_e} = \frac{1}{C}$ where $C \gg 1$.
- (2) For all link-flows f_l where l is connected to a special node, we set n, r : the number of packets sent/dropped in such a way that if l is connected to a special node, then $1 + \alpha_{f_l} = 1 + C$. This ensures that O^{AI} will always assign a label of 1 to a link connected to a special node (recall that 1 denotes the link being up).
- (3) For all link-flows f_l where both endpoints of l are connected to vertex nodes, we set n, r : the number of packets sent/dropped such that $1 + \alpha_{f_l} = 1 + \epsilon$ where ϵ is a small number > 0 . This assigns a small cost of assigning a label 0 to a link in \mathcal{T} whose both endpoints of l are connected to vertex nodes.

The vertex cover is obtained by simply picking vertices in G corresponding to links in \mathcal{T} that are deemed as failed by O^{AI} . Conditions (1) and (2) above ensure that only vertex-links will be classified as failed by O^{AI} . Conditions (2) and (3) ensure that the result set of vertices will cover all edges. Condition (3) ensures that the resultant vertex cover will be of the smallest size. \square

Agent/Collector scalability : Our agent's CPU usage for an end-host sending traffic grew linearly with the data rate (see Fig 7b) and was $<2\%$ of one core for a 1Gbps uplink and 10-15% of a core for a 10 Gbps uplink. As can be seen from figure 7c, the resource usage was independent of the number of flows. We verified that our multicore collector can handle 8K connections/sec from agents (see Fig 7). We tested this by launching several agent processes that generate dummy flow reports to send to the collector. These results show that the passive monitoring can be handled by an end-host agent and a centralized collector.

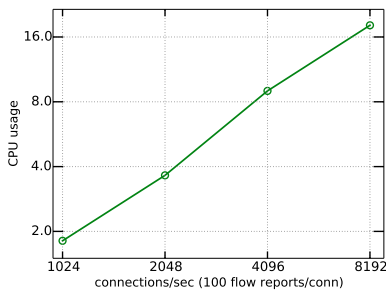


Scheme	Predicted failed links
007	(I1, I2)
NetBouncer	(S2, I1), (I2, D2)
Flock	(I2, D2)

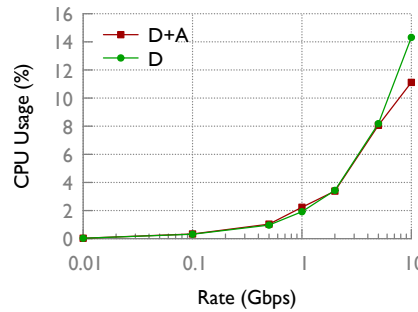
(a) Example: failed link shown via the red cross

(b) Output of various schemes

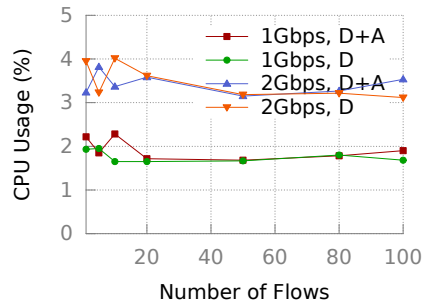
Figure 6: Example: (a) 5 links in the network. 5 flows shown in 5 different colors, annotated with packets dropped/packets sent. (b) Flock correctly localizes the failed link.



(a) Collector scaling

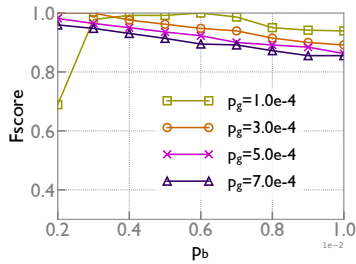


(b) agent CPU usage: single flow

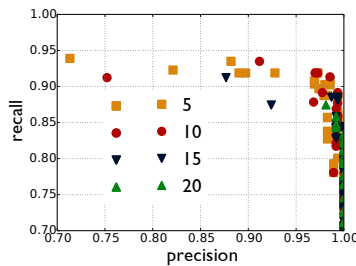


(c) several flows

Figure 7: (a) CPU usage at the collector, varying number of agents (b) CPU usage on the agent for handling a single flow. D: cost of packet header dumping, A: cost of compiling flow reports from packet headers. (c) agent CPU usage with many concurrent flows



(a) Parameter sensitivity



(b) Effect of changing priors

Figure 8: (a) Effect of changing p_g and p_b : intuitively, we expect precision to increase when p_g or p_b is increased, at the cost of reduced recall, which is confirmed by the figure (b) Higher priors result in points to the right. Priors ρ resulted in a significant reduction in false positives.

THEOREM 4. For any topology with $(1/\alpha)$ -skewed traffic, Flock’s inference returns the set of all failed links if the number of failures is $\leq \alpha/2$ with high probability, the number of packets T_{min} crossing every link is larger than a certain threshold, and the drop probabilities are $< p_g$ on all good links and $> p_b$ on all failed links where (p_g, p_b) satisfy the condition $5p_g < p_b < 0.05$.

The condition $T_{min} > T_0$ ensures that the effect of variance in the number of packets dropped by a link is small. Theorem 2 holds for a wide range of values of p_g and p_b (see lemma 1).

PROOF. If γ_f^H be the 0/1 status of the path taken by flow f as per hypothesis H (0 being that the path has at least one failed link as per H and hence is labeled as failed), we can express the (normalized) log likelihood for a flow given a hypothesis as follows

$$\begin{aligned}
& \log P[f|H = \{l_1, l_2, \dots, l_n\}] - \log P[f|l_1 = 1, l_2 = 1, \dots, l_n = 1] \\
&= \log \left(\gamma_f^H p_g^r (1 - p_g)^{n-r} + (1 - \gamma_f^H) p_b^r (1 - p_b)^{n-r} \right) \\
&\quad - \log \left(p_g^r (1 - p_g)^{n-r} \right) \\
&= (1 - \gamma_f^H) \left(r \log \frac{p_b}{p_g} + (n - r) \log \frac{(1 - p_g)}{(1 - p_b)} \right) \\
&= (1 - \gamma_f^H) \lambda (r - n\mu) = (1 - \gamma_f^H) \lambda \sum_{i=0}^n (b_i - \mu) \\
&= (1 - \gamma_f^H) X(f, H)
\end{aligned}$$

where $\lambda = \log \frac{p_b(1-p_g)}{p_g(1-p_b)}$ and $\mu = \log \frac{(1-p_g)}{(1-p_b)} / \log \frac{p_b(1-p_g)}{p_g(1-p_b)}$ are constants and b_i is a binary variable which is 1 if the i^{th} packet of the flow was dropped and 0 otherwise. One can check that $p_g < \mu < p_b$ for any $0 < p_g < p_b < 1$ by taking partial derivatives or using a plotting tool. Note that we can simply ignore the constant λ for the purpose of log likelihood maximization. If the maximum allowed drop rate on a correctly working link be p^* and the maximum path length for the given topology be k , then we set $p_g \geq kp^*$ so that for any flow f that does not go through any of the failed links, $(p_f - \mu) < 0$, where p_f is the packet drop probability of the path taken by f . We show the following lemma which holds for any reasonable settings for p_g and p_b for the purpose of detecting packet drops. The expected value of the $X(f, H)$ is given as

$$\begin{aligned}
E[X(f, H)] &= E \left[\lambda \sum_{i=0}^n (b_i - \mu) \right] \\
&= \lambda \sum_{i=0}^n E[(b_i - \mu)] = n\lambda(p_f - \mu)
\end{aligned}$$

LEMMA 1. If $\mu = \frac{\log \frac{(1-p_g)}{(1-p_b)}}{\log \frac{p_b(1-p_g)}{p_g(1-p_b)}}$ and $5p_g < p_b \leq 0.05$, then $0 \leq p_g < \mu < 2\mu < p_b$

Lemma 1 can be seen by taking partial derivatives or via a numerical plotting tool.

LEMMA 2. If p_l denotes the drop probability of link l , $L_f = \{l_1, l_2, \dots, l_k\}$ be the links in the path taken by flow f , p_f denotes the drop probability of the path L_f and H^* denote the set of failed links, then $(p_f - \mu) \leq \sum_{l \in H^* \cap L_f} p_l$

PROOF. For any $l_i \in L_f$, we have-

$$\begin{aligned}
p_f &= 1 - \left((1 - p_{l_1}) \dots (1 - p_{l_i}) \dots (1 - p_{l_k}) \right) \\
&= 1 - (1 - p_{l_i}) \left((1 - p_{l_1}) \dots (1 - p_{l_{i-1}}) (1 - p_{l_{i+1}}) \dots (1 - p_{l_k}) \right) \\
&= 1 - \left((1 - p_{l_1}) \dots (1 - p_{l_{i-1}}) (1 - p_{l_{i+1}}) \dots (1 - p_{l_k}) \right) \\
&\quad + p_{l_i} \left((1 - p_{l_1}) \dots (1 - p_{l_{i-1}}) (1 - p_{l_{i+1}}) \dots (1 - p_{l_k}) \right) \\
&\leq 1 - \left((1 - p_{l_1}) \dots (1 - p_{l_{i-1}}) (1 - p_{l_{i+1}}) \dots (1 - p_{l_k}) \right) + p_{l_i}
\end{aligned}$$

Applying the same argument recursively for all links in $H^* \cap L_f$, we get:

$$\begin{aligned}
p_f - \mu &= 1 - (1 - p_{l_1})(1 - p_{l_2}) \dots (1 - p_{l_k}) - \mu \\
&\leq \left(1 - \prod_{l \in L_f \setminus H^*} (1 - p_l) \right) - \mu + \sum_{l \in H^* \cap L_f} p_l \leq \sum_{l \in H^* \cap L_f} p_l
\end{aligned}$$

The last inequality follows from the fact that $(p_f - \mu) < 0$ for a path consisting of only good links. This completes the proof of lemma 2. \square

Consider the set of hypotheses with single link failures- say S_1 . We first show that $H_{opt} = \arg \max_{H \in S_1} L(H)$, corresponds to a failed link as $T_{min} \rightarrow \infty$ which in turns implies that the greedy algorithm succeeds in picking a failed link in the first iteration. Let $H_l \in S_1$ denote the hypothesis $\{l\}$ where l is a good link in the topology, $F(l)$ is the set of flows that go through the link l , n_f is the number of packets sent by flow f and p_f as before is the ground truth drop probability of the path taken by f . We have,

$$\begin{aligned}
E[LL(H_l)] &= \sum_{f \in F(l)} E[X(f, H)] = \sum_{f \in F(l)} n(p_f - \mu) \\
&\leq \sum_{l^* \in H^*} \sum_{f \in F(l) \cap F(l^*)} n(p_f - \mu) \\
&\leq \sum_{l^* \in H^*} \sum_{f \in F(l) \cap F(l^*)} n_f p_{l^*} \leq \sum_{l^* \in H^*} p_{l^*} T(l, l^*) \\
&\leq \sum_{l^* \in H^*} \frac{1}{\alpha} p_{l^*} T(l^*) \leq \sum_{l^* \in H^*} \frac{1}{\alpha} \sum_{F(l^*)} n_f p_{l^*} \\
&\leq \sum_{l^* \in H^*} \frac{1}{\alpha} \sum_{F(l^*)} n_f 2(p_{l^*} - \mu) < \frac{2}{\alpha} \sum_{l^* \in H^*} E[LL(H_{l^*})] \\
&< \arg \max_{l \in H^*} E[LL(H_l)]
\end{aligned}$$

Note that $p_f - \mu < 0$ if flow f takes a path with all good links and $|H^*| \leq \alpha/2$.

This shows that the greedy inference algorithm will pick a failed link in the first iteration, say l_1 . Greedily picking the link that maximizes the log likelihood of the hypothesis $\{l_1\}$ is equivalent to deleting all flows crossing l_1 and the link l_1 itself from the input for analysis for subsequent iterations. Hence, the same proof about iteratively picking a failed link works for subsequent iterations if we ensure that the traffic-skew is maintained after deleting l_1 and all flows crossing it.

For any pair of links $l_2, l_3 \neq l_1$, we have $T(\{l_2, l_3\})/T(\{l_2\}) \leq \frac{1}{\alpha}$. Let's say that the number of packets crossing link l_2 is $T'(\{l_2\})$ after we delete l_1 and all

Algorithm 1 Flock inference: Greedy search with JLE (crux)

```

1: procedure GREEDYSEARCH()
2:   current_hypothesis ← []
3:   Δ = ComputeInitialDelta()
4:   while max(Δ) > 0 do
5:     link = argmax(Δ)
6:     Δ = UpdateDeltaArr(Δ, current_hypothesis, link)
7:     current_hypothesis.add(link)
8:   return current_hypothesis
9: procedure UPDATEDELTAARR(Δ, hypothesis, link)
10:  new_hypothesis ← hypothesis + [link]
11:  for F in FlowsIntersectingWithLink(link) do
12:    for l in F.links do
13:      Δ[l] += GetFlowDelta(new_hypothesis, l, F)
14:      Δ[l] -= GetFlowDelta(hypothesis, l, F)
15:  return Δ

```

flows crossing l_1 . Then we have,

$$\begin{aligned} \frac{T'(\{l_2, l_3\})}{T'(\{l_2\})} &\leq \frac{T(\{l_2, l_3\})}{T(\{l_2\})} = \frac{T(\{l_2, l_3\})}{T(\{l_2\}) - T(\{l_1, l_2\})} \\ &\leq \frac{T(\{l_2, l_3\})}{(1 - 1/\alpha)T(\{l_2\})} = \frac{1}{\alpha - 1} \end{aligned}$$

Thus, for subsequent iterations, after deleting l_1 , traffic is $1/(\alpha - 1)$ -skewed and the number of failures is $\alpha/2 - 1$ in the deleted graph. The same argument as before shows that the greedy algorithm will pick a failed link in every iteration as long as there is at least one failed link not in the current hypothesis.

Stopping Criteria: Finally we need to show that once all failed links are picked, the greedy algorithm will halt. This happens when $LL(H_l) < 0$ for all l not in the current hypothesis. Note that the input to log likelihood computations in the

current iteration is the topology obtained after deleting all links in the current hypothesis and all flows crossing any of those links. As before we have,

$$E[LL(H_l)] = \sum_{f \in F(l)} E[X(f, H)] = \sum_{f \in F(l)} n(p_f - \mu)$$

Note that for a path with all good links, $p_f - \mu \leq p_g - \mu < 0$. Hence, $E[LL(H_l)]$ is bounded away from 0 towards $-\infty$. Since, $LL(H_l)$ is the sum of independent binary variables corresponding to packets each of whose expectation is $\leq (p_g - \mu) < 0$, applying Chernoff bounds followed by a union bound for all links shows that $LL(H_l) < 0$ for all links l with high probability. This complete the proof of Theorem 2. \square

A.1 Defining precision/recall

Precision is the fraction of predicted failed links that had actually failed and **recall** is the fraction of failed links that were correctly reported as failed. A faulty device or any of its links are considered to be correct for calculating precision. For calculating recall, including the faulty device itself in H counts as 100% recall, and including $x\%$ of the device links in H counts as $x\%$ recall, where H is the set of failed links/devices predicted by the algorithm. More precisely, if H is the set of failed links predicted by the algorithm and H^* is the actual set of failed links, then precision = $|H \cap H^*|/|H|$ and recall = $|H \cap H^*|/|H^*|$.

We define precision to be 1 if the algorithm returns the empty hypothesis. For 0 actual failures, precision represents the fraction of examples where the algorithm returns a wrong answer and recall is 1 since there are no failures to detect.

Algorithm 2 Flock’s Hypotheses search: Greedy with Joint Likelihood Exploration

```

1: procedure GREEDYSEARCH()
2:   current_hypothesis  $\leftarrow$  []
3:    $\Delta$  = ComputeInitialDelta()
4:   while max( $\Delta$ ) > 0 do
5:     link = argmax( $\Delta$ )
6:      $\Delta$  = UpdateDeltaArr( $\Delta$ , current_hypothesis, link)
7:     current_hypothesis.add(link)
8:   return current_hypothesis
9:
10: procedure UPDATEDELTAARR( $\Delta$ , hypothesis, link)
11:   new_hypothesis  $\leftarrow$  hypothesis + [link]
12:   for F in FlowsIntersectingWithLink(link) do
13:      $\triangleright$  these counters are a simple data structure
14:      $\triangleright$  trick to speed-up the subsequent for loop
15:     old_counters = GetCounters(hypothesis, flow)
16:     new_counters = GetCounters(new_hypothesis, flow)
17:     for l in F.links do
18:        $\Delta[l]$  += GetFlowDelta(l, *new_counters)
19:        $\Delta[l]$  -= GetFlowDelta(l, *old_counters)
20:   return  $\Delta$ 
21:
22: procedure GETCOUNTERS(hypothesis, flow)
23:   paths_failed  $\leftarrow$  0
24:   num_paths = dict()
25:   for path in flow.paths do
26:     if (PathFailedAsPerHypothesis(path, hypothesis)) then
27:       paths_failed++
28:     else
29:       for link in path do
30:         num_paths[link]++
31:   return (paths_failed, num_paths[link], flow.paths.size(),
32:         flow.packets_sent, flow.bad_packets)
33:
34: procedure GETFLOWDELTA(link, paths_failed, num_paths,
35:   n_flow_paths, packets_sent, bad_packets)
36:   bad_paths = paths_failed + num_paths[link]
37:   return GetLogLikelihood(bad_paths, n_flow_paths, packets_sent,
38:     packets_dropped)
39:
40: procedure COMPUTEINITIALDELTA()
41:   for l in links do
42:      $\Delta[l]$   $\leftarrow$  0
43:   for flow in flows do
44:     counters = GetCounters([], flow)
45:     for l in flow.links do
46:        $\Delta[l]$  += GetFlowDelta(l, *counters)
47:   return  $\Delta$ 
48:
49: procedure GETLOGLIKELIHOOD(bad_paths, n_flow_paths,
50:   bad_packets, packets_sent)
51:   good_packets = packets_sent - bad_packets
52:   log_likelihood = bad_paths * pow( $p_b$ , bad_packets) *
53:     pow(1 -  $p_b$ , good_packets)
54:   good_paths = n_flow_paths - bad_paths
55:   log_likelihood += good_paths * pow( $p_g$ , bad_packets)
56:     * pow(1 -  $p_g$ , good_packets)
57:   log_likelihood /= n_flow_paths
58:   return log_likelihood

```

Algorithm 3 JLE can be used to speedup Sherlock’s inference, with max concurrent failures = K

```

1: best_hypothesis = None
2: max_likelihood =  $-\infty$ 
3: procedure SHERLOCKWITHJLE()
4:    $\Delta$  = ComputeInitialDelta()
5:   ExploreBranch([],  $\Delta$ , 0.0)
6:   return best_hypothesis
7:
8: procedure EXPLOREBRANCH(current_hypothesis,  $\Delta$ ,
9:   current_likelihood)
10:   if current_likelihood > max_likelihood then
11:     max_likelihood = current_likelihood
12:     best_hypothesis = current_hypothesis
13:   if current_hypothesis.size < K then
14:     for l in links do
15:       new_hypothesis = current_hypothesis.add(l)
16:       new_likelihood = current_likelihood +  $\Delta[l]$ 
17:        $\Delta_{new}$  = UpdateDeltaArr( $\Delta$ , current_hypothesis, l)
18:       ExploreBranch(new_hypothesis,  $\Delta_{new}$ , new_likelihood)

```

B Joint Likelihood Exploration: pseudocode

Refer to Algorithm 1 for a short summary of Flock’s inference algorithm and Algorithm 2 for full pseudocode.

C Full Runtime analysis

Let n be the number of links, m be the number of flows, T be an upper bound on the number of links that any flow intersects with, D be an upper bound on the number of flows that any link intersects with and K be the maximum number of concurrent failures (note that Flock’s algorithm doesn’t need to know K).

We describe the components in the running time of Flock’s overall inference, including Greedy and JLE:

- Before the first greedy iteration, the Δ array is computed once, in a linear pass over all flows and their path sets, incurring $O(n + mT)$ time.
- After each greedy iteration, updating the Δ array via JLE requires iterating over all flows that intersect with the newly added link. For each such flow F , let L_F be the links that F intersects with. Updating all entries $\Delta_{H'}(l, F)$ for all $l \in L_F$ requires a couple of passes over L_F . Thus, the execution time for subsequent $(K - 1)$ greedy iterations, barring the first, is $O((K - 1)DT)$.

Hence, the running time of Greedy inference with JLE is $O(n + mT + (K - 1)DT)$. If we had used just Greedy without JLE (computing likelihood of each hypothesis individually), the runtime would be $O(n + mT + (K - 1)nDT)$.

We compare runtime with Sherlock. To compute the MLE, Sherlock scans all $O(n^K)$ hypotheses with $\leq K$ failures. Sherlock is better than brute force, however: it uses $LL(H)$, for an explored hypothesis H , to compute $LL(H \oplus l)$ by updating the flow contributions $LL_F(H)$ for all flows F that intersect with link l (since their likelihoods would have changed after flipping the status of link l), giving $O(n^K DT)$ runtime. We can apply JLE to accelerate Sherlock’s inference by evaluating n

neighbor hypotheses at once (Algorithm 3 in appendix), improving Sherlock's runtime, by a factor of n , to $O(n^{K-1}DT)$.

However, this is exponential in K and too slow for our purposes. From the analysis above (and our experiments later), it can be seen that Greedy + JLE is dramatically faster.