

m3: Accurate Flow-Level Performance Estimation using Machine Learning



Chenning Li*, Arash Nasr-Esfahany*, Kevin Zhao, Kimia Noorbakhsh, Prateesh Goyal, Mohammad Alizadeh, and Thomas Anderson

Predict Network Performance

- Data center network operators need to predict the impact of design choices on network performance (e.g., tail latency, throughput, etc)
- Simulation** is a common tool to predict network performance

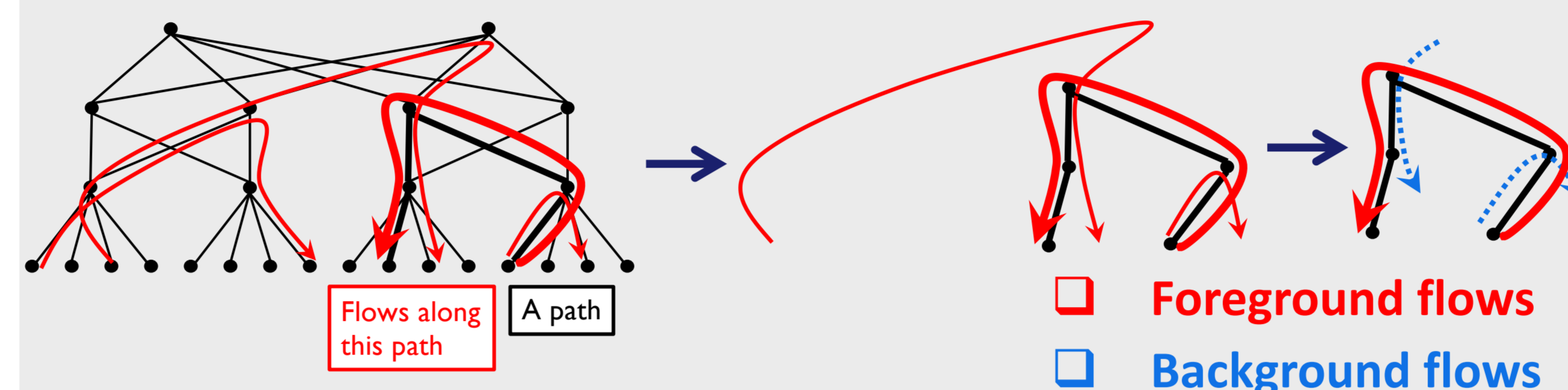


- Recent work on fast simulation: DeepQueueNet, Parsimon, ...
- All are **packet-level simulators** -> **slow** for large-scale networks especially as the networks become larger and faster

Path-level Decomposition

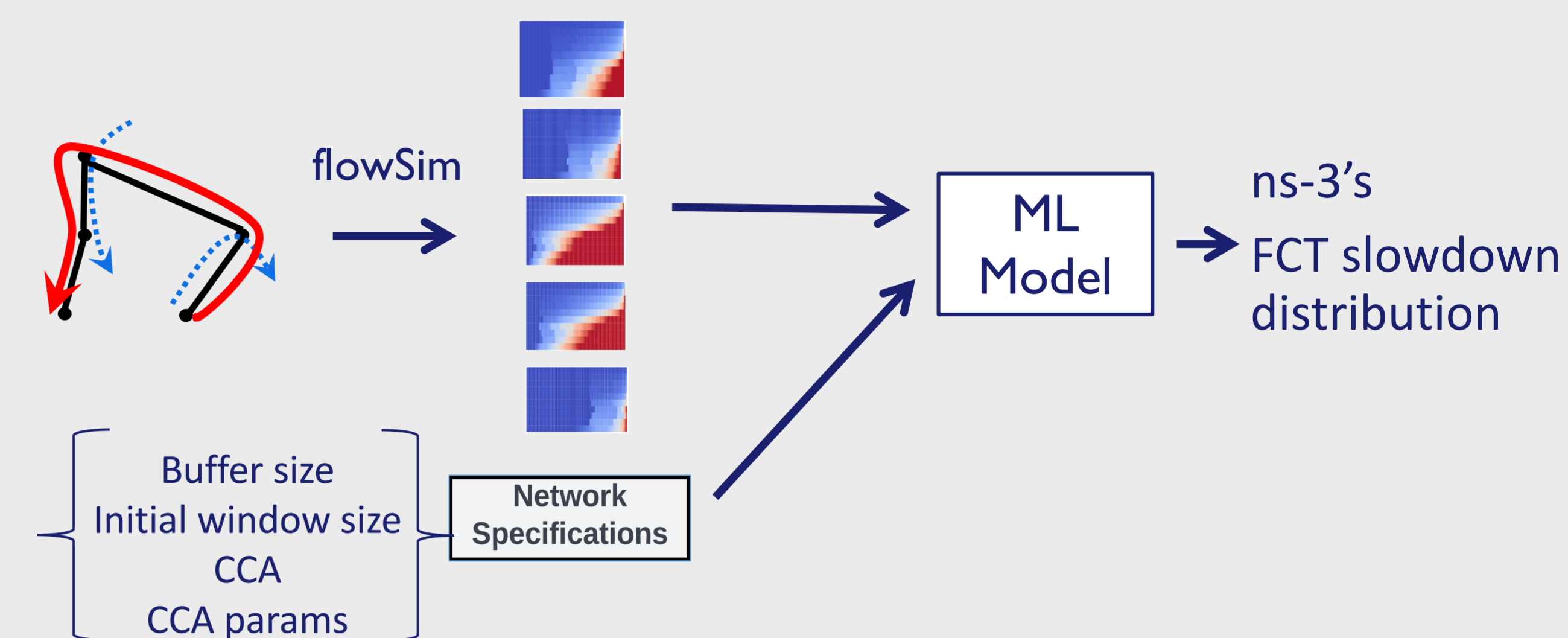
m3 decomposes the network topology into **independent paths**, predicts performance on sampled paths, aggregates results

- Assumption:** Flows that do not intersect a path have a second-order effect on the behavior of flows entirely along that path



- Goal:** for each path-level simulation, m3 estimates the flow completion time (FCT) distribution of the **foreground flows**
- Benefits:** Path-level sims. are easier to learn & enable parallelism, yet produce accurate estimates of network-wide behavior

m3's fast path-level simulation using ML

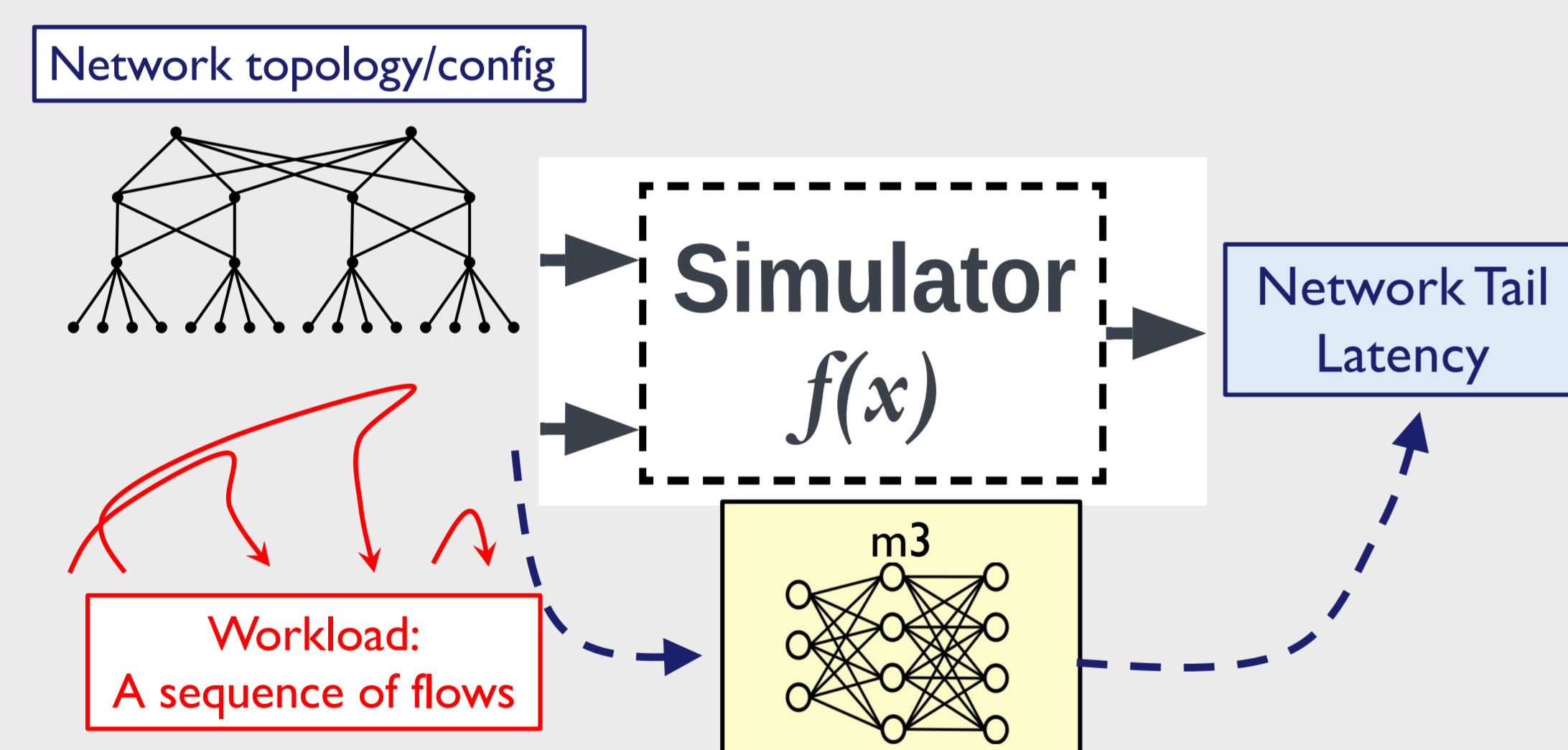


- m3 uses ML to correct flowSim, "translating" between flowSim and ns-3's output

General principle: use a simple reference system to extract features to learn a model of a complex system

Abstract network simulator as a function

- Learn a model approximating the simulator function mapping a network scenario to **aggregate performance statistics**
- Example: **network tail latency**



m3: ~1200X speedup with only ~10% estimation error in a 384-rack, 6144-host topology (vs. ns-3)

From ~10 hours (ns-3) to less than 1 min (m3)

Two main challenges:

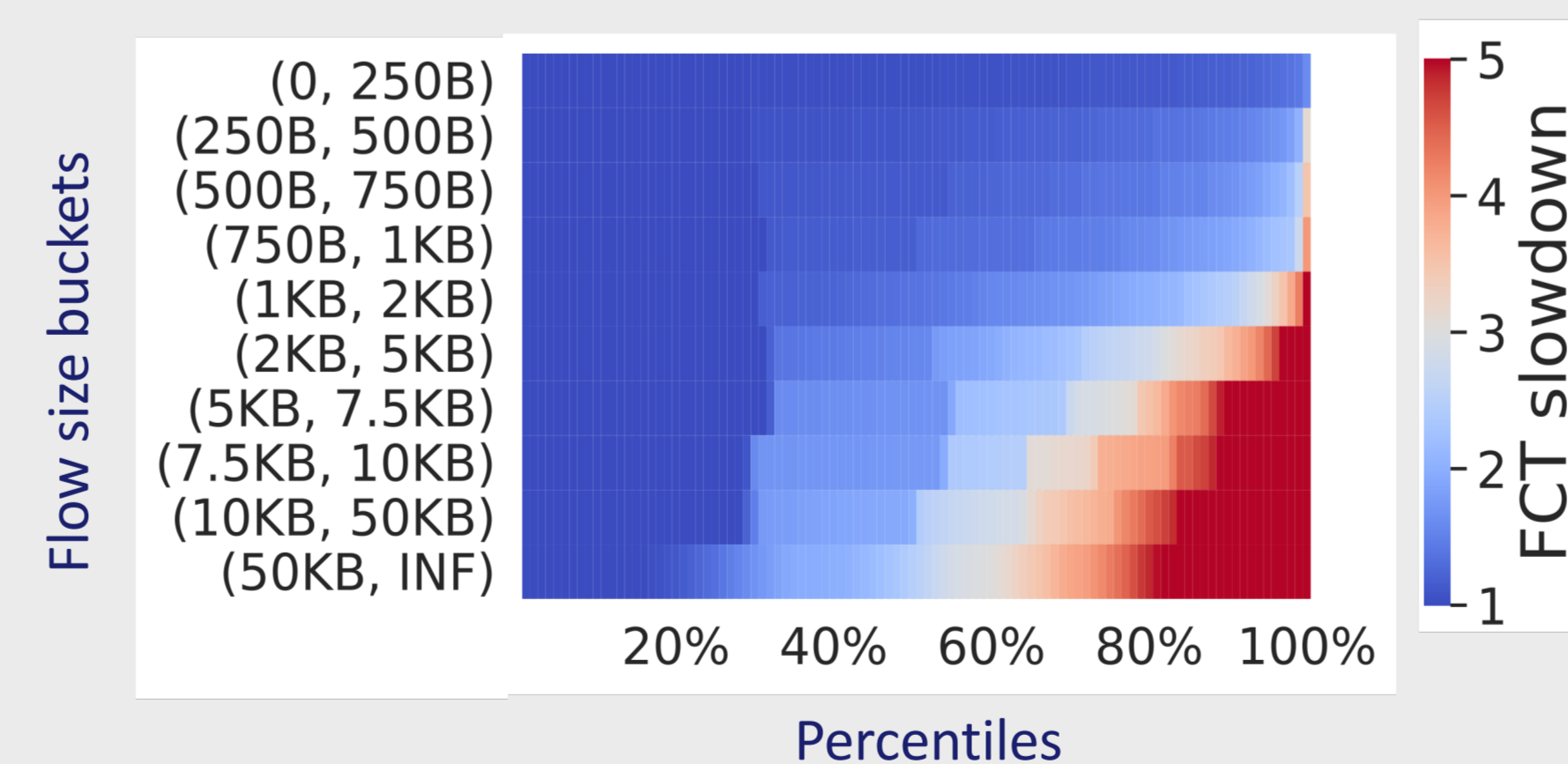
- Hard to represent the function in a compact way
- Slow to generate the training dataset for the ML model

m3 uses (i) path decomposition and (ii) feature extraction from a flow-level simulator to tackle these challenges

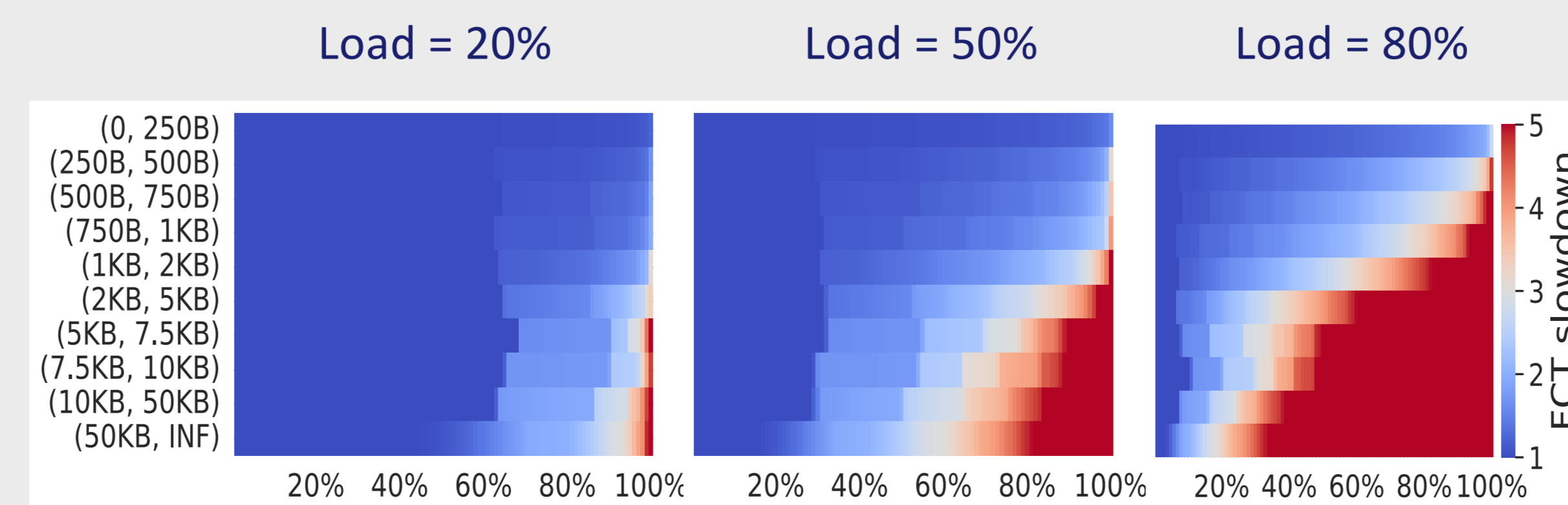
Workload Featurization

m3 uses ML to estimate the path-level perf. quickly and uses a **flow-level simulator** to extract compact features for its ML models

- flowSim: max-min flow-level simulation**
 - Fast:** <1 sec for a path-level sim with 1 million flows
 - Not accurate:** no queuing -> underestimates short flow FCT
- Extract a **compact feature map** from the complex workload:
 - Run path-level simulation with flowSim
 - Summarize** per flow FCT slowdown into a feature map



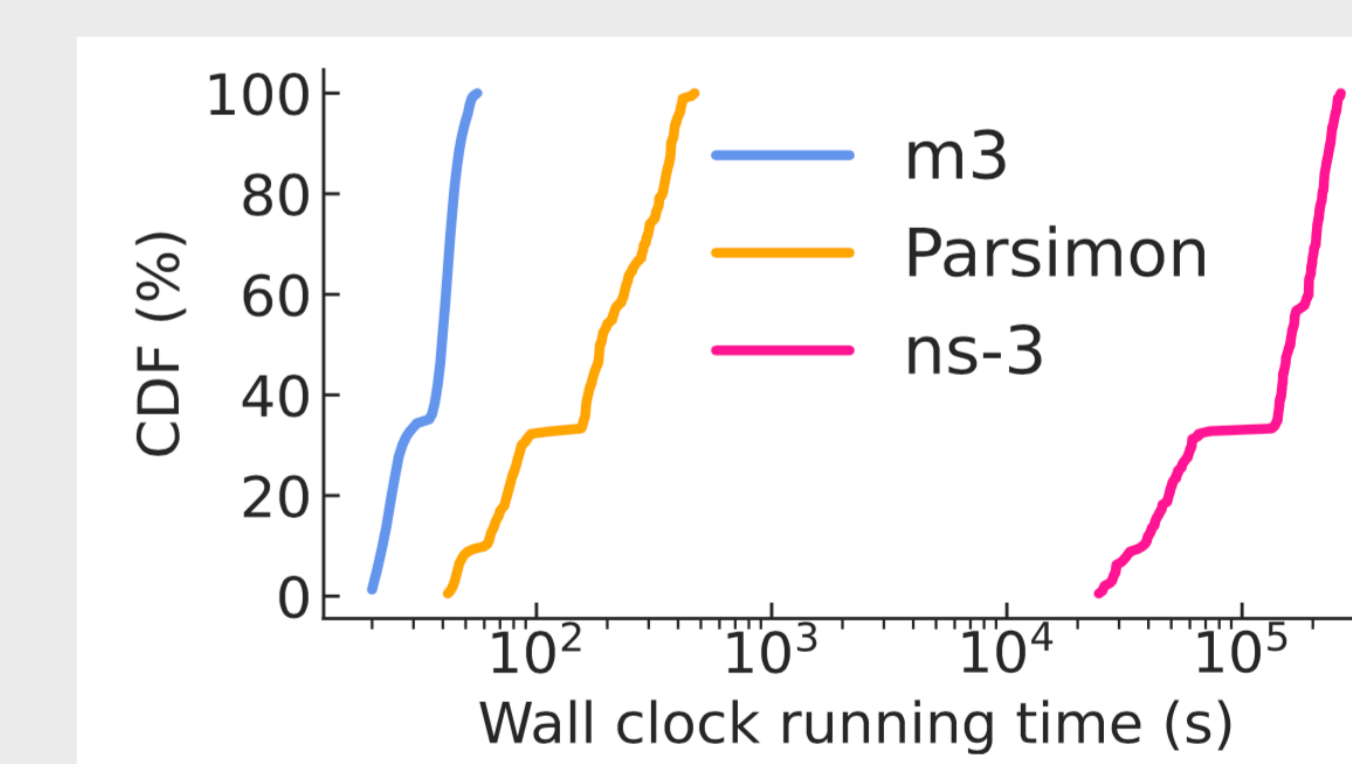
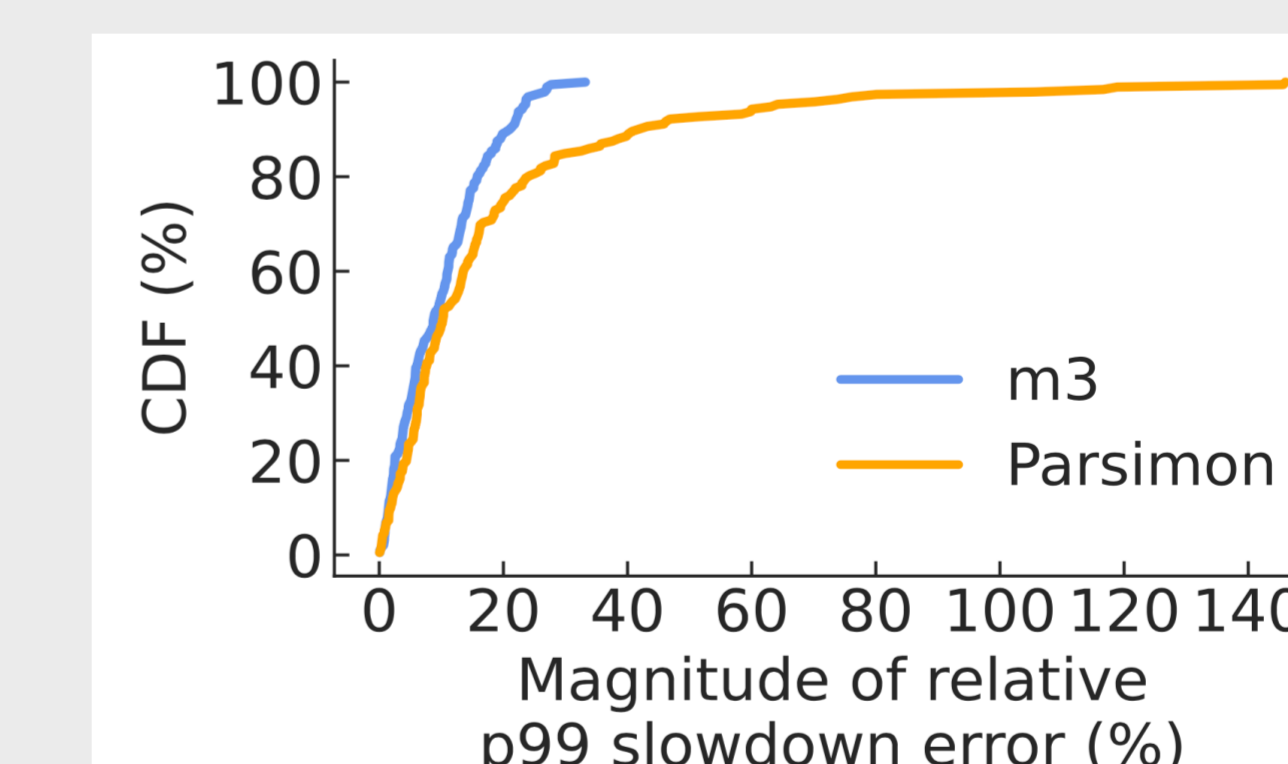
- The feature map distinguishes between workloads in a logical way



Results

Extensive simulation

- Various production workloads based on Meta's data center network
- Baseline: Parsimon^[1]
- m3 reduces Parsimon's mean error (relative to ns-3) from **18.3% to 9.9%**
- m3 has 5.7x speed-up over Parsimon, and is 3 orders **faster** than ns-3



Large-scale simulation

- Meta production workload
- A 384-racks, 6,144-host Meta's data center fabric
- Baseline: Parsimon^[1]

Init. Window	Methods	p99	Error	Time	Speedup
10KB	ns-3	2.05	-	13.5h	-
	Parsimon	4.29	+109%	1m29s	546x
	m3	2.10	+2.44%	37s	1314x
18KB	ns-3	2.44	-	11.9h	-
	Parsimon	2.73	+11.9%	1m24s	510x
	m3	2.30	-5.74%	40s	1071x

[1] Zhao, Kevin, et al. "Scalable tail latency estimation for data center networks." In proceedings of USENIX Symposium on Networked Systems Design and Implementation (NSDI). 2023.