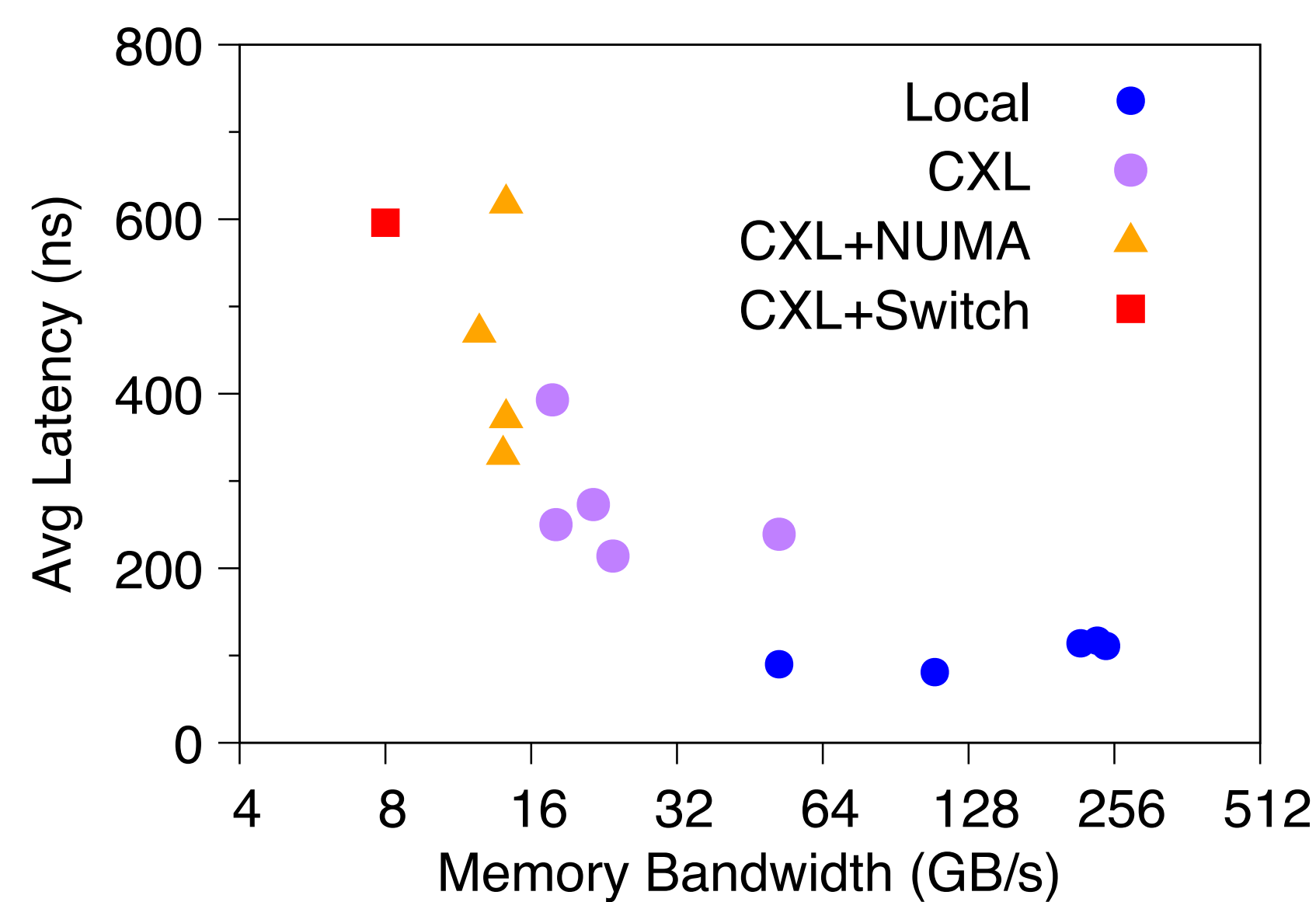


Heterogenous CXL Latency and Bandwidth



CXL introduces **diverse and higher** latency, what is the performance implication of CXL memory across **CXL devices, processors, and workloads at scale?**

Research Questions:

Is CXL latency as **stable/predictable** as regular DRAM?

How does CXL latency affect **workload performance**?

How does CXL latency affect CPU **pipeline** (e.g., **prefetching**)?

Melody Overview

A comprehensive framework for CXL characterization and analysis

265 workloads across 4 CXL devices under 7 memory latency levels on 5 processors

1. Unstable and unpredictable latency introduced by CXL

μs-scale memory tail latency even when bandwidth is not saturated

2. Extensive CXL characterization across diverse workloads

Quantitative slowdowns due to latency or bandwidth boundness

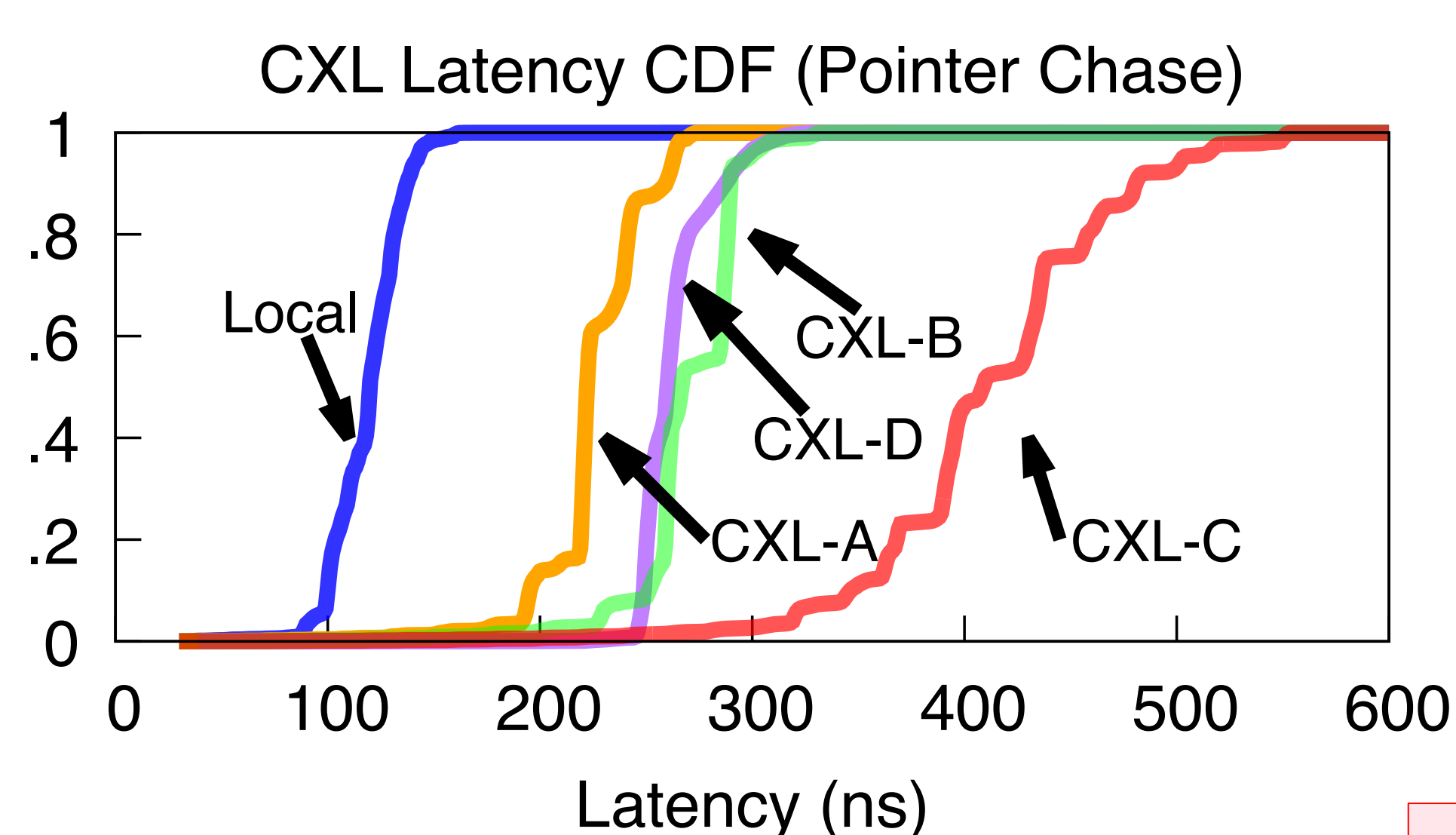
3. SPA: A simple and accurate performance analysis approach

9 CPU counters for **accurate** slowdown estimation (<5% inaccuracy for over 95% workloads)

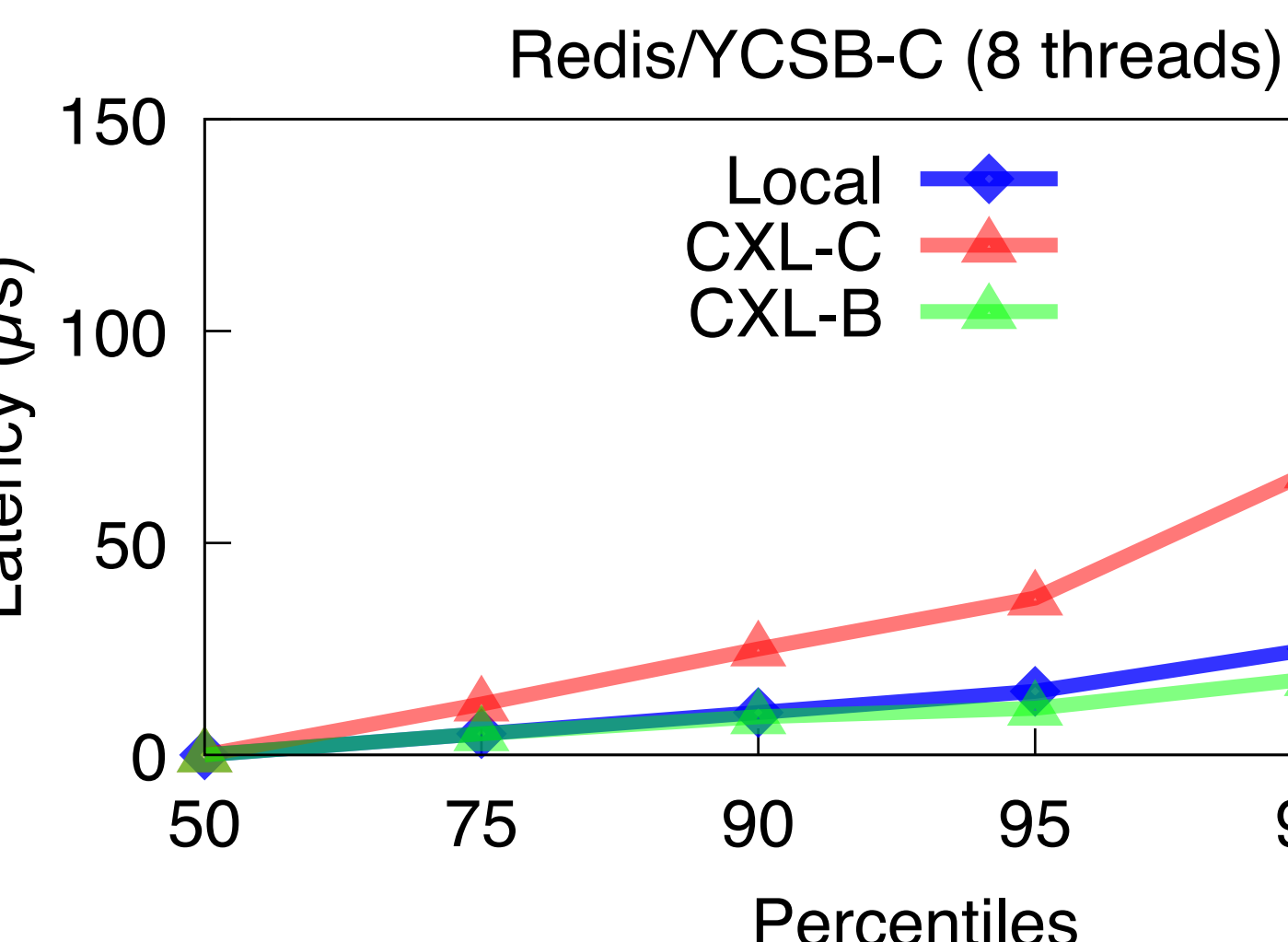
Dissect the **root causes** of CXL slowdown

Disclose CPU **prefetching** inefficiency

CXL Tail Latency



CXL devices exhibit **unstable/unpredictable** latency compared to regular DRAM.



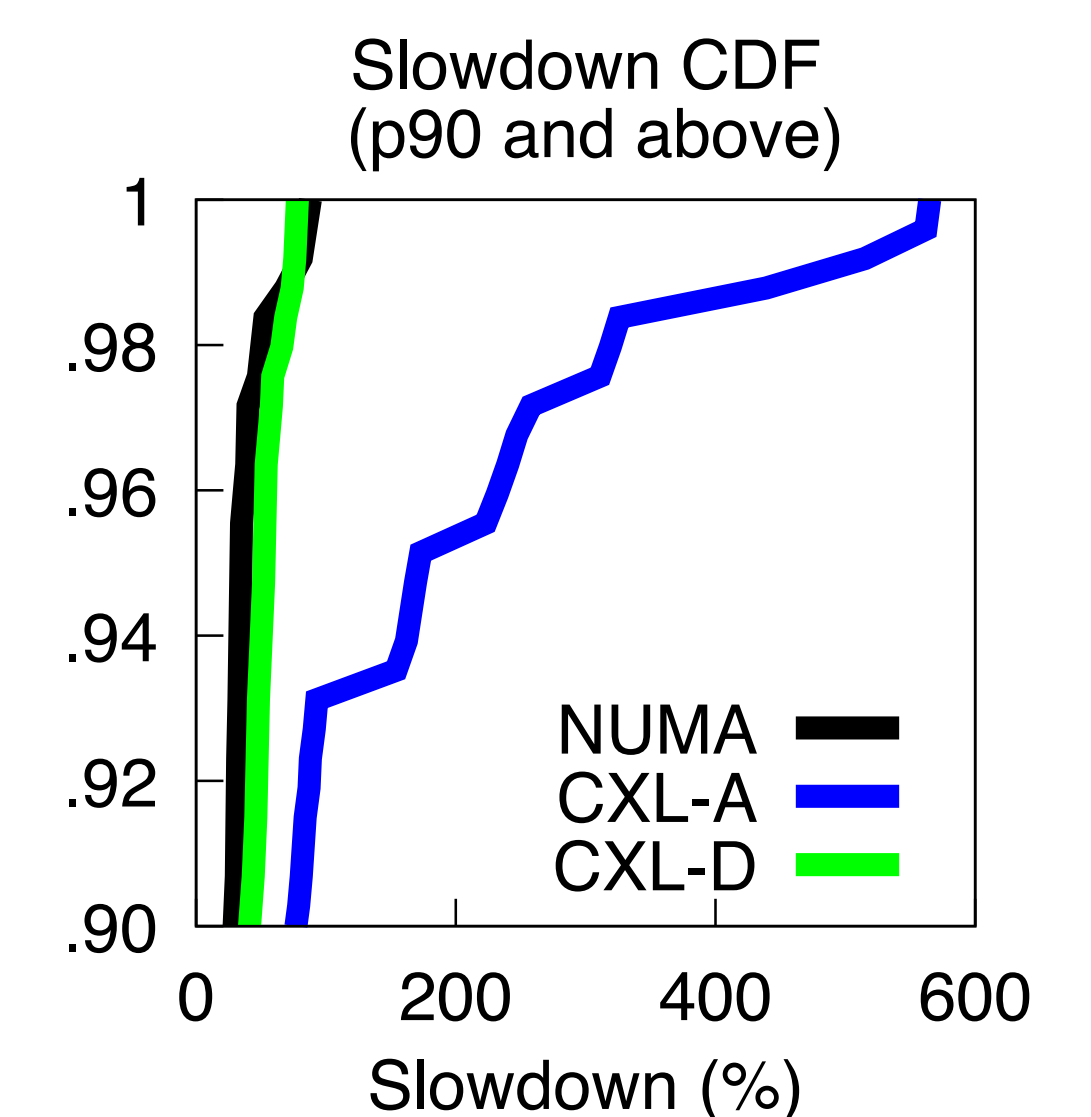
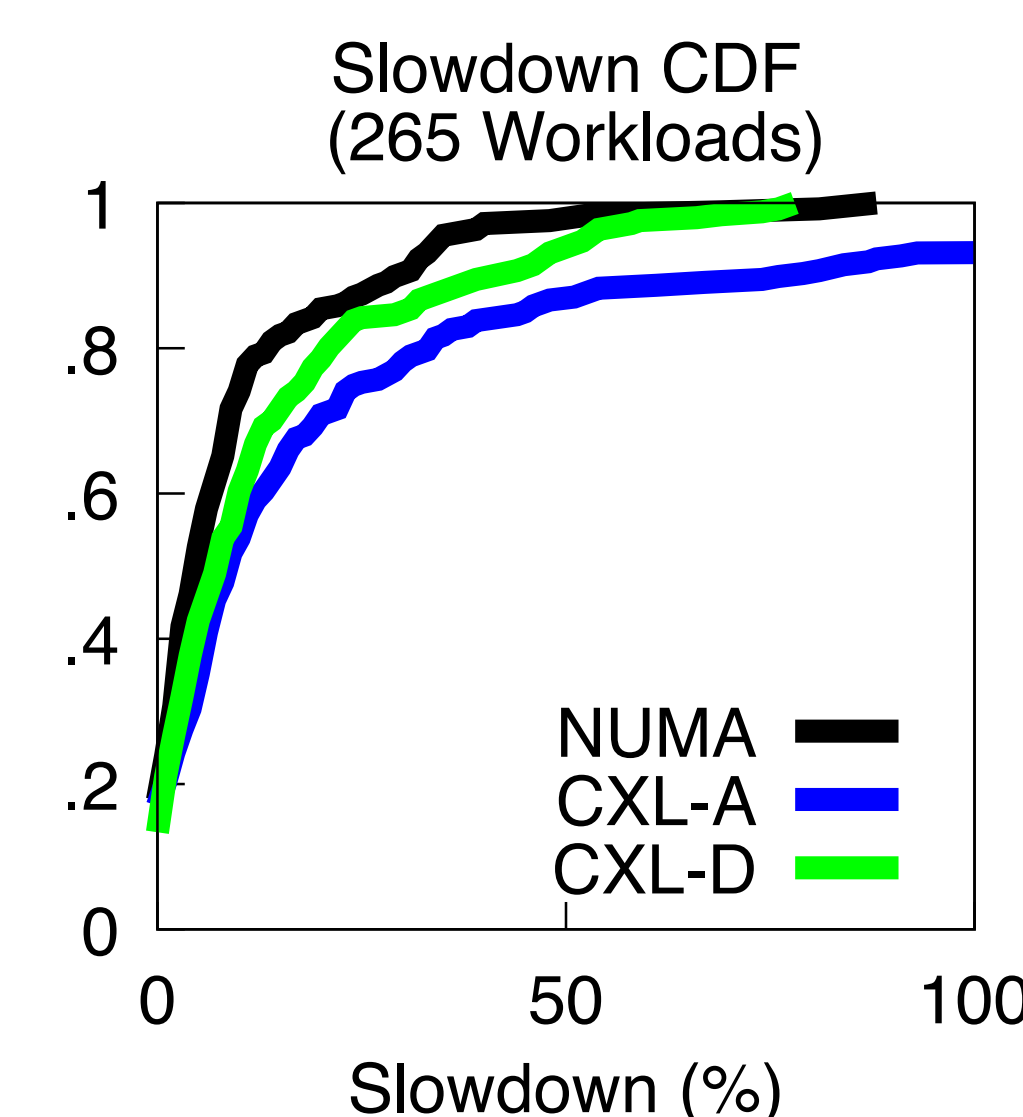
p99.9 – p50 = 133μs

p99.9 – p50 = 40μs

CXL tail latency can lead to unpredictable application performance.

Workload Characterization on CXL

CXL	#ch	Size	Lat	BW
	DDR	GB	ns	GB/s
CXL-A	2×DDR4	128	214	24
CXL-B	1×DDR5	128	271	22
CXL-C	1×DDR4	16	394	18
CXL-D	2×DDR5	756	239	52



The bandwidth limitation on CXL causes some workloads with high slowdown

The performance gap between CXL(-D) and NUMA diminishes due to its higher bandwidth even though its latency is worse

CXL memory can be used as a viable alternative to NUMA memory

SPA: Stall-based CXL Performance Analysis

1. Workload slowdown breakdown

$$\Delta \text{CPU Cycles} \approx \Delta \text{DRAM-Stalls} + \Delta \text{Cache-Stalls} + \Delta \text{Store-Stalls}$$

$$\text{Slowdown (S)} = \Delta \text{CPU Cycles} / \text{Cycles on Local}$$

$$\text{Slowdown (S)} = S_{\text{DRAM}} + S_{\text{Cache}} + S_{\text{Store}}$$

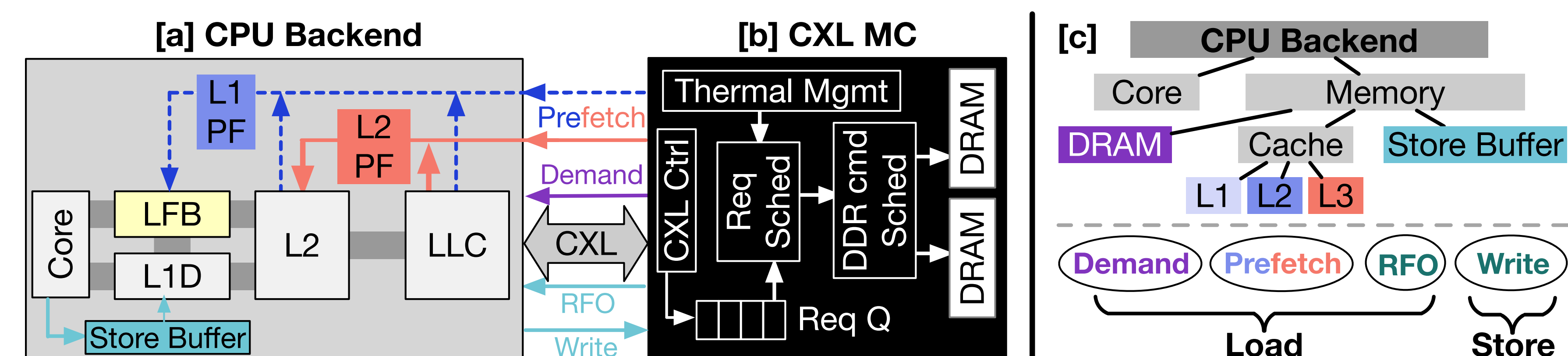
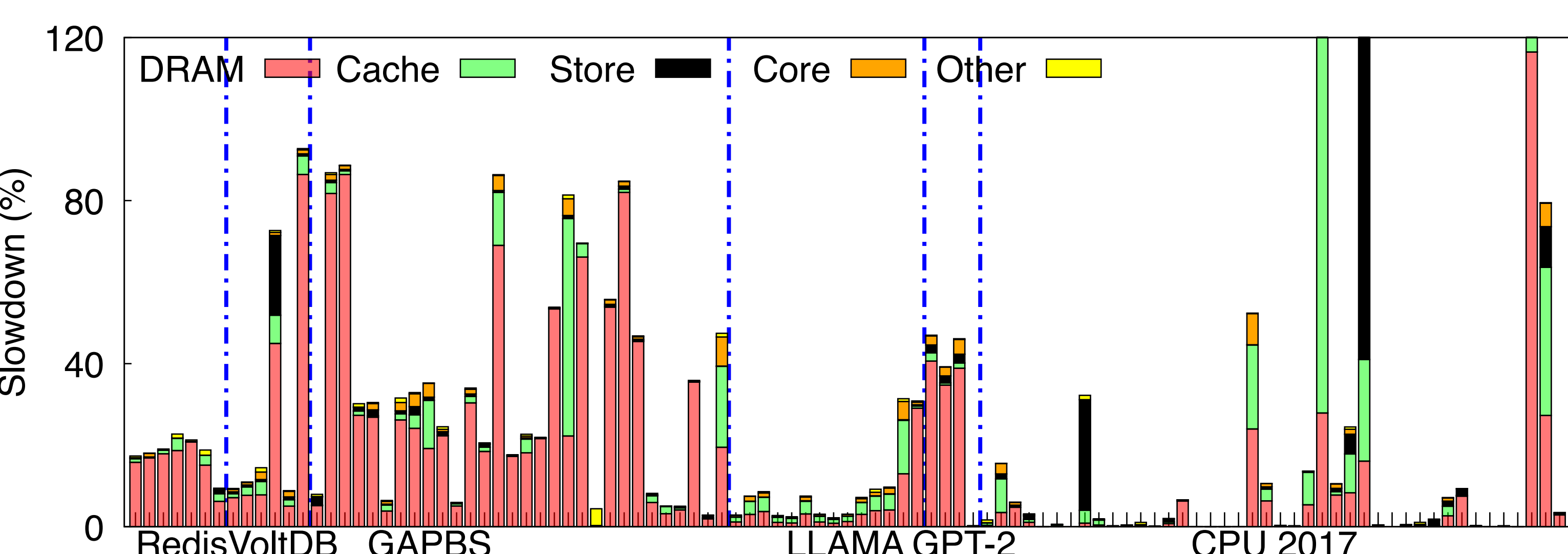
S_{DRAM} Demand read miss on L3

S_{Cache} Less efficient prefetching under longer memory latency

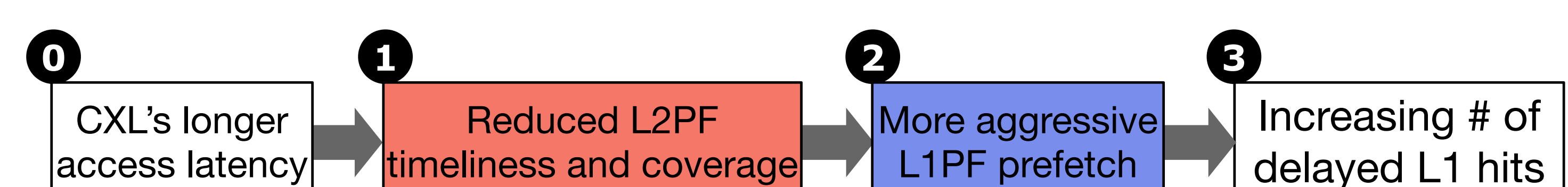
S_{Store} Intensive RFO on Store Buffer with limited size

2. CXL slowdown for real-world workloads

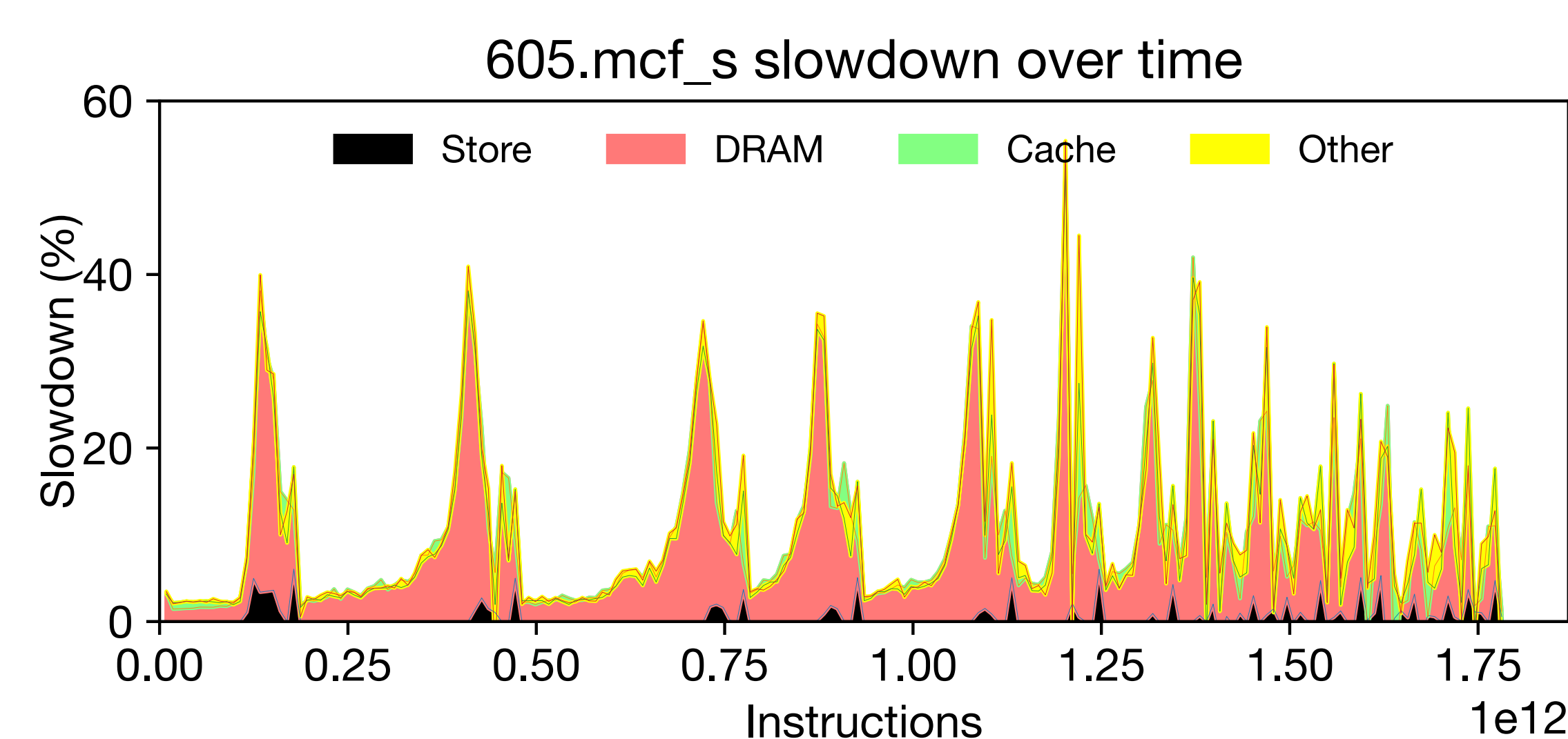
The sources of slowdown vary across workloads



3. Cache slowdown reasoning



4. Dynamic slowdown



More in the paper

