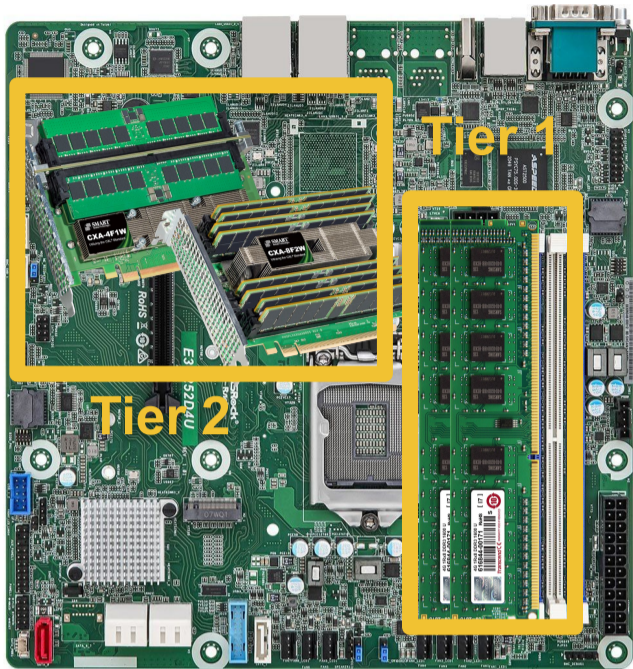# Handling Latency in Tiered Memory with Prefetching

Musa Unal, Vishal Gupta, Yueyang Pan, Yujie Ren, Sanidhya Kashyap
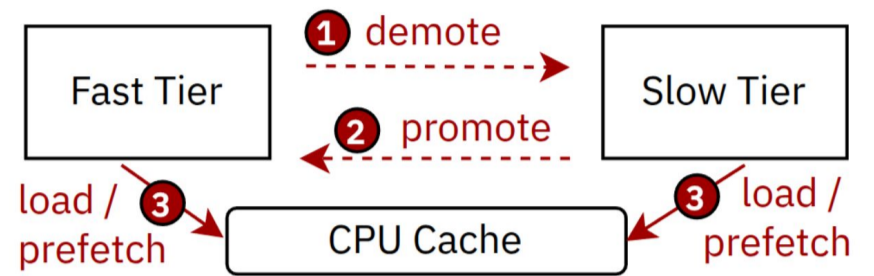
## Problem: Tiering systems do not consider how application is accessing the data

Tier 1

Tier 2

Data centers use different types of memories which have different characteristics. Trade-off between
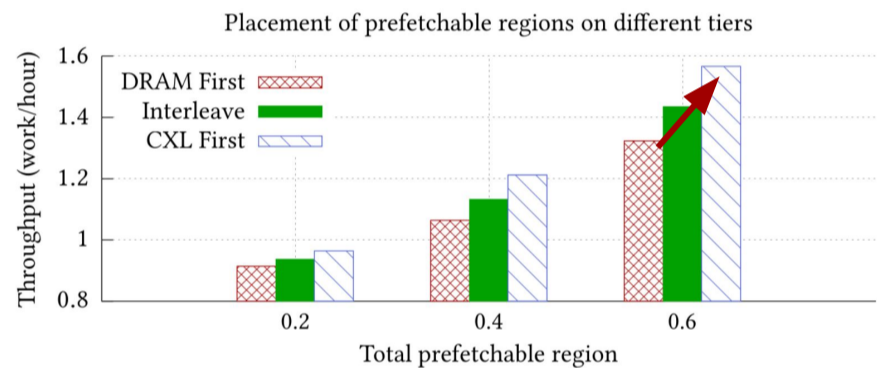
- Latency/Bandwidth
- Price
- Energy

Fast Tier ①demote→ Slow Tier
②promote←
load / prefetch ③ → CPU Cache ← ③ load / prefetch

Current tiering systems are trying to optimize load time (3) by demoting (1) and promoting (2) between the nodes.

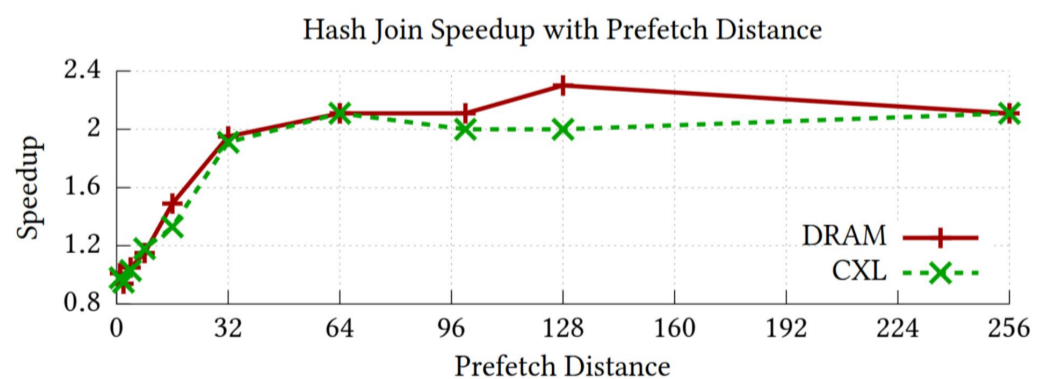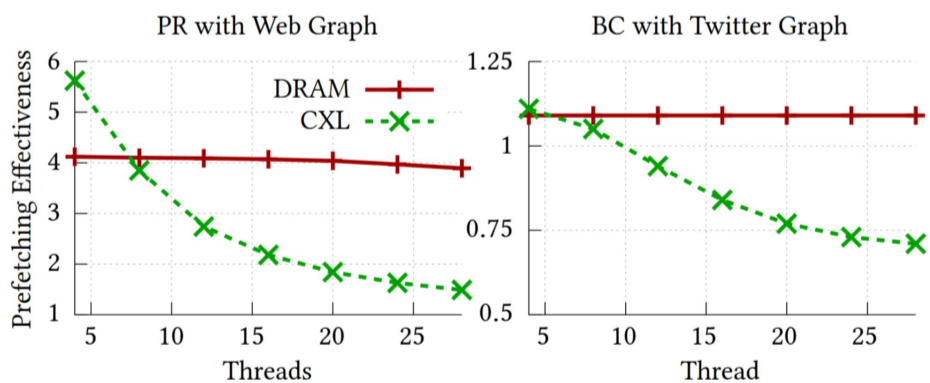## Insight: Data prefetchers can tolerate CXL latency.

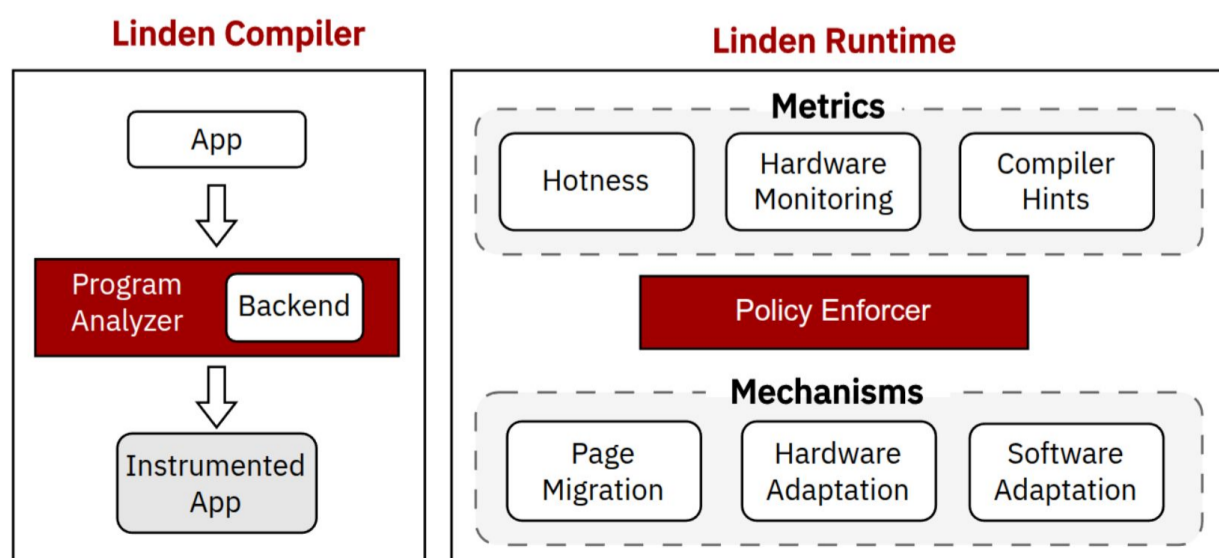Hardware prefetchers and software prefetching can *hide* the CXL latency

Placement of prefetchable regions on different tiers
(bar chart: Throughput (work/hour) vs Total prefetchable region; DRAM First, Interleave, CXL First at 0.2, 0.4, 0.6)

*Hardware prefetchers* can degrade performance under <u>high load</u>.

PR with Web Graph (Prefetching Effectiveness vs Threads; DRAM, CXL)
BC with Twitter Graph (vs Thread; DRAM, CXL)

*Software prefetches* must use <u>tiering-aware prefetch distances</u> to be effective.

Hash Join Speedup with Prefetch Distance (Speedup vs Prefetch Distance; DRAM, CXL)

## Proposal: A runtime system to effectively use prefetchers on tiered memory

**Linden Compiler**

App → Program Analyzer / Backend → Instrumented App

**Linden Runtime**

Metrics: Hotness, Hardware Monitoring, Compiler Hints

Policy Enforcer

Mechanisms: Page Migration, Hardware Adaptation, Software Adaptation

Reduce: Minimize latency by increasing the locality.
Tolerate: Hide the latency by prefetching.

1. Detect prefetchable regions
2. Monitor memory bandwidth
3. Enforce different policies regarding to application's needs