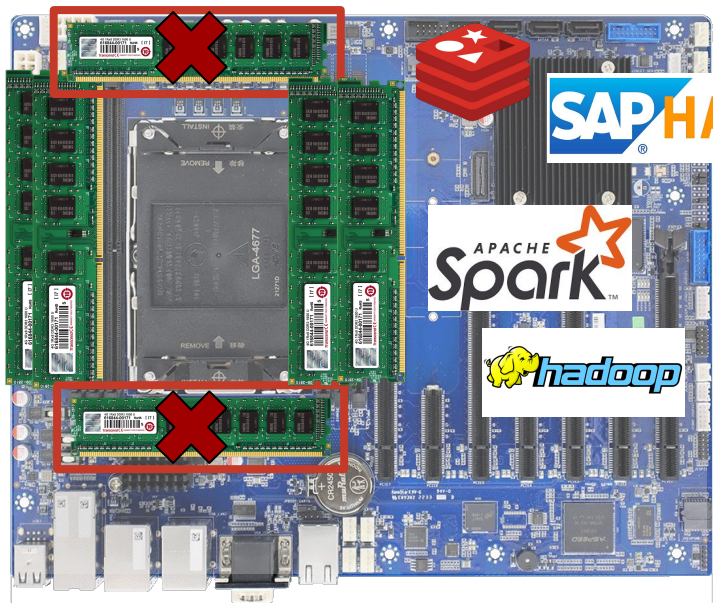# Tolerate It if You Cannot Reduce It: Handling Latency in Tiered Memory

Musa Unal, Vishal Gupta, Yueyang Pan, Yujie Ren, Sanidhya Kashyap

EPFL

1

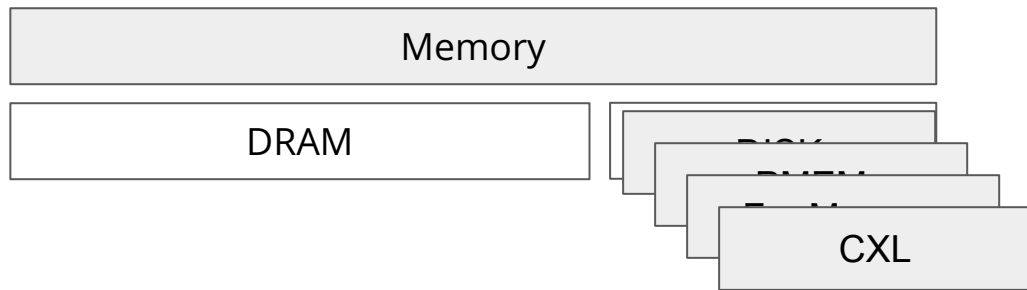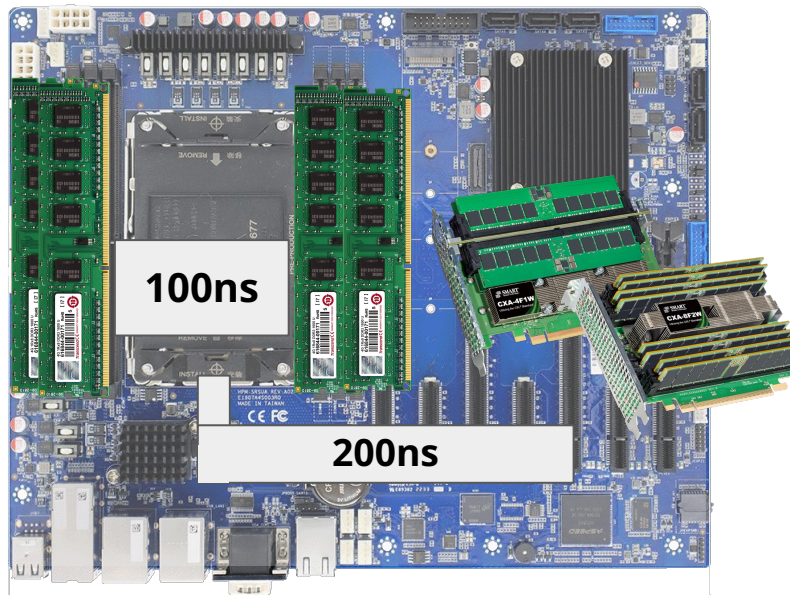# Tiered memory



**Why ?**

Applications demand for memory.
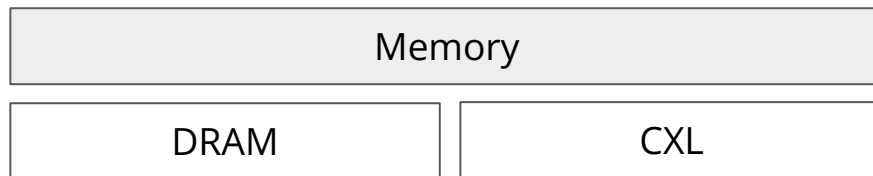
We want to expand main memory.

We can not add more DRAMs.

**How ?**

| Memory |
|---|

| DRAM | | DISK |
|---|---|---|
| | | PMEM |
| | | CXL |

# CXL expands main memory



100ns

200ns

## CXL Attached Memory

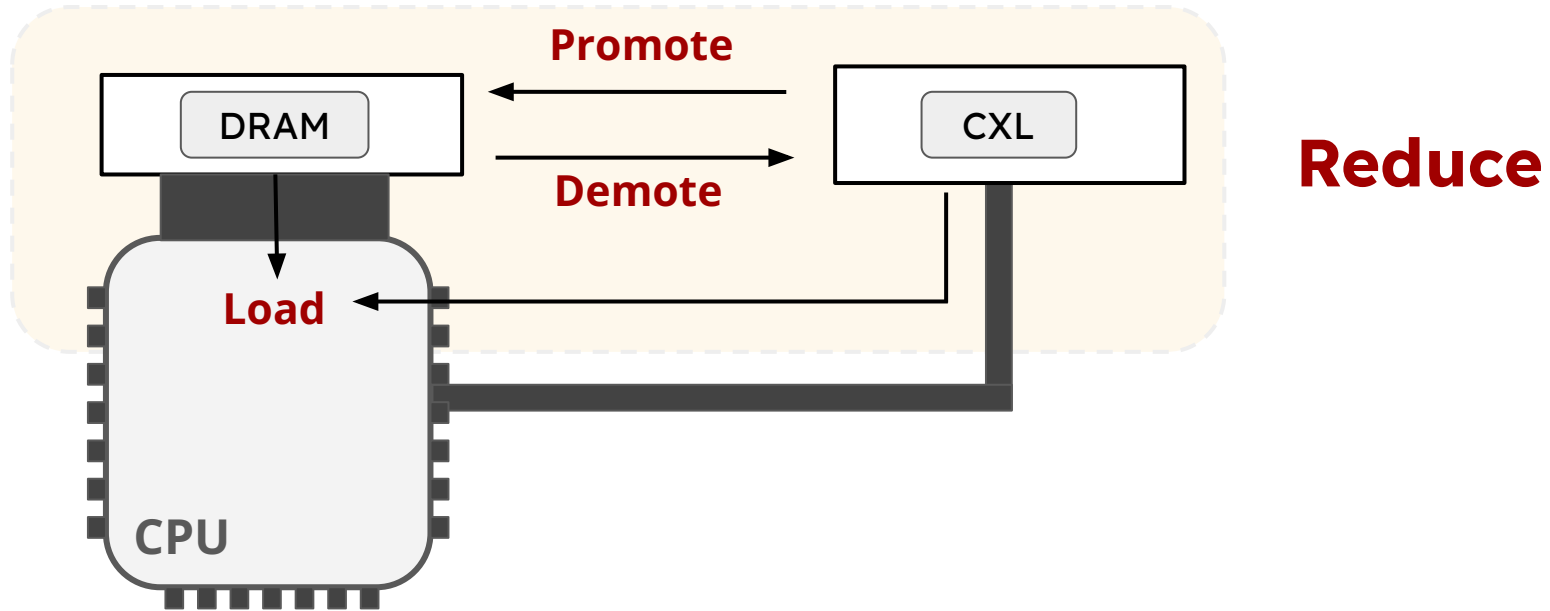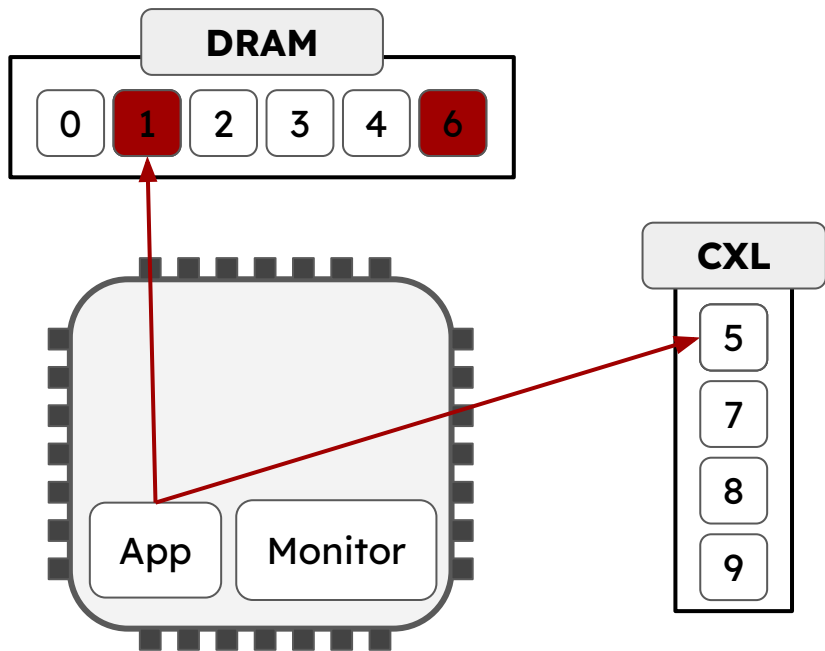| Memory |
|---|
| DRAM | CXL |

- Protocol for fast memory access over PCIe

- Latency is 2x higher (or higher with different topologies)

**How can we handle the latency in tiered memory?**

# Conventional approach to handle latency

# Page migration can reduce the latency
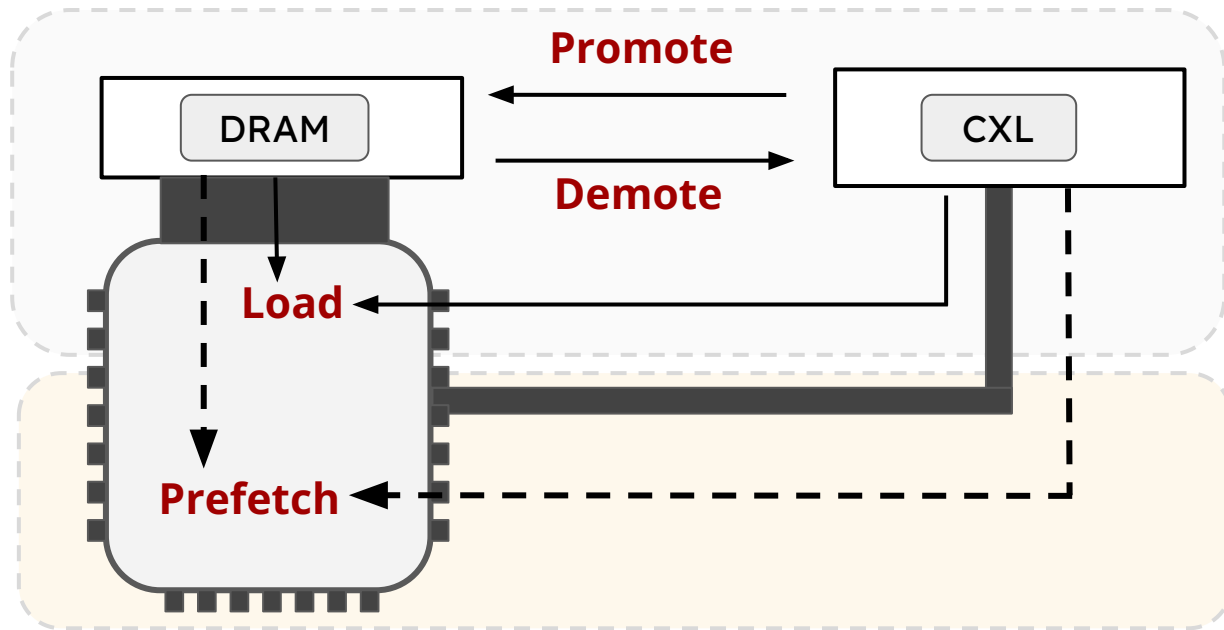


**DRAM**

| 0 | 1 | 2 | 3 | 4 | 6 |

**CXL**

5

7

8

9

App    Monitor

1. Tiering systems monitor pages hotness information.

2. Promote hot pages
   Demote cold pages
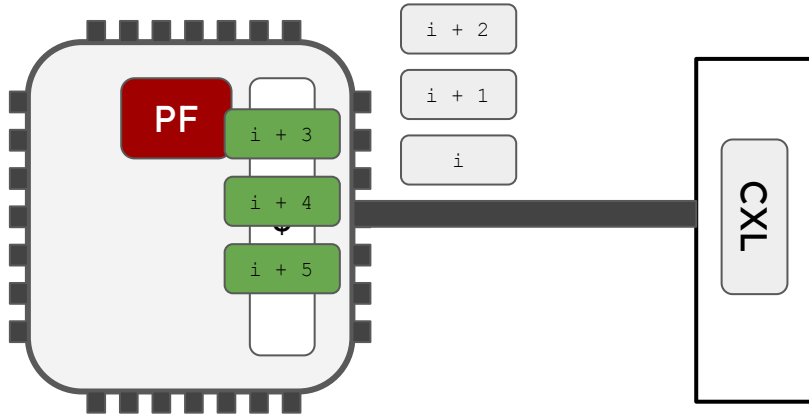
3. Increase **locality**, and **reduce** the latency

5

# Missing aspect to handle latency

# Hardware prefetchers can tolerate latency



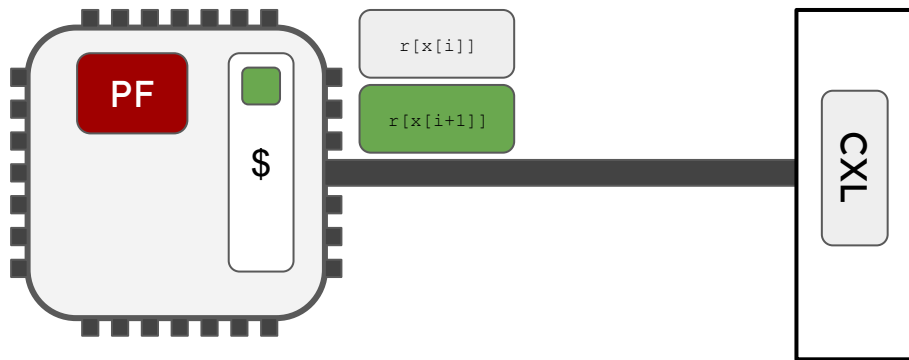| Sequential Pointer Chase | |
| --- | --- |
| **Stride Size** | **CXL Lat / DRAM Lat** |
| 4096 | ~2.01 |
| 128 | ~1.79 |
| 64 | ~1.55 |

**Hardware prefetchers can help to hide CXL's latency.**

# Hardware prefetchers can be harmful



**Prefetching can degrade the performance under high loads.**

# Software prefetching can help as well



**Timeliness is important!**
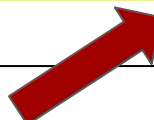
Early prefetches can be evicted.

Late prefetches are not helpful.

```
for i in range(1M):
  e = records[x[i]]
  work(e)
```

```
for i in range(1M):
  e = records[x[i]]
  prefetch(records[x[i+1]])
  work(e)
```
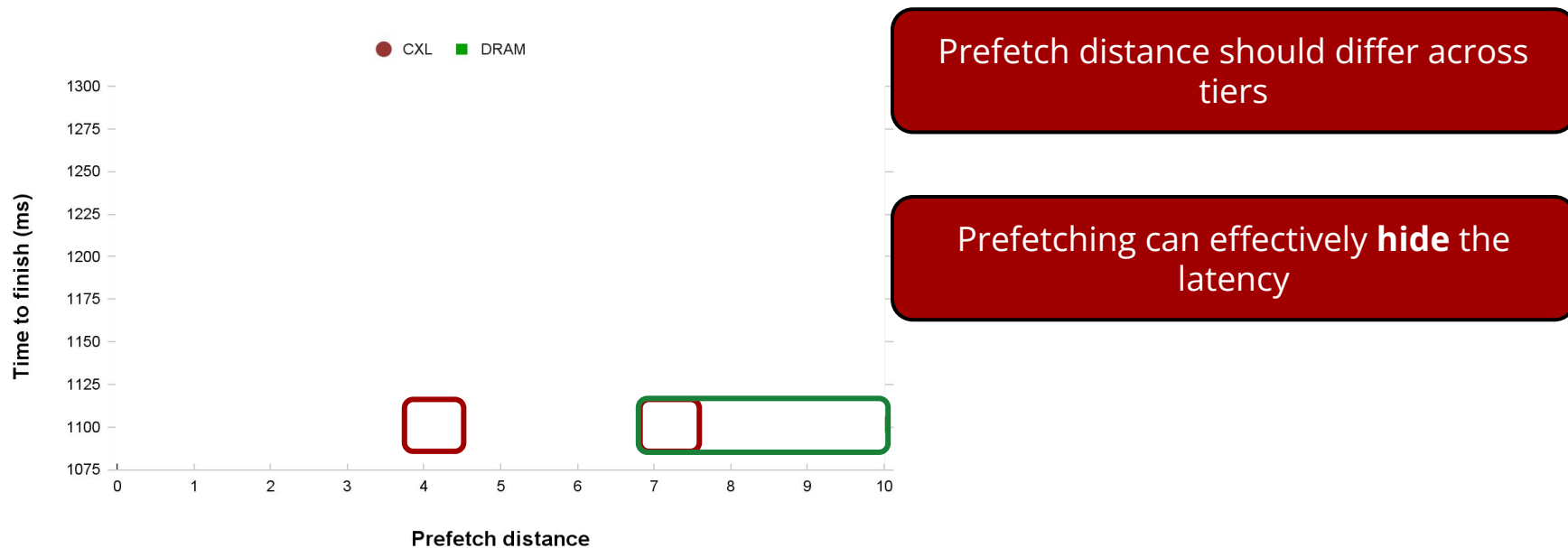
```
for i in range(1M):
  e = records[x[i]]
  prefetch(records[x[i+ k]])
  work(e)
```
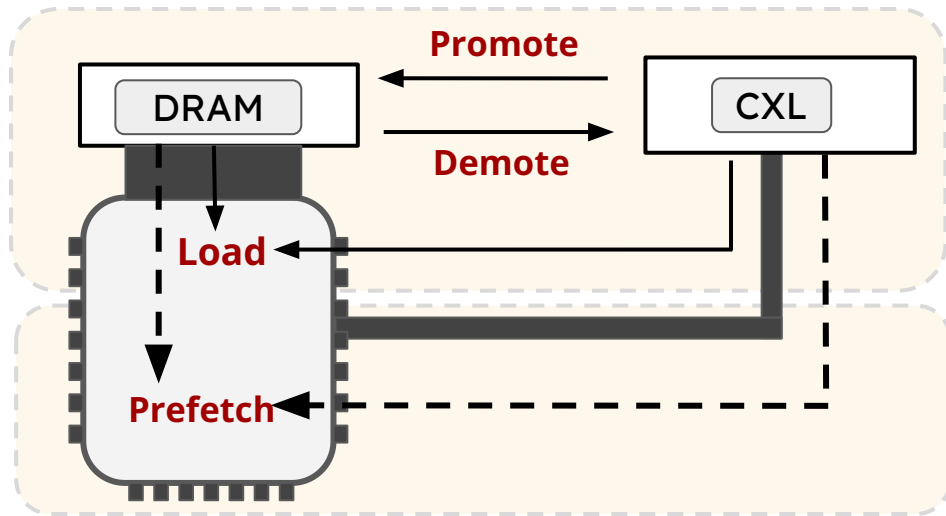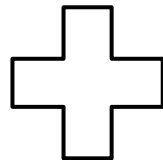
**Prefetch distance**

# Prefetch distance matters

Scan microbenchmark



Prefetch distance should differ across tiers

Prefetching can effectively **hide** the latency

# Linden: Combining reduce and tolerate

# Prefetchability as a tolerate metric

Can we identify memory regions that will get benefit from prefetching ?
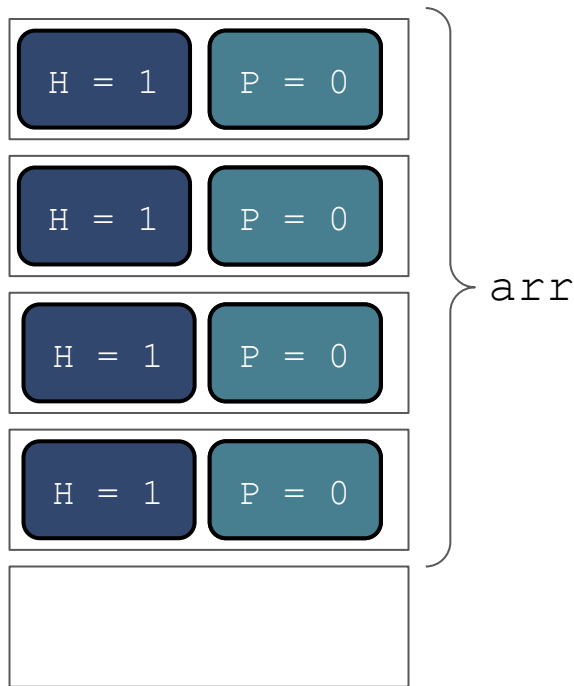
**Prefetchability**:

Percentage of the loads to particular region that can be hidden by prefetching.

$$\frac{\textit{Number of loads hidden with prefetching}}{\textit{Issued loads}}$$

| Thread | Region | Access Pattern | Prefetchability Ratio |
|--------|--------|----------------|----------------------|
| T_1 | arr | Sequential | 1 |
| T_2 | arr_2 | Random | 0 |

# Combining hotness and prefetchability



Memory Pages

```
// rand access

arr = [...]

for j in range(10):

    for i in range(1000):

        sum += arr[rand]
```

H=0-1 Hotness

P=0-1 Prefetchability

# Linden's objectives and workflow

**Prefetching can hide CXL's latency.**

**Under high contentions prefetching can be harmful**

**Prefetch distances matter in different tiers**

**Compiler**

Identify prefetchable access patterns

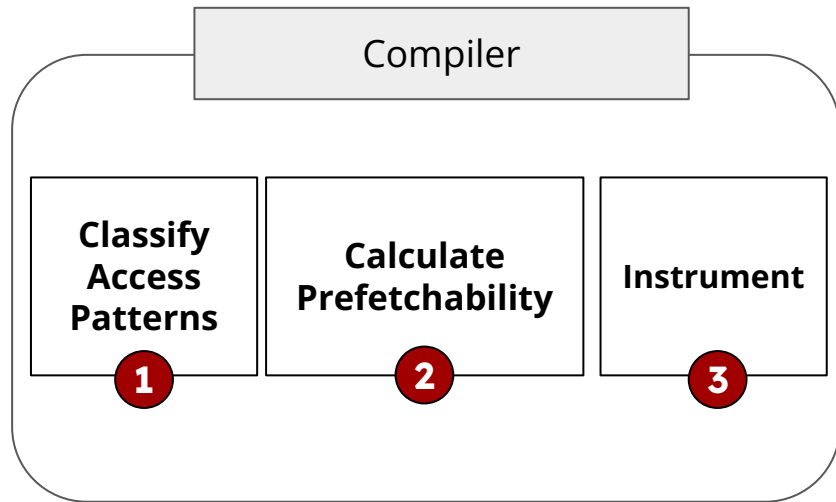Insert software prefetch hints.

**Runtime**

Observe dynamic behaviours

Adapts to the current state
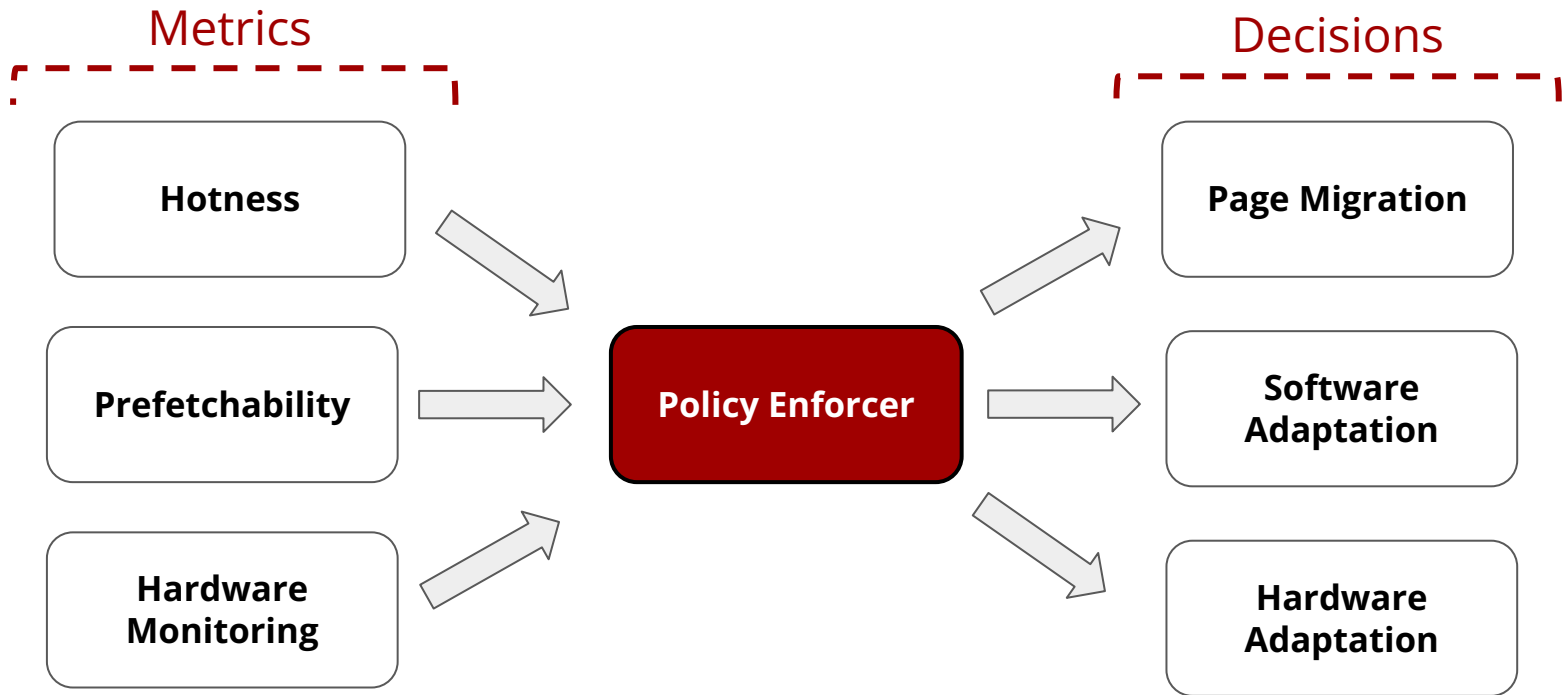
# Linden compiler

**Unmodified Program**

```
arr = [...]

for i in range(1000):

    sum += arr[i]
```

**Compiler**

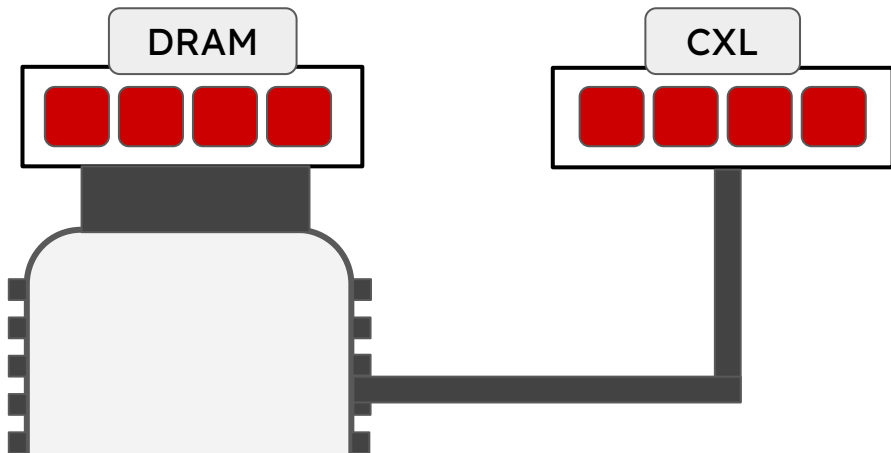| Classify Access Patterns ① | Calculate Prefetchability ② | Instrument ③ |

**Instrumented Program**

```
arr = [...]

register_region(arr, pf_rat);

for i in range(1000):

    prefetch(arr[i+pf_dist])

    sum += arr[i]
```

# Linden's runtime responsibilities

**Hotness**

**Prefetchability**

**Hardware Monitoring**

**Policy Enforcer**

**Page Migration**

**Software Adaptation**

**Hardware Adaptation**
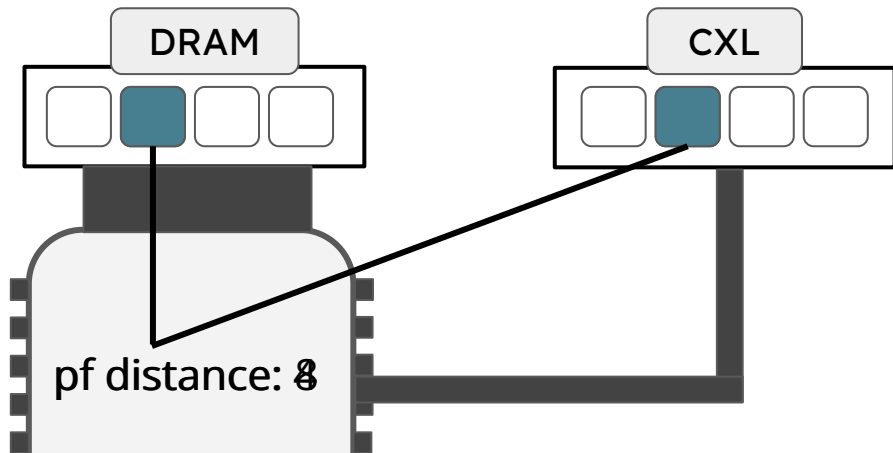
16

# Page migration

What happens if all regions are "hot" ?



- *Hotness + Prefetchability*
- Find memory regions that can **tolerate** the latency.
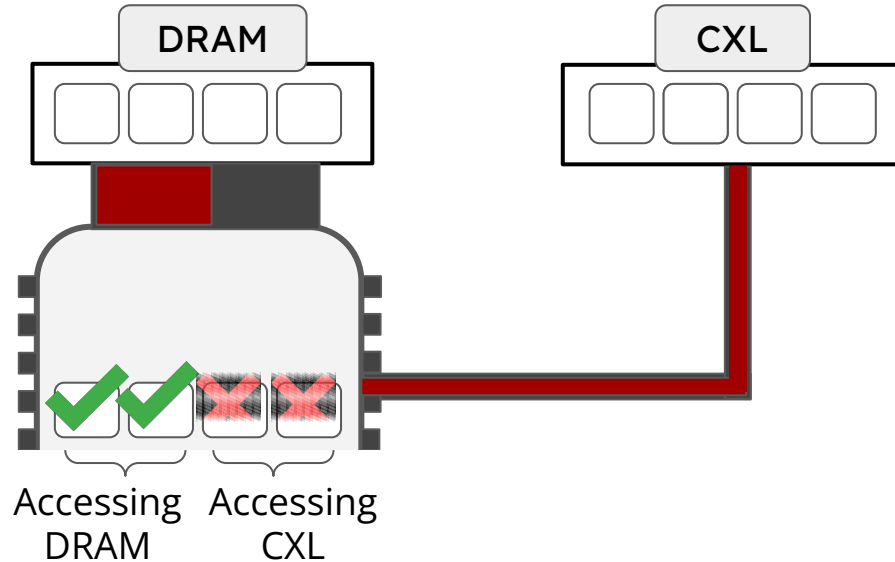- Then migrate these pages to the CXL.

# Software adaptation

What if a region with prefetch migrated?



- Prefetch distance should be updated!
- Dynamically adjust the prefetch distances with considering the tier information.

# Hardware adaptation

What happens under high contention.



Accessing DRAM    Accessing CXL

- Disable hardware prefetching!
- Check which threads are accessing the CXL memory.

# Linden: Tolerate It if You Cannot Reduce It

- Two techniques for tiered memory management:
  - Reduce understands the hotness.
  - Tolerate understands the prefetchability.

- Why toleration matters ?
  - Can help to hide CXL's latency
  - Careful monitoring is essential
  - Hardware and software adaptations needed

## Thanks