

# CS698Y: Modern Memory Systems Lecture-10 (Cache Hierarchy)

# **Biswabandan Panda**

biswap@cse.iitk.ac.in

https://www.cse.iitk.ac.in/users/biswap/CS698Y.html

#### Logistics

#### Quiz 1.0: Sept. 23, 13:00 to 15:00 hrs

## P.A. 1: Sept. 10, 17:00 hrs, code + report @Canvas

#### Presentation: Sept. 14, 10:30 hrs, (9 + 1) mins

#### **Insights: Why and how of what? For example:**

What: astar, Hawkeye performs x% worse than LRU. Why: x% of blocks have feature y, which is not captured by Hawkeye. How: We can augment with a blackbox that does some black magic to improve Hawkeye.

#### **Cache Hierarchy**



## **Revisiting Cache Hierarchy (The Organization)**



Modern Memory Systems

#### **Cache Hierarchies**

L1 Inclusive: L1 cache blocks are present in L3 LLC Exclusive: L1 cache blocks are absent in L3 Non-inclusive: L1 cache blocks may/may not be in L3

#### **Inclusive Hierarchy**



BackInval

#### **Non-inclusive**



#### **Exclusive Hierarchy**



## **Back-invalidation [TLA, MICRO '10]**



**Reference 'e' misses and evicts 'a' from hierarchy** 

Next Reference to 'a' misses 😕

## **On-chip Traffic?**



#### **Non-Inclusive Hierarchy**

**Exclusive Hierarchy** 

#### **Shared Last-level Cache in Multi-core**



#### **Application Behavior @LLC**



#### **Application Behavior @LLC**



Biswabandan Panda, CSE@IITK

#### **Shared LLC**

Shared LLC provides a good tradeoff for all kinds of apps.

Space unutilized by one app. can be utilized by other apps.

However bandwidth is an issue  $\mathfrak{S}$ 

1000 monkeys: one banana

#### **Sliced Shared LLC**



## Intel's Sandybridge [Courtesy: Intel]



#### Intel Xeon [Courtesy: Intel]



#### **IBM POWER 4**



Biswabandan Panda, CSE@IITK

#### **Read More about**

Ring, bus, and crossbar

## Non-uniform Cache Architecture [ASPLOS '02]



#### Bigger and slower caches into smaller and faster slices.

## Shared to Distributed [ISCA '09]



#### **STATIC and DYNAMIC NUCA**

Static NUCA: Fixed mapping of cache lines into banks. Sets spread over banks.

Long wire latency to detect a cache hit.

Dynamic NUCA: Dynamic mapping of cache lines into banks. Ways spread over banks.

Blocks migrate from one bank to another. Move frequently used cache lines closer to CPU

#### Simple Mapping – Bank sets to Banks



# Fair Mapping



# **Shared Mapping**



#### How to search a block?

Incremental Search: starting from closest bank till you find the block or miss.

Multicast search: Send requests to few or all banks of a bank set.

Limited search: Send requests to M banks and search in parallel and then incremental search.

#### What about Cache Access Time

$$T (Hit) = T_{search} + T_{bank_{hit}} + T_{transer}$$

What about address mapping ? Instead of sets – bank sets and bank id

For many-core, T(hit) can become > DRAM-access