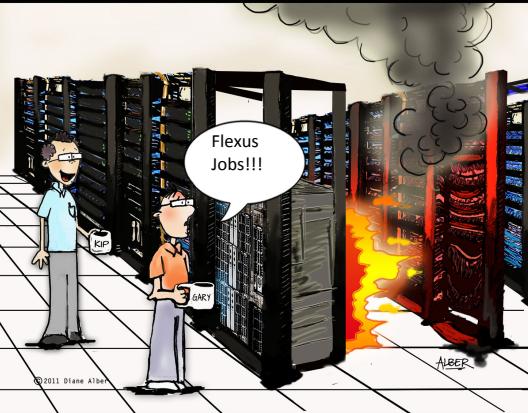


**ACA**

**Servers**

**Fall 2016**



**Pejman Lotfi-Kamran**

Adapted from slides originally developed by Profs. Hill, Hoe, Falsafi and Wenisch of CMU, EPFL, Michigan, Wisconsin

Lec.23 - Slide 1

## Where Are We?

- This Lecture
  - Server
- Next Lecture:
  - Data Center

Fr	Sa	Su	Mo	Tu
27-Shahrivar		29-Shahrivar		
3-Mehr		5-Mehr		
10-Mehr		12-Mehr		
17-Mehr		19-Mehr		
24-Mehr		26-Mehr		
1-Aban		3-Aban		
8-Aban		10-Aban		
15-Aban		17-Aban		
22-Aban		24-Aban		
29-Aban		1-Azar		
6-Azar		8-Azar		
13-Azar		15-Azar		
20-Azar		22-Azar		
27-Azar		29-Azar		
4-Dey		6-Dey		

Fall 2016

Lec.23 - Slide 2

### Modern Computing



Database systems      Web servers

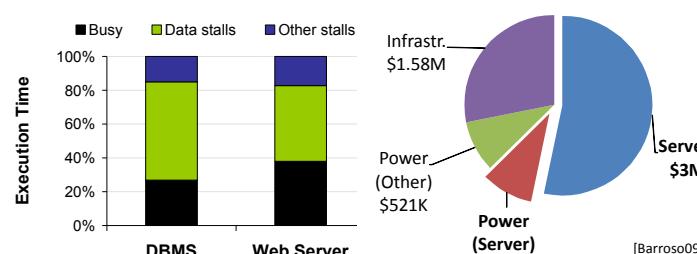
► Server applications: the backbone of computing

Fall 2016

Lec.23 - Slide 3

### But, Servers Are Wasting \$\$\$

Inefficient computation      Servers are expensive



Execution Time

■ Busy    ■ Data stalls    ■ Other stalls

System	Busy (%)	Data stalls (%)	Other stalls (%)
DBMS	~20	~45	~35
Web Server	~35	~35	~30

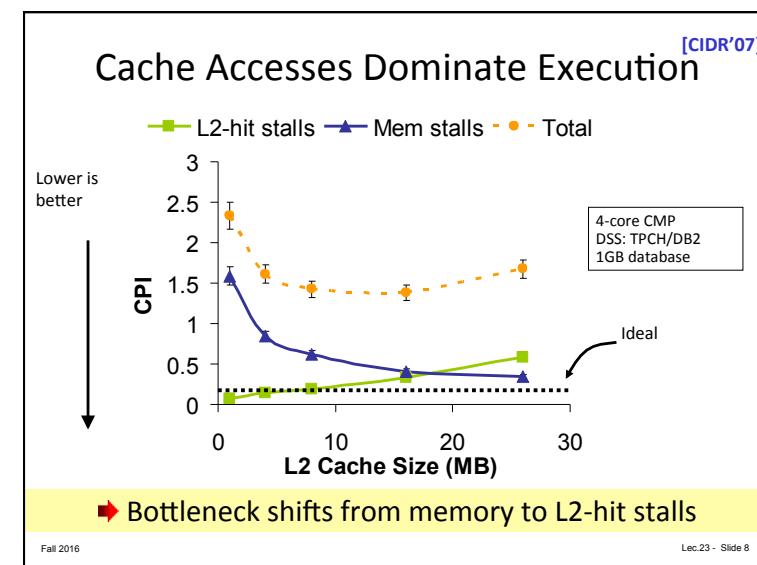
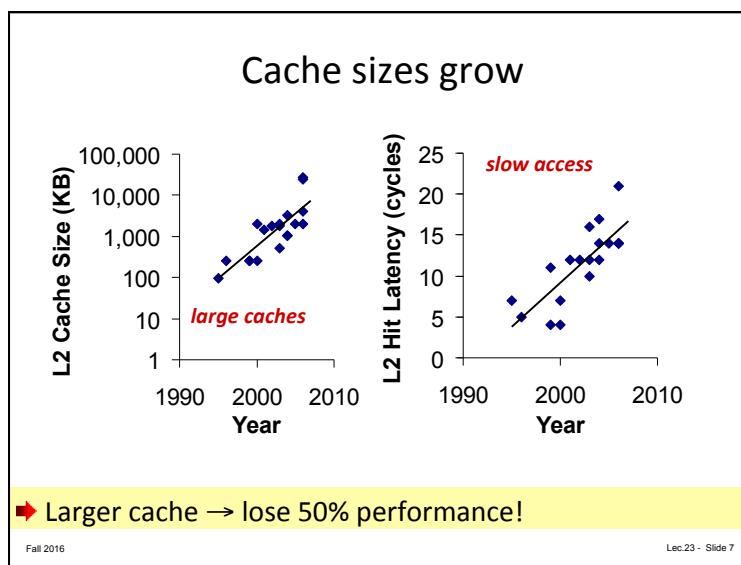
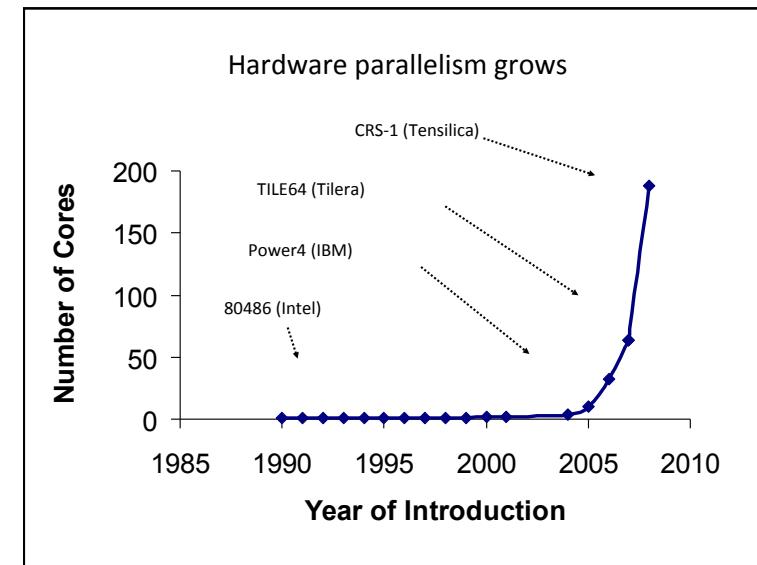
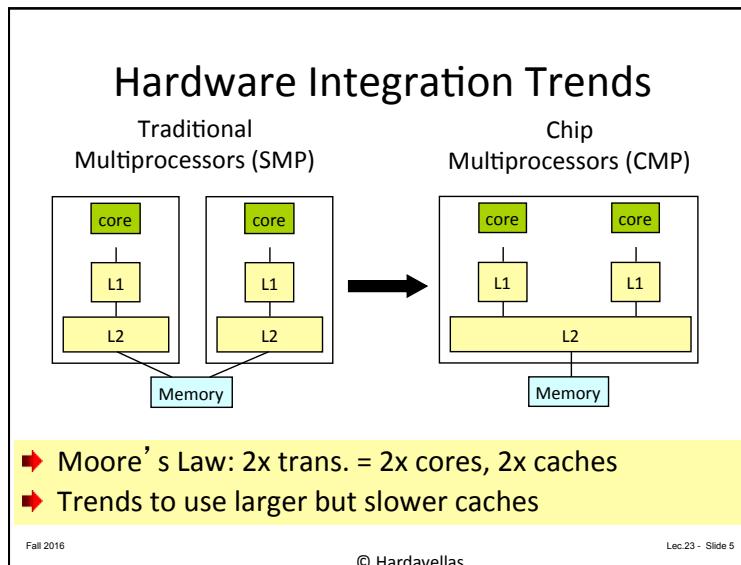
Infrast. \$1.58M  
Power (Server) \$521K  
Power (Other) \$521K  
Servers \$3M

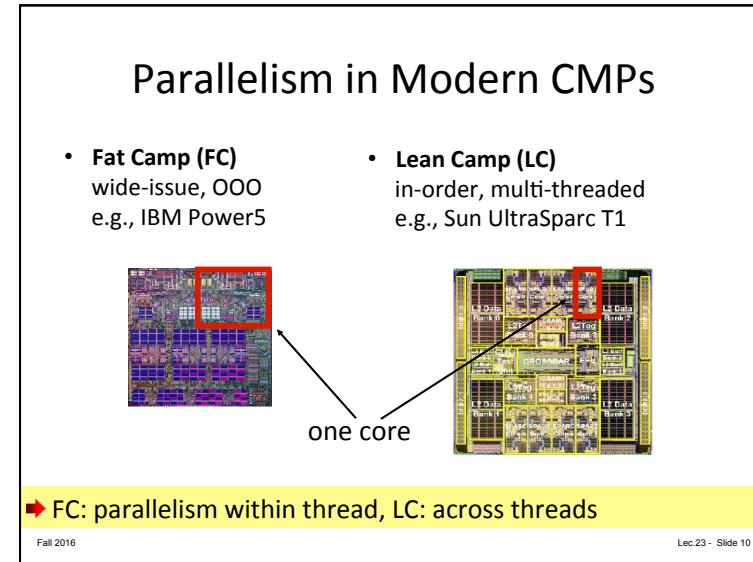
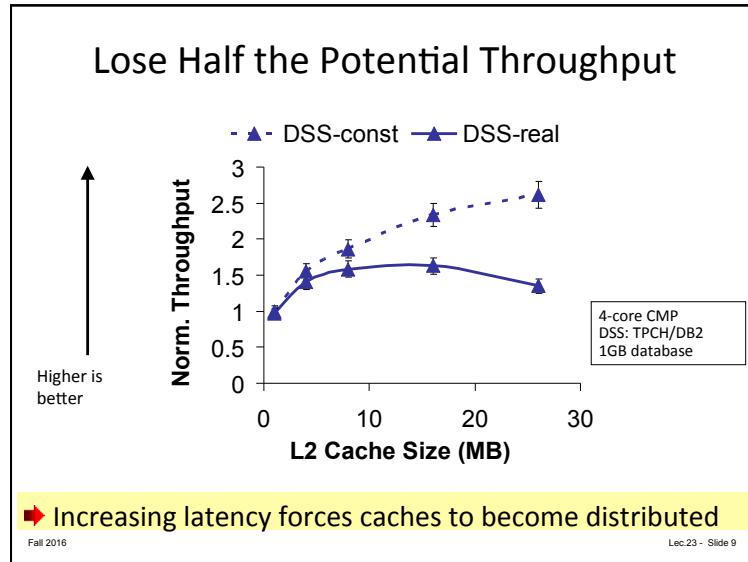
[Barroso09] [Hamilton09]

► Must optimize processors for server workloads

Fall 2016

Lec.23 - Slide 4

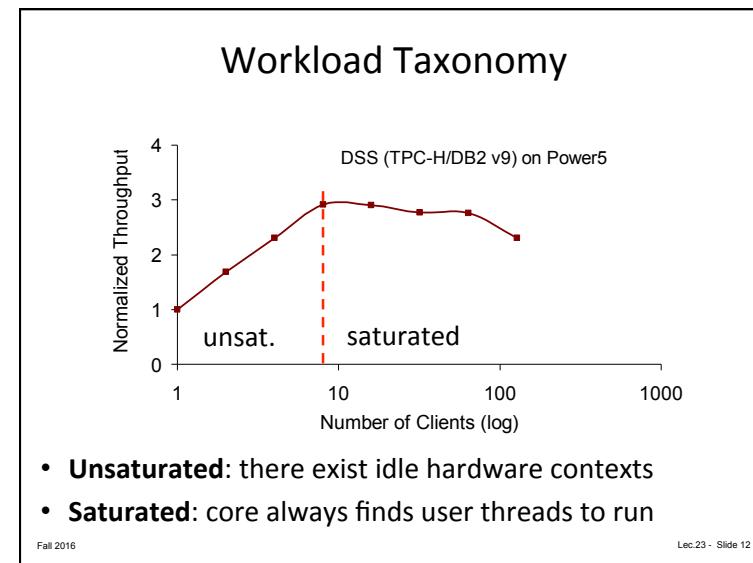


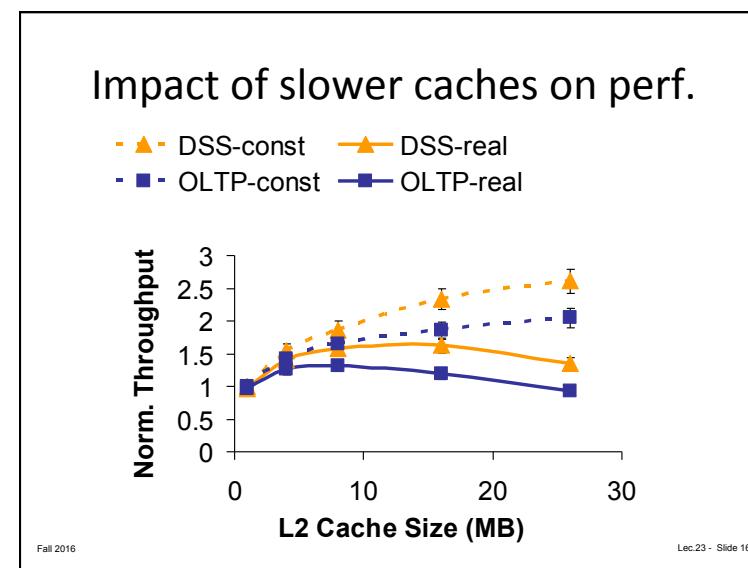
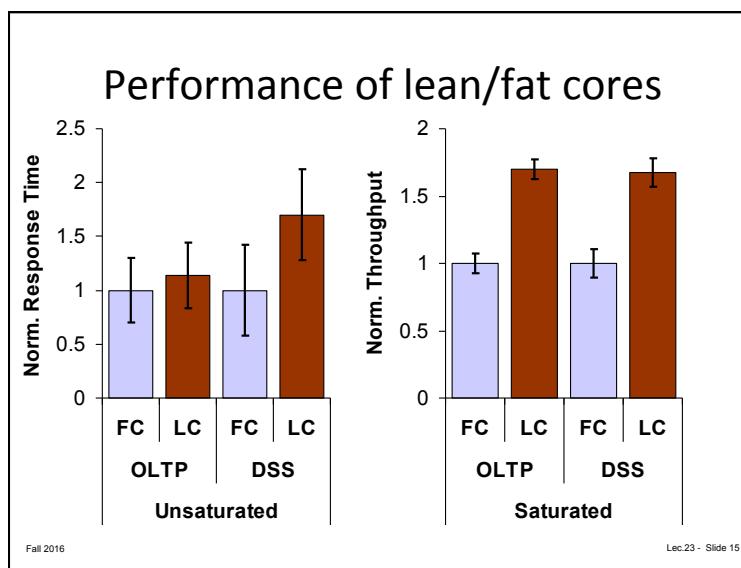
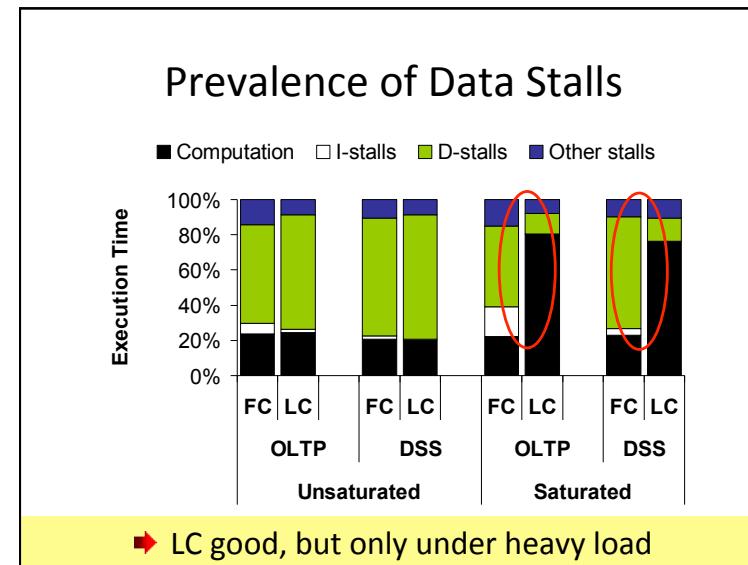
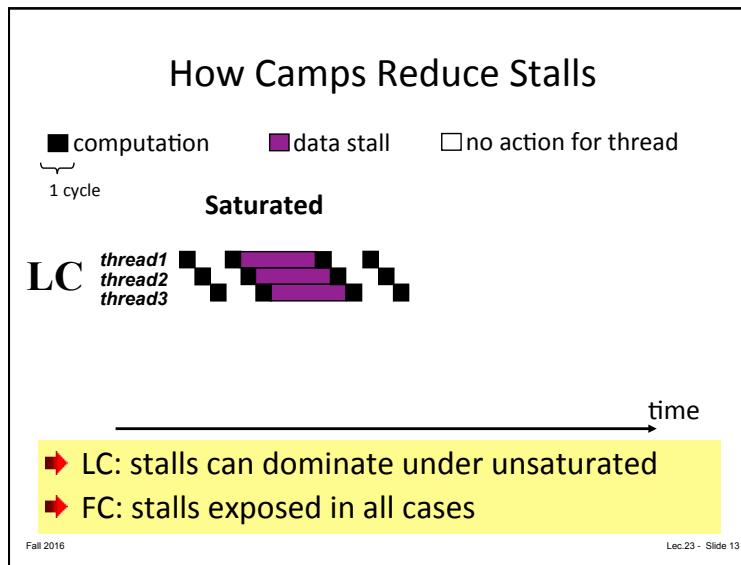


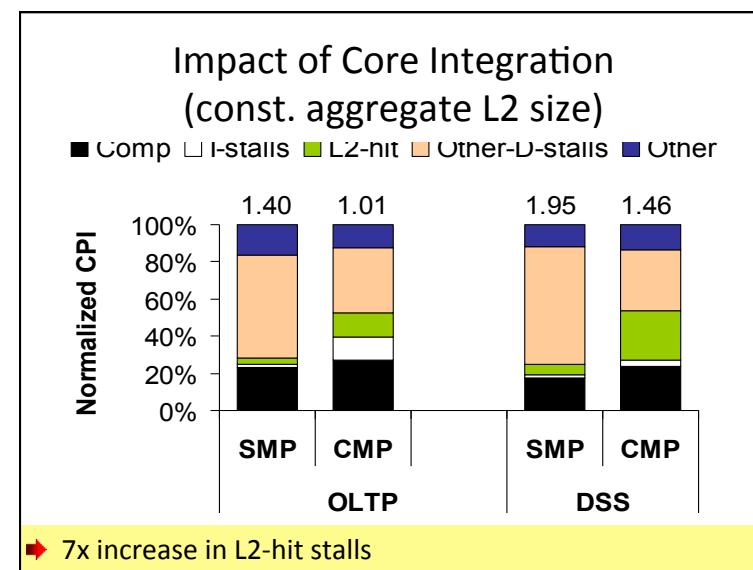
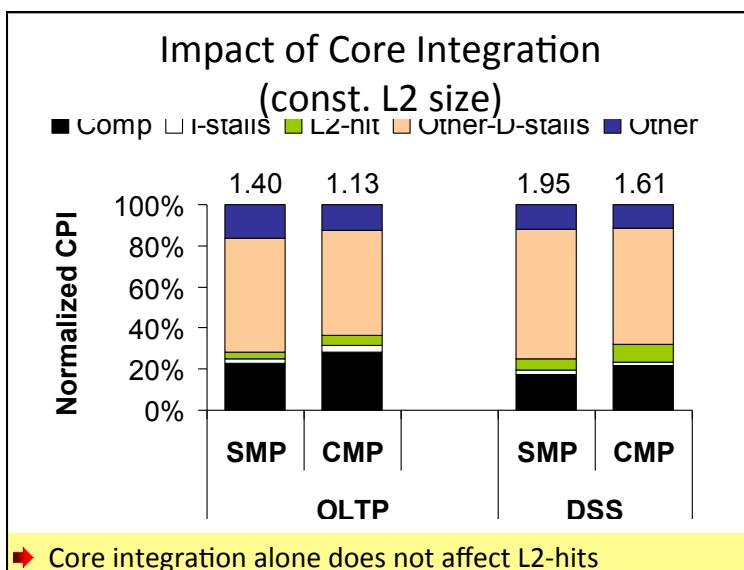
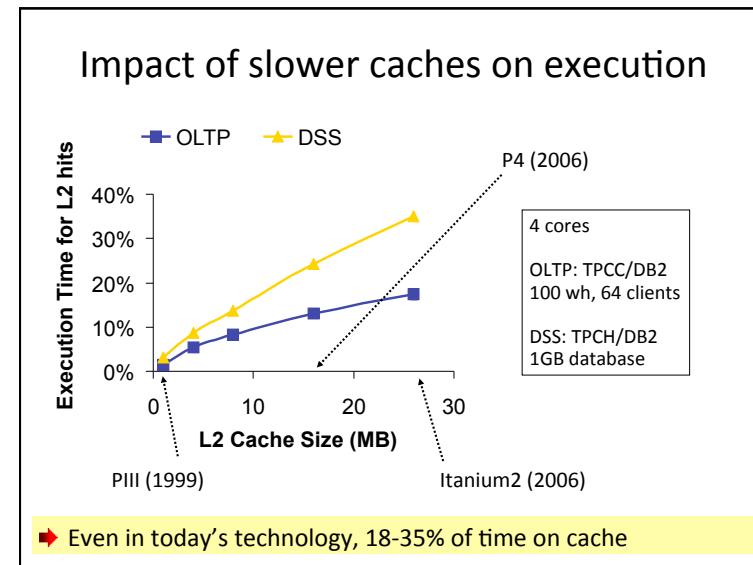
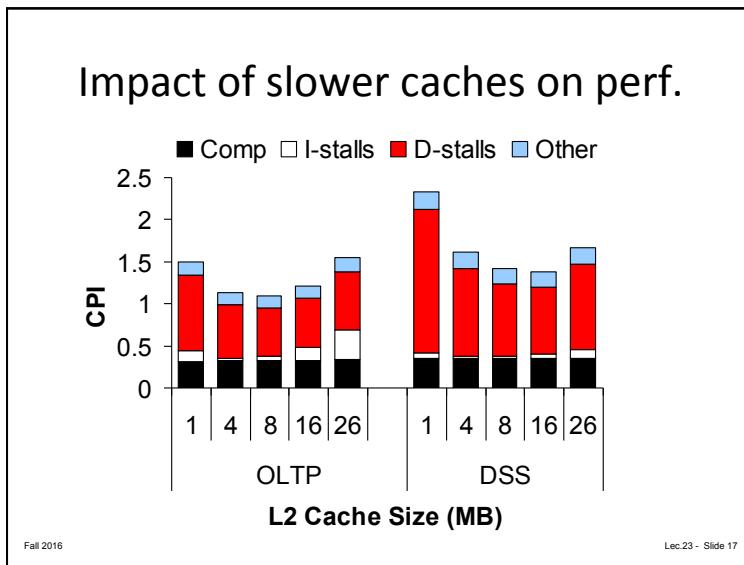
### Processor Designs Taxonomy

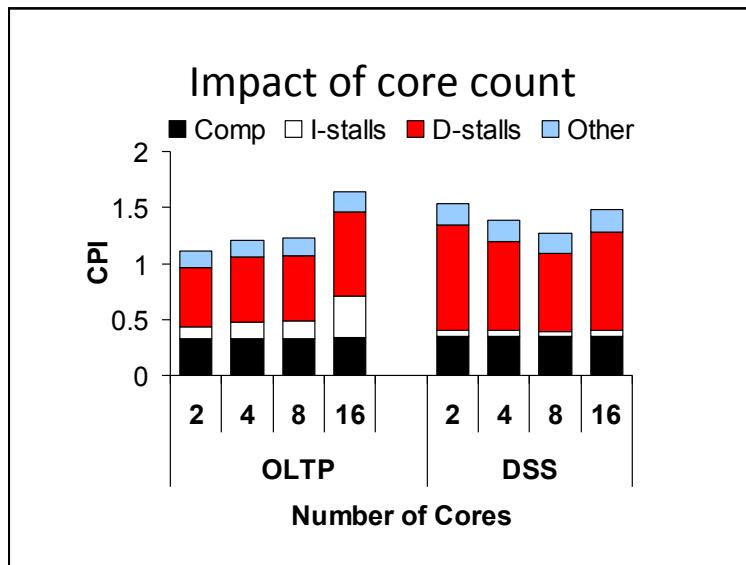
Technology	Fat Camp (FC)	Lean Camp (LC)
Issue Width	Wide (4+)	Narrow (1-2)
Execution Order	Out-of-Order	In-order
Pipeline Size	Deep (14+)	Shallow (5-6)
HW Threads	Few (1-2)	Many (4+)
Core Size	Large (3xLCsize)	Small (LCsize)

Fall 2016 Lec.23 - Slide 11

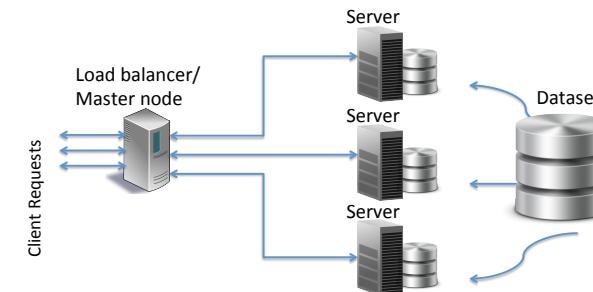








### Characteristics of Scale-Out Apps



- Many independent requests/tasks
- Huge dataset split into shards
- Minimal communication among servers

Fall 2016

Lec.23 - Slide 22

### Scale-Out Datacenter Internals

The diagram illustrates the process of populating a datacenter with servers. It starts with a single server icon, followed by an arrow labeled "Populate a datacenter with servers", and ends with a large rack of servers in a datacenter aisle.

- Tens of thousands of servers
- Large memory per server
- High efficiency = high utilization

**Maximize parallelism for efficiency**

Fall 2016

Lec.23 - Slide 23

### How Efficient are Today's Servers?

[“Clearing the Clouds”, ASPLOS ‘12]

- Created benchmark suite
  - Diverse set of cloud workloads
  - Quantified high-level behavior
- Studied off-the-shelf hardware
  - Used performance counters
  - Identified needs of cloud apps



**Modern CPUs don't match needs of cloud apps**

## Cloud Suite 1.0

(released @ [parsa.epfl.ch/cloudsuite](http://parsa.epfl.ch/cloudsuite))

The dashboard displays the following services:

- Data Serving:** Cassandra NoSQL (Icon: Facebook logo, Cassandra logo)
- MapReduce:** Machine learning on Hadoop (Icon: Mahout logo)
- Linux 2.6.32:** CentOS (Icon: Penguin logo)
- Media Streaming:** Apple Quicktime Server (Icon: Apple logo, Quicktime logo)
- SAT Solver:** Symbolic VM constraint solver (Icon: LVM logo)
- Web Frontend:** Nginx, PHP server (Icon: NGINX logo, PHP logo)
- Web Search:** Apache Nutch (Icon: Nutch logo)

Fall 2016      Lec.23 - Slide 25

## Hardware

The hardware components shown are:

- Dell PowerEdge M1000e
- Dell Blades M610
- Two Intel x5670 2.9GHz 6-core, 12MB LLC (Processor and RAM highlighted with a red circle)
- 24GB RAM

Fall 2016      Lec.23 - Slide 26

## Methodology: “Server-grade” CPU

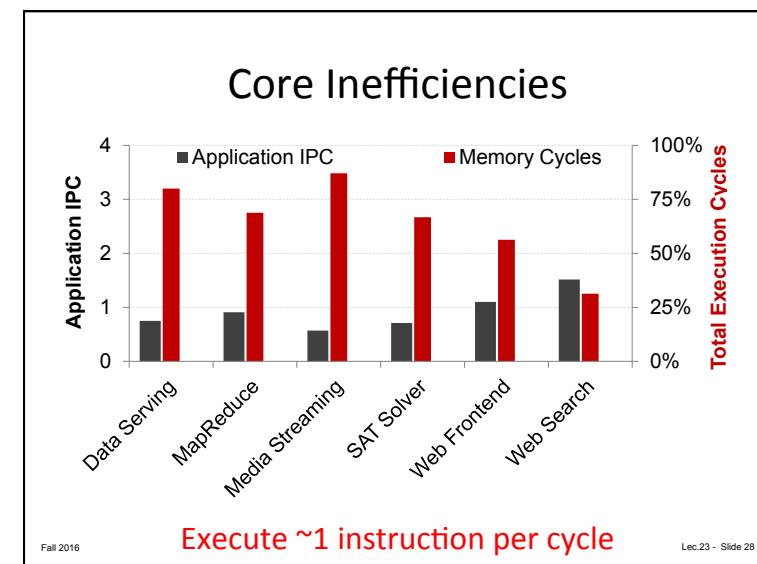
- Aggressive OoO cores
  - Run up to 4 instructions per cycle
- Large instruction window
  - 128 instructions in flight
  - 48 loads in flight
- L2 and large L3 caches
- Vast off-chip b/w

The diagram illustrates the architecture of a "Server-grade" CPU, showing the following layers from inner to outer:

- Core:** Contains **Exec. Units** (represented by four downward-pointing arrows) and an **Inst. Window** (represented by a vertical stack of horizontal bars).
- Memory Accesses:** Handles **L1-I** and **L1-D** requests.
- L2:** A 256KB cache layer.
- L3:** A 12MB cache layer.
- Off-chip bandwidth:** 32GB/s.

Timing metrics shown: 4 cycles for L1 access, 10 cycles for L2 access, 39 cycles for L3 access, and 200+ cycles for off-chip access.

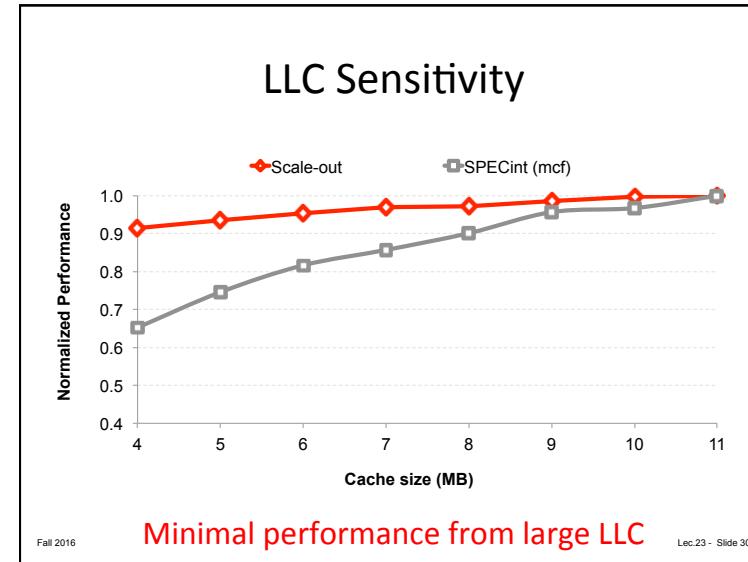
Fall 2016      Lec.23 - Slide 27



## Core Inefficiencies

- Underutilized complexity
- Scale-out requirements low
  - couple parallel memory ops.
  - one execution unit

Fall 2016 Lec.23 - Slide 29



## LLC and Bandwidth Inefficiencies

- Scale-out needs modest LLC
  - Beyond 3-4MB useless
  - Area & latency w/o payoff
- Low per-core BW needs
  - <15% utilization
  - Too many channels
  - Too high frequency

Fall 2016 Lec.23 - Slide 31

