# ACA

# Servers

**Fall 2016**

Flexus Jobs!!!

©2011 Diane Alber
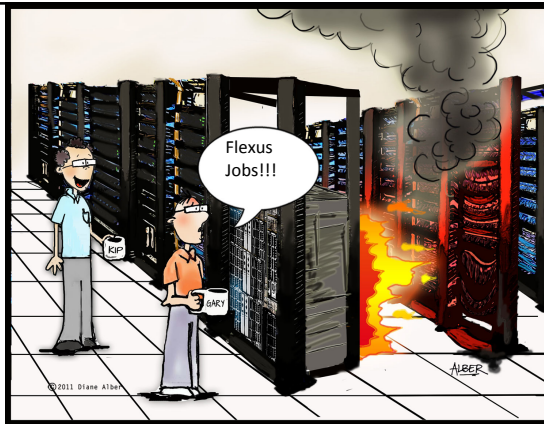
**Pejman Lotfi-Kamran**

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

Lec.24 - Slide 1

---

# Where Are We?

| 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 |  |

- This Lecture
  – Server

- Next Lecture:
  – Data Center

Fall 2016

Lec.24 - Slide 2

---

# What do Scale-Out Apps Need? [ISCA 2012] [IEEE Micro 2012]

Cores share instructions
– Large code footprint fits in LLC

requests

Data is in memory
– Data footprint dwarfs LLC

Code access

Data access

Server

Cores don't communicate
– Independent requests
– Mostly Read-only accesses

Memory

Fall 2016

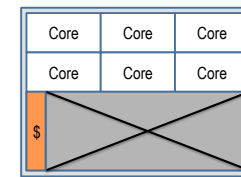**Poor match to existing server processors** Lec.24 - Slide 3

---

# Conventional Processors: Major Mismatch

✘ Few fat cores
– Limited on-chip parallelism

| Core | Core | Core |
|------|------|------|
| Core | Core | Core |

$

✘ Large LLC
– Dwarfed by vast datasets
– Slow to access
– Takes up much die area

e.g., *Intel Westmere*

**Large cores & LLC wasteful on scale-out apps**

---

1

## Small-chip Processors: Not Enough

✘ Few lean cores
– Low parallelism

✘ Modestly-sized cache
– Fast to access, but…
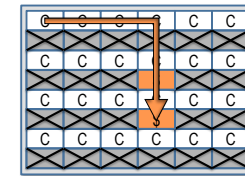–    takes up large fraction of die area

e.g., *Calxeda*

## Emerging Tiled Processors: Not Optimal

✓ Many lean cores
– High parallelism

✘ Large distributed cache
– Dwarfed by vast datasets
– Slow to access
– Takes up much die area

e.g., *Tilera Tile64*

✘ Tiled organization
– Distance hurts performance

**Large LLC & distance prevent scalability**

## Outline

• Why not today's processors?
• What do scale-out workloads need?
• Scale-Out Processors
– Overview
– Design considerations

## Roadmap for Scale-Out Processor

1. Eliminate wasteful capacity in the LLC
– Shrink the cache, add cores in freed space

2. Reduce delay to LLC
– Partition the die into independent multi-core *pods*

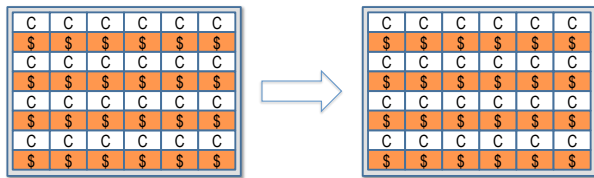3. Enable technology scalability
– Do not interconnect the pods

## Step 1: Less cache, more cores

✓ Small LLC
  – Capture instructions
  – More area for the cores



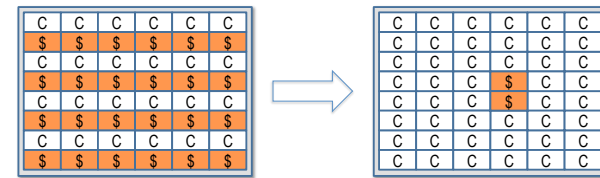Fall 2016                                           Lec.24 - Slide 9

## Step 1: Less cache, more cores

✓ Small LLC
  – Capture instructions
  – More area for the cores
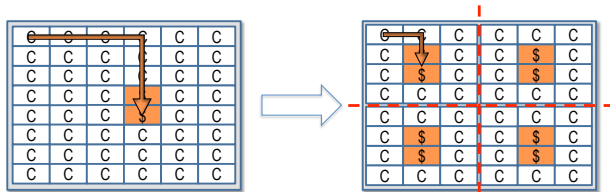✓ More cores for higher throughput



Fall 2016                                           Lec.24 - Slide 10

## Step 2: Form "pods" to reduce distance

✓ Small LLC
✓ More cores for higher throughput
✗ Fast path to LLC for instructions



How to choose the optimal pod?

Fall 2016                                           Lec.24 - Slide 11

## Sizing Pods: Not So Simple

Too few cores ➜ limited parallelism

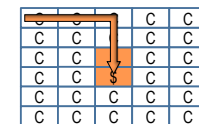Too many cores ➜ long distance to LLC
  – Slower instruction fetch

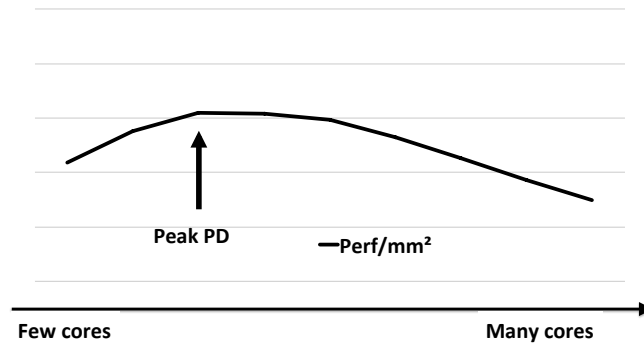*Cache dominates die area*                          *Distance hurts performance*



Optimal pod?

**Few cores**                                       **Many cores**

How do we characterize optimality?

Fall 2016                                           Lec.24 - Slide 12

## Our Optimization Criterion: Performance Density (PD)



**Peak PD**  —**Perf/mm²**

**Few cores**                    **Many cores**

Fall 2016

PD pinpoints optimal use of silicon
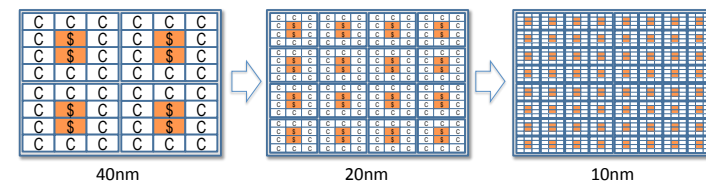
Lec.24 - Slide 13

## Step 3: Scale-Out Processors

One or more pods

Each pod is a standalone server
– Runs a full software stack

No inter-pod connectivity or coherence
– Enhances scalability



40nm          20nm          10nm

Fall 2016

Inherently optimal & scalable

Lec.24 - Slide 14

## Methodology

Flexus simulation infrastructure [Wenisch '06]
CMP analytic model [Hardavellas '09]
TCO model [Hardy '11]

Core Parameters
– Out-of-order, 3-way, 2 GHz
– L1 (I & D): 32KB, 2-way

Chip components
– Technology node: 20 nm
– Core: 1.1 mm²
– 1MB cache: 1.2 mm²
– Mem channel: 12 mm²

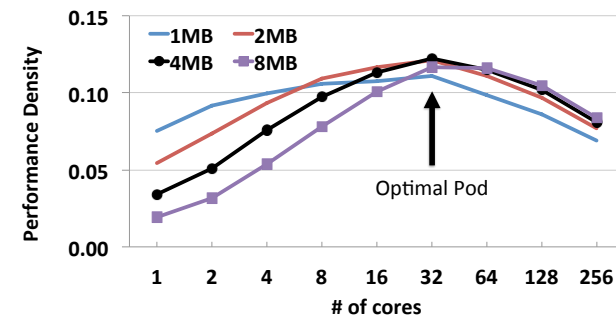Chip budget
– Die area: 280 mm²
– Power: 95W
– BW: 6 x 3.2 GT/s DDR4

Datacenter
– Power: 20 MW
– Rack: 42 U, 17 kW
– Server: 64 GB, N sockets

Fall 2016
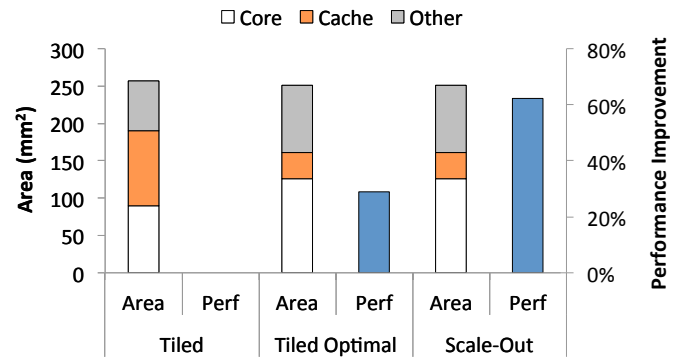
Lec.24 - Slide 15

## Pod Design Space Exploration



Optimal Pod

Fall 2016

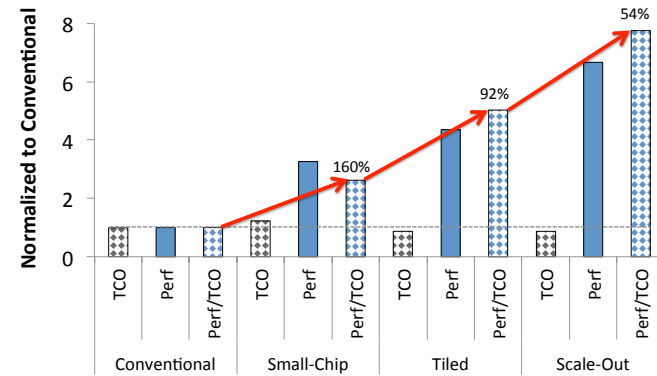Optimal Pod: 32 cores & 4MB of cache

Lec.24 - Slide 16

## Processor Efficiency



Core  Cache  Other

Area (mm²) vs Performance Improvement

Tiled | Tiled Optimal | Scale-Out
(Area / Perf for each)

Fall 2016

**SOP: highest performance per mm²**

Lec.24 - Slide 17

## Datacenter Efficiency



Normalized to Conventional

Conventional | Small-Chip | Tiled | Scale-Out
(TCO / Perf / Perf/TCO for each)

160%    92%    54%

Fall 2016

**SOP: highest performance/TCO**

Lec.24 - Slide 18

## Outline

- Why not today's processors?
- What do scale-out workloads need?
- Scale-Out Processors
  - Overview
  - Design considerations
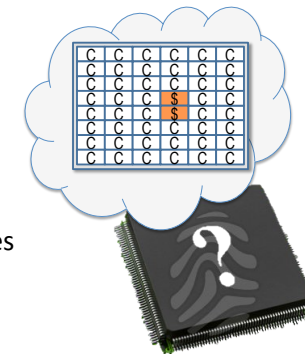
Fall 2016

Lec.24 - Slide 19

## How do we architect a pod? [MICRO 2012]

Pod characteristics:
- Several MB of LLC
- Many cores

Challenge:
- Efficiently connect the pieces
- Maximize PD



Fall 2016

Lec.24 - Slide 20

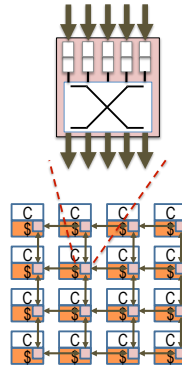## Tiled Pod

Mesh NOC
- Nearest-neighbor connectivity
- 5-ported routers

Distributed LLC

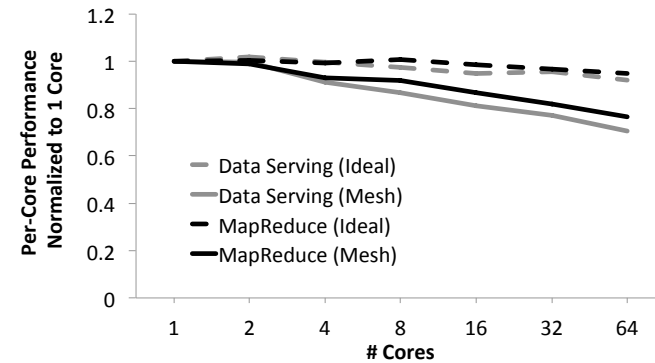✓ Low complexity

✓ Small NOC area

✗ High LLC access latency



Low cost, low performance

## Performance: Tiled Pod vs Ideal



Legend:
- Data Serving (Ideal)
- Data Serving (Mesh)
- MapReduce (Ideal)
- MapReduce (Mesh)

Y-axis: Per-Core Performance Normalized to 1 Core

X-axis: # Cores (1, 2, 4, 8, 16, 32, 64)

## Fast Tiled Pod

Flattened Butterfly NOC *Kim, MICRO '07*
- Rich connectivity
- Many-ported routers

Distributed LLC

✗ High complexity
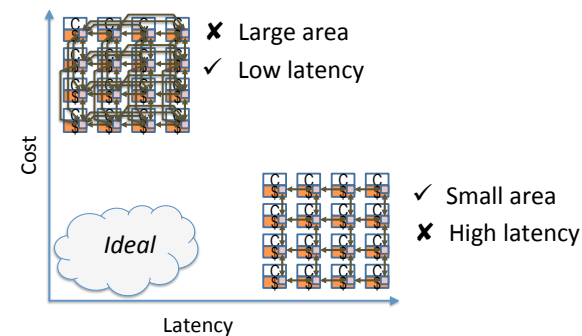
✗ Large NOC area (7x of mesh!)

✓ Low LLC access latency



row   column

column   row

$k^2/2$ links per row/column

High cost, high performance

## Off-the-shelf designs not ideal



✗ Large area

✓ Low latency

✓ Small area

✗ High latency

*Ideal*

Cost

Latency

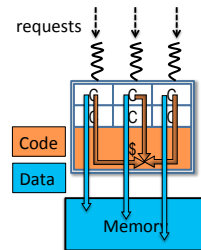Exploit workload characteristics for efficiency

## How do threads access data?

**Where is the code and data?**
- Code: in LLC
- Data: in memory
- ➢Nothing useful in remote L1s!

**Data access pattern:** *bilateral*
- L1 ◀▶ LLC
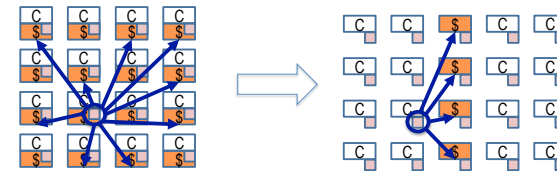- L1 ✖ L1

requests

Code
Data
Memory

Fall 2016 — Lec.24 - Slide 25

## What is the actual traffic pattern?

**Distributed LLC: all-to-all traffic**
- Mismatch with underlying bilateral access pattern
- Requires expensive connectivity for low latency

**Turn all-to-all traffic into bilateral to reduce cost!**

## Roadmap for an "Ideal" Pod

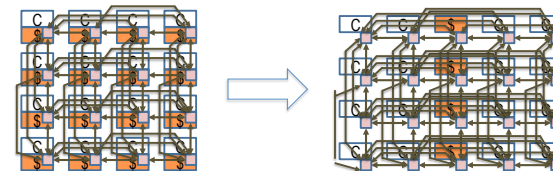Starting point: fast tiled pod

1. LLC in the center
   - Decouple cores and LLC banks
2. Remove links from flattened butterfly
   - Leverage the *bilateral* traffic pattern to lower cost
3. Share and specialize core-to-LLC interconnect

Fall 2016 — Lec.24 - Slide 27

## Step 1: LLC in the Center

✓ Decouple core and LLC tiles
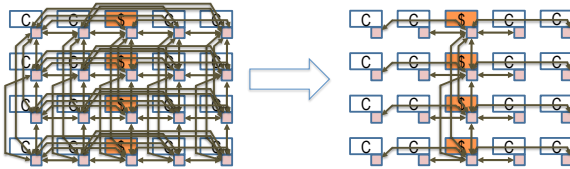   - Natural fit for the bilateral traffic pattern

Fall 2016 — Lec.24 - Slide 28

7

## Step 2: Remove unneeded links

✓ Core-to-core connectivity not needed
  – Less cost, same performance
✓ Rich intra-LLC connectivity helps performance
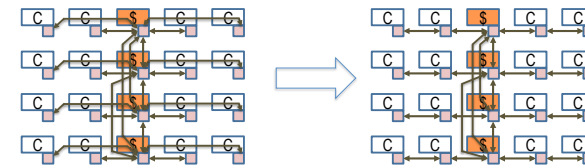  – Expense limited to a fraction of the die

## Step 3: Share links and specialize

✓ Cores share links to LLC
  – Dedicated core-to-LLC links have poor cost/perf
✓ Specialize *request*/*reply* (to/from LLC) networks
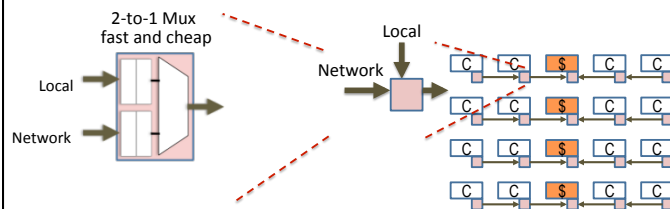  – Maximize cost/perf

## NOC-Out: near-ideal pod design

*Request* network
  – Tree topology
  – Each node: 2-to-1 flow-controlled mux



2-to-1 Mux
fast and cheap
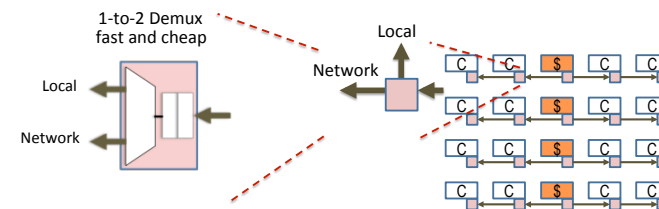
Local

Network

Local

Network

## NOC-Out: near-ideal pod design

*Reply* network
  – Tree topology
  – Each node: 1-to-2 flow-controlled demux



1-to-2 Demux
fast and cheap
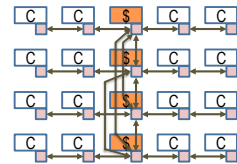
Local

Network

Local

Network

## NOC-Out: near-ideal pod design

LLC network: flattened butterfly topology
– Expense limited to a fraction of the die
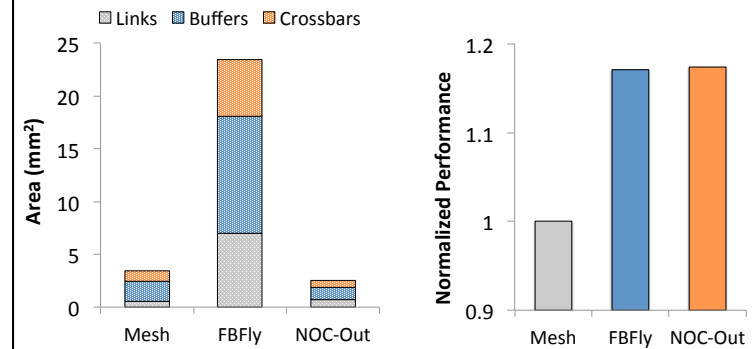
*Request* & *Reply* networks: tree topology
– Limited connectivity for efficiency
– SW stack unaffected



Fall 2016

Lec.24 - Slide 33

## Area Breakdown



NOC-Out: FBfly's performance at 1/10$^{th}$ cost

## Summary

Scaling Trends
– More lean cores
– Beware of larger caches -> slow you down

Scale-out workloads
– Lots of request-level parallelism
– Large instruction footprints
– Little core-to-core communication

Chip organization & design
– Can improve Performance/TCO

Fall 2016

Lec.24 - Slide 35