

February 25, 2025

Announcements

- paper responses: due
Monday
- talk slides coll & perm.
to publish

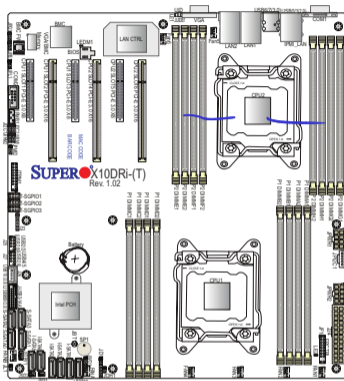
Goals

- no class Tue Mar 4
- ←

Review

- multi processor memory

Topology and NUMA



[SuperMicro Inc. '15]

Demo:

- ▶ Show `lstopo` on porter, from [hwloc](#).
- ▶ [lstopo on MI300](#)

Machine (31GB total)

Package L#0

NUMANode L#0 P#0 (31GB)

L3 (12MB)

L2 (1280KB)

L2 (1280KB)

L2 (2048KB)

L2 (2048KB)

L1d (48KB)

L1d (48KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1d (32KB)

L1i (32KB)

L1i (32KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

L1i (64KB)

Core L#0

Core L#1

Core L#2

Core L#3

Core L#4

Core L#5

Core L#6

Core L#7

Core L#8

Core L#9

PU L#0

P#0

PU L#2

P#2

PU L#4

P#4

PU L#5

P#5

PU L#6

P#6

PU L#7

P#7

PU L#8

P#8

PU L#9

P#9

PU L#10

P#10

PU L#11

P#11

PU L#1

P#1

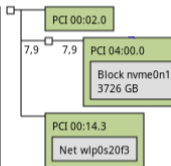
PU L#3

P#3

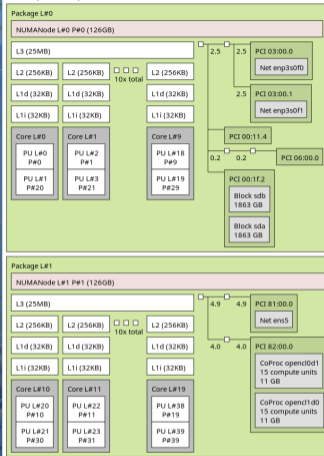
SMT

Host: arc

Date: Di 25 Feb 2025 11:14:47 CST



Machine (252GB total)



Host: koelsch

Date: Tue 25 Feb 2025 11:16:51 AM CST

Placement and Pinning

Who decides on what core my code runs? How?

OMP_PLACES = cores
pthread - set - affinity - up

Who decides on what NUMA node memory is allocated?

First touch, libnuma

[Demo: intro/NUMA and Bandwidths](#)

What is the main expense in NUMA?

Cache Coherence

What is *cache coherence*?

A large, empty rounded rectangular box with a thin black border, intended for the user to provide an answer to the question above.

How is cache coherence implemented?

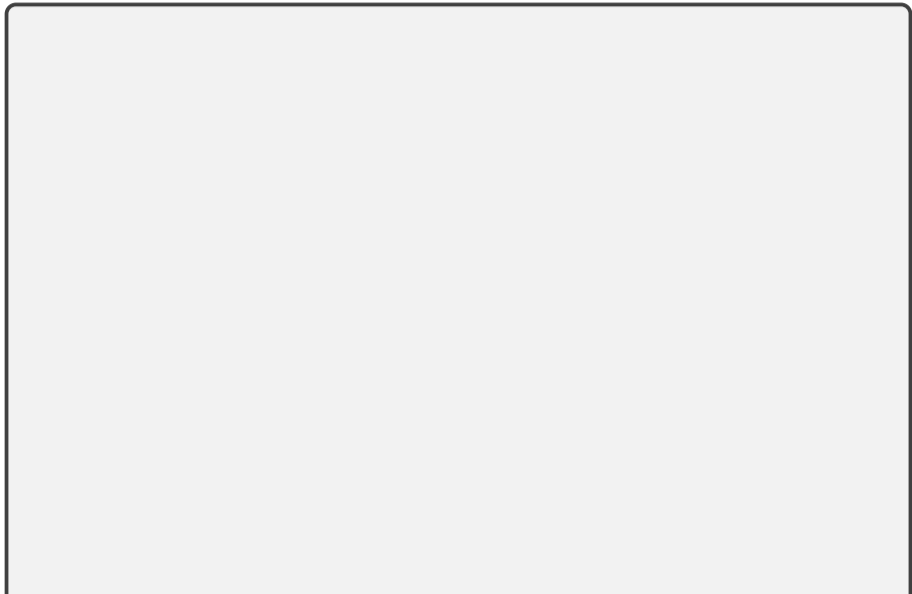
A large, empty rounded rectangular box with a thin black border, intended for the user to provide an answer to the question above.

[Cache coherence simulation](#)

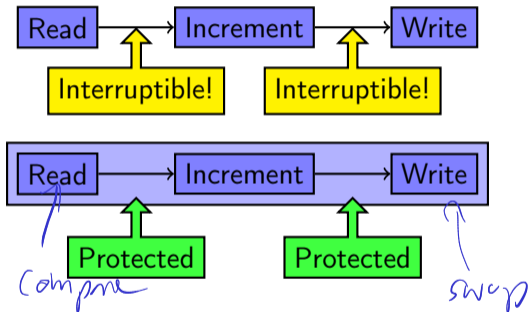
What are the performance impacts?

- ▶ [Demo: intro/Threads vs Cache](#)
- ▶ [Demo: intro/Lock Contention](#)

MESI Sequence to Try



'Conventional' vs Atomic Memory Update



Outline

Introduction

Machine Abstractions

C

OpenCL/CUDA

Convergence, Differences in Machine Mapping

Lower-Level Abstractions: SPIR-V, PTX

Performance: Expectation, Experiment, Observation

Performance-Oriented Languages and Abstractions

Polyhedral Representation and Transformation

Outline

Introduction

Machine Abstractions

C

OpenCL/CUDA

Convergence, Differences in Machine Mapping

Lower-Level Abstractions: SPIR-V, PTX

Performance: Expectation, Experiment, Observation

Performance-Oriented Languages and Abstractions

Polyhedral Representation and Transformation

Atomic Operations: Compare-and-Swap

```
#include <stdatomic.h>
_Bool atomic_compare_exchange_strong(
    volatile A* obj, C* expected, C desired );
```

What does `volatile` mean?

What does this do?

How might you use this to implement atomic FP multiplication?

read, op, CAS

Memory Ordering

Why is Memory Ordering a Problem?

e.g. lock vs compiler/processor reordering

What are the different memory orders and what do they mean?

- seq_cst

- acquire : no r/w in this thread reorder before
all releasing writes elsewhere are visible!

- release : no r/w in this thread reorder after
all writes here are visible to acquirers elsewhere

Example: A Semaphore With Atomics

```
#include <stdatomic.h> // mo_→memory_order, a_→atomic
typedef struct { atomic_int v;} naive_sem_t;
void sem_down(naive_sem_t *s)
{
    while (1) {
        while (a_load_explicit(&(s->v), mo_____) < 1)
            spinloop_body();
        int tmp=a_fetch_add_explicit(&(s->v), -1, mo_____rel);
        if (tmp >= 1)
            break; // we got the lock
        else // undo our attempt
            a_fetch_add_explicit(&(s->v), 1, mo_____);
    }
}
void sem_up(naive_s_t *s) {
    a_fetch_add_explicit(&(s->v), 1, mo_____);
}
```