



(HPC)
High Performance
Coding & Computing
Alberto Dassatti - 2023

WHAT'S NEXT

You now know how to optimize your code.

Bit Hacks, caching, branch prediction and SIMD have no secrets for you.
You know how to write code in such a way your compiler enables the most
advanced optimizations.

WHAT'S NEXT

You now know how to optimize your code.

Bit Hacks, caching, branch prediction and SIMD have no secrets for you. You know how to write code in such a way your compiler enables the most advanced optimizations.

What if all of this is not enough?

WHAT'S NEXT

You now know how to optimize your code.

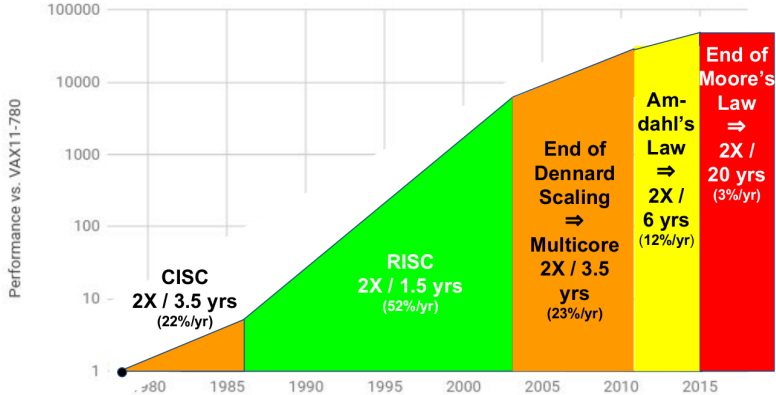
Bit Hacks, caching, branch prediction and SIMD have no secrets for you. You know how to write code in such a way your compiler enables the most advanced optimizations.

What if all of this is not enough?

Parallel hardware is ubiquitous: let's go parallel!

End of Growth of Performance?

40 years of Processor Performance



Based on SPECintCPU. Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

PARALLELISM VS CONCURRENCY

Tasks are parallel when they execute at the same time.

Tasks are concurrent when their execution order is not predetermined.

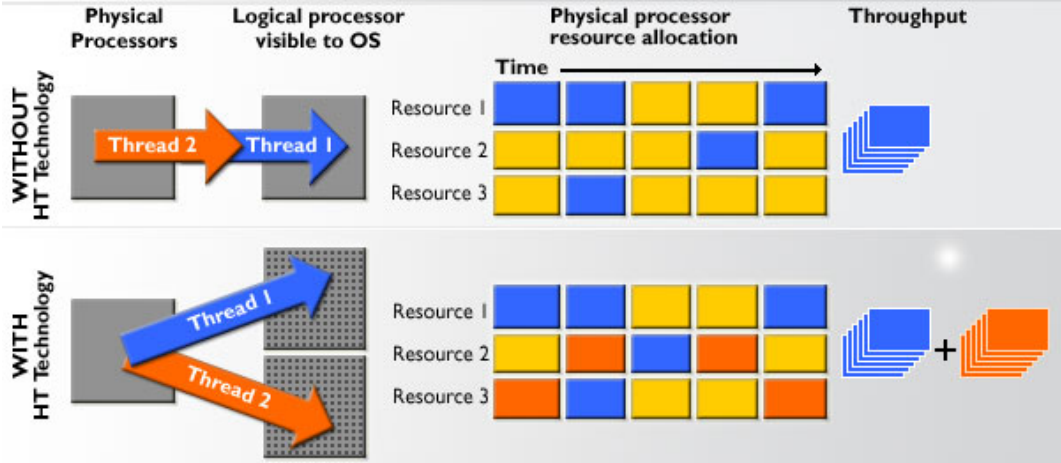
PARALLEL MACHINES

Shared-Memory Hardware
or
Distributed-Memory Hardware

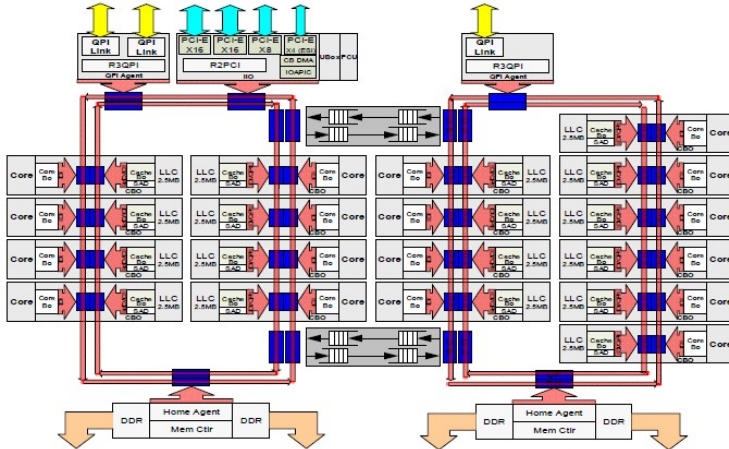
Shared-Memory machines have HW Cache Coherence, the semantic and implementation of which is extremely complex and out of scope for our purposes.

Distributed-Memory machines use Message passing (MPI) to transfer data and distribute job on several nodes.

How Hyper-Threading Technology Works



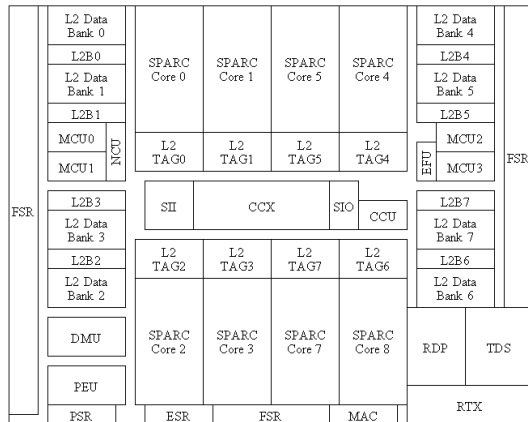
XEON E7-8890 v4



up to 24 cores.

[source]

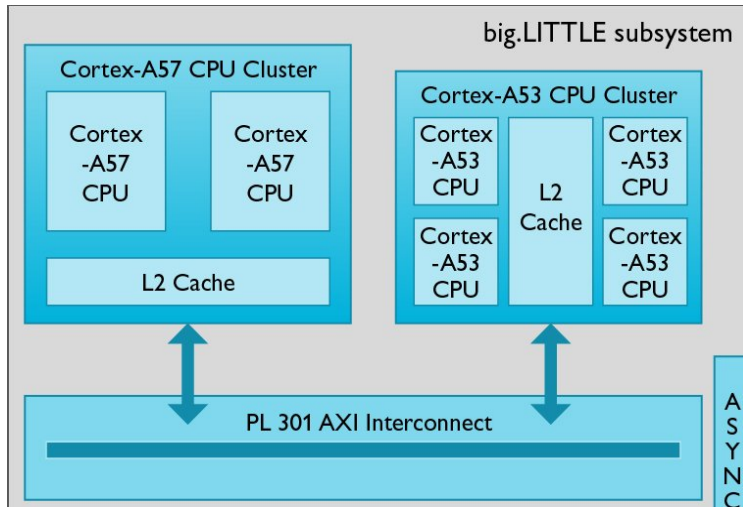
ULTRASPARC T2



Niagra 2 / UltraSPARC T2 / OpenSPARC T2 - Die Micrograph Diagram (davidhalko)

up to 8 cores, 8 threads each.

CORTEX A72

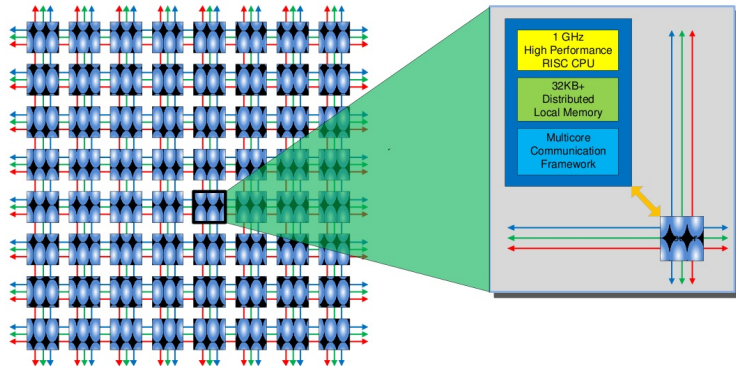


up to 8 cores.

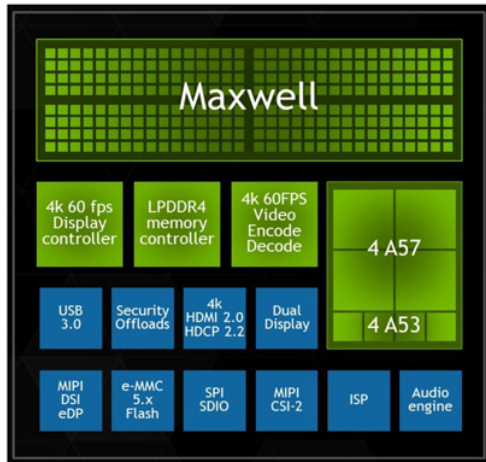
PARALLELA BOARD



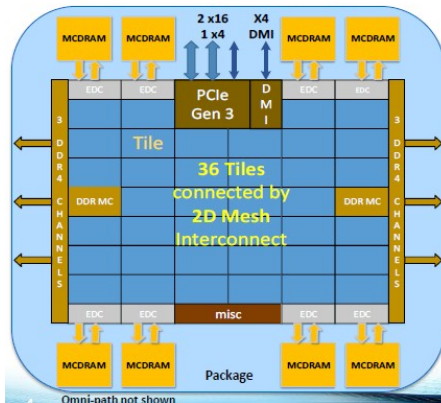
PARALLELA BOARD



NVIDIA X1



Knights Landing Overview



TILE

2 VPU	CHA	2 VPU
Core	1MB L2	Core

Chip: 36 Tiles interconnected by **2D Mesh**

Tile: 2 Cores + 2 VPU/core + 1 MB L2

Memory: MCDRAM: 16 GB on-package; High BW

DDR4: 6 channels @ 2400 up to 384GB

IO: 36 lanes PCIe Gen3. 4 lanes of DMI for chipset

Node: 1-Socket only

Fabric: Omni-Path on-package (not shown)

Vector Peak Perf: 3+TF DP and 6+TF SP Flops

Scalar Perf: ~3x over Knights Corner

Streams Triad (GB/s): MCDRAM : 400+; DDR: 90+

Source Intel: All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. Binary Compatible with Intel Xeon processors using Haswell Instruction Set Extension (Xeon Phi). Bandwidth numbers are based on STREAM-like memory access pattern when used in conjunction with Intel Memory. Results have been estimated based on internal Intel analysis and are not intended for commercial purposes only. Any difference in system configuration or software settings may affect actual performance.

HARDWARE

There are many parallel machines and many more are coming, but the question is still there:

HARDWARE

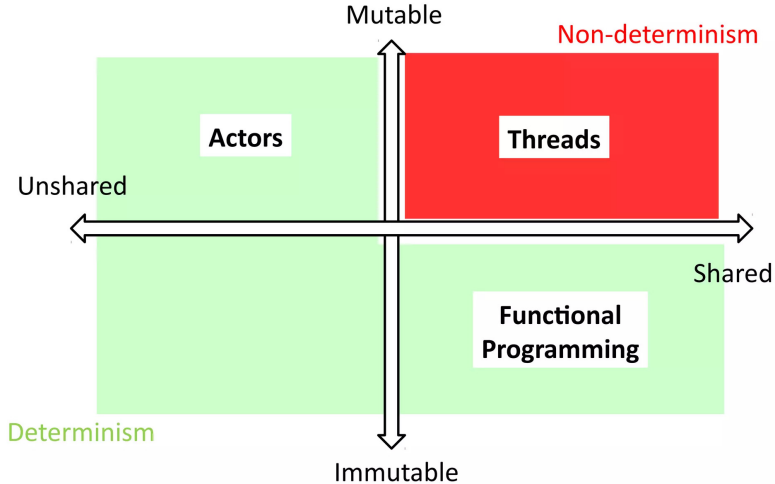
There are many parallel machines and many more are coming, but the question is
still there:
how we program them?

HARDWARE

There are many parallel machines and many more are coming, but the question is
still there:
how we program them?

You already know `pthread`s. Can they be the answer? Why?

The state quadrants



PARALLELISM

Using parallelism means doing more tasks (work) in parallel.

Main difficulties:

1. data dependency (linked list traversal).
2. synchronization (mutex, semaphores, data races, deadlocks, etc).
3. problem formulation (a less performant serial algorithm may be easier to parallelize).
4. compilers (and in general tools) do not help much (as of now).
5. memory sub-system is critical (SMP, NUMA, rDMA).
6. DSL and languages can have a better impact on parallel programming, C/C++ do not shine here.

Suggested reading: *"Is Parallel Programming Hard, And, If So, What Can You Do About It?"*, by Paul E. McKenney ([pfd](#)).

PARALLELISM: GOOD NEWS

There are many frameworks (languages, libraries, etc.) that helps the developer.

Our objective is knowing the existing options and making some experience with
(some of) them.

PARALLELISM: GOOD NEWS

- pthread (all low level details are exposed to the programmer).
- OpenMP (#pragma based).
- OpenACC (#pragma based).
- Threading Building Blocks (Intel TBB C++).
- Cilk Plus (fork-join parallelism, now Intel, Tapir MIT open source).
- OpenCL.
- CUDA (only NVIDIA).
- SYCL.
- DSL (break computation and schedule apart)

- MPI (non shared memory machines).

PARALLELISM: BAD NEWS

There are many frameworks (languages, libraries, etc.) that helps the developer,
but

PARALLELISM: BAD NEWS

There are many frameworks (languages, libraries, etc.) that helps the developer,

but

they are incompatible, closed source, quickly evolving or abandoned, complex, miss tooling, are not portable, etc. We will only focus on:

- `pthread`: POSIX compliant, high availability.
- OpenMP: mature and stable, well supported by open tooling.

BEFORE PROGRAMMING SOME THEORY

what are the limits to parallelization?

BEFORE PROGRAMMING SOME THEORY

what are the limits to parallelization?

Most programs have a sequential part and a parallel part.

AMDAHL'S LAW

S : fraction of serial runtime in a serial execution.

P : fraction of parallel runtime in a serial execution.

$$S + P = 1 \quad (1)$$

T_s : time for the program to run in serial

N : number of processors/parallel executions

T_p : time for the program to run in parallel

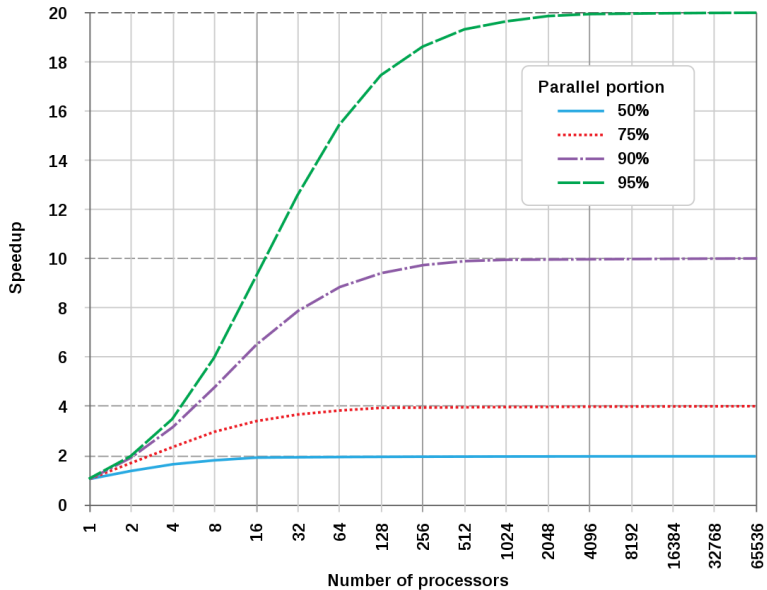
$$T_p = T_s \left(S + \frac{P}{N} \right) \quad (2)$$

AMDAHL'S LAW

If no overhead (never the case, but could be low) is needed to parallelize, the speedup can be expressed as:

$$\begin{aligned} \text{speedup} &= \frac{T_s}{T_p} \\ &= \frac{T_s}{T_s \cdot (S + \frac{P}{N})} \\ &= \frac{1}{S + \frac{P}{N}} \end{aligned}$$

Amdahl's Law



AMDAHL'S LAW

$$\begin{aligned} \text{speedup} &= \frac{T_s}{T_p} \\ &= \frac{T_s}{T_s \cdot (S + \frac{P}{N})} \\ &= \frac{1}{S + \frac{P}{N}} \\ &= \frac{1}{(1 - P) + \frac{P}{N}} \end{aligned}$$

or

$$\text{max_speedup} = \frac{1}{(1 - P)} \quad (3)$$

THREAD & PROCESS

Modern OSes give us two ways of describing parallelism and concurrency: threads and processes.

Process

- Private address space (virtual memory).
- One private stack.
- Private resources.
- High overhead (`fork` and `switch`).

Thread

- Shared address space (parent virtual memory).
- Multiple private stacks.
- Shared resources.
- Low overhead (`pthread_create` and `switch` without MMU involvement).

THREAD

There are three main threading models:

- Hardware Thread (hyperthreading or cores).
- Kernel Thread (POSIX threading model). Available in C++11 and C11, [1].
- Green Thread or user level threading.
 - ▶ Quick context switches, no need for system call.
 - ▶ Cannot use multiple processors, only for concurrency.
 - ▶ GNU Portable Threads, GO, Ruby, etc., [list](#).

Windows 7 use an hybrid model. It has the benefit of the two worlds, but has increased complexity.

Each task (process/thread) can be associated with a set of processors (Processor Affinity, `taskset`).

PTHREAD

Some aspects of `pthread` you may be interested in with respect of performance:

- Detached or joinable state
- Scheduling inheritance
- Scheduling policy
- Scheduling parameters
- Scheduling contention scope
- Stack size
- Stack address
- Stack guard (overflow) size

Make sure the libraries you use are thread-safe! (`glibc` reentrant functions are.)

IO

If you are blocked by *I/O*, may be files, may be sockets, Asynchronous *I/O* could be the answer.

It's not easy to implement correctly and usually the OS does a great job for you, but nobody knows your task better than you do.

```
fd = open (... , O_NONBLOCK );  
read (...); // returns instantly !  
close ( fd );
```

Modern linux includes [io uring](#), go and check it out!

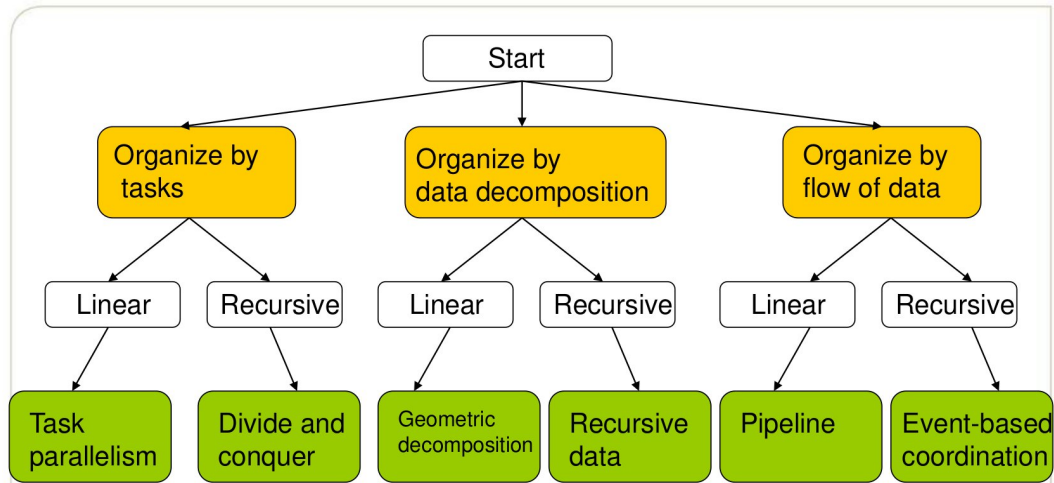
FEW WORDS ON SYNCHRONIZATION

You know what a `mutex` is and how to use it. There are some useful functions if your focus is performance:

- `try-lock`.
- `pthread_spin_lock`: repeatedly try the lock and will not put the thread to sleep.
- if reads can happen in parallel, as long as there's no write, `pthread_rwlock_rdlock` and `pthread_rwlock_wrlock`.
- `pthread_barrier`: Allows you to ensure that (some subset of) a collection of threads all reach the barrier before returning.
- `Atomics!` [1],[2].
- (not really synchronization but) do not forget `volatile`.

But much better than all of them... Lock-Free Algorithms!

PARALLELIZATION PATTERNS



source.

THREAD POOLS

Instead of creating threads, destroying them and recreating them, you can use a thread pool. It creates n threads; you just push work onto them.

Only question is: How many threads should you create? (Experiment to find out). For an example implementation look at [GLib](#).

FOR MORE...

- Memory consistency models.
- Work stealing scheduler.
- Lock free algorithms and data structures.
- Heterogeneous systems.