# Introduction to Parallel Programming with MPI

Dr. Alireza Ghasemi, Dr. Georg Hager

Erlangen National High Performance Computing Center

Odds and Ends – what we have left out

# What we have left out

- **Point-to-point** bells and whistles
  - Persistent communication
  - Message probing: MPI_Probe,…
  - One-sided communication: MPI_Put, MPI_Get, MPI_Accumulate,…
  - Partitioned communication
- **Collectives** bells and whistles
  - MPI_Reduce_scatter, MPI_Scan,…
- **MPI I/O**
- **Virtual topologies**
- **MPI shared memory** communication

# Introduction to Parallel Programming with MPI

Dr. Alireza Ghasemi, Dr. Georg Hager

Erlangen National High Performance Computing Center

Computer Architecture and Performance issues
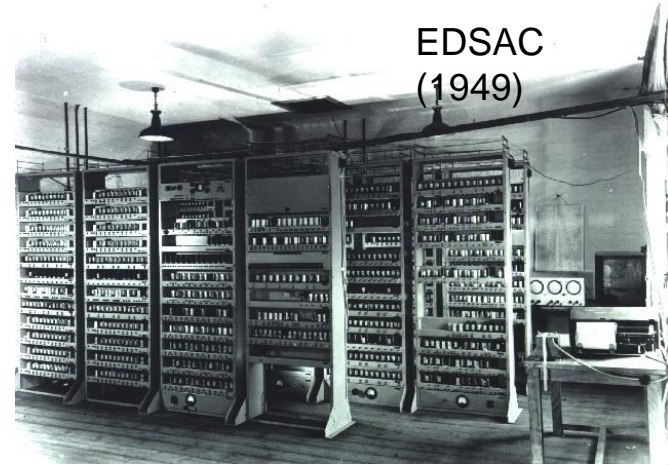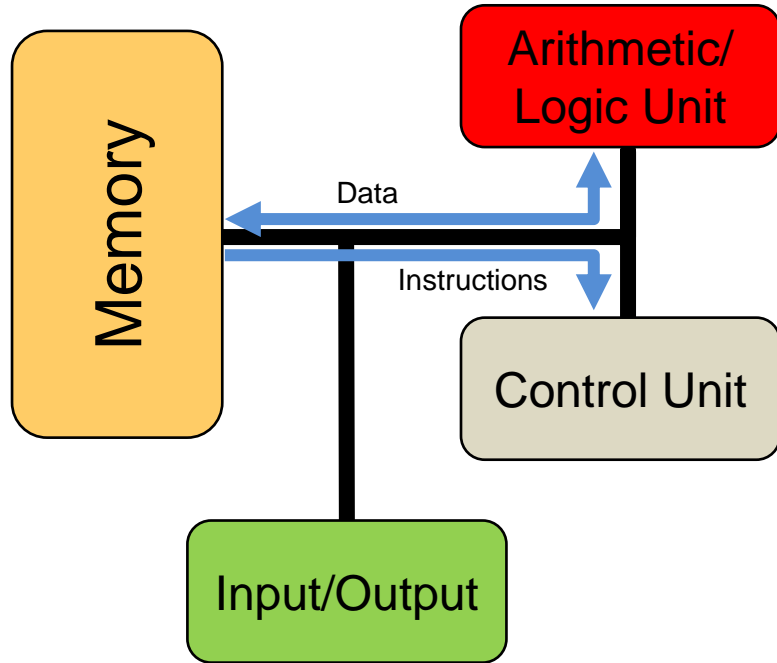
In MPI programming

# Performance issues – overview

- Basics of parallel computer architecture

- Affinity and pinning

- Simple scaling laws

- Benchmarking and performance assessment

- Tracing tools

# Basics of parallel computer architecture

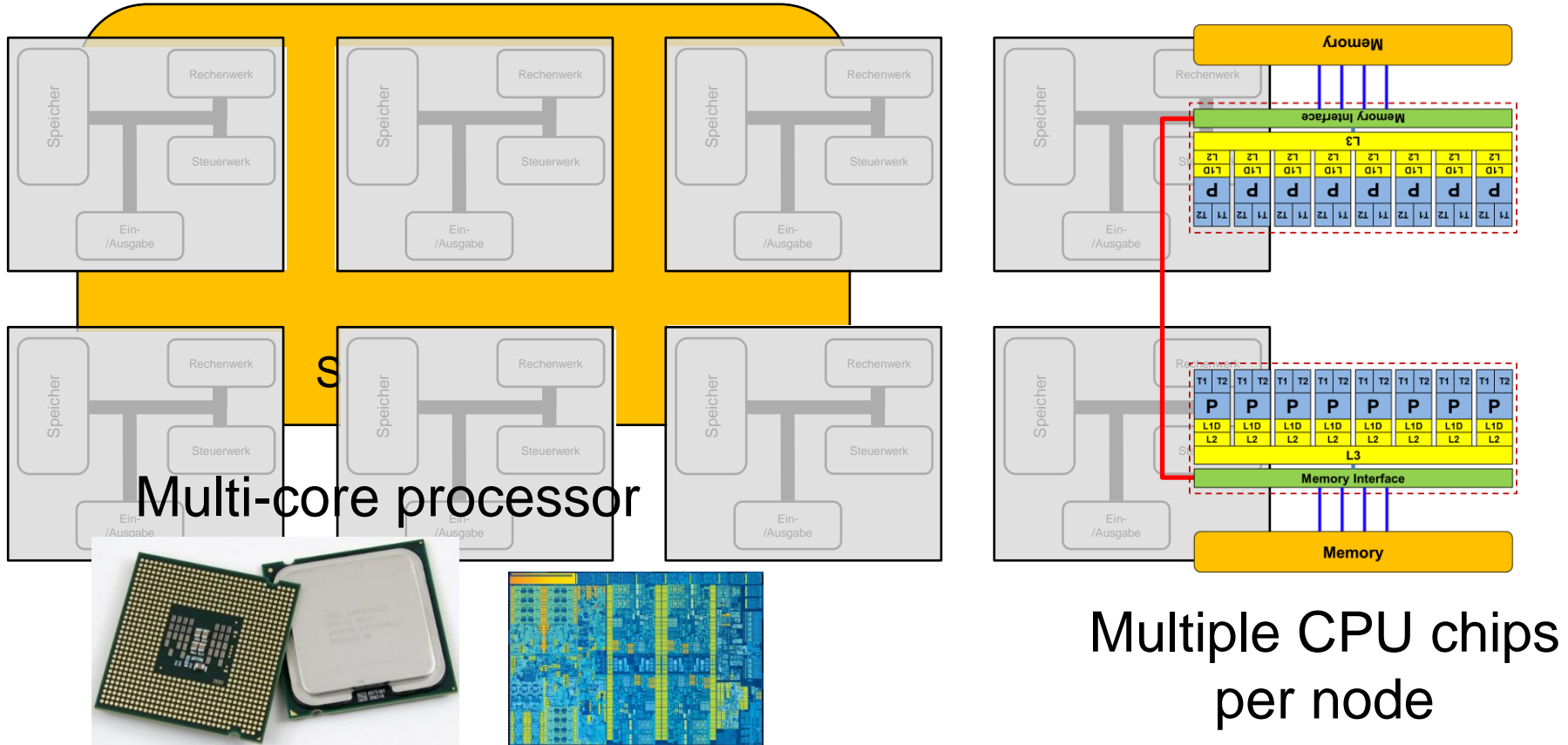# At the core: the stored-program computer



EDSAC (1949)

CC BY 2.0, https://commons.wikimedia.org/w/index.php?curid=432935

TRS-80 Model 4 (1983)

iPhone 7 (2016)

By Rafael Fernandez - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=512180 06
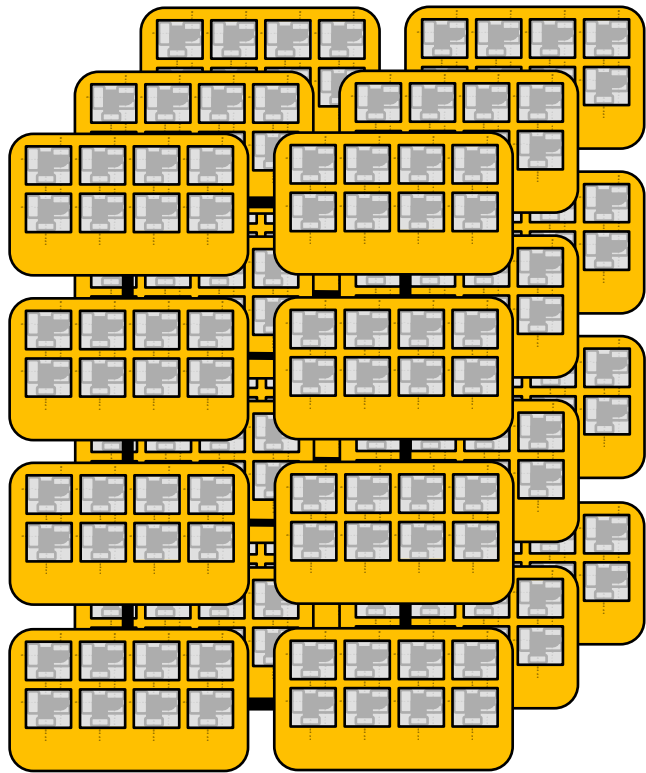
# Shared memory: a single cache-coherent address space



Multi-core processor
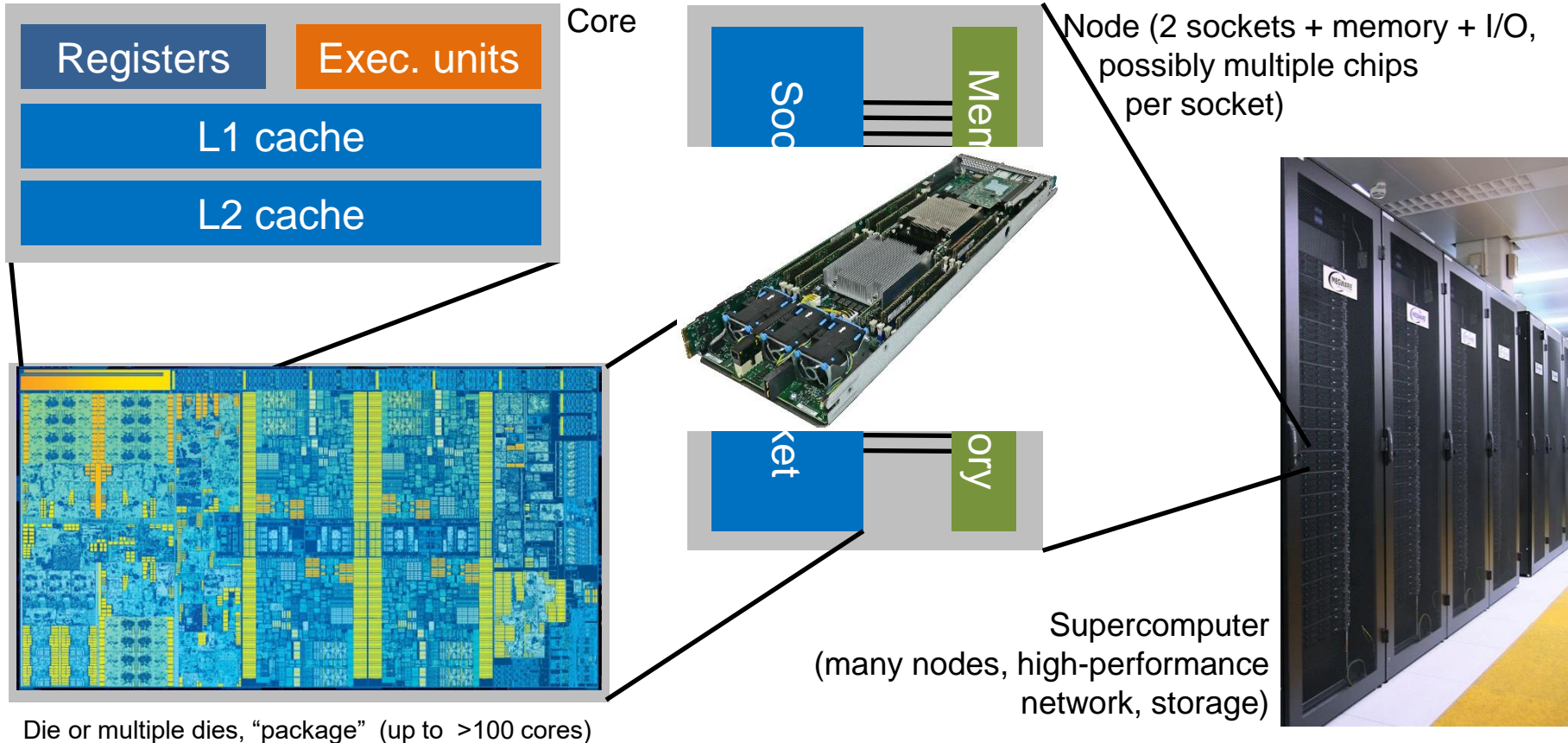
Multiple CPU chips per node

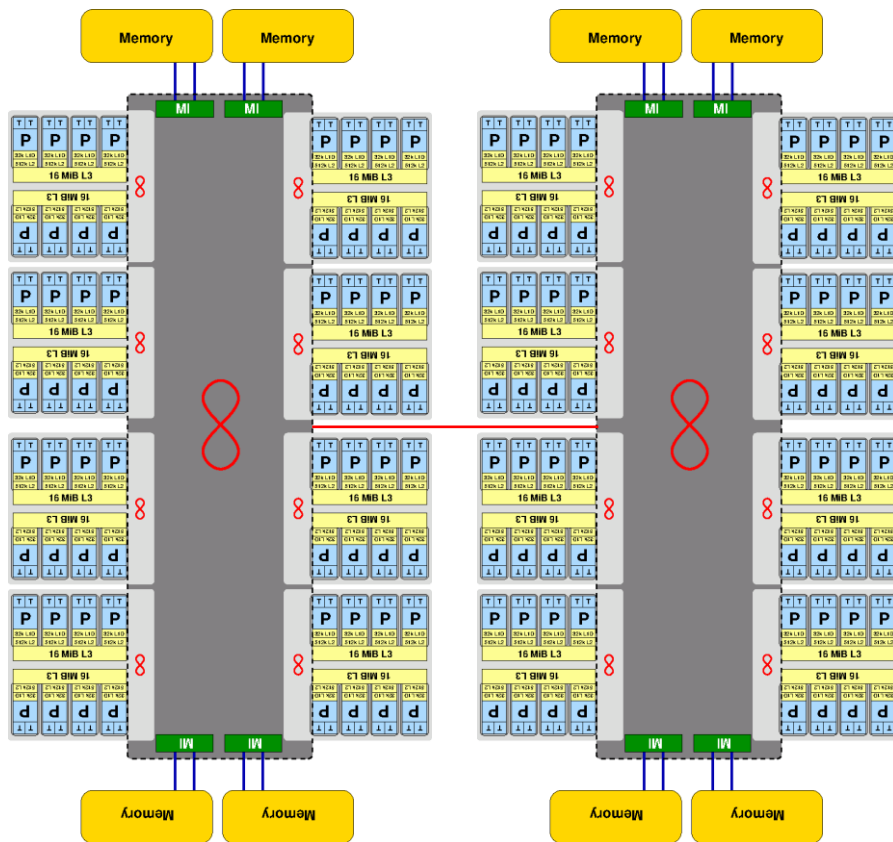# Distributed memory: no cache-coherent single address space



Cluster/
supercomputer

Modern supercomputers are
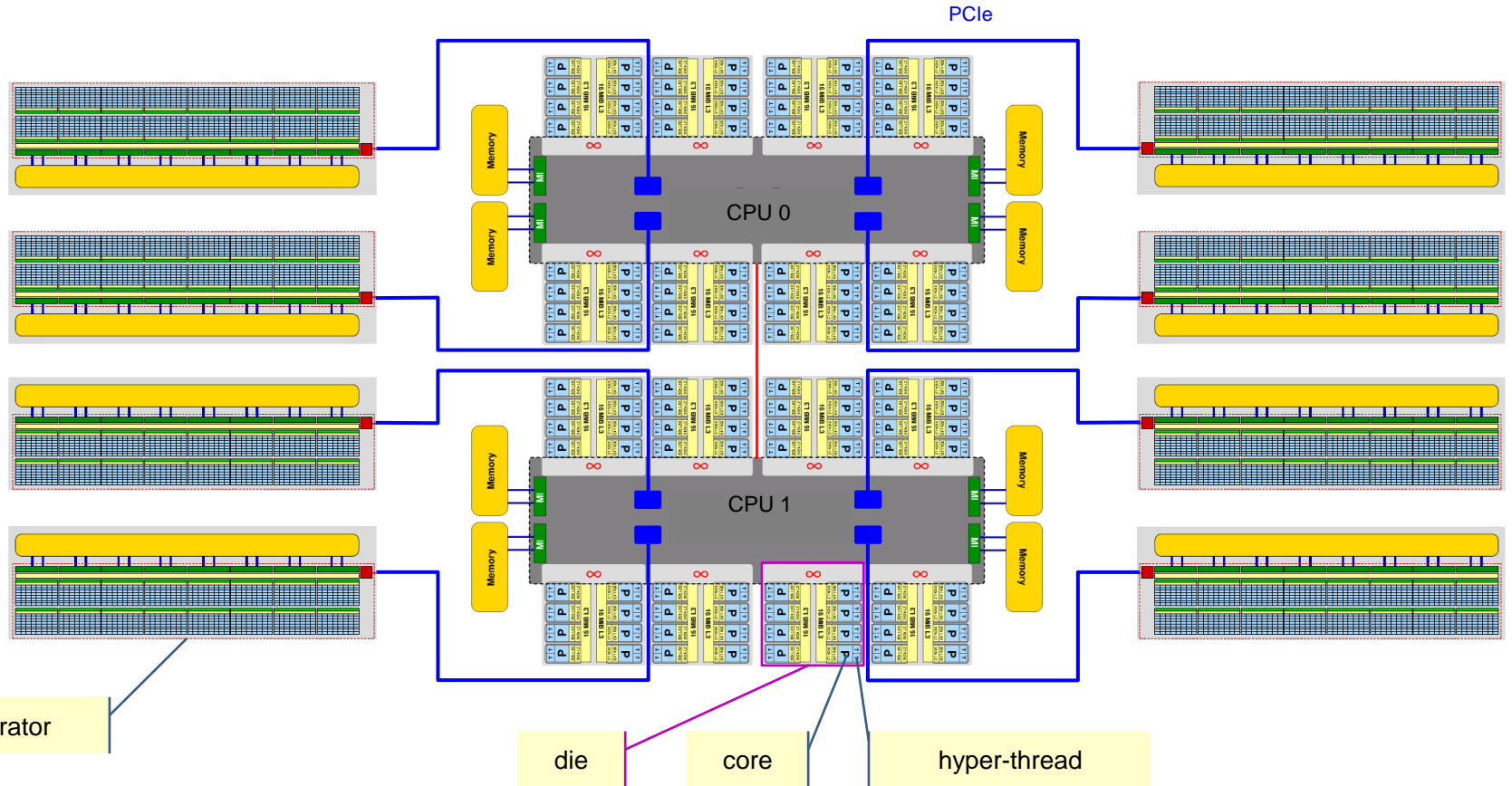shared-/distributed-memory hybrids

# Parallelism in modern computers

Core

| Registers | Exec. units |
|-----------|-------------|

L1 cache

L2 cache

Node (2 sockets + memory + I/O, possibly multiple chips per socket)

Soc... Mem...

...ket ...ory

Die or multiple dies, "package"  (up to  >100 cores)

Supercomputer
(many nodes, high-performance network, storage)

# A modern CPU compute node (AMD Zen2 "Rome")

# Adding accelerators to the node



PCIe

CPU 0

CPU 1

Memory

accelerator

die    core    hyper-thread

# Turning it into a cluster
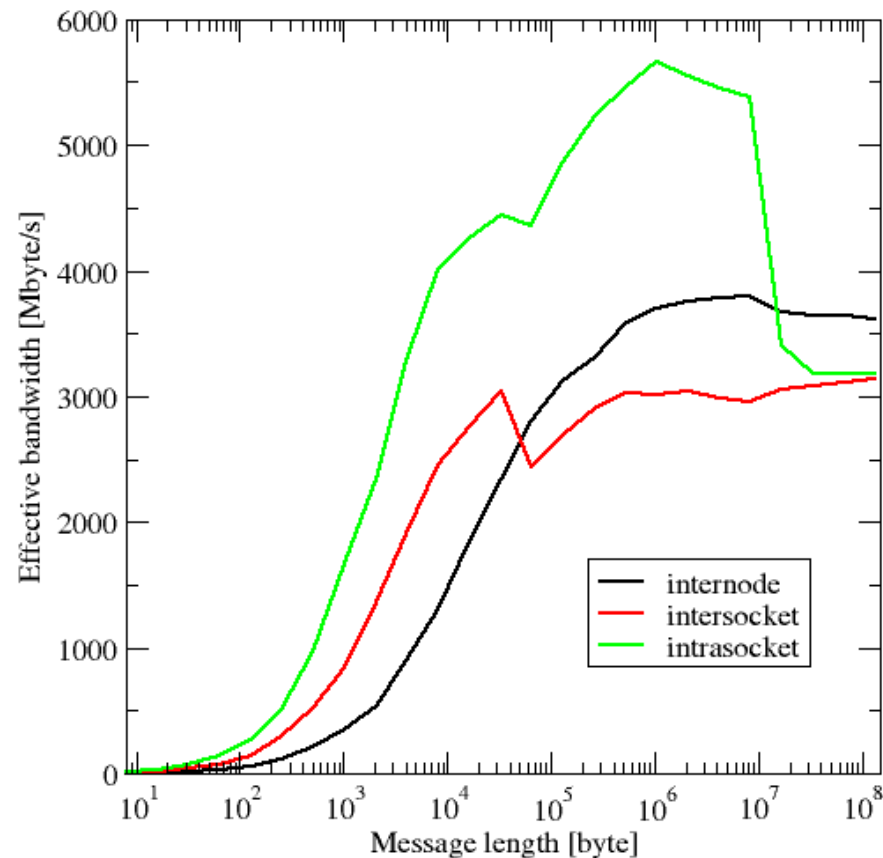


communication network

# Adding permanent storage

# Point-to-point data transmission performance

- Simple "Hockney model" for data transfer time

$$T_{comm} \; = \; \lambda + \frac{V}{b}, \; \; B_{\text{eff}} = \frac{V}{T_{comm}}$$
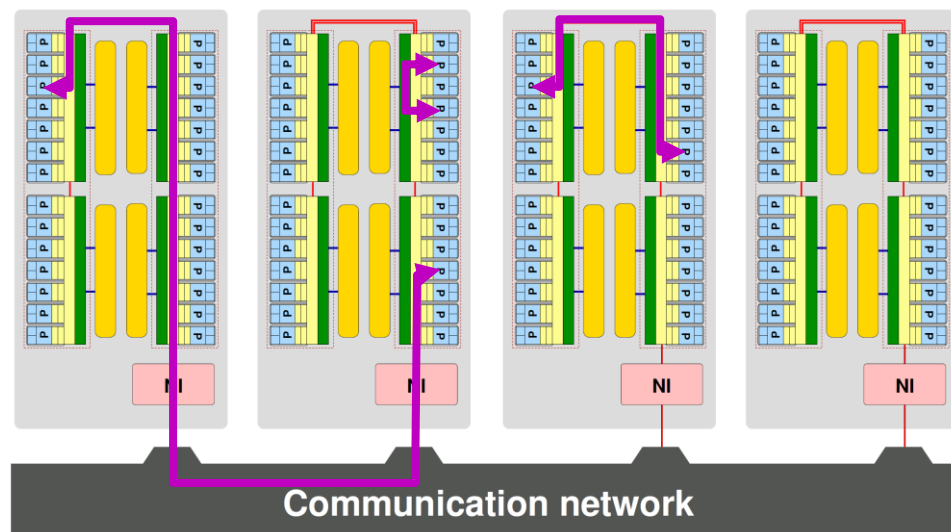
$\lambda$: latency, $b$: asymptotic BW

- Reality is more complicated
  - System topology
  - Caching effects
  - Contention effects
  - Protocol switches
  - Collective communication

# Distributed-memory systems today

"Hybrid" distributed-/shared-memory systems

- Cluster of networked shared-memory nodes

- ccNUMA architecture per node

- Multiple cores per ccNUMA domain



Communication network

- Expect strong topology effects in communication performance
  - Intra-socket, inter-socket, inter-node, all have different $\lambda$ and $b$
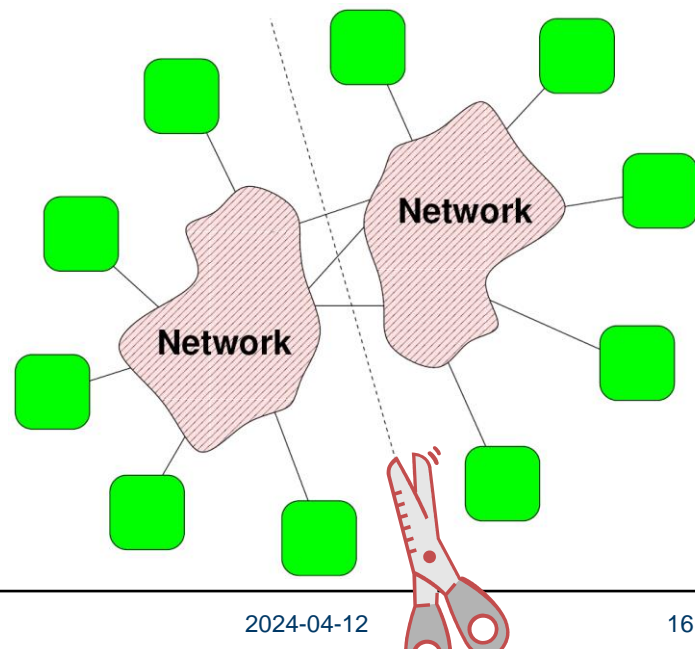  - On top: Effects from network structure

# Characterizing communication networks

- Network bisection bandwidth $B_b$ is a general metric for the data transfer "capability" of a system:

  Minimum sum of the bandwidths of all connections cut when splitting the system into two equal parts
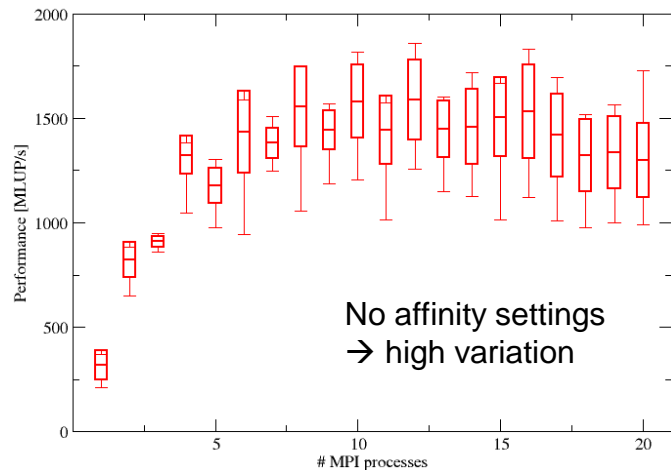


- More meaningful metric for system scalability: bisection BW per node: $B_b/N_{nodes}$

- Bisection BW depends on
  - Bandwidth per link
  - Network topology

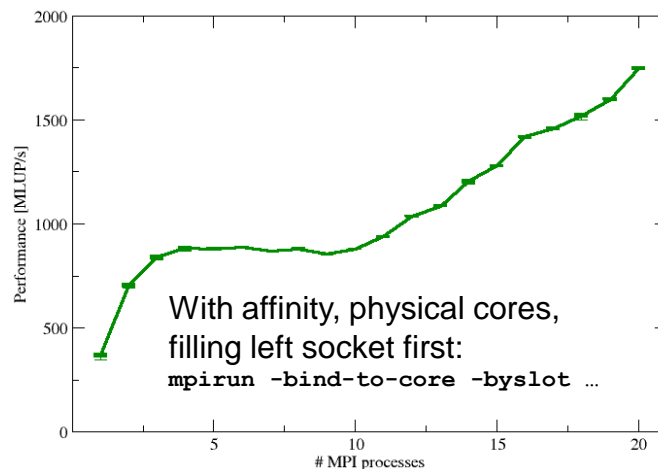# Affinity control (pinning) of processes

# Anarchy vs. affinity with a heat equation solver



No affinity settings
→ high variation

## Reasons for caring about affinity:

- Eliminating performance variation

- Making use of architectural features

- Avoiding resource contention



2x 10-core Intel Ivy Bridge, OpenMPI



With affinity, physical cores,
filling left socket first:
`mpirun -bind-to-core -byslot` …

# Pinning of MPI processes

- Highly implementation and system dependent!

- Intel MPI: env variable `I_MPI_PIN_PROCESSOR_LIST` (MPI only) or `I_MPI_PIN_DOMAIN` (MPI+OpenMP)

- OpenMPI: choose between several mpirun options, e.g.,
  -bind-to-core, -bind-to-socket, -bycore, -byslot …

- Cray's aprun

- Platform-independent tools: `likwid-mpirun`
  (`likwid-pin`, `numactl`)

# Which way to pin

- MPI-only code: `I_MPI_PIN_PROCESSOR_LIST`
- Many options

- Straightforward use:

  `$ mpirun -genv I_MPI_PIN_PROCESSOR_LIST=0-71  -np 144 ./a.out`

  pins one process on each physical core
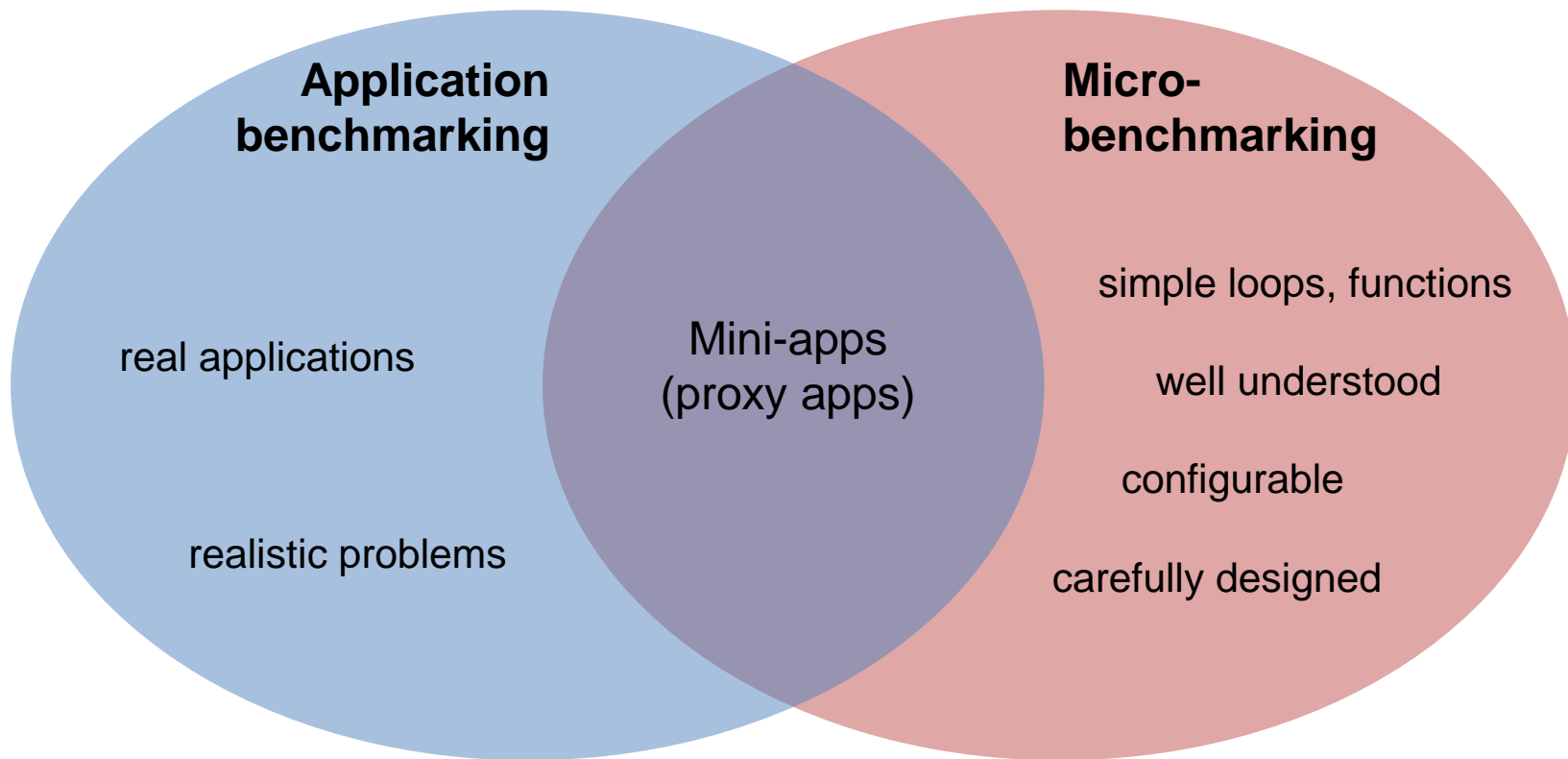
# Benchmarking and performance assessment

More info:

Lecture "Experiments and Data Presentation in High Performance Computing"

https://youtu.be/y1n0IJZiPuw

# Benchmarking: two kinds (and a half)



**Application benchmarking**

real applications

realistic problems

Mini-apps
(proxy apps)

**Micro-benchmarking**

simple loops, functions

well understood

configurable

carefully designed

# Proper definition of benchmark cases

Benchmarking is a vital part of development and performance analysis

1. Define proper benchmark case(s) (input data sets)
   - Reflect(s) "production" workload
   - Tolerable runtime (minutes at most)
2. Document system settings and execution environment
   - Software: compilers, compiler options, library versions, OS version, …
   - Hardware: CPU type, network, [… many more …]
   - Runtime options: Threads/processes per node, affinity, large pages, [… many more …]
3. Document measurement methodology
   - Number of repetitions, statistical variations, …

# Performance and time

- **Performance** is a "higher is better" metric: $P(N) = S(N) \times P(1)$
  - How much work can be done per time unit?
- Work: flops, iterations, "the problem," …
- Time: wall-clock time

- Measuring performance:

```
double s = get_walltime();
// do your work here
double e = get_walltime();
double p = work/(e-s);
```

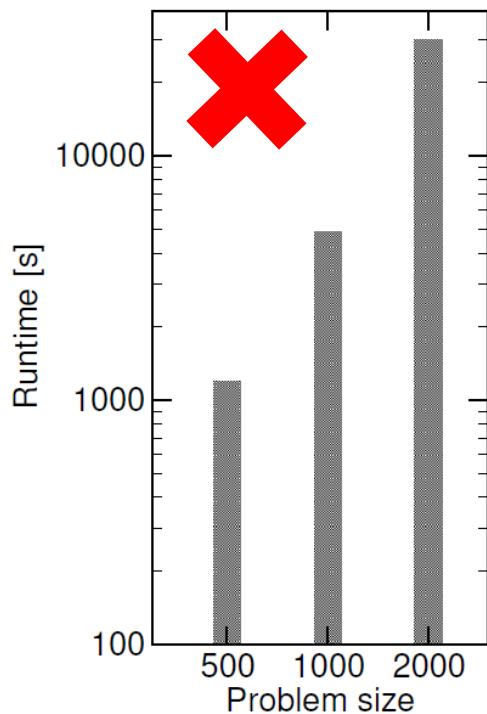- Caveat:
  Timer resolution is finite!

Return time stamp

For Fortran

```
#if !defined(_POSIX_C_SOURCE)
#define _POSIX_C_SOURCE 199309L
#endif

#include <time.h>


double get_walltime() {
  struct timespec ts;
  clock_gettime(CLOCK_MONOTONIC, &ts);
    return (double)ts.tv_sec +
           (double)ts.tv_nsec * 1.e-9;
}


double get_walltime_() {
  return get_walltime();
}
```
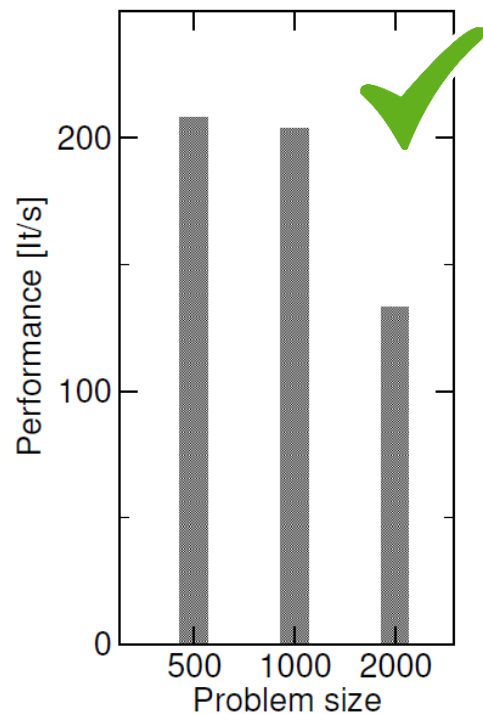
# Popular blunders: runtime != performance

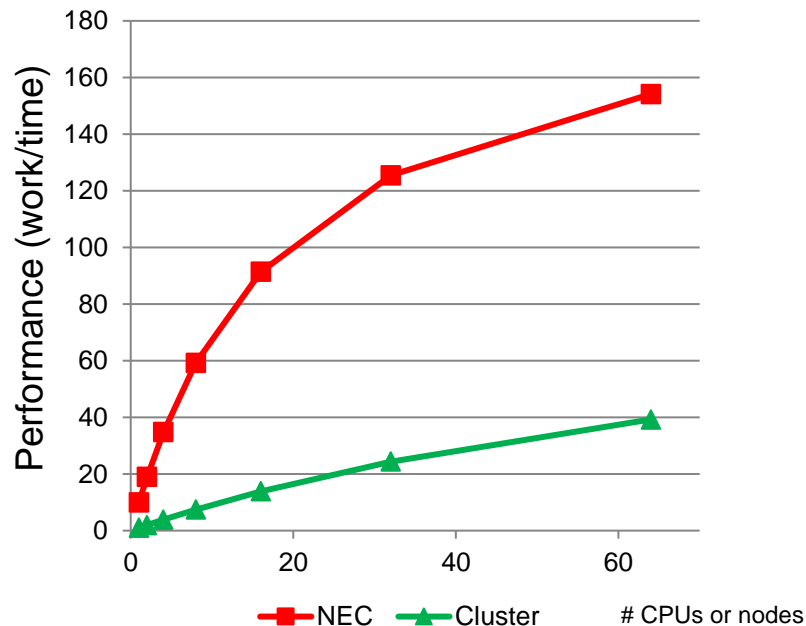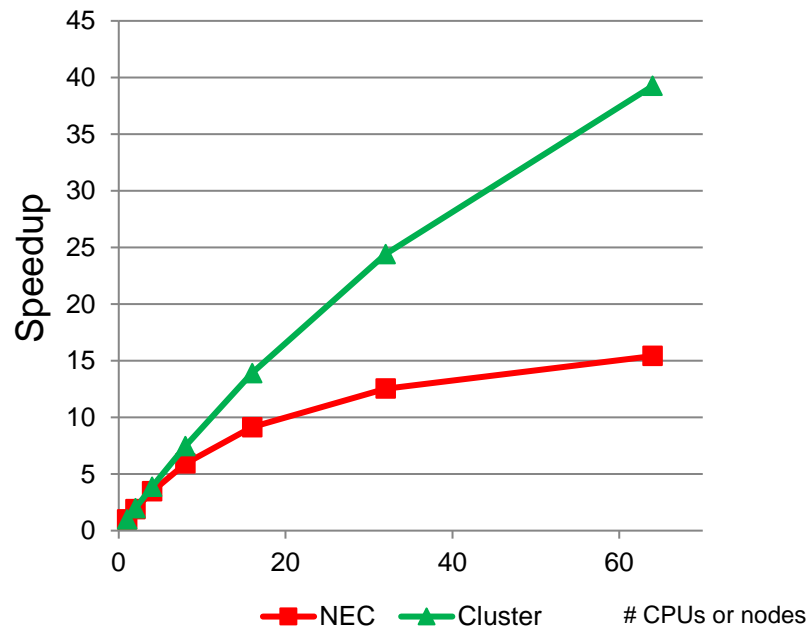- Just presenting runtime is almost always a bad idea!



Insights hidden by trivial dependency: "larger problems need more time"

Performance metric reveals interesting behavior worth investigating!

# Popular blunders: speedup != performance

Speedup hides the "higher is better" quality when comparing different systems or cases

# Limits of parallelism: simple scaling laws

# Metrics to quantify the efficiency of parallel computing

- $T(N)$: execution time of some fixed workload with $N$ workers
- How much faster than with a single worker?

$$\rightarrow \text{parallel speedup: } S(N) = \frac{T(1)}{T(N)}$$

- How efficiently do those $N$ workers do their work?

$$\rightarrow \text{parallel efficiency: } \varepsilon(N) = \frac{S(N)}{N}$$

Can we predict $S(N)$? Are there limits to it?

- Warning: These metrics are not performance metrics!

# Assumptions for basic scalability models

- Scalable hardware: $N$ times the iron can work $N$ times faster
- Work is either fully parallelizable or not at all
- For the time being, assume a constant workload

Ideal world:
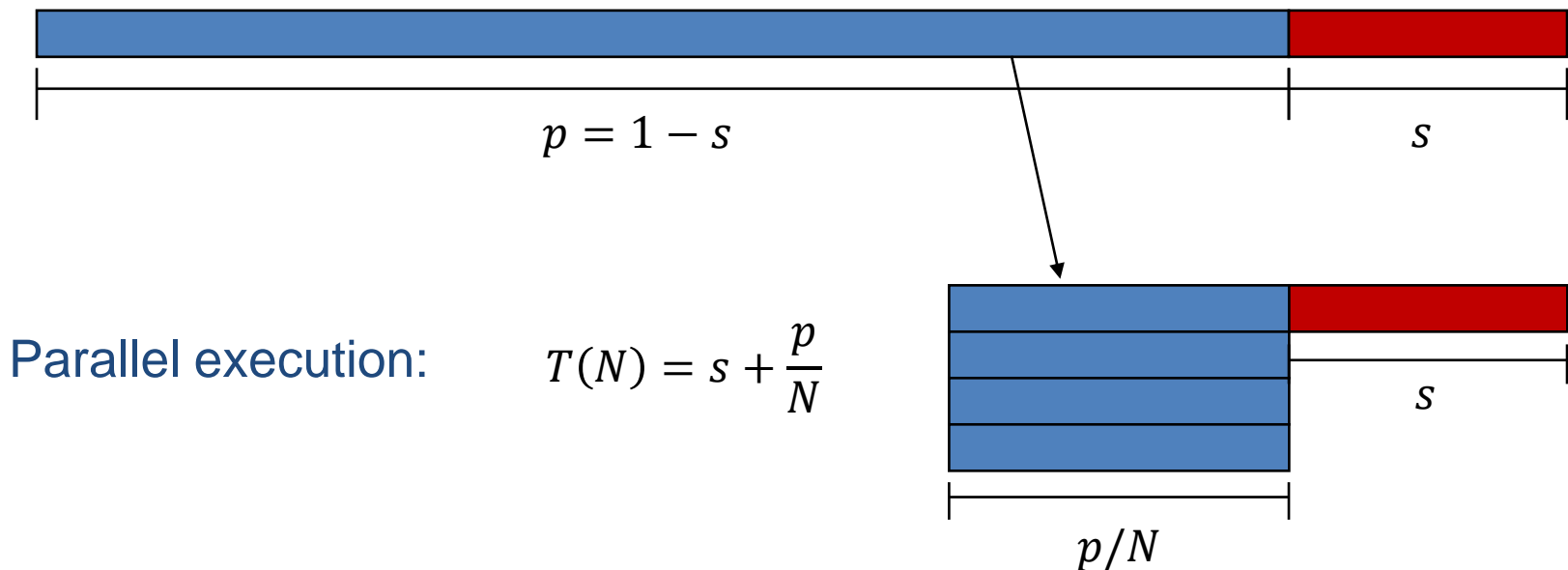All work is perfectly parallelizable
$$S(N) = N, \qquad \varepsilon = 1$$

# A simple speedup model for fixed workload

One worker normalized execution time: $T(1) = s + p = 1$
$s$: runtime of purely serial part
$p$: runtime of perfectly parallelizable part

$$p = 1 - s$$

$$s$$

Parallel execution:

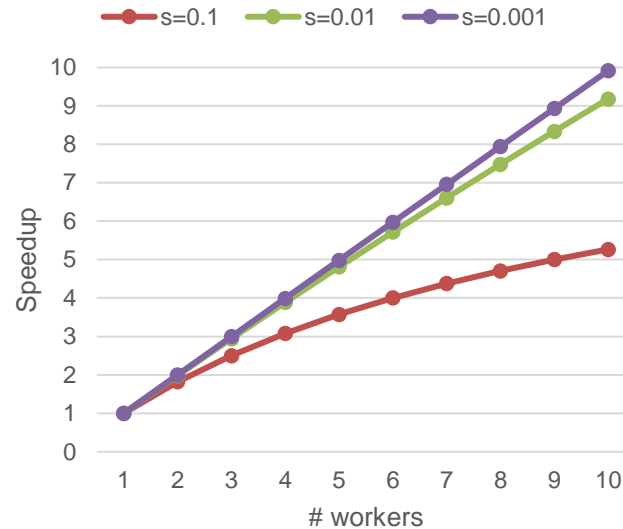$$T(N) = s + \frac{p}{N}$$

$$s$$

$$p/N$$

# Amdahl's Law (1967) – "Strong Scaling"

- Fixed workload speedup with $s$ being the fraction of nonparallelizable work

$$S(N) = \frac{T(1)}{T(N)} = \frac{1}{s + \frac{1-s}{N}}$$

- Parallel efficiency: $\varepsilon(N) = \frac{1}{s(N-1)+1}$



Gene M. Amdahl: *Validity of the single processor approach to achieving large scale computing capabilities*. In Proceedings of the April 18-20, 1967, spring joint computer conference (AFIPS '67 (Spring)). Association for Computing Machinery, New York, NY, USA, 483–485. DOI:10.1145/1465482.1465560

# Fundamental limits in Amdahl's Law

- Asymptotic speedup

$$\lim_{N \to \infty} S(N) = \frac{1}{s}$$

- Asymptotic parallel efficiency

$$\lim_{N \to \infty} \varepsilon(N) = 0$$

→ Asymptotically, nobody is doing anything except the worker that gets the serial work!

- In reality, it's even worse…

# Strong scaling plus overhead

- Let $c(N)$ be an overhead term that may include communication and/or synchronization

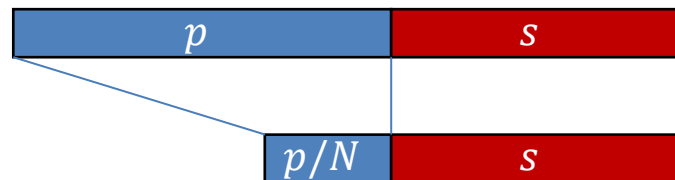$$\rightarrow \quad T(N) = s + \frac{p}{N} + c(N)$$

- What goes into $c(N)$?
  - Communication pattern
  - Synchronization strategy
  - Message sizes
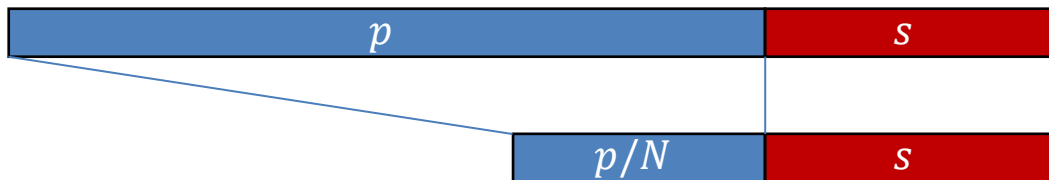  - Network structure
  - …

Typical examples: $c(N) =$

- $kN^2$           (all-to-all on bus network)
- $k \log N$       (optimal synchronization)
- $kN$             (one sends to all)
- $\lambda + kN^{-\frac{2}{3}}$     (Cartesian domain decomposition, nonblocking network)

# A simple speedup model for scaled workload

- What if we could increase the parallel part of the work only?
  → the larger $p$, the larger the speedup



- This is not possible for all applications, but for some



- "Weak scaling"

# A simple speedup model for scaled workload

- Parallel workload grows linearly with $N$

$$\rightarrow T(N) = s + \frac{pN}{N} = s + p, \text{ i.e., runtime stays constant}$$

- Scalability metric?
  $\rightarrow$ How much more work per second can be done with $N$ workers than with one worker?

$$S(N) = \frac{(s + pN)/(s + p)}{(s + p)/(s + p)} = s + (1 - s)N$$

Gustafson's Law ("weak scaling")

John L. Gustafson: *Reevaluating Amdahl's law*. Commun. ACM 31, 5 (May 1988), 532–533.
DOI:10.1145/42411.42415

# Gustafson's Law for weak scaling

- Linear speedup (but not proportional unless $s = 0$) with $N$:

$$S(N) = s + (1-s)N \quad \rightarrow \quad \text{unbounded speedup!}$$

- Weak scaling is the solution to the Amdahl dilemma: Why should we build massively parallel systems if all parallelism is limited by the serial fraction?

- Extension to communication?

$$\rightarrow T(N) = s + \frac{pN}{N} + c(N) = 1 + c(N)$$

$$\rightarrow S(N) = \frac{(s+pN)/(1+c(N))}{(s+p)/1} = \frac{s+(1-s)N}{1+c(N)}$$

Much more relaxed conditions on $c(N)$

# How can we determine the model parameters?

- Manual analysis: Requires in-depth knowledge of hardware and program
- Curve fitting: Less insight, but also less cumbersome

- Example: Strong scaling of hypothetical code on "Meggie" node @FAU (10 cores per socket, 2 sockets per node)
- Use "extended Amdahl's" with $kN$ overhead
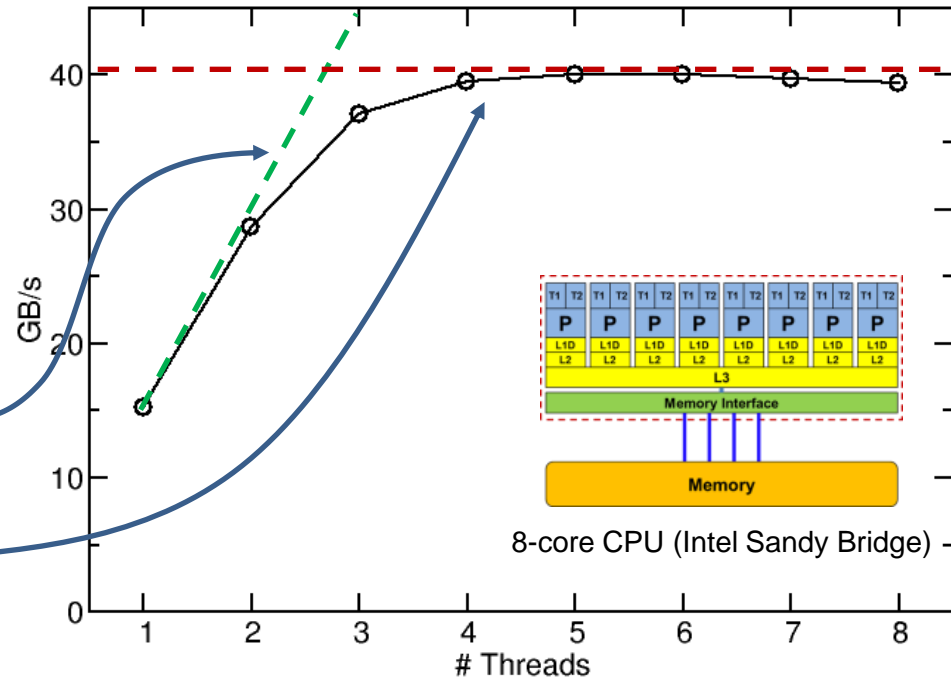
- Result:
  Best fit is not a good fit at all

# Resource bottlenecks

- Amdahl's Law assumes perfect scalability of resources
- Reality: Computer architecture is plagued by bottlenecks!
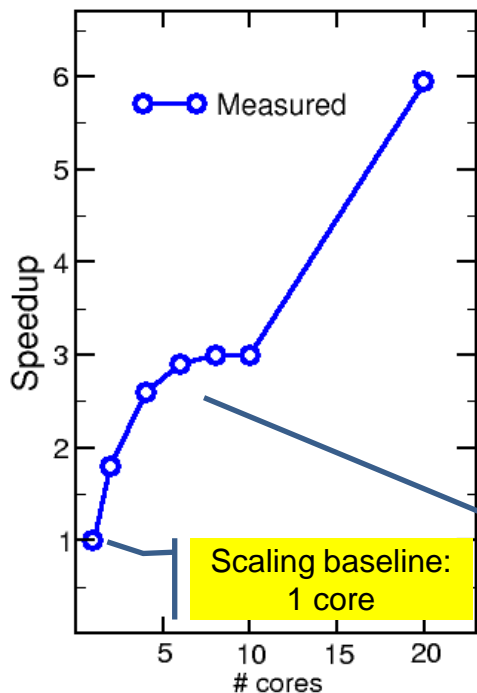- Example: array update loop

```
// MPI-parallel

for(i=0; i<10000000; ++i)

  a[i] = a[i] + s * c[i];
```

- Amdahl's: $s = 0, c(N) = 0$
  - Perfect speedup? No!
  - Saturation because of memory bandwidth exhaustion



8-core CPU (Intel Sandy Bridge)

# Separation of scaling baselines is key!

- **Intra-socket scaling is often not covered by the model**
  - Model assumes "scalable resources"
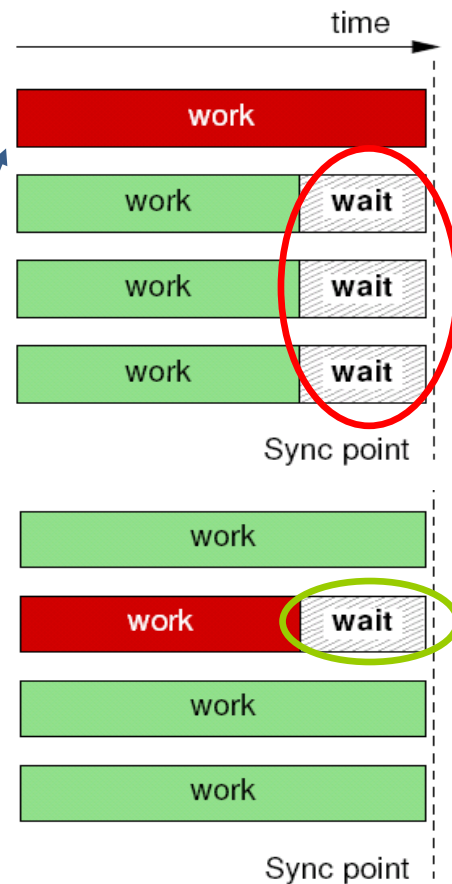


Model well suited for internode scaling!

Separating scaling baselines is important in modeling!

Socket saturation due to memory bandwidth

Scaling baseline: 1 core

Scaling baseline: 1 node

# Amdahl generalized: load imbalance

- Load imbalance at sync points
  - More specifically, execution time imbalance
  - $p/N$ assumption no longer valid in general

- Hard to model in general, but two corner cases:
  - A few "laggers" waste lots of resources
    - Single lagger → Amdahl's Law
  - A few "speeders" might be harmless

- Tuning advice
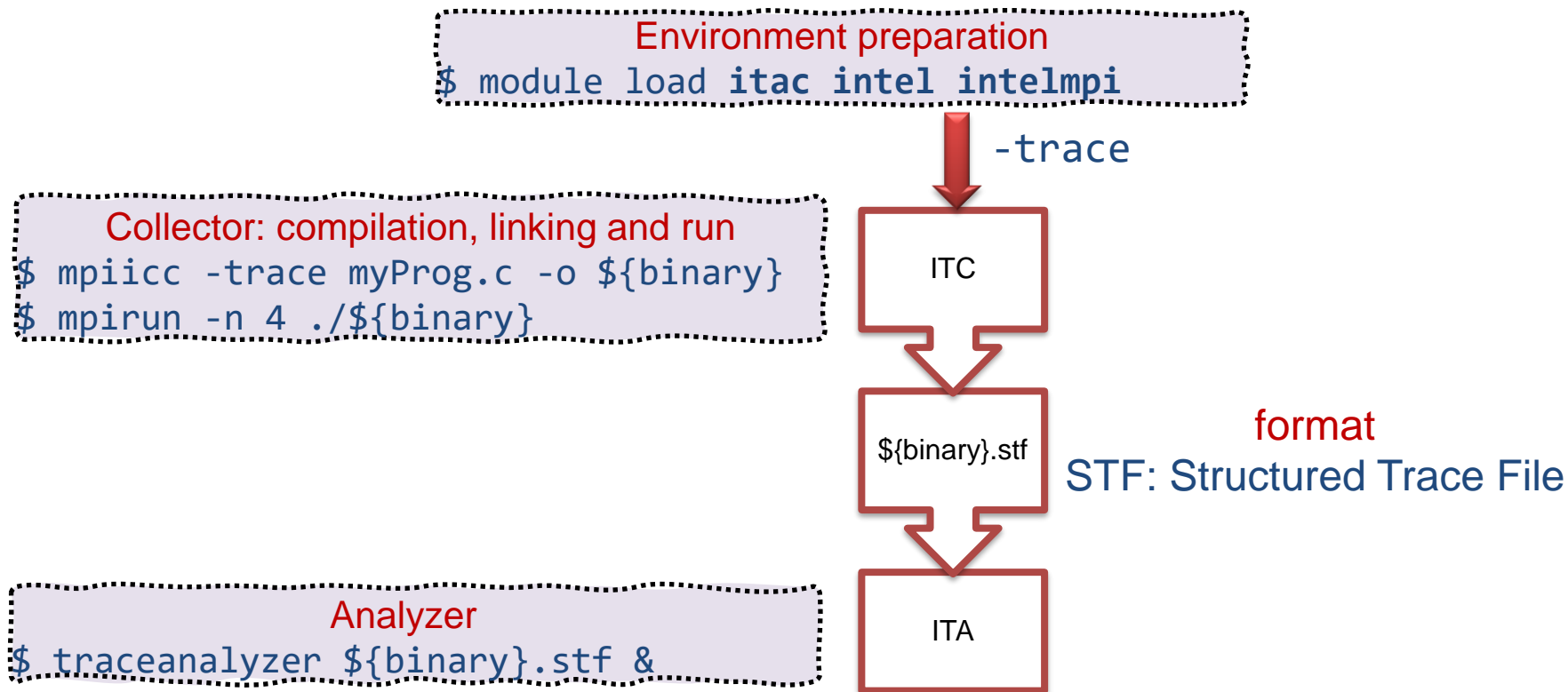  - Avoid sync points
  - Turn laggers into speeders

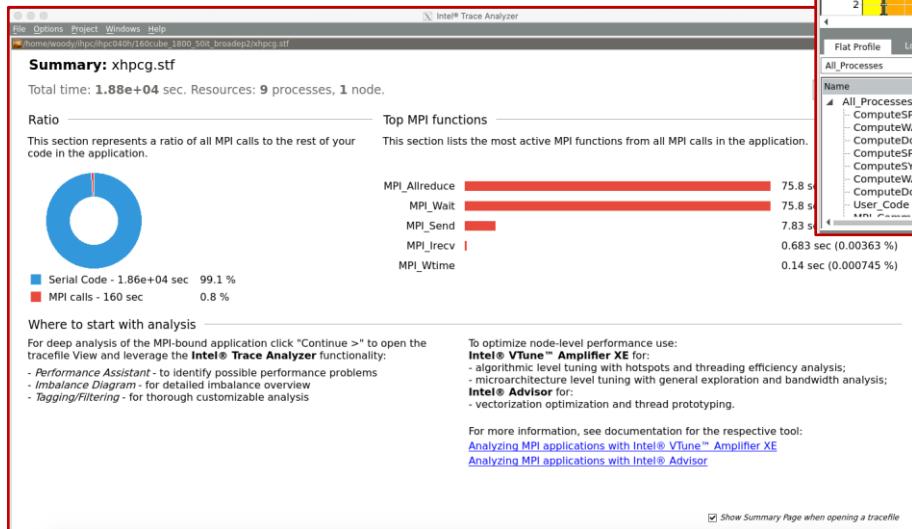# MPI tracing tools

# MPI tracing tools

- Allow the user to track events and statistics pertaining to MPI communication and code execution
- Popular tools
  - Intel Trace Analyzer and Collector (ITAC)
  - VAMPIR (commercial)
  - Paraver

- Powerful tools
- Potential to produce massive amounts of data
- Danger of "drowning in data"

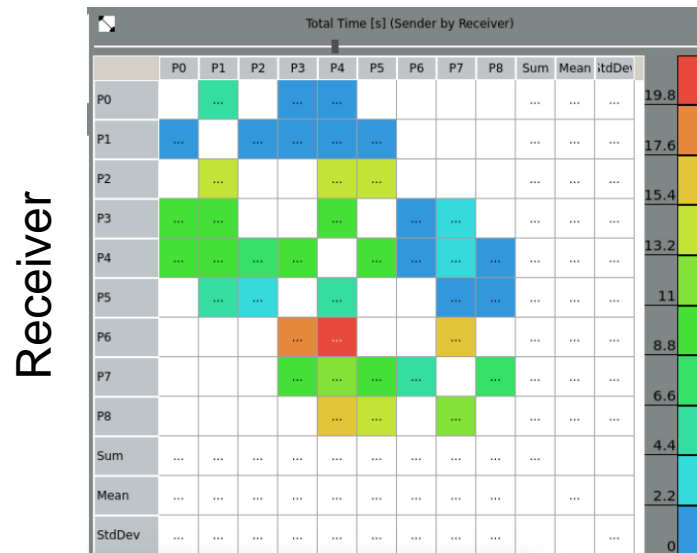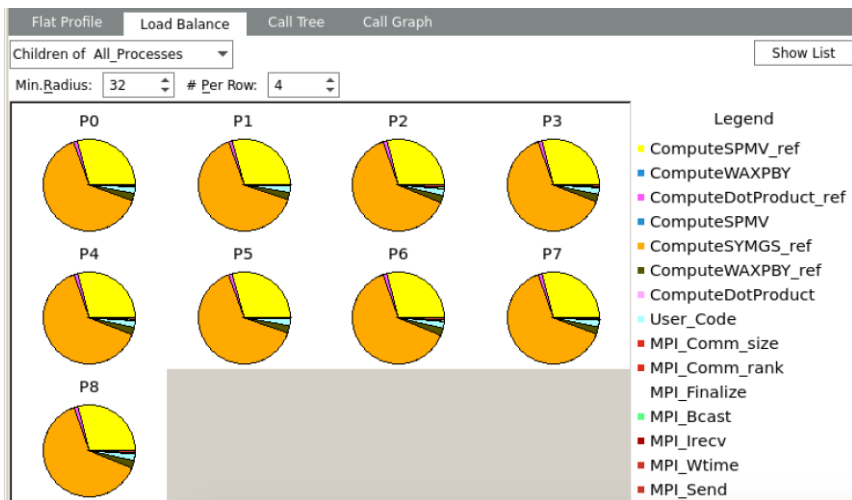# Intel Trace Analyzer and Collector

Environment preparation
```
$ module load itac intel intelmpi
```

-trace

Collector: compilation, linking and run
```
$ mpiicc -trace myProg.c -o ${binary}
$ mpirun -n 4 ./${binary}
```

ITC

${binary}.stf

format
STF: Structured Trace File

Analyzer
```
$ traceanalyzer ${binary}.stf &
```

ITA

# Basic features of ITAC

Event-based approach that record
- user function calls
- MPI communication calls

# Some features of ITAC

# Timeline view