# High Performance Computing
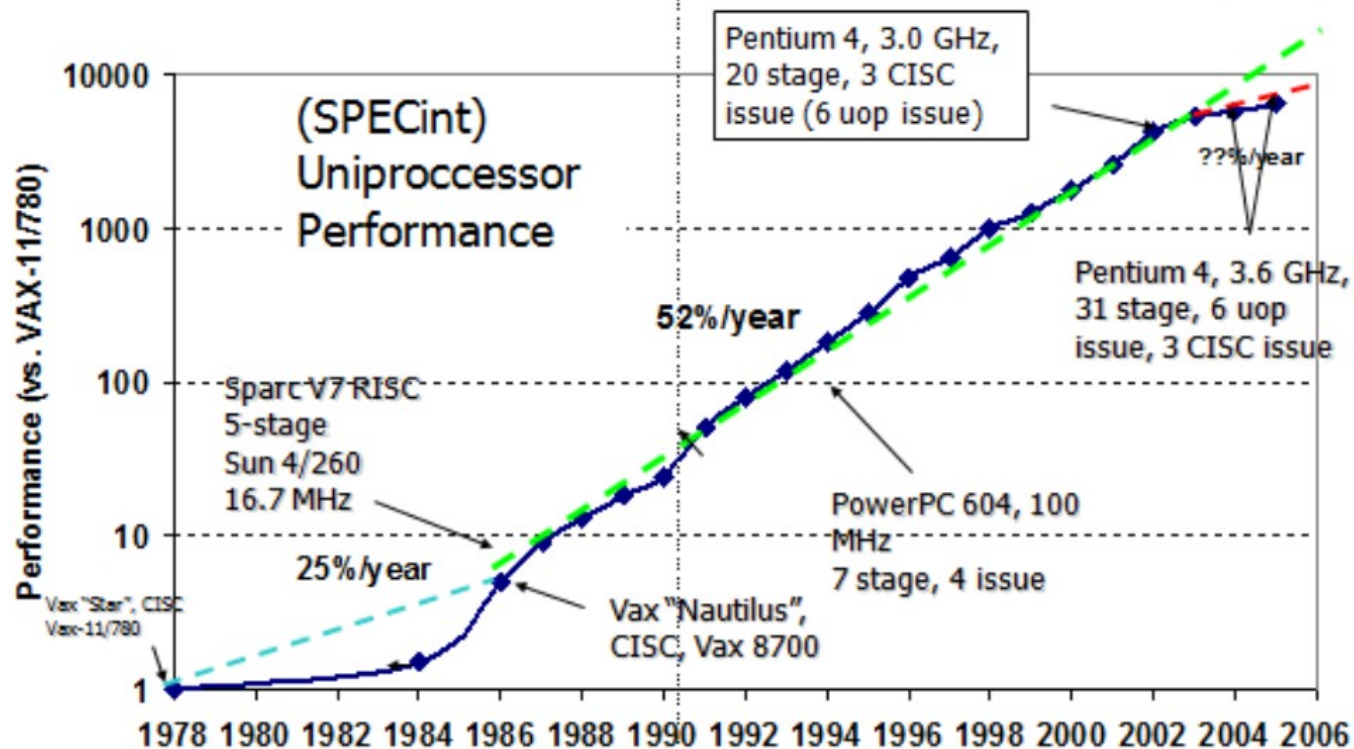
# Shared Memory Parallel CPU & OpenMP
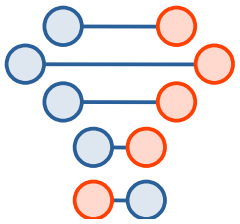
Computer Eng., KMITL

Assoc. Prof. Dr. Surin. K.

The good old days …

Moore's Law

(SPECint) Uniproccessor Performance

Pentium 4, 3.0 GHz, 20 stage, 3 CISC issue (6 uop issue)

Pentium 4, 3.6 GHz, 31 stage, 6 uop issue, 3 CISC issue

??%/year

52%/year

Sparc V7 RISC 5-stage Sun 4/260 16.7 MHz

25%/year

PowerPC 604, 100 MHz 7 stage, 4 issue

Vax "Star", CISC Vax-11/780

Vax "Nautilus", CISC, Vax 8700

From Hennessy and Patterson, *Computer Architecture: A Quantitative Approach*, 4th edition, Sept. 15, 2006
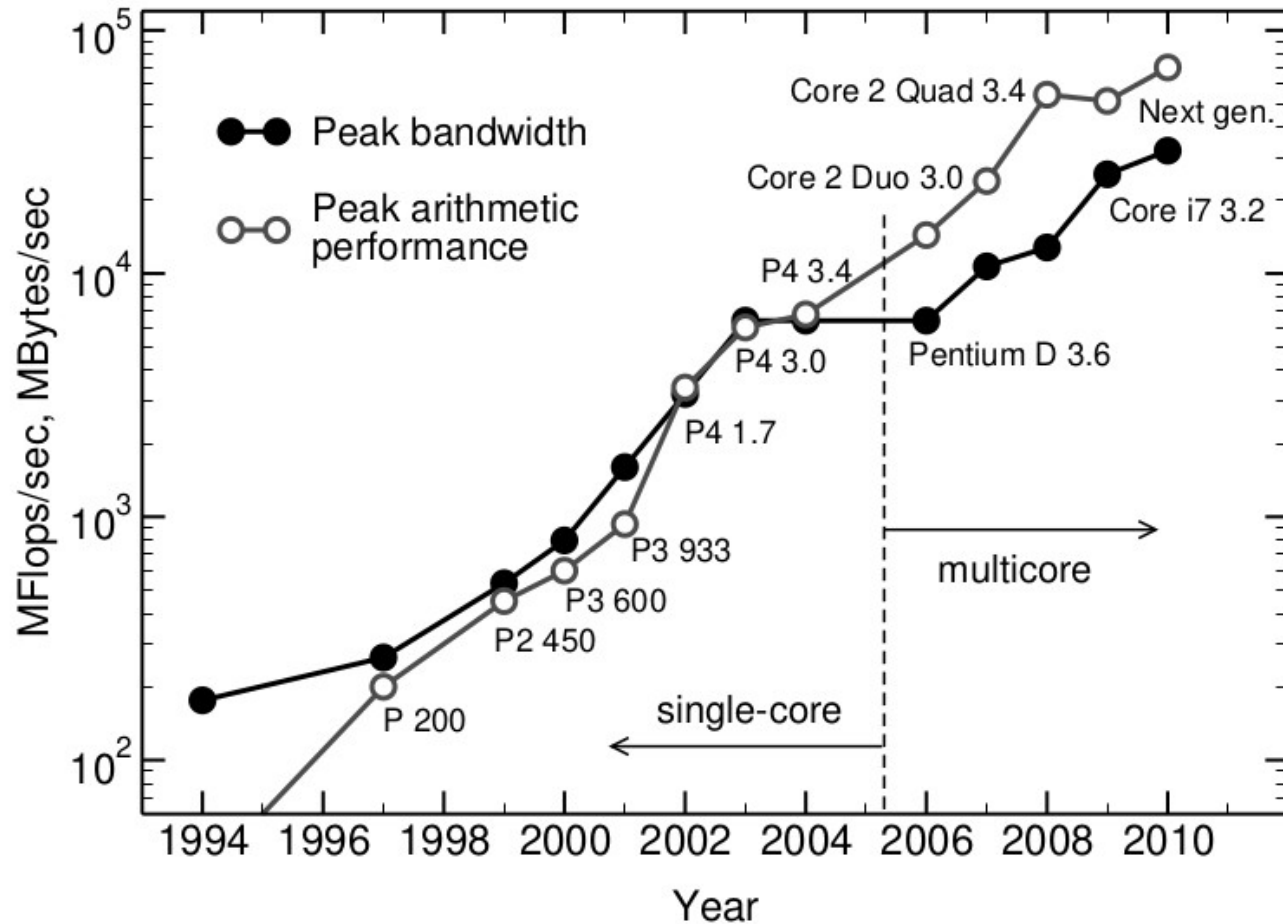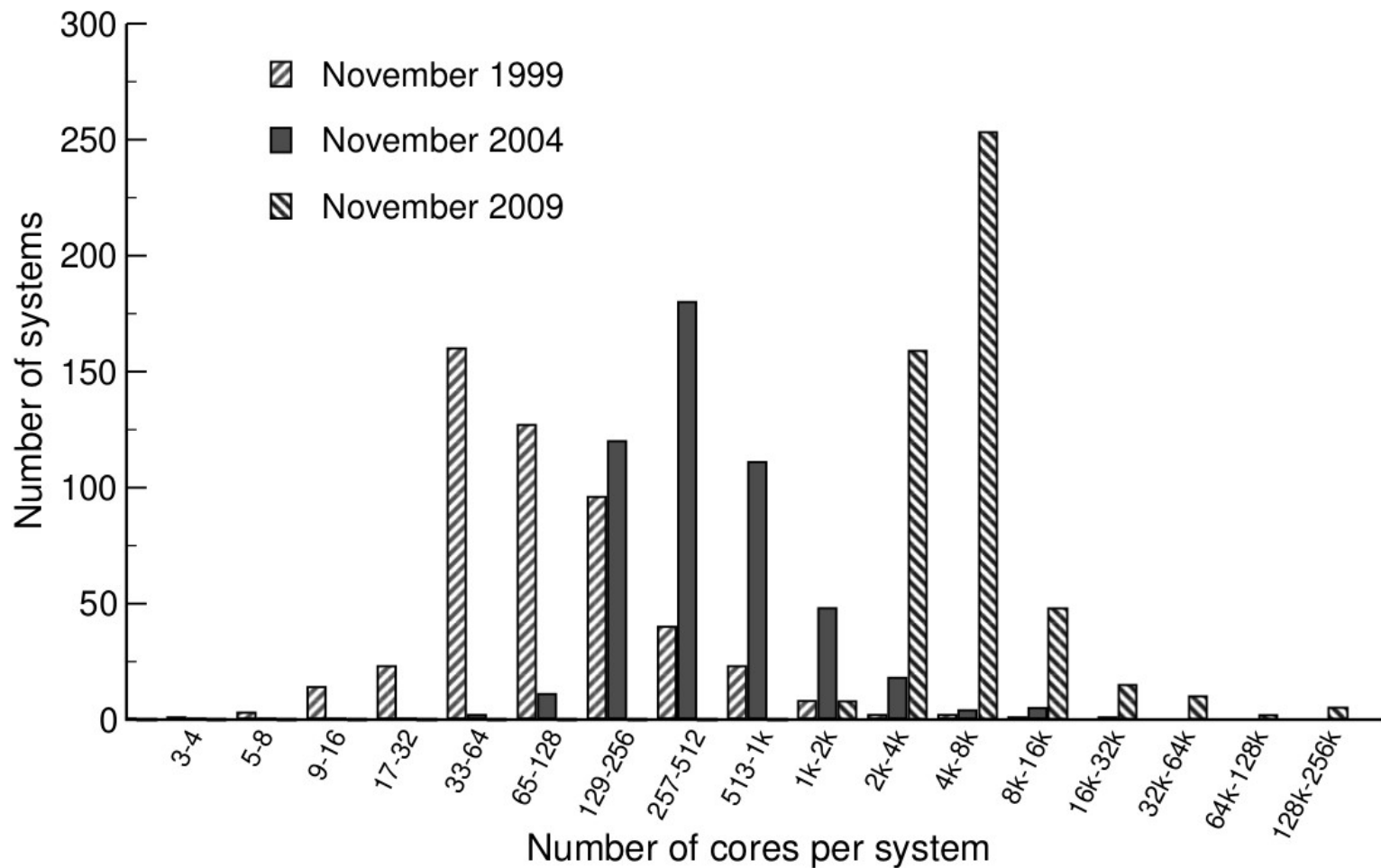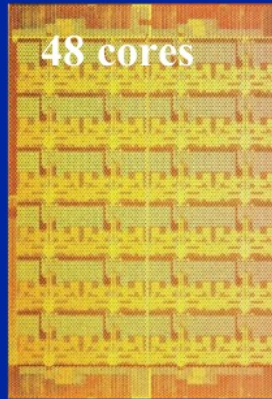
2

# Moore's Law



**Figure 3.2:** Progress of maximum arithmetic performance (open circles) and peak theoretical memory bandwidth (filled circles) for Intel processors since 1994. The fastest processor in terms of clock frequency is shown for each year. (Data collected by Jan Treibig.)

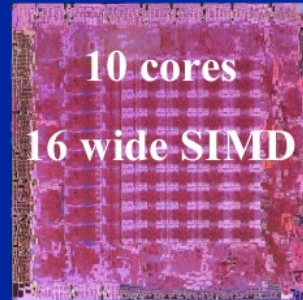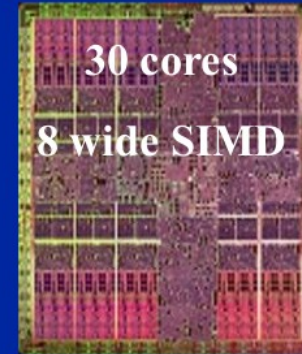# Number of systems versus core count

# Microprocessor trends

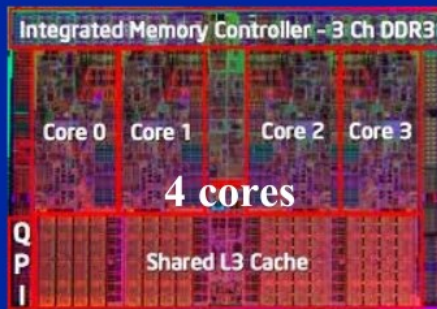**Individual processors are many core (and often heterogeneous) processors.**

48 cores

**Intel SCC Processor**

10 cores

16 wide SIMD

**AMD ATI RV770**

30 cores

8 wide SIMD

**NVIDIA Tesla C1060**

Integrated Memory Controller - 3 Ch DDR3

Core 0   Core 1      Core 2   Core 3

4 cores

QPI   Shared L3 Cache

**Intel® Xeon® processor**

1 CPU + 6 cores

**IBM Cell**

4 cores

**ARM MPCORE**

3<sup>rd</sup> party names are the property of their owners.

Source: OpenCL tutorial, Gaster, Howes, Mattson, and Lokhmotov, HiPEAC 2011
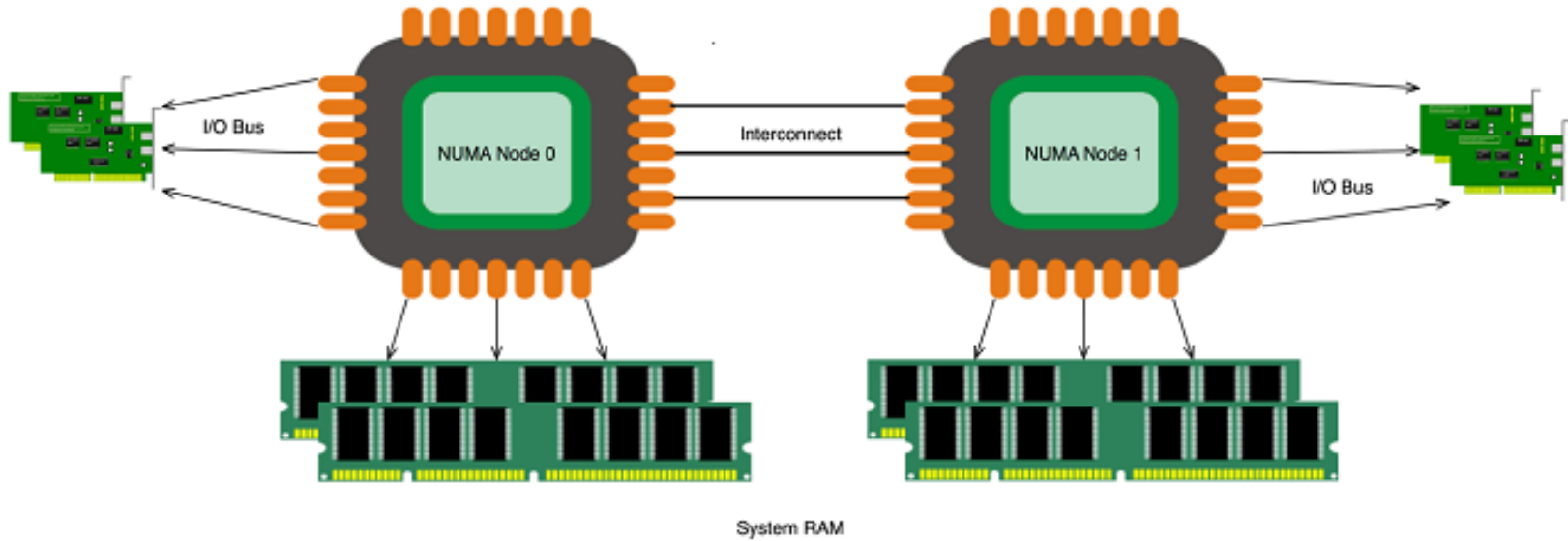
# NUMA: Intel Sandy Bridge



CPU architecture (Intel Sandy Bridge)

# NUMA: Intel Quick Assist Technology



I/O Bus

NUMA Node 0
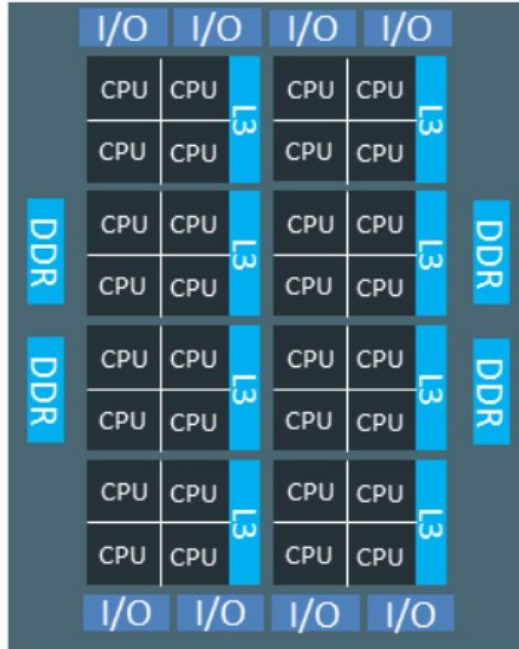
Interconnect

NUMA Node 1

I/O Bus

System RAM

# NUMA: Intel Quick Assist Technology
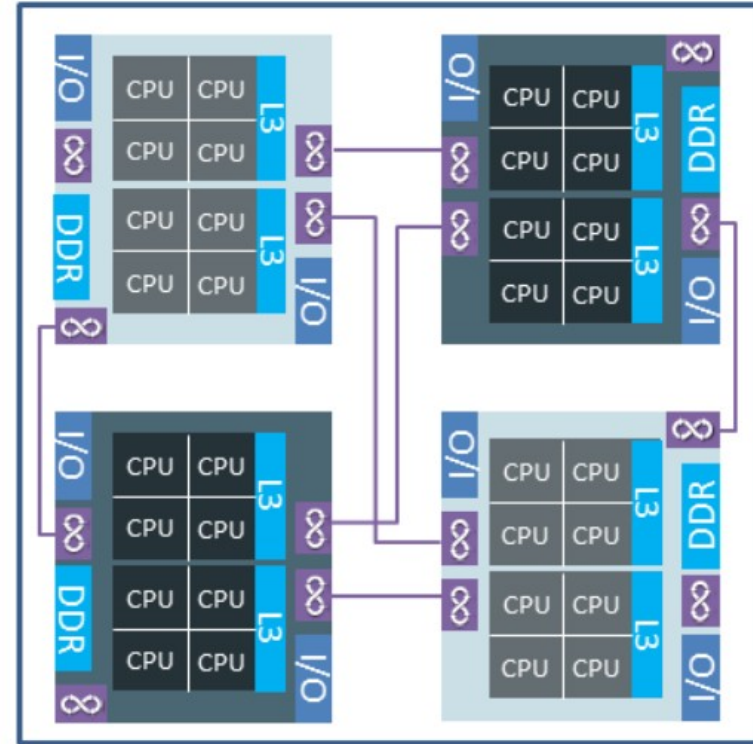
8

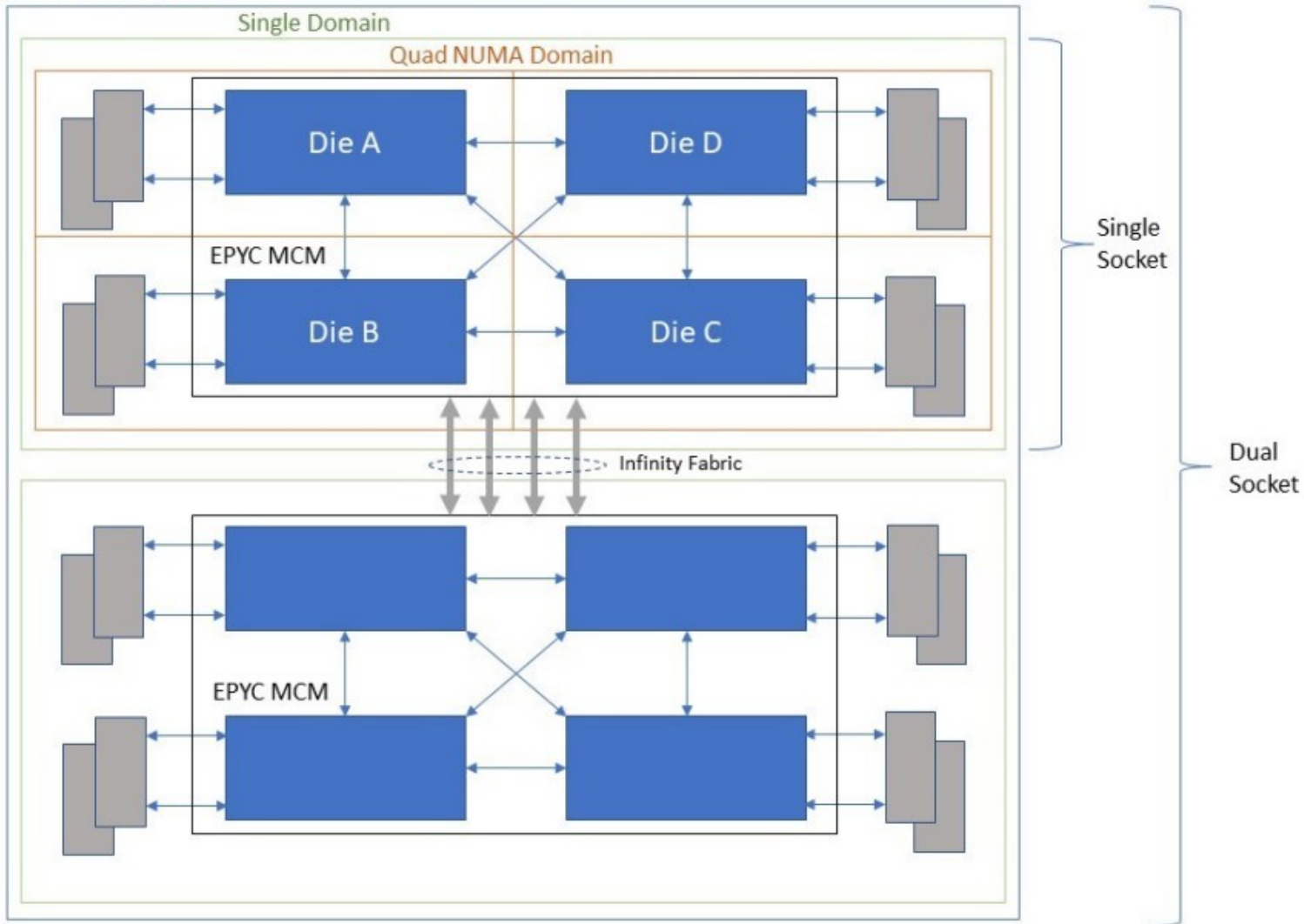# NUMA: AMD Epyc



Monolithic Die

EPYC MCM

32C Die Cost
1.0X

4 x 8C Die Cost
0.59X[1]

NUMA: AMD Epyc

# OpenMP

## Compiler notes: Other

- Linux and OS X with gcc:

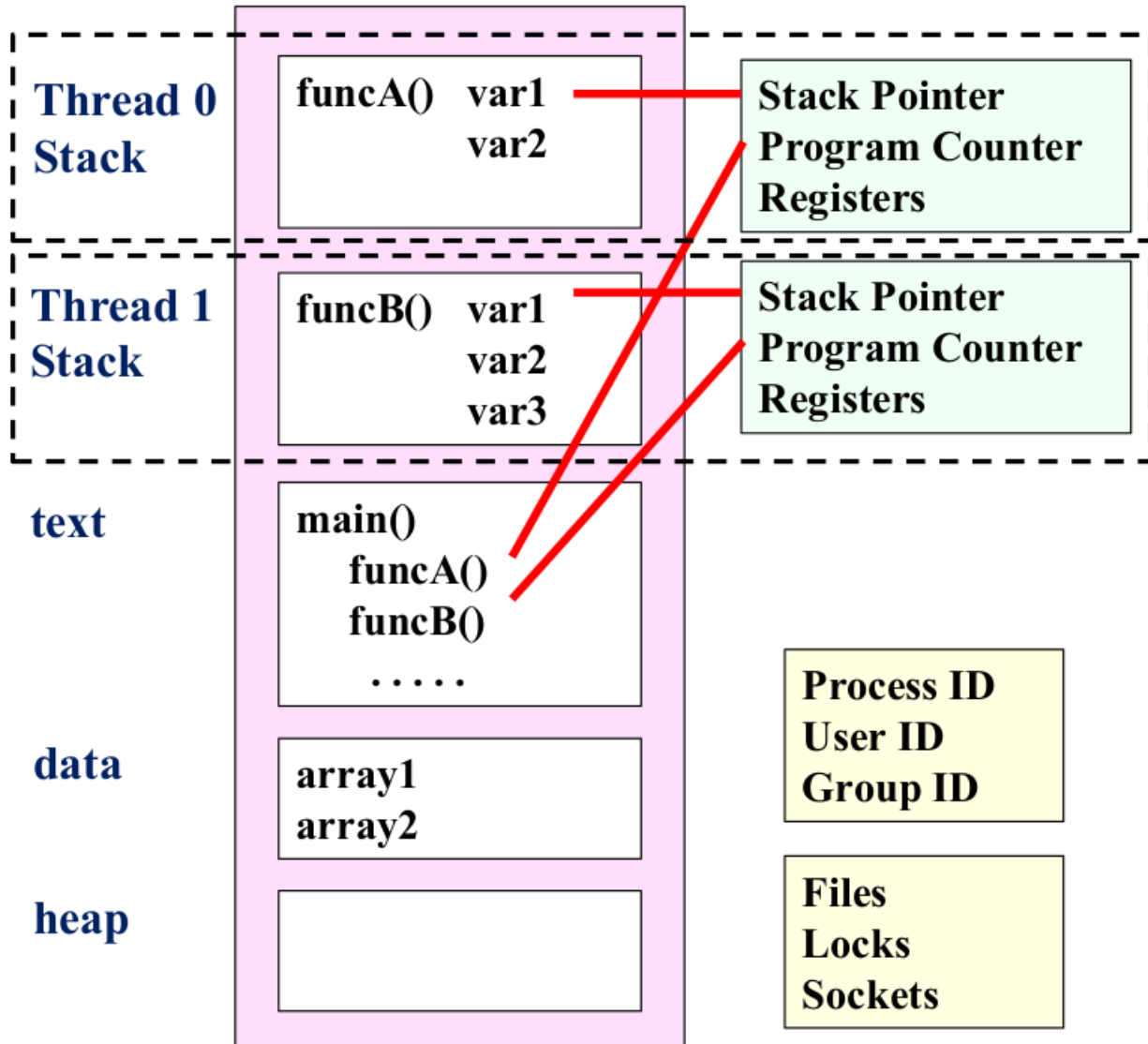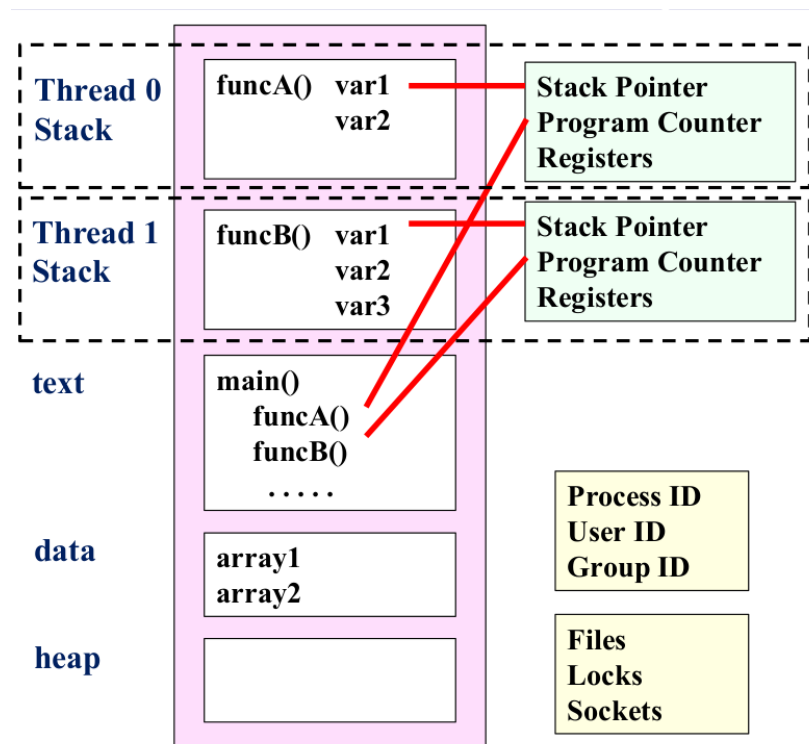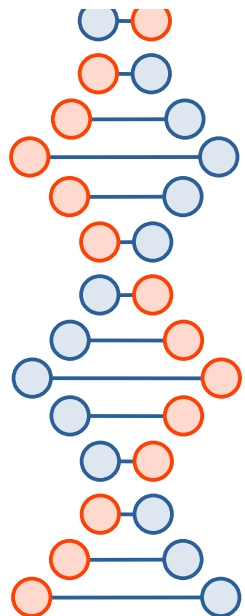  > gcc -fopenmp foo.c

  > export OMP_NUM_THREADS=4

  > ./a.out

for the Bash shell

# OpenMP & NUMA



| Thread 0 Stack | funcA()  var1 var2 | Stack Pointer Program Counter Registers |

| Thread 1 Stack | funcB()  var1 var2 var3 | Stack Pointer Program Counter Registers |

text
main()
    funcA()
    funcB()
    . . . . .

data
array1
array2

heap

Process ID
User ID
Group ID

Files
Locks
Sockets

# OpenMP & NUMA



NUMA Node 0

NUMA Node 1

Interconnect

I/O Bus

I/O Bus

System RAM

| | | |
|---|---|---|
| **Thread 0 Stack** | funcA() var1 var2 | **Stack Pointer Program Counter Registers** |
| **Thread 1 Stack** | funcB() var1 var2 var3 | **Stack Pointer Program Counter Registers** |
| text | **main() funcA() funcB() . . . . .** | |
| data | **array1 array2** | **Process ID User ID Group ID** |
| heap | | **Files Locks Sockets** |

13

# Example: Medium-grained loop parallelism
## Number of systems versus core count



```
do i=1,500
    a(i)=c*b(i)
enddo
```
P1

```
do i=501,1000
    a(i)=c*b(i)
enddo
```
P2

```
do i=1,1000
    a(i)=c*b(i)
enddo
```

**Figure 5.1:** An example for medium-grained parallelism: The iterations of a loop are distributed to two processors P1 and P2 (in shared memory) for concurrent execution.

# 5.2.2 Functional parallelism
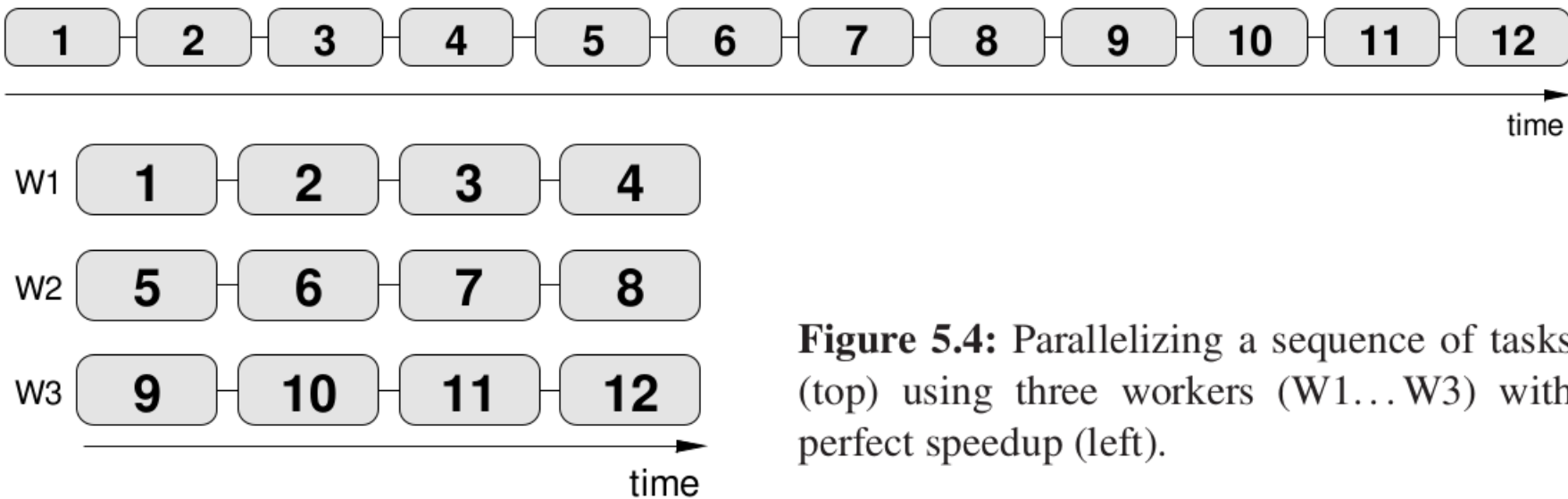# Example: Master-worker scheme



**Figure 5.4:** Parallelizing a sequence of tasks (top) using three workers (W1…W3) with perfect speedup (left).
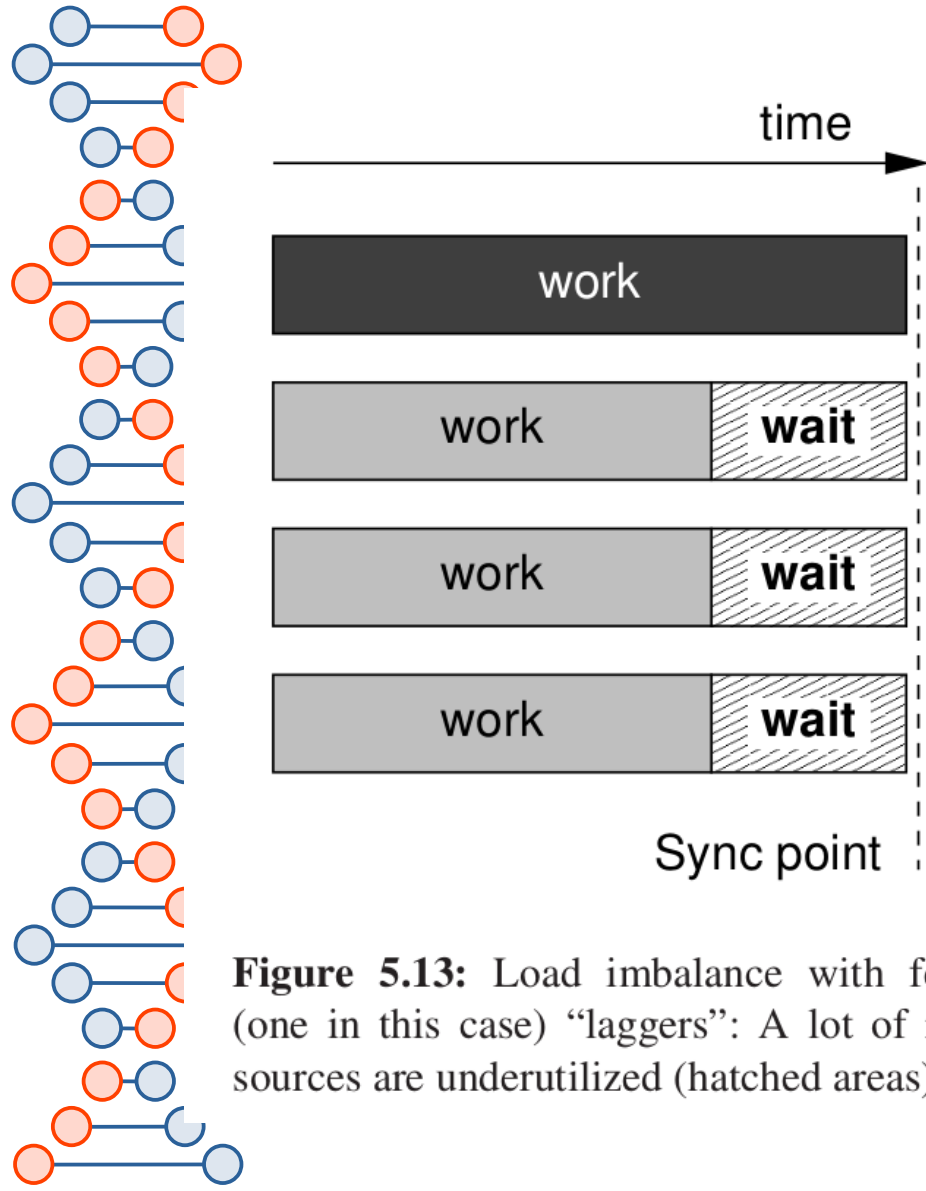
15

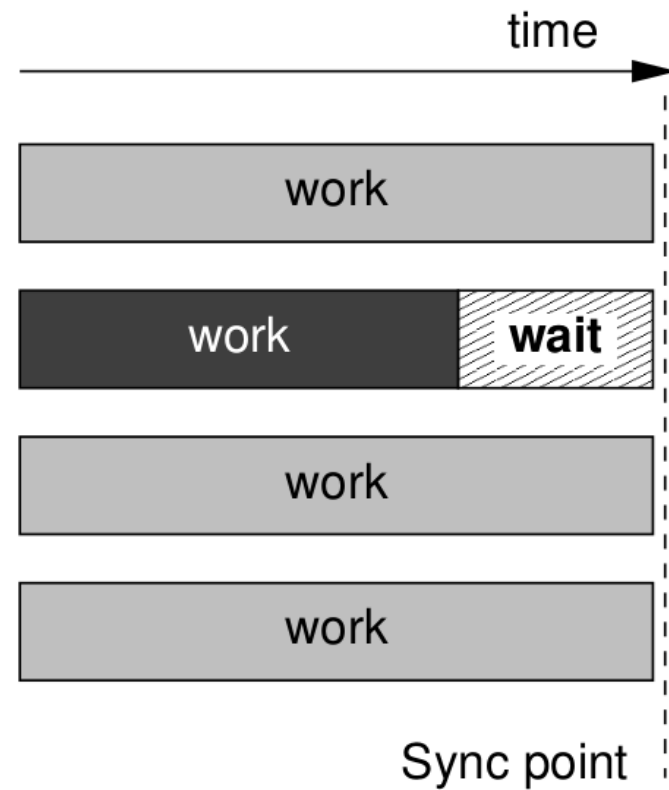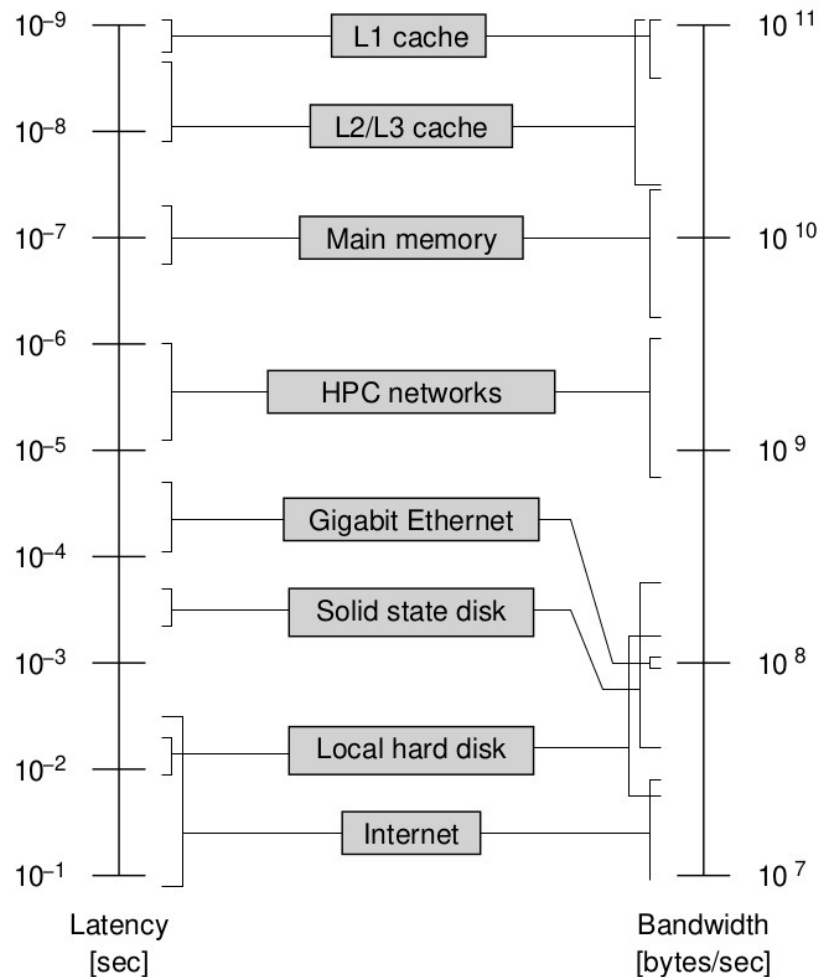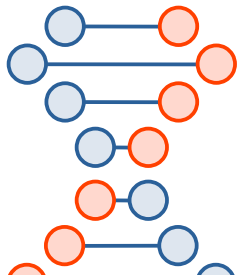**Figure 5.13:** Load imbalance with few (one in this case) "laggers": A lot of resources are underutilized (hatched areas).

**Figure 5.14:** Load imbalance with few (one in this case) "speeders": Underutilization may be acceptable.

16

# Latency vs Bandwidth

$$B_{\mathrm{m}} = \frac{\text{memory bandwidth [GWords/sec]}}{\text{peak performance [GFlops/sec]}} = \frac{b_{\max}}{P_{\max}}$$

# Latency vs Bandwidth

| data path | balance [W/F] |
|---|---|
| cache | 0.5–1.0 |
| **machine (memory)** | 0.03–0.5 |
| interconnect (high speed) | 0.001–0.02 |
| interconnect (GBit ethernet) | 0.0001–0.0007 |
| disk (or disk subsystem) | 0.0001–0.01 |

**Table 3.1:** Typical balance values for operations limited by different transfer paths. In case of network and disk connections, the peak performance of typical dual-socket compute nodes was taken as a basis.
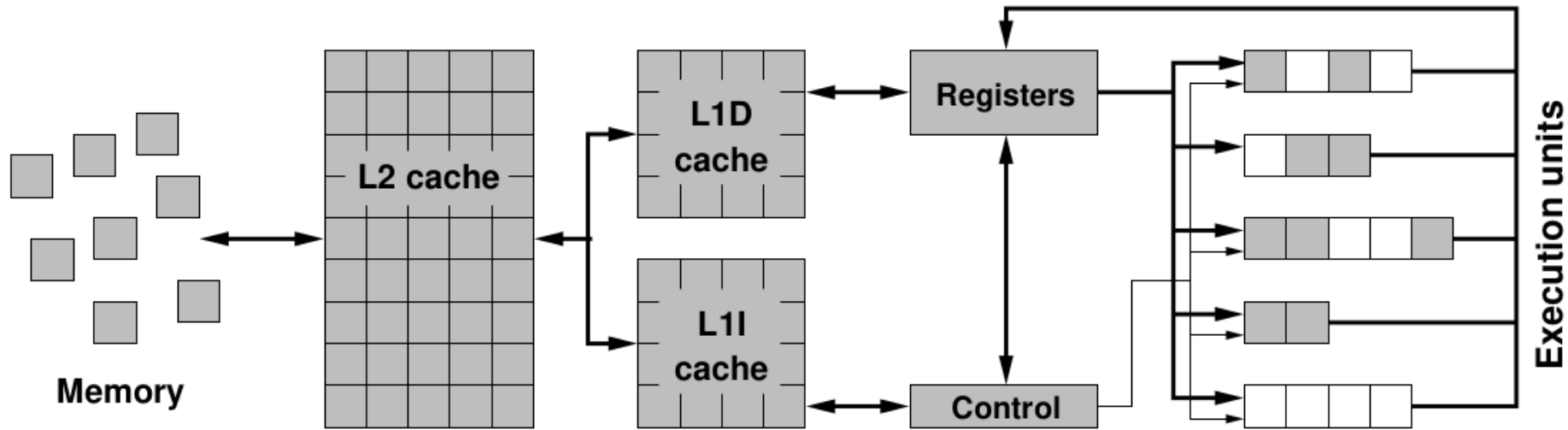
# 1 CPU core



**Figure 1.19:** Simplified diagram of control/data flow in a (multi-)pipelined microprocessor without SMT. White boxes in the execution units denote pipeline bubbles (stall cycles). Graphics by courtesy of Intel.

# Hardware performance counters Per Core

```
 1  CPU Cycles........................................ 8721026107
 2  Retired Instructions.............................. 21036052778
 3  Average number of retired instructions per cycle........ 2.398151
 4  L2 Misses......................................... 101822
 5  Bus Memory Transactions........................... 54413
 6  Average MB/s requested by L2...................... 2.241689
 7  Average Bus Bandwidth (MB/s)...................... 1.197943
 8  Retired Loads..................................... 694058538
 9  Retired Stores.................................... 199529719
10  Retired FP Operations............................. 7134186664
11  Average MFLOP/s................................... 1225.702566
12  Full Pipe Bubbles in Main Pipe.................... 3565110974
13  Percent stall/bubble cycles...................... 40.642963
```
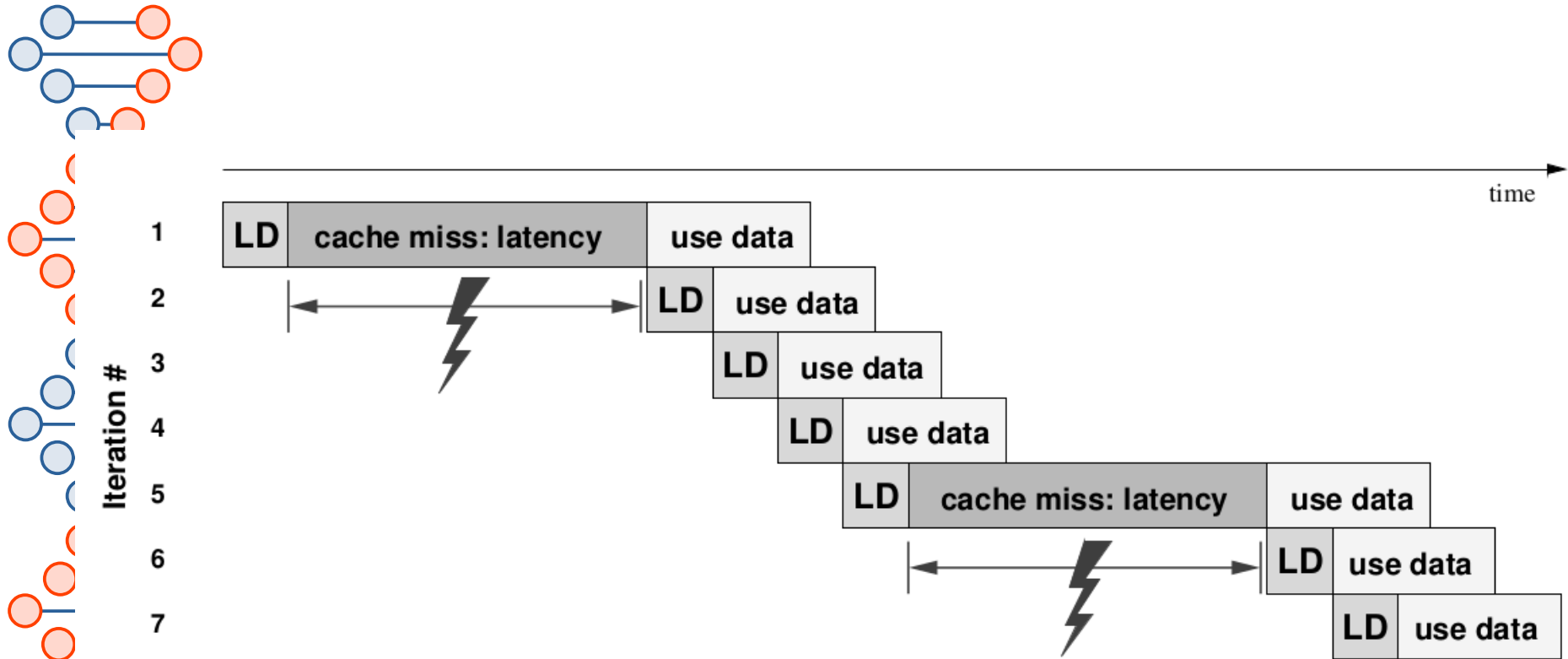
**Figure 1.12:** Timing diagram on the influence of cache misses and subsequent latency penalties for a vector norm loop. The penalty occurs on each new miss.
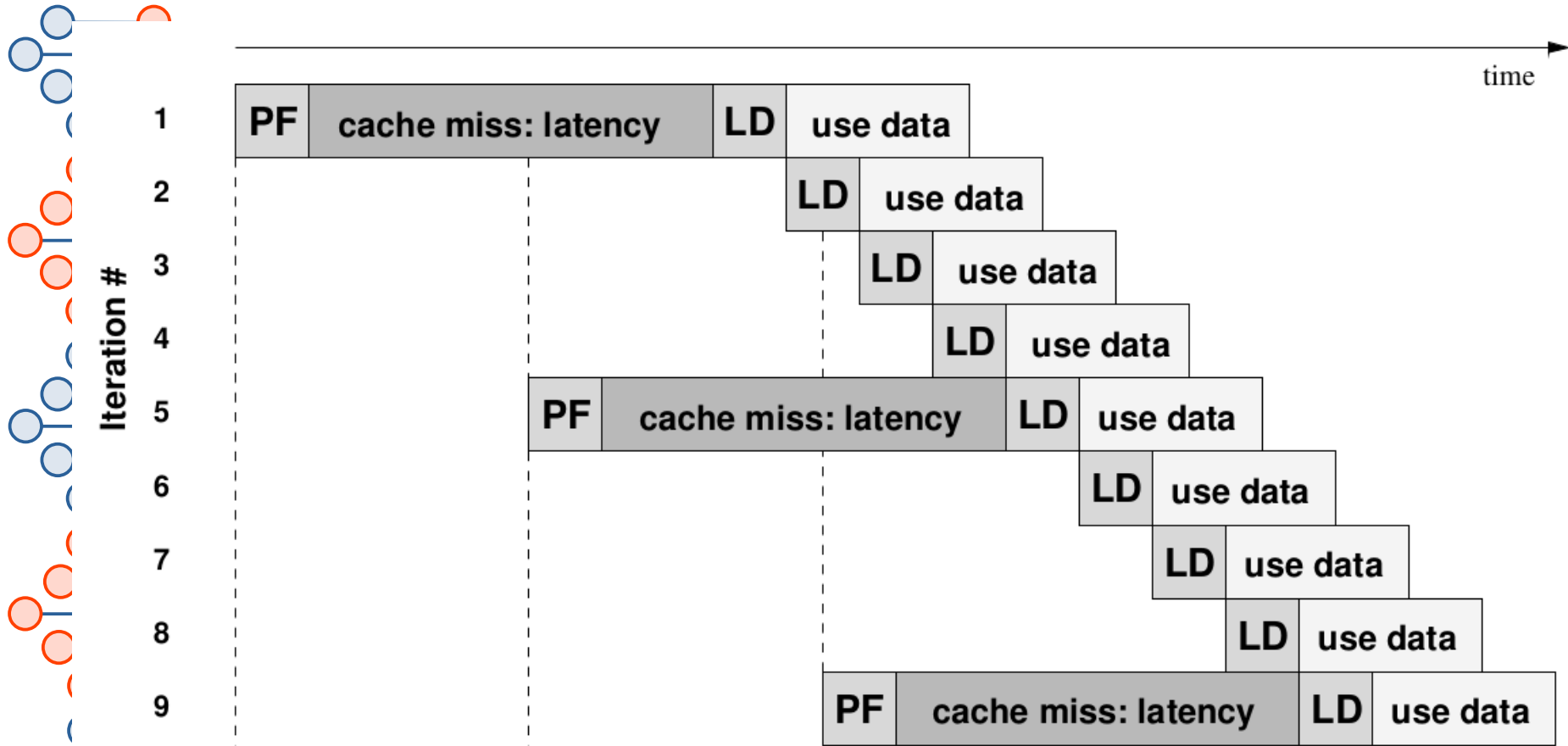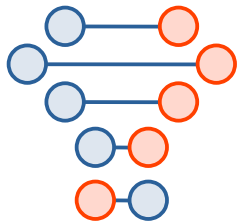
21

**Figure 1.13:** Computation and data transfer can be overlapped much better with prefetching. In this example, two outstanding prefetches are required to hide latency completely.
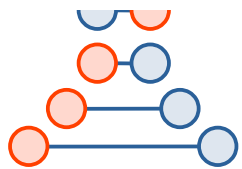
22

# 4.5 Networks

## Point-to-point connections

Whatever the underlying hardware may be, the communication characteristics of a single point-to-point connection can usually be described by a simple model: Assuming that the total transfer time for a message of size $N$ [bytes] is composed of latency and streaming parts,

$$T = T_\ell + \frac{N}{B} \qquad (4.1)$$

and $B$ being the maximum (asymptotic) network bandwidth in MBytes/sec, the effective bandwidth is

$$B_{\text{eff}} = \frac{N}{T_\ell + \frac{N}{B}} \ . \qquad (4.2)$$
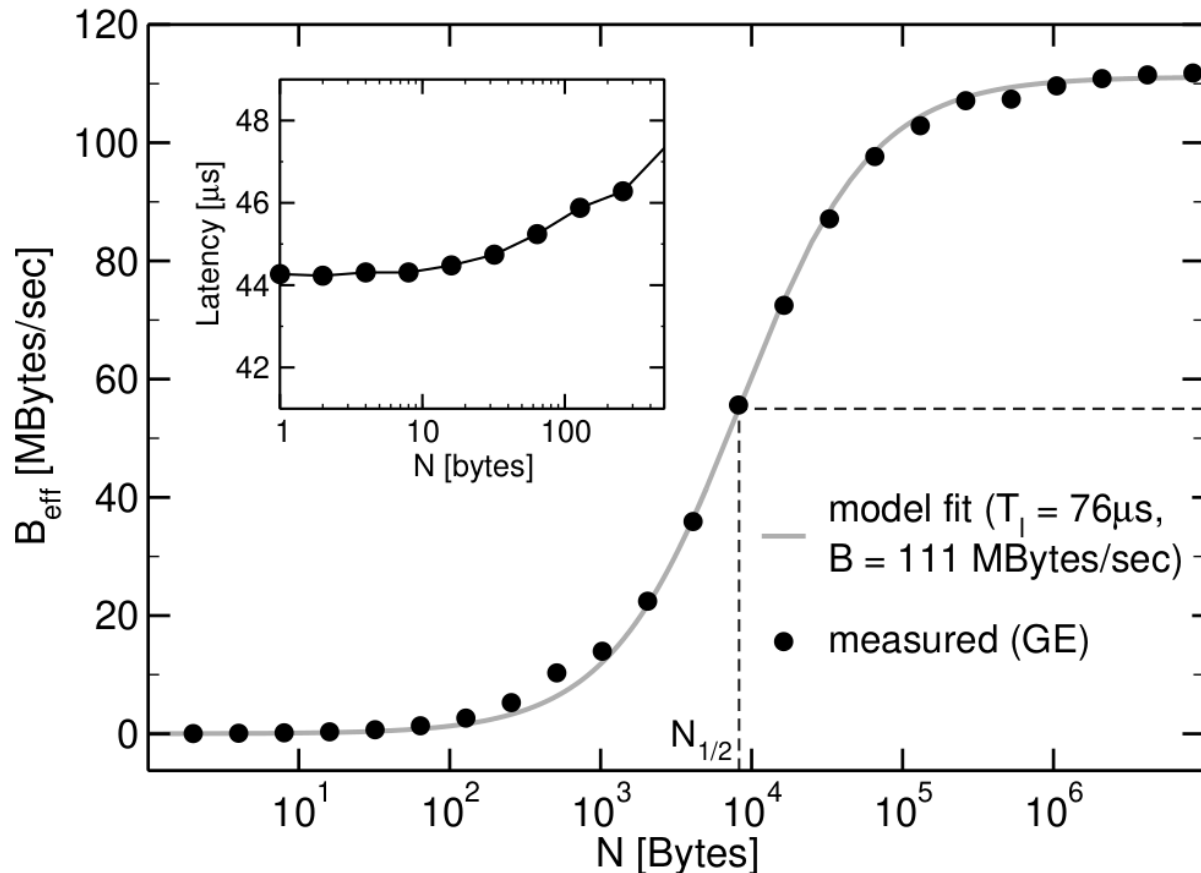
**Figure 4.10:** Fit of the model for effective bandwidth (4.2) to data measured on a GigE network. The fit cannot accurately reproduce the measured value of $T_\ell$ (see text). $N_{1/2}$ is the message length at which half of the saturation bandwidth is reached (dashed line).
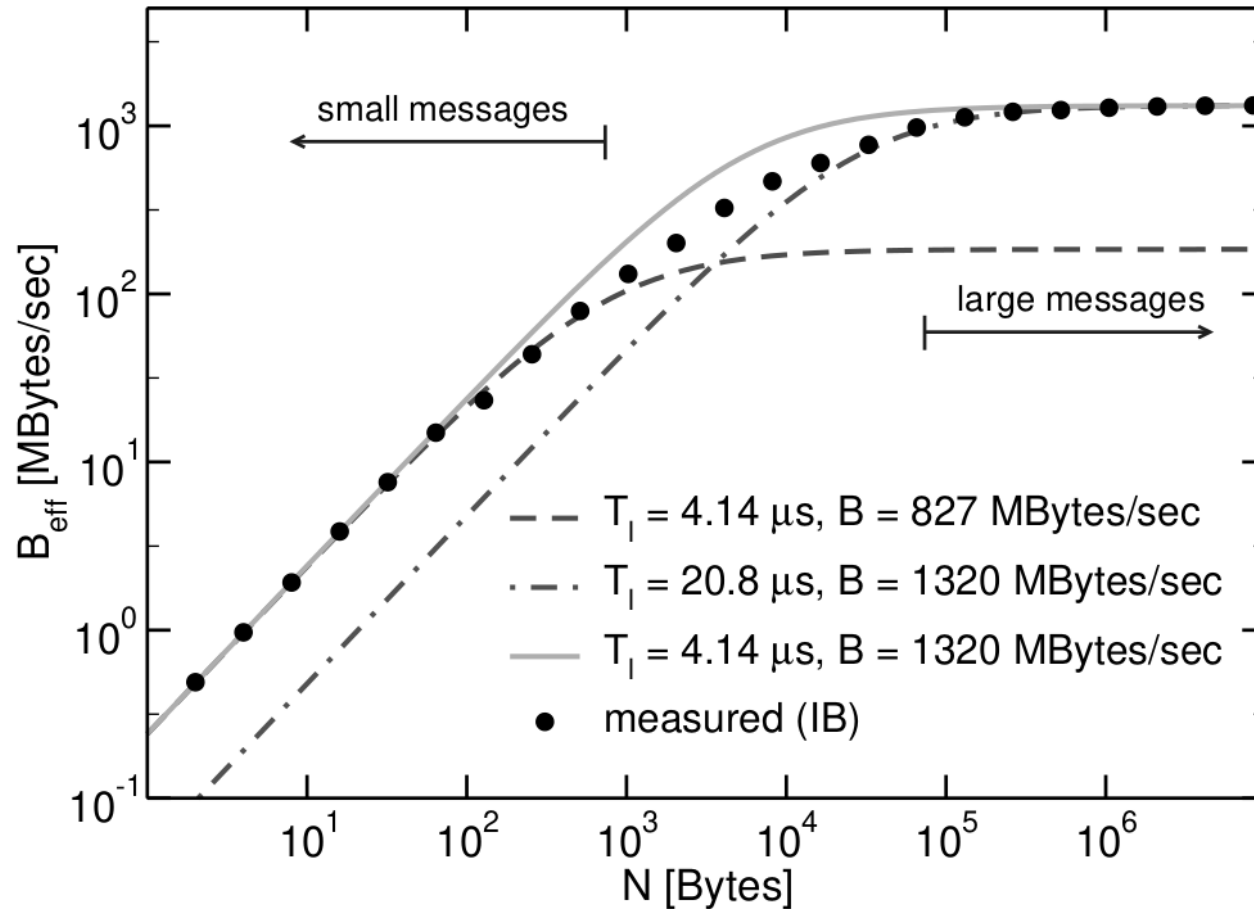
24

**Figure 4.11:** Fits of the model for effective bandwidth (4.2) to data measured on a DDR In-finiBand network. "Good" fits for asymptotic bandwidth (dotted-dashed) and latency (dashed) are shown separately, together with a fit function that unifies both (solid).
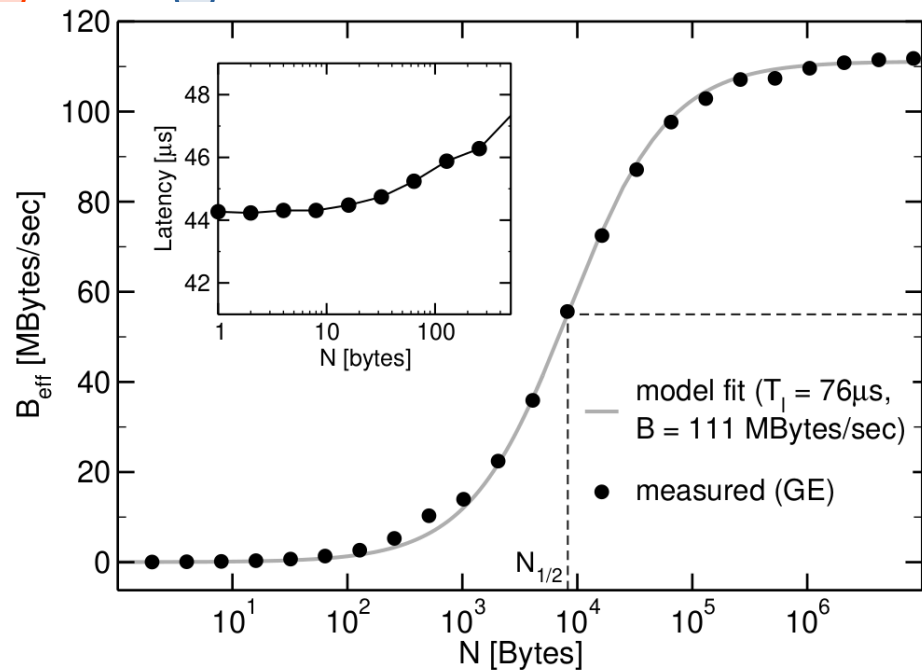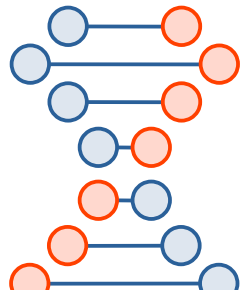
# 4.5 Networks



**Figure 4.10:** Fit of the model for effective bandwidth (4.2) to data measured on a GigE work. The fit cannot accurately reproduce the measured value of $T_\ell$ (see text). $N_{1/2}$ message length at which half of the saturation bandwidth is reached (dashed line).
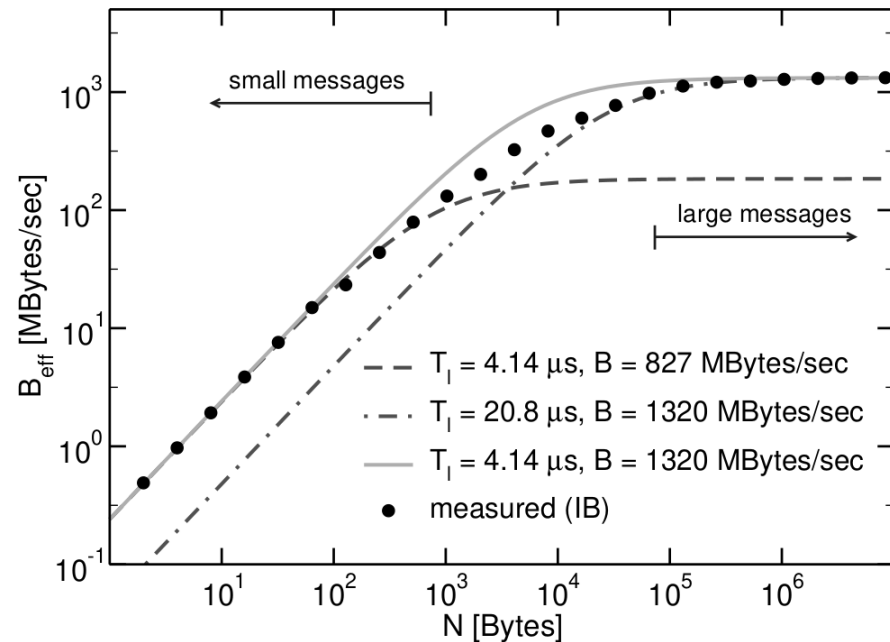
**Figure 4.11:** Fits of the model for effective bandwidth (4.2) to data measured on a DDR In-finiBand network. "Good" fits for asymptotic bandwidth (dotted-dashed) and latency (dashed) are shown separately, together with a fit function that unifies both (solid).