

PARALLEL PROGRAMMING

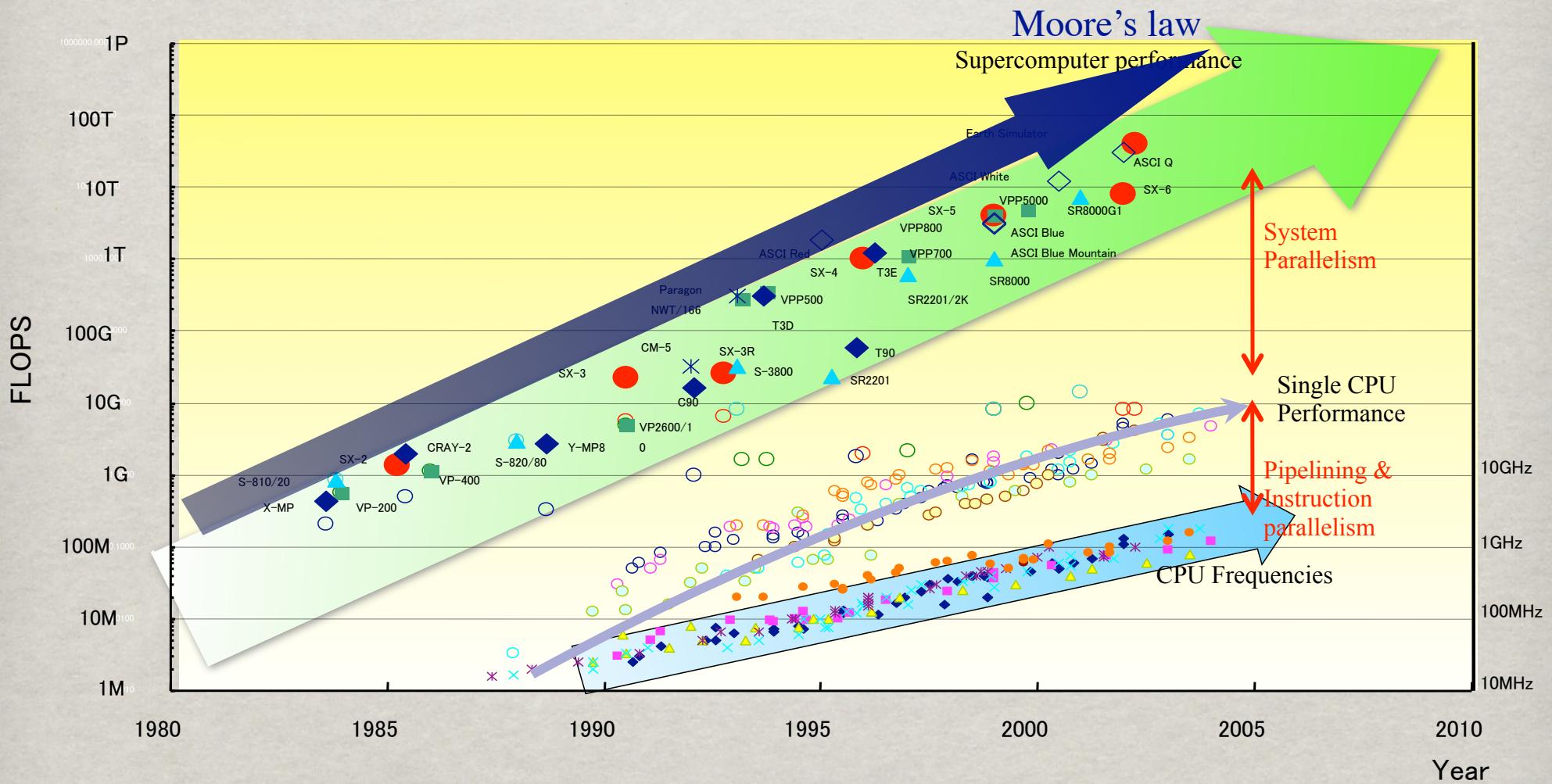
DAY 1

DISTRIBUTED COMPUTING

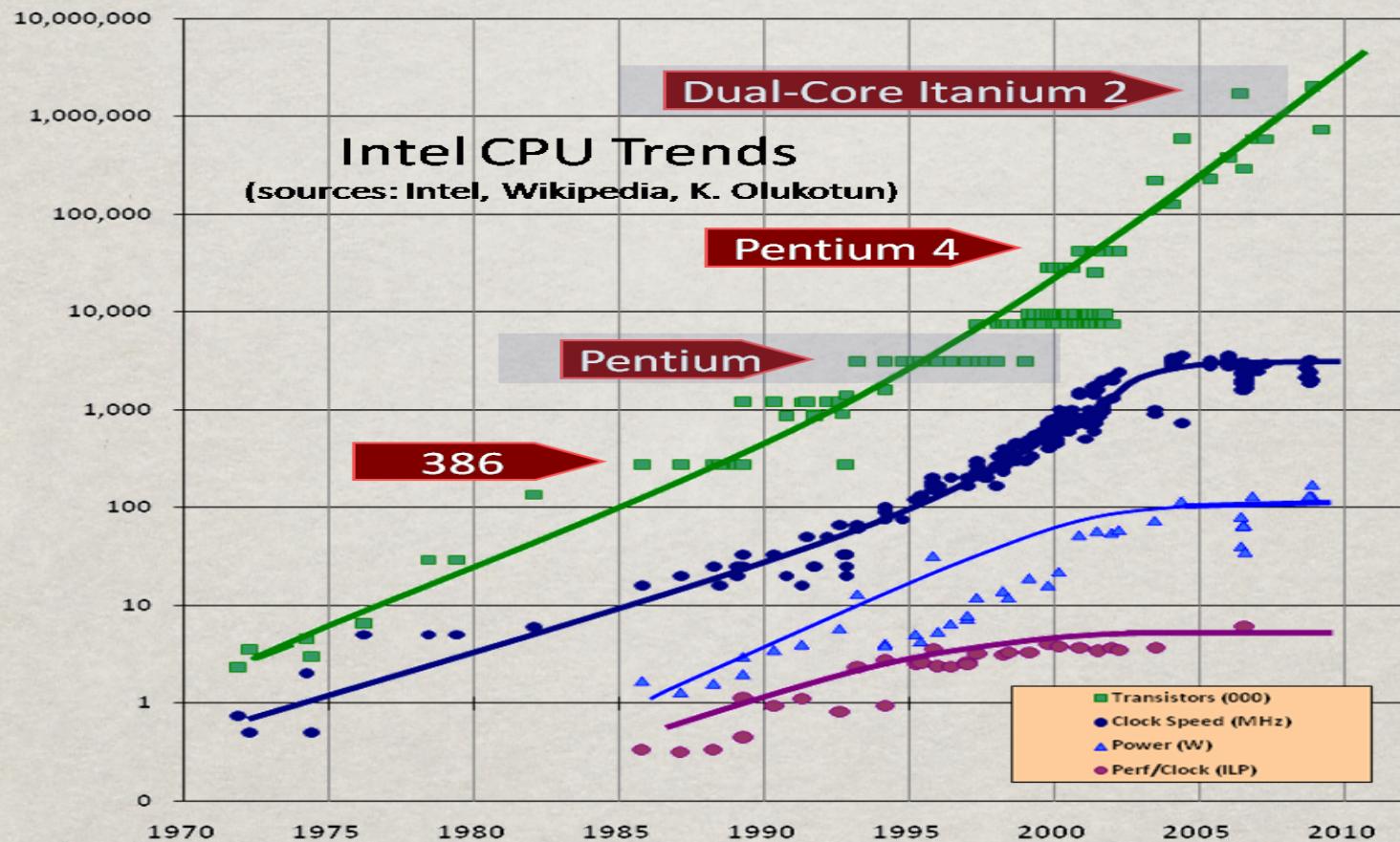
LOGISTICS

- ✿ 10:00 - 12:00 : lectures & questions
- ✿ 13:00 - 17:00 :
practicals & questions & discussion
- ✿ Assistants: Roeland Douma & Roy Bakker

TRENDS



TRENDS



WHY PARALLEL PROGRAMMING?

- ✿ Computing needs: web apps, weather, medical research, gaming, astrophysics, etc.
- ✿ Single processor computers do not get faster fast enough to match the needs
- ✿ Goals:
 - ✿ **Maximize throughput** : computations / second
 - ✿ Minimize overall cost :
maximize computations / euro
 - ✿ NB: energy is a form of cost

PERFORMANCE VS THROUGHPUT

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation
- ✿ **throughput** = computations / second

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation
- ✿ **throughput** = computations / second
- ✿ What's the difference?

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation
- ✿ **throughput** = computations / second
- ✿ What's the difference?
- ✿ Performance very constrained by hardware

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation
- ✿ **throughput** = computations / second
- ✿ What's the difference?
 - ✿ Performance very constrained by hardware
 - ✿ Throughput is more flexible, but requires that **problems can be decomposed** into independent computations

PERFORMANCE VS THROUGHPUT

- ✿ **performance** = seconds / computation
- ✿ **throughput** = computations / second
- ✿ What's the difference?
 - ✿ Performance very constrained by hardware
 - ✿ Throughput is more flexible, but requires that **problems can be decomposed** into independent computations
 - ✿ Limits: Amdahl's and Gustafson's laws
(check them out!)

PART 1: PLATFORMS

COMPUTERS

This is a computer:



This too:



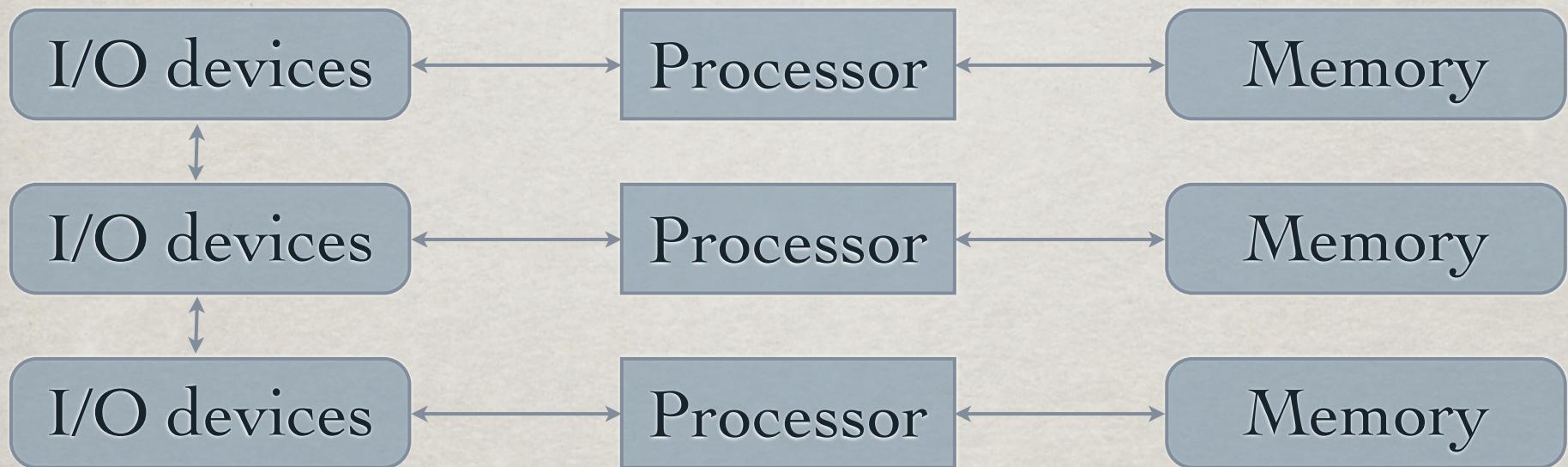
This too:



Can you recognize examples of each?

COMPUTERS

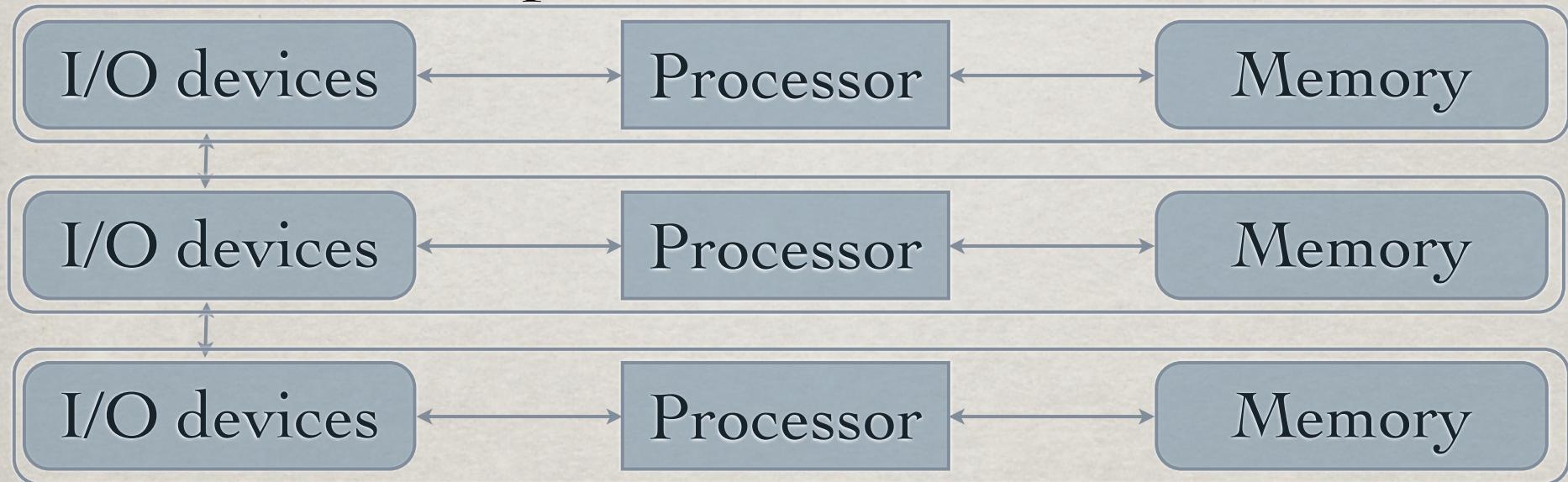
This is also a “computer”:



... but is commonly called a **distributed system**

COMPUTERS

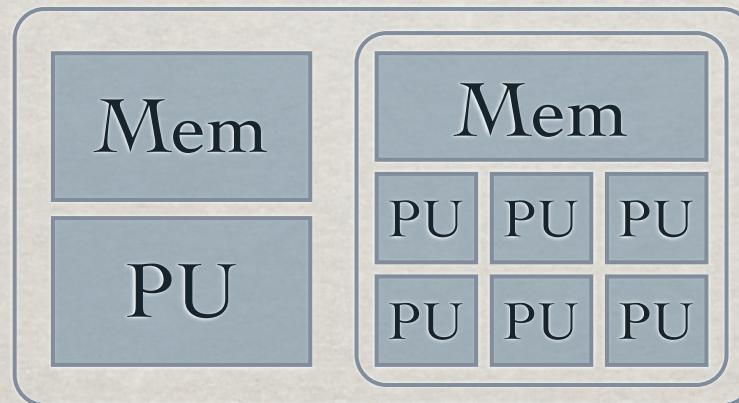
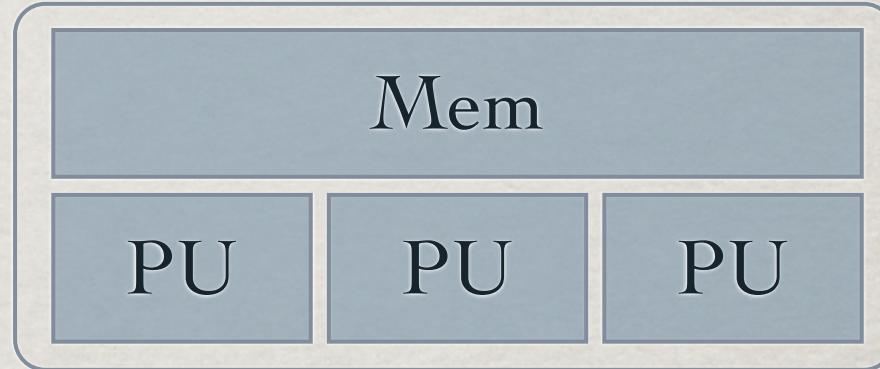
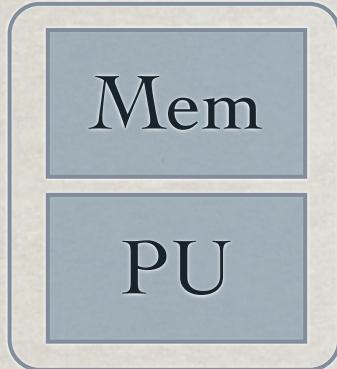
This is also a “computer”:



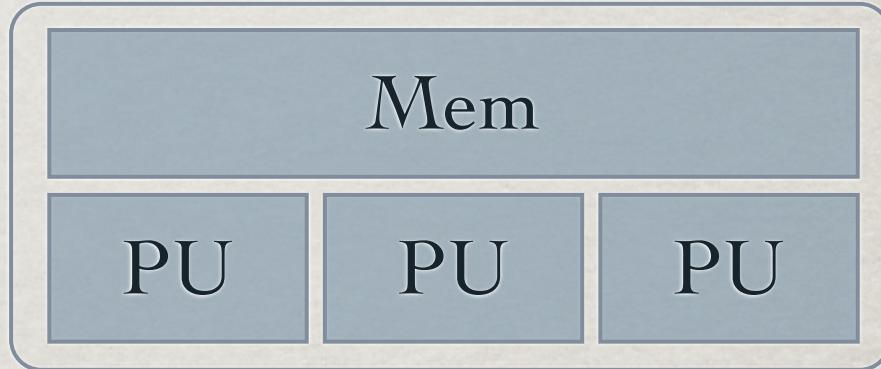
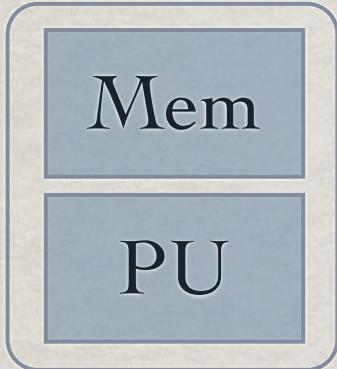
... but is commonly called a **distributed system**

“Node” =
group of PU(s) + Memory(ies) sharing an I/O interface

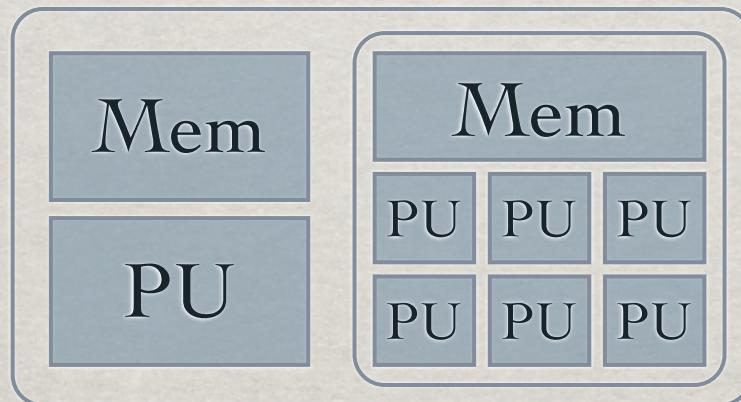
TOPOLOGIES



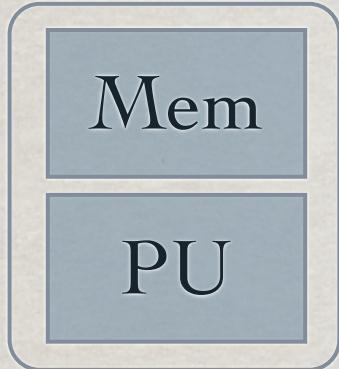
TOPOLOGIES



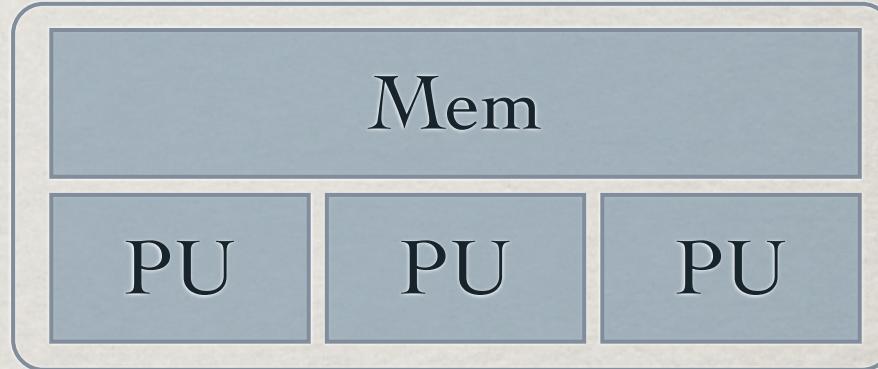
Uni-processor



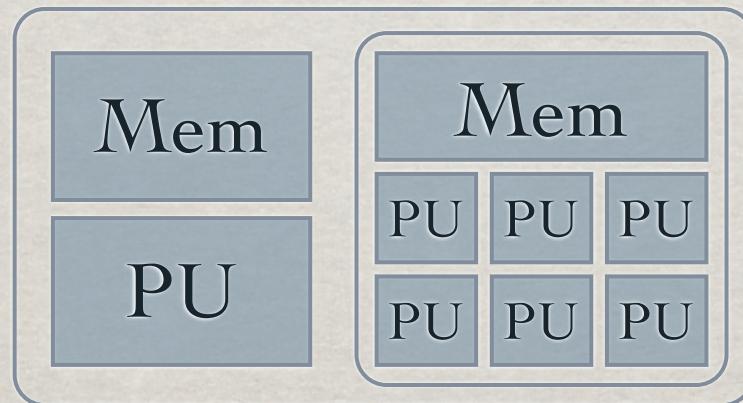
TOPOLOGIES



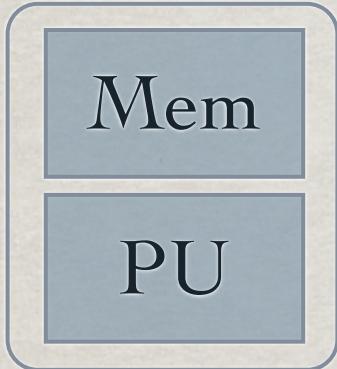
Uni-processor



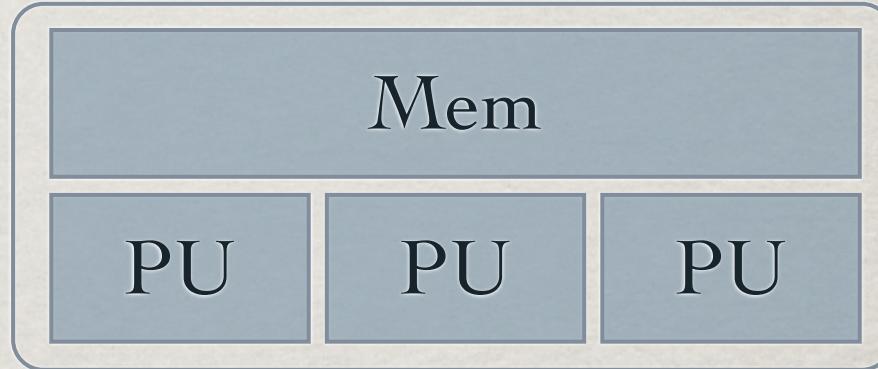
Multi-processor
Multi-core



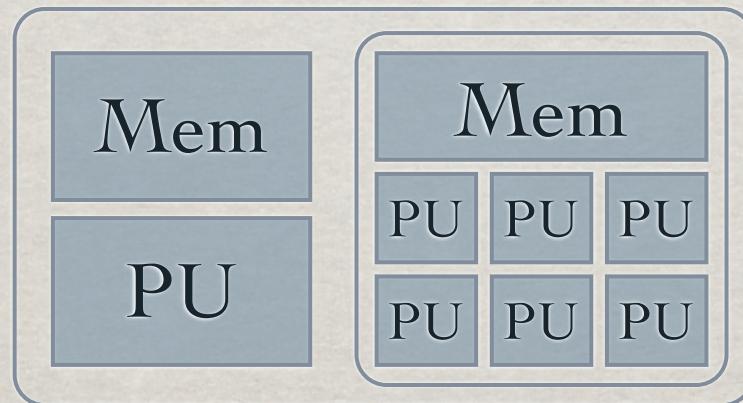
TOPOLOGIES



Uni-processor



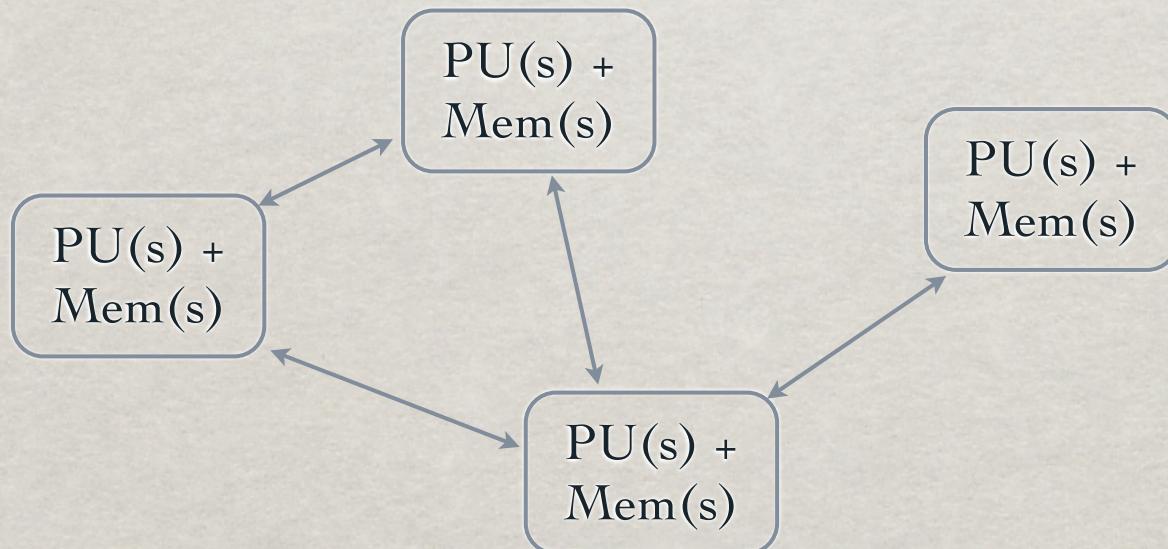
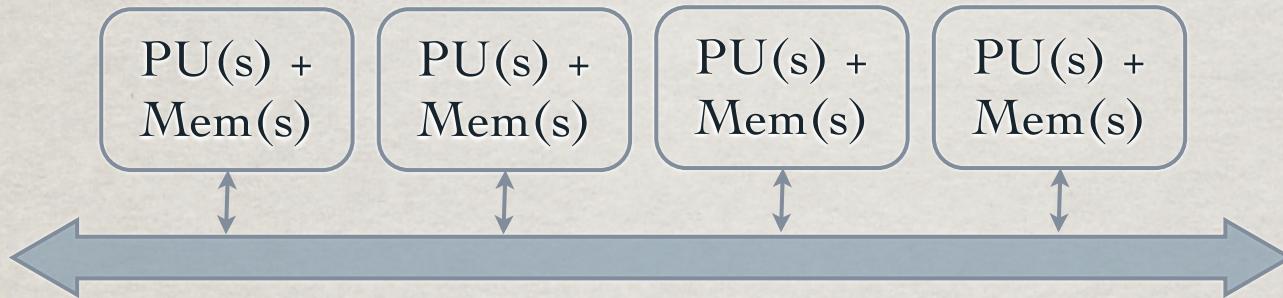
Multi-processor
Multi-core



Heterogeneous multi-cores
Accelerators

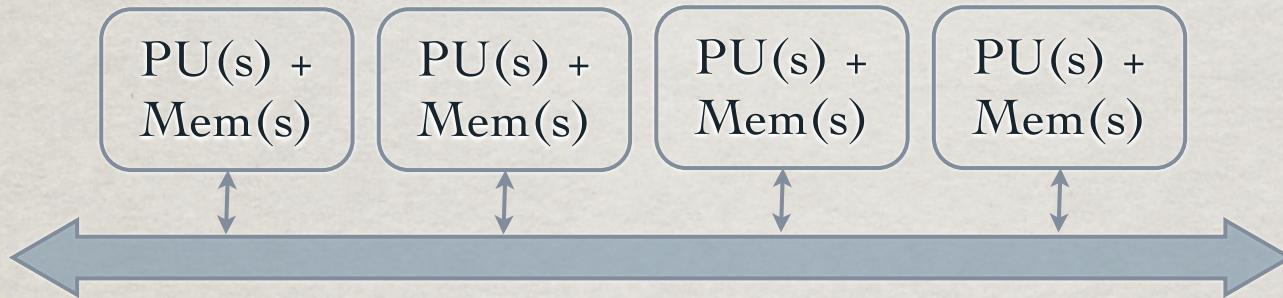
TOPOLOGIES FOR DISTRIBUTED SYSTEMS

“nodes”

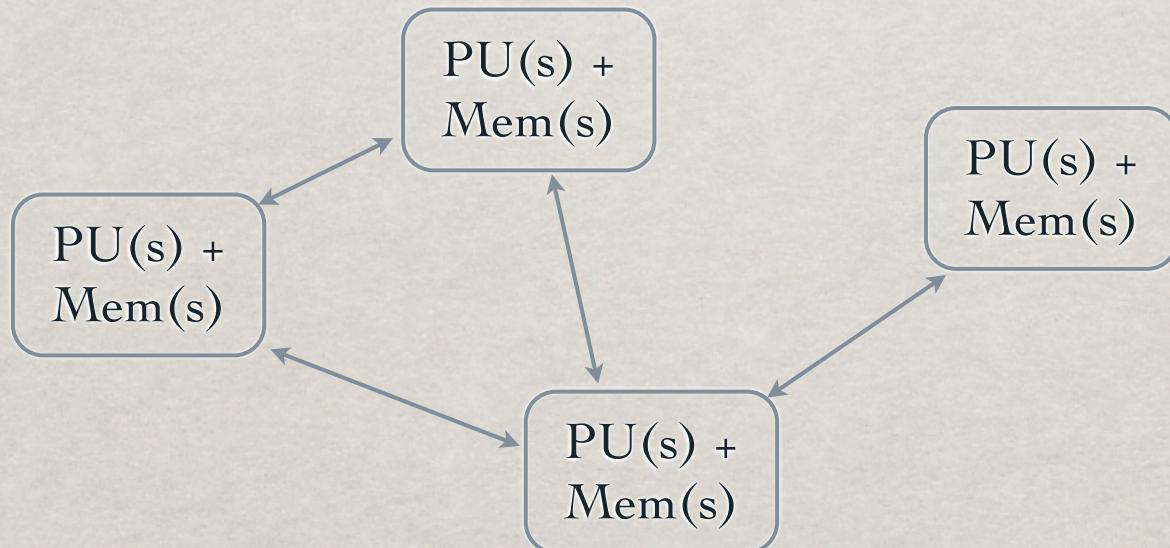


TOPOLOGIES FOR DISTRIBUTED SYSTEMS

“nodes”

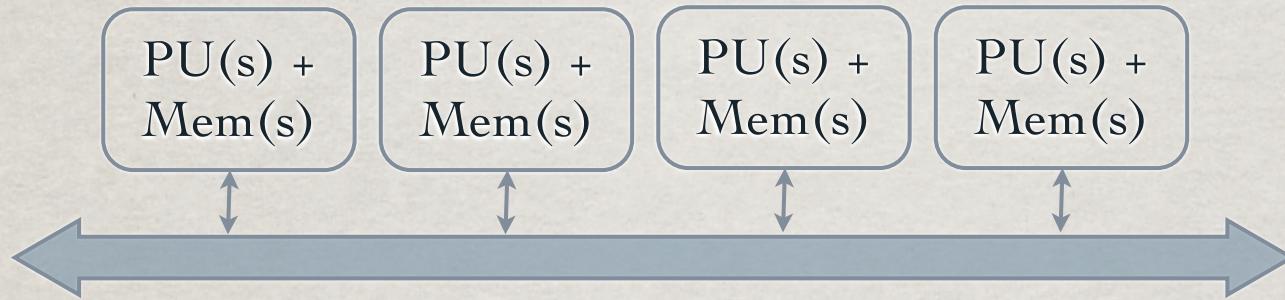


Regular interconnect: clusters, grids

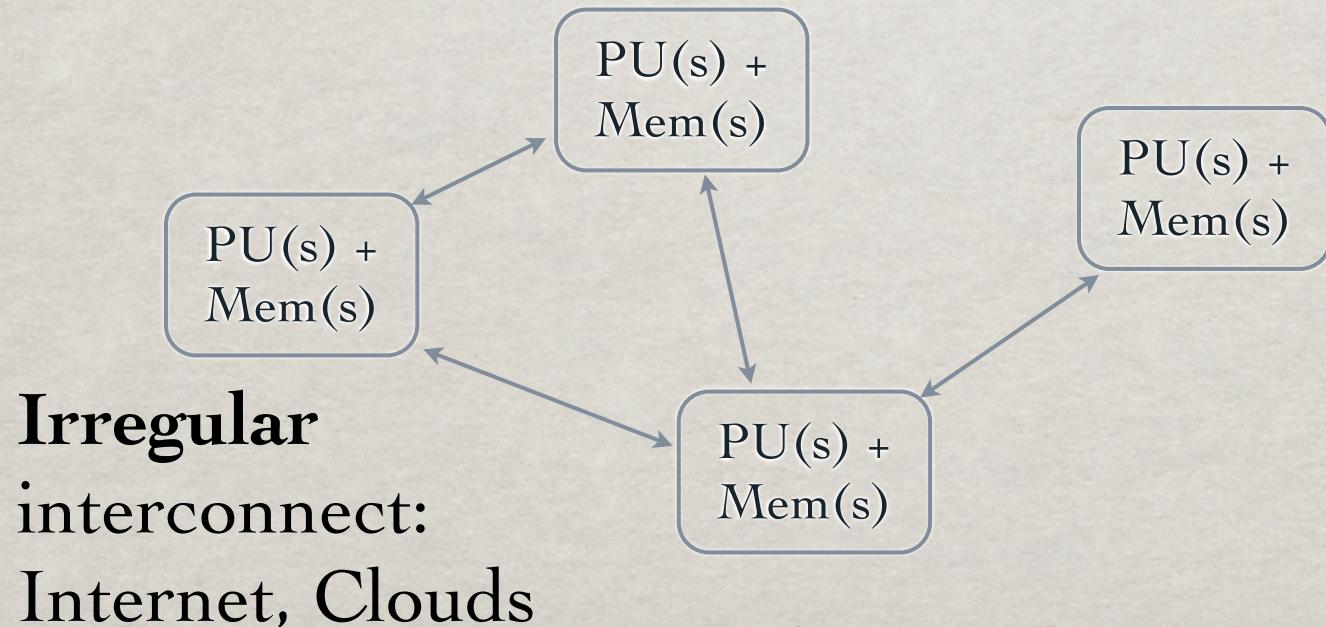


TOPOLOGIES FOR DISTRIBUTED SYSTEMS

“nodes”



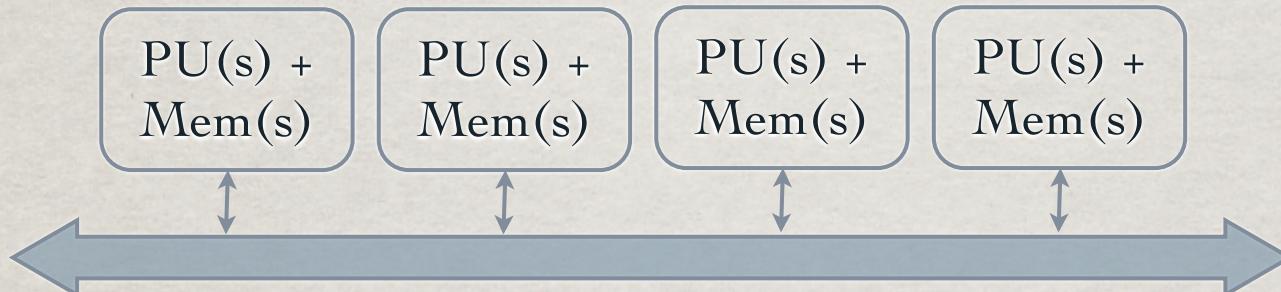
Regular interconnect: clusters, grids



Irregular
interconnect:
Internet, Clouds

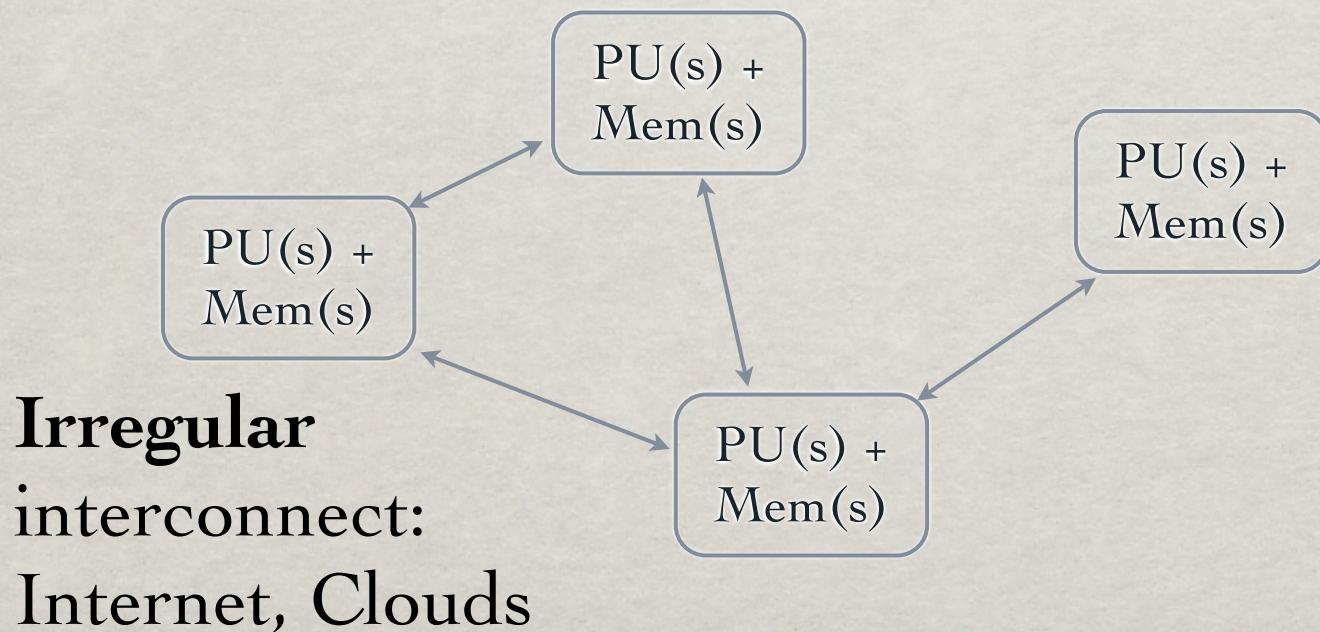
TOPOLOGIES FOR DISTRIBUTED SYSTEMS

“nodes”



Regular interconnect: clusters, grids

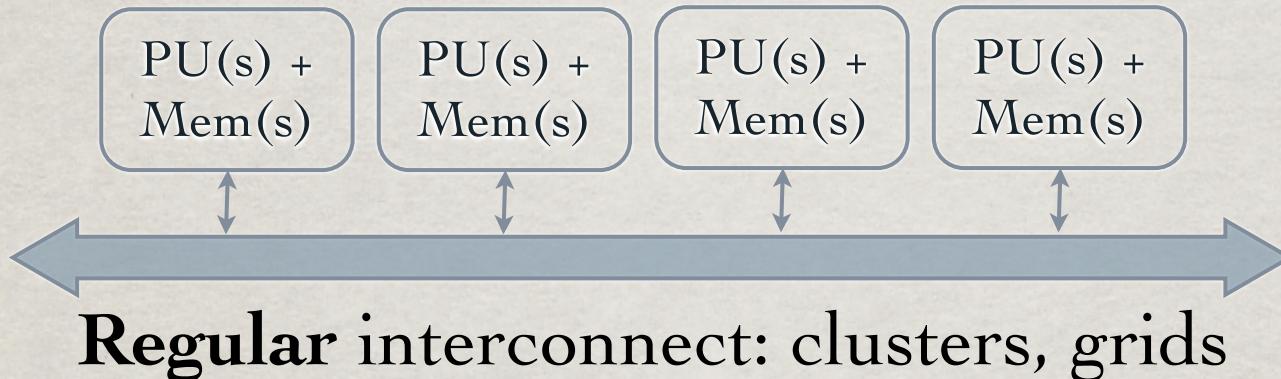
Easier to
program



Irregular
interconnect:
Internet, Clouds

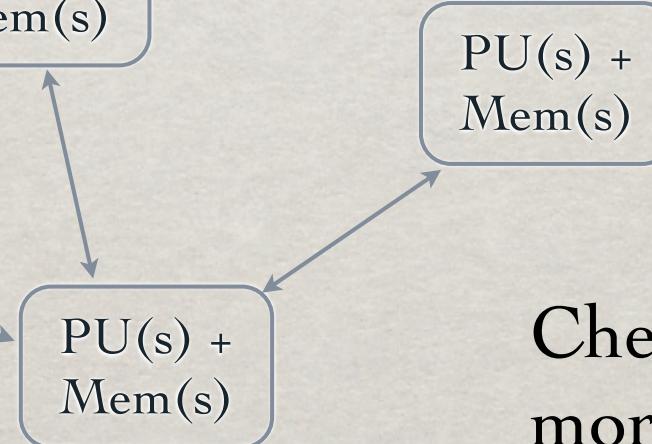
TOPOLOGIES FOR DISTRIBUTED SYSTEMS

“nodes”



Easier to
program

Irregular
interconnect:
Internet, Clouds



Cheaper,
more robust to faults

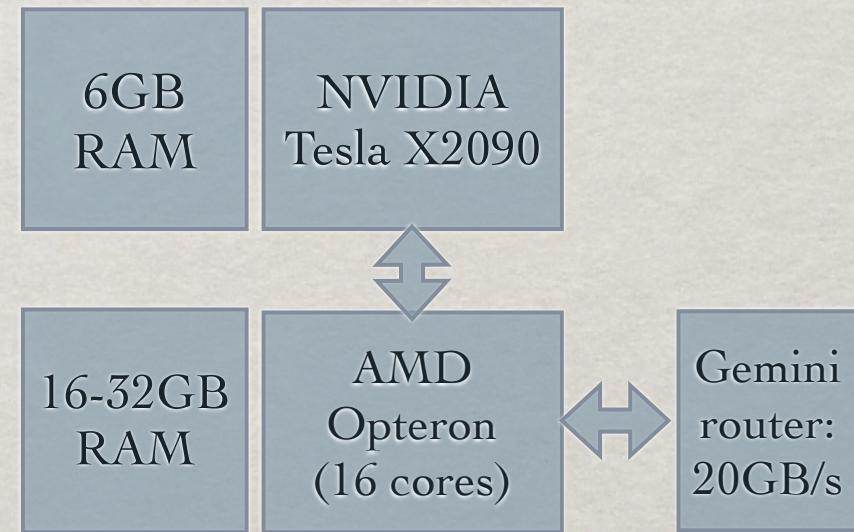
EXAMPLES



XE6 blade

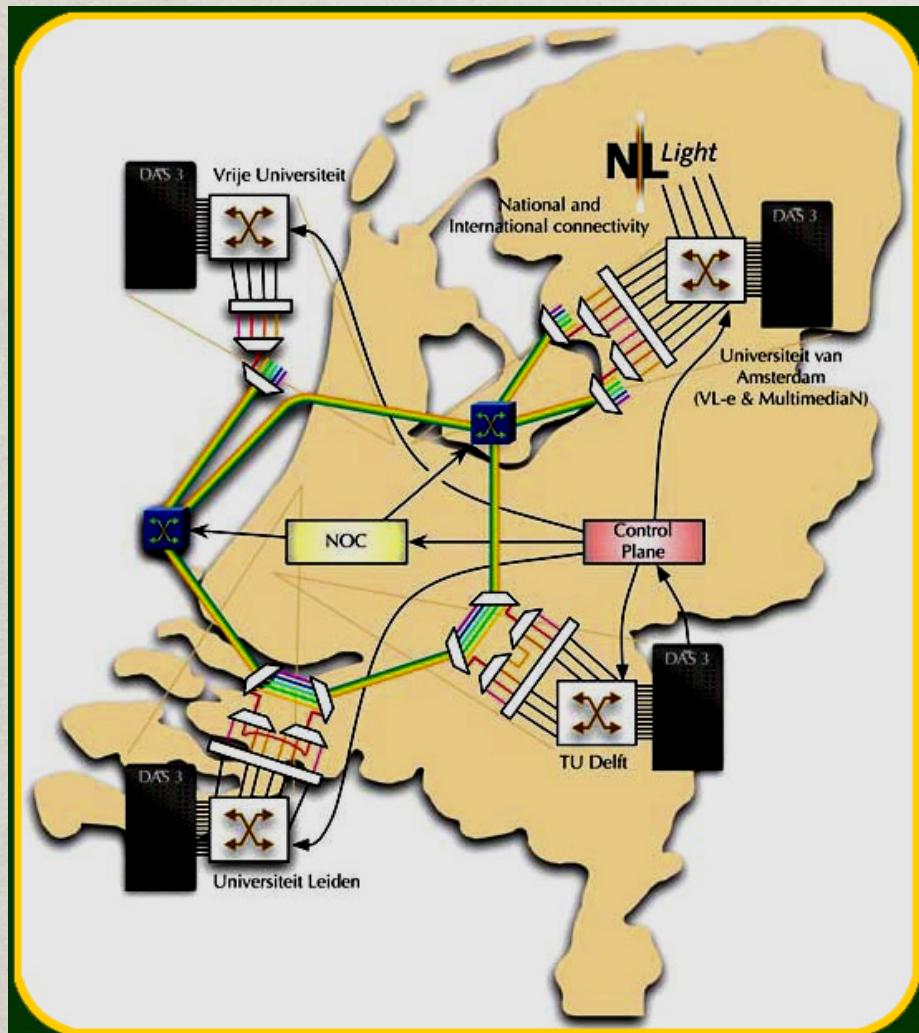
Cray XE6:

- 24 blades per cabinet
- 4 nodes per blade
- 16 cores + GPU per node
- 3D regular torus interconnect



XE6 node
inside bandwidth: 52GB/s

EXAMPLES



DAS: Distributed ASCI Computer

DAS-4:

- 4 sites
- 16-74 nodes per site
- 8 cores + accels per node
- 24-48GB RAM per node
- GbE network for control
- InfiniBand network for throughput: 20GB/s

SUMMARY

SUMMARY

- ✿ Characteristics of a distributed system?

SUMMARY

- ✿ Characteristics of a distributed system?
- ✿ Two levels: “local” (nodes) and “network”

SUMMARY

- ✿ Characteristics of a distributed system?
- ✿ Two levels: “local” (nodes) and “network”
- ✿ Lots (GBs/TBs) of local memory/storage per node

SUMMARY

- ✿ Characteristics of a distributed system?
 - ✿ Two levels: “local” (nodes) and “network”
 - ✿ Lots (GBs/TBs) of local memory/storage per node
 - ✿ Access times: fast to local node, longer to other nodes, typically 2x-100x difference

SUMMARY

- ✿ Characteristics of a distributed system?
 - ✿ Two levels: “local” (nodes) and “network”
 - ✿ Lots (GBs/TBs) of local memory/storage per node
 - ✿ Access times: fast to local node, longer to other nodes, typically 2x-100x difference
 - ✿ Variable size by adding/removing nodes

PART 2: BASICS OF PARALLEL PROGRAMMING

ESSENCE OF PARALLEL PROGRAMMING

- ✿ Maximize throughput : computations / second
- ✿ Requires to **decompose problems** into separate sub-computations
- ✿ Two separate **phases**:
 - ✿ **Exposing concurrency**:
find and describe the sub-problems
(done by the human)
 - ✿ **Parallel execution**:
distribute the sub-problems on the platform
(shared responsibility between human and computer)

EXAMPLE: DATA PARALLELISM

- ✿ Multiply two vectors:

$$S = A[0]*B[0] + A[1]*B[1] + \dots + A[n]*B[n]$$

- ✿ Sequential program:

# Python	/* C */
s = 0	s = 0;
for i in range(n):	for (i = 0; i < n; ++i)
s += a[i] * b[i]	s += a[i] * b[i];

- ✿ Can we not do better?

EXAMPLE: DATA PARALLELISM

- ❖ Decomposition: each $A[i]^*B[i]$ can be **computed independently**
- ❖ Solution:
 - ❖ decompose in p partial sums; each partial sum k has n/p steps to compute:

```
def partialsum(k, p, n, A, B):  
    s = 0  
    for i in range(k * n/p, (k + 1) * n/p):  
        s = A[i] * B[i]  
    return s
```

- ❖ all partial sums can execute in parallel; the partial results can be summed at the end

GRANULARITY

GRANULARITY

- n = problem size

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations
- ✿ n/p = “how much work to do per sub-computation”
= **granularity**

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations
- ✿ n/p = “how much work to do per sub-computation”
= **granularity**
- ✿ Trade offs?

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations
- ✿ n/p = “how much work to do per sub-computation”
= **granularity**
- ✿ Trade offs?
 - ✿ **p small: less concurrency**
less opportunity for speed-up, **lower throughput**

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations
- ✿ n/p = “how much work to do per sub-computation”
= **granularity**
- ✿ Trade offs?
 - ✿ **p small: less concurrency**
less opportunity for speed-up, **lower throughput**
 - ✿ **p large: more communication and synchronization** may introduce overhead that **limit performance**

GRANULARITY

- ✿ n = problem size
- ✿ p = number of independent sub-computations
- ✿ n/p = “how much work to do per sub-computation”
= **granularity**
- ✿ Trade offs?
 - ✿ **p small: less concurrency**
less opportunity for speed-up, **lower throughput**
 - ✿ **p large: more communication and synchronization** may introduce overhead that **limit performance**
- ✿ In practice: choose p to match the number of PUs

DATA PARALLELISM: SIMD / SPMD

DATA PARALLELISM: SIMD / SPMD

- ❖ What defines data parallelism?

DATA PARALLELISM: SIMD / SPMD

- ✿ What defines data parallelism?
- ✿ The same function (instruction/program) is applied to each item in the input data: SI/SP

DATA PARALLELISM: SIMD / SPMD

- ✿ What defines data parallelism?
- ✿ The same function (instruction/program) is applied to each item in the input data: SI/SP
- ✿ The large number of items in the input data provides concurrency: MD

DATA PARALLELISM: SIMD / SPMD

- ✿ What defines data parallelism?
- ✿ The same function (instruction/program) is applied to each item in the input data: SI/SP
- ✿ The large number of items in the input data provides concurrency: MD
- ✿ General form: `for i: Output[i] = f(Input[i])`
Also noted: `Out = parallel_map(n, p, f, Input)`

REDUCTIONS

- ✿ Typically a data parallel computation
is followed by a reduction

- ✿ In the previous example:

```
S = parallel_map(n, p, partialsum, (A, B))
```

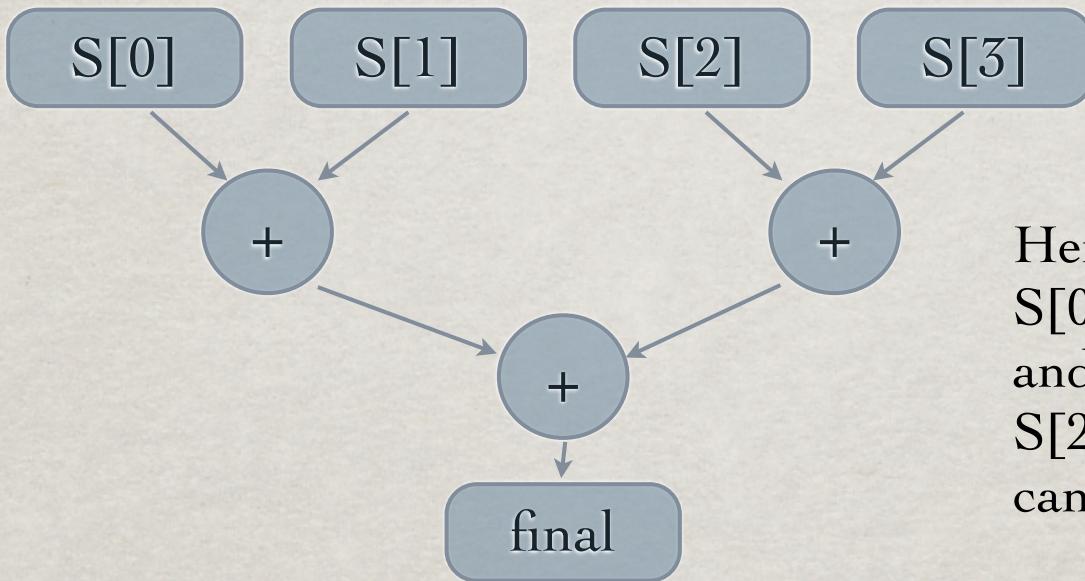
```
final = 0
for k in range(p): # reduce
    final += S[k]
```

- ✿ General form of a reduction:

$$X = A[0] \odot A[1] \odot \dots \odot A[n]$$

REDUCTIONS

- Reductions can be partially parallelized too!



Here
 $S[0]+S[1]$
and
 $S[2]+S[3]$
can be computed in parallel

- Condition: the operation must be **associative** and **commutative**, eg: min, max, +, *, and, or
- Noted also: `Out = parallel_reduce(n, p, op, Input)`

EXAMPLE MAP/REDUCE (SIMPLIFIED)

EXAMPLE MAP/REDUCE (SIMPLIFIED)

- Sequential:

```
def prod(n, A, B):  
    sum = 0  
    for i in range(n):  
        sum += A[i] * B[i]  
    return sum
```

EXAMPLE MAP/REDUCE (SIMPLIFIED)

- ✿ Sequential:

```
def prod(n, A, B):  
    sum = 0  
    for i in range(n):  
        sum += A[i] * B[i]  
    return sum
```

- ✿ Parallel:

```
def parallel_prod(n, A, B, p):  
    return parallel_reduce(n, p, __add__,  
                           parallel_map(n, p, partialsum, (A, B)))
```

EXAMPLE MAP/REDUCE (SIMPLIFIED)

- ✿ Sequential:

```
def prod(n, A, B):  
    sum = 0  
    for i in range(n):  
        sum += A[i] * B[i]  
    return sum
```

- ✿ Parallel:

```
def parallel_prod(n, A, B, p):  
    return parallel_reduce(n, p, __add__,  
                           parallel_map(n, p, partialsum, (A, B)))
```

- ✿ How many steps? eg with $n=10^6$ and $p=1000$

EXAMPLE MAP/REDUCE (SIMPLIFIED)

- ✿ Sequential:

```
def prod(n, A, B):  
    sum = 0  
    for i in range(n):  
        sum += A[i] * B[i]  
    return sum
```

- ✿ Parallel:

```
def parallel_prod(n, A, B, p):  
    return parallel_reduce(n, p, __add__,  
                           parallel_map(n, p, partialsum, (A, B)))
```

- ✿ How many steps? eg with $n=10^6$ and $p=1000$

- ✿ Sequential: **n steps**, eg 10^6 steps

EXAMPLE MAP/REDUCE (SIMPLIFIED)

- ✿ Sequential:

```
def prod(n, A, B):  
    sum = 0  
    for i in range(n):  
        sum += A[i] * B[i]  
    return sum
```

- ✿ Parallel:

```
def parallel_prod(n, A, B, p):  
    return parallel_reduce(n, p, __add__,  
                           parallel_map(n, p, partialsum, (A, B)))
```

- ✿ How many steps? eg with $n=10^6$ and $p=1000$

- ✿ Sequential: **n steps**, eg 10^6 steps

- ✿ **p parallel PUs:**

- $\frac{n}{p}$ steps for map + $\log(p)$ steps for reduce = 1006 steps

EXPOSING CONCURRENCY WITH DATA PARALLELISM

- ✿ **Role of the programmer:**
 - ✿ **Find and describe the independent sub-problems**
 - ✿ **Describe the dependencies** for reductions
 - ✿ If possible, leave the parameter **p configurable** at run-time
- ✿ **The programming environment helps:**
 - ✿ Usually provides an **API for map and reduce**
(no need to manually code the distribution and communication)
 - ✿ Sometimes also determines **p automatically**

SUMMARY

- ✿ You: **decompose** the problem, try to express using **n/p complexity**; computer: run work over parallel resources
- ✿ Limits: p too small = less throughput, p too big = more overhead; try $p \approx$ number of PUs
- ✿ **map/reduce** suitable for **data-parallelism followed by reductions**

PART 3: PATTERNS OF DISTRIBUTED PROGRAMMING

LIMITS OF DATA-PARALLELISM

LIMITS OF DATA-PARALLELISM

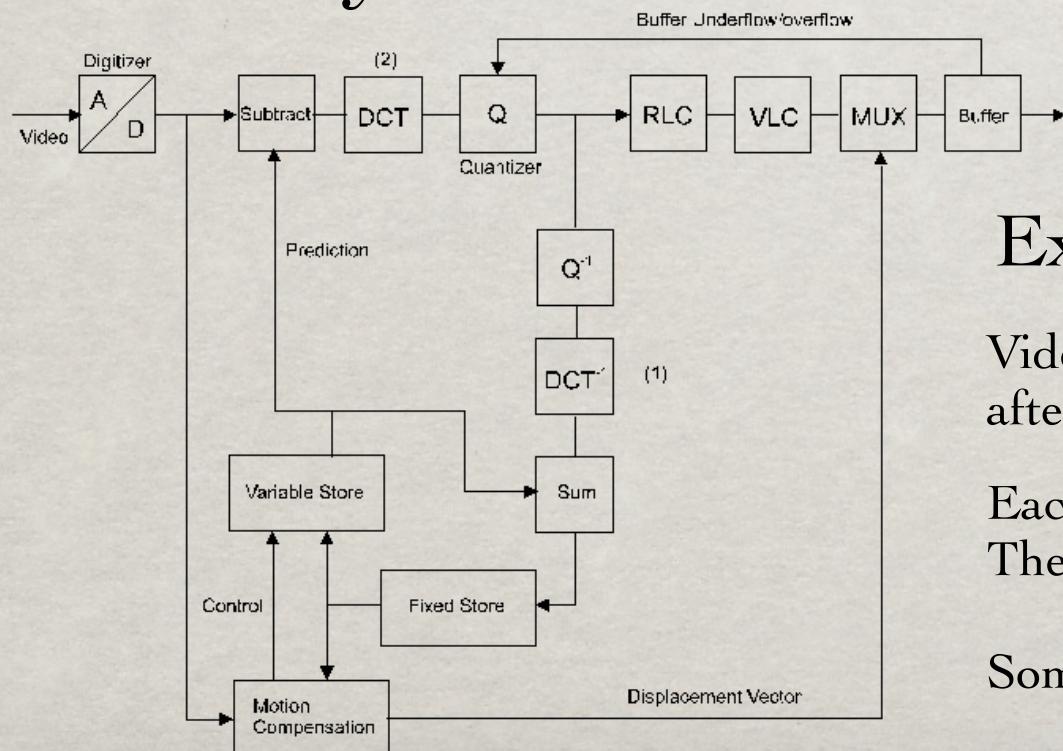
- ✿ Unfortunately not all problems are purely data-parallel

LIMITS OF DATA-PARALLELISM

- ✿ Unfortunately not all problems are purely data-parallel
- ✿ Can you think of counter examples?

LIMITS OF DATA-PARALLELISM

- ✿ Unfortunately not all problems are purely data-parallel
- ✿ Can you think of counter examples?



Example: MPEG encode

Video stream: many frames(picture) one after another

Each frame decomposed in blocks
The workload per frame varies

Some dependencies across frames

LIMITS OF DATA- PARALLELISM

LIMITS OF DATA-PARALLELISM

- Also, SIMD/SPMD case yields best throughput
only if all p PUs run at equal speed

LIMITS OF DATA-PARALLELISM

- Also, SIMD/SPMD case yields best throughput
only if all p PUs run at equal speed
- Is this always guaranteed?

LIMITS OF DATA-PARALLELISM

- Also, SIMD/SPMD case yields best throughput
only if all p PUs run at equal speed
- Is this always guaranteed?
- heat / cooling induces **transient failures/ slowdowns**

LIMITS OF DATA-PARALLELISM

- Also, SIMD/SPMD case yields best throughput
only if all p PUs run at equal speed
- Is this always guaranteed?
 - heat / cooling induces transient failures/ slowdowns
 - interference from other computations** running on the same resources

OTHER PATTERNS: PIPELINE

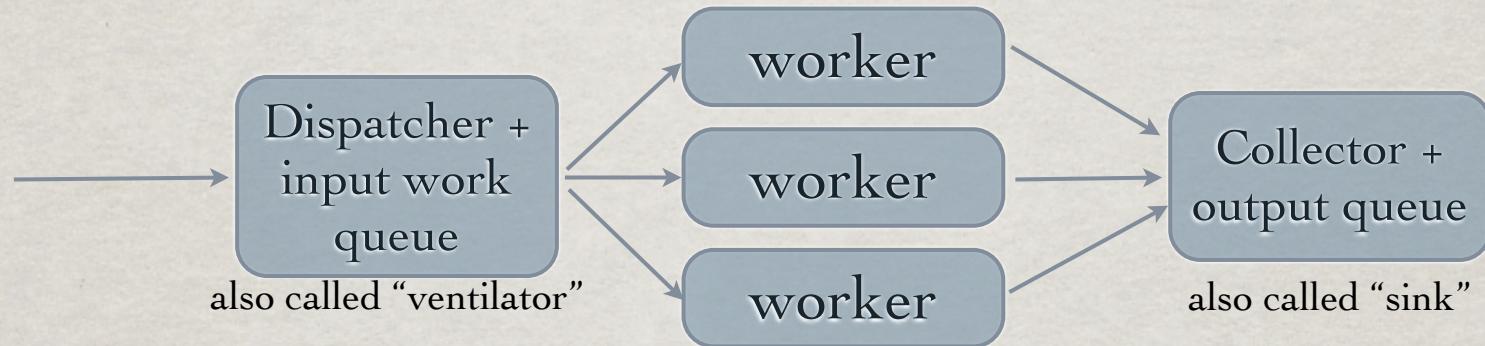
- ✿ concept: “assembly line”



- ✿ Works best when each sub-task takes approximately the same time
- ✿ Typical example: signal processing, video

OTHER PATTERNS: FARMING

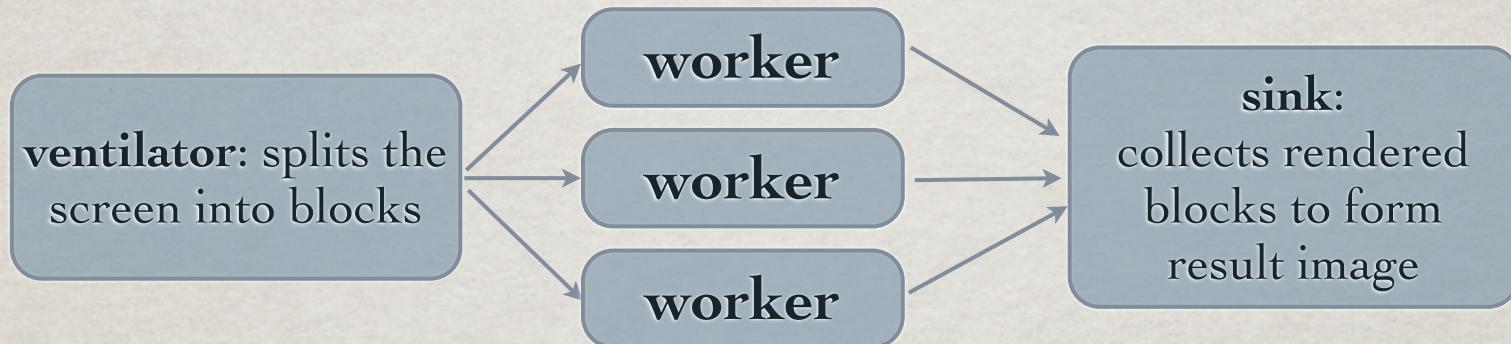
- ✿ concept: “clerks at the city hall”



- ✿ Works best when each task work is independent from the others + separate data
- ✿ (Can be used to simulate SIMD/SPMD on heterogeneous platforms)
- ✿ Typical example: web server, database

PRACTICALS

EXAMPLE: RAYTRACER



- ✿ We provide you Python source code for the 3 components
- ✿ It uses $\emptyset\text{MQ}$ for communication: a library for inter-process messaging
- ✿ Your task: try it out, optimize: how many workers and which block size are best?
- ✿ Run on your local workstation / local network

BASICS WITH MPI

- ✿ MPI is an “old school” API for parallel programming - more mature, widely used
- ✿ Like ØMQ, MPI provides send/recv
- ✿ MPI also provides automated reduction, barriers and other abstractions to simplify the programmer’s job

BASICS WITH MPI

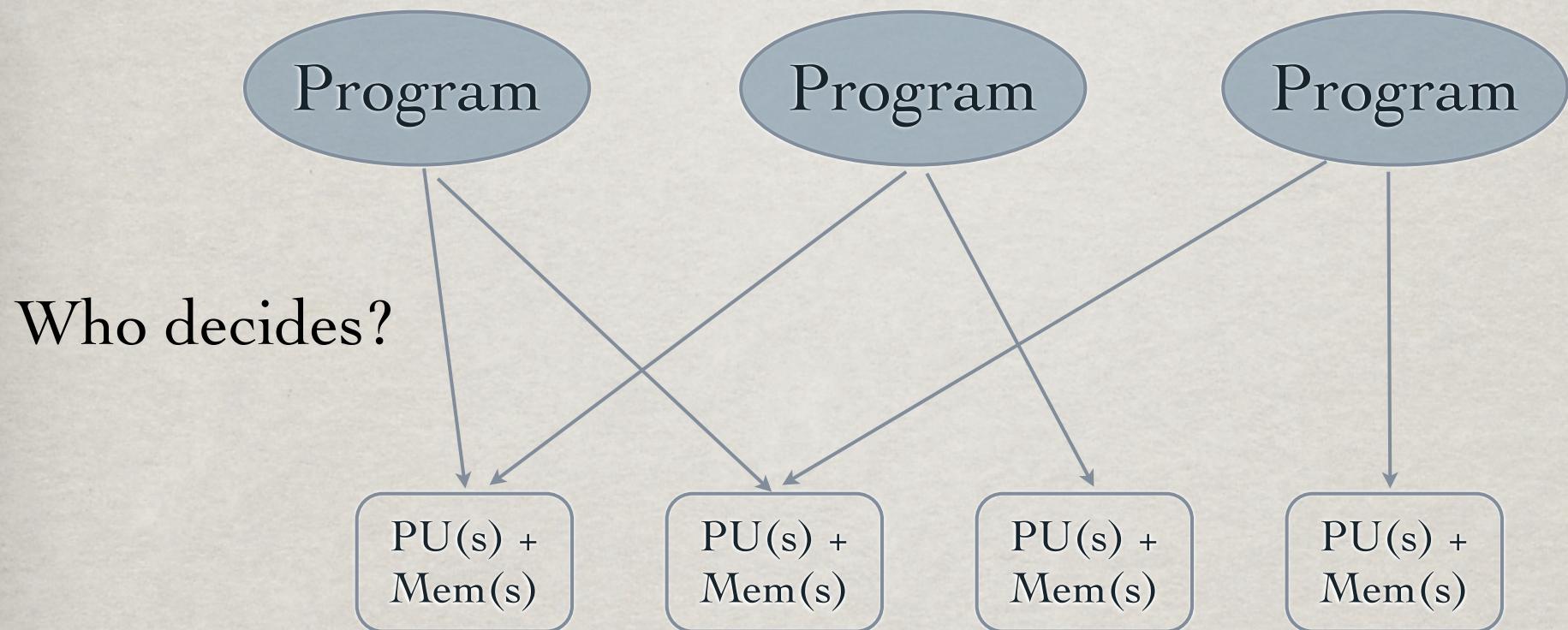
- ✿ Basic principle: MPI automatically starts **the same C program multiple times** on multiple processes/nodes
- ✿ Example:

```
#include <stdio.h>
#include <mpi.h>
int main(int argc, char **argv){
    int k, p;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &k);
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    printf("hello world k=%d p=%d\n", k, p);
    // (place computation here)
    MPI_Finalize();
    return 0;
}
```

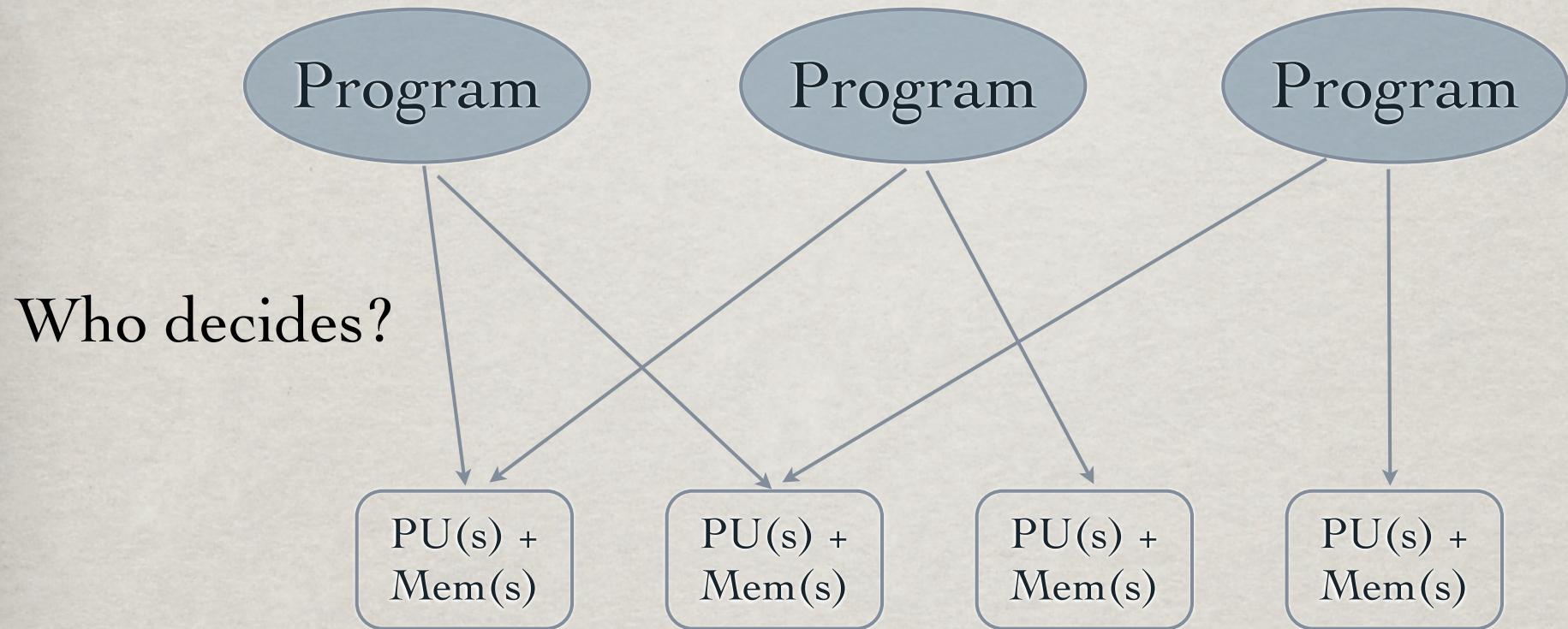
BASICS WITH MPI

- ✿ SIMD/SPMD “map”: in MPI, no special primitive needed
(the entire program is already duplicated)
- ✿ “reduce”: use `MPI_Reduce(...)`
- ✿ Check with the assistants for details!

SHARING RESOURCES



SHARING RESOURCES



Typically:

- programmers define **jobs**
- **job = program + resource request**
- a common **job manager** maps jobs to the compute nodes
- Example: <http://www.cs.vu.nl/das4/jobs.shtml>