

CSC367 Parallel computing

Lecture 8: Parallel Architectures and Parallel Algorithm Design Continued!

Methods for containing interaction overheads

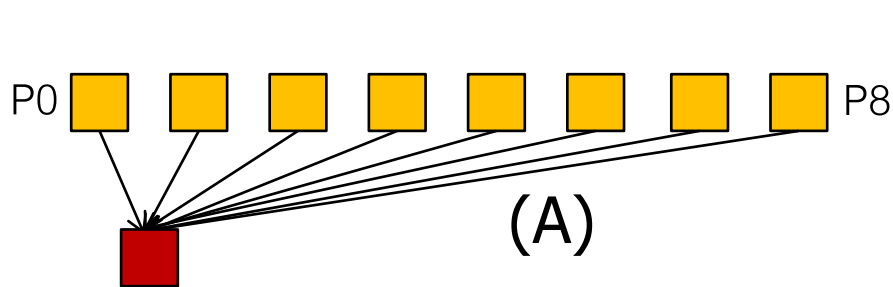
- Maximizing data locality
- Minimizing contention and hot spots
- Overlapping computations with interactions
- Replicating data or computations

Maximize data locality

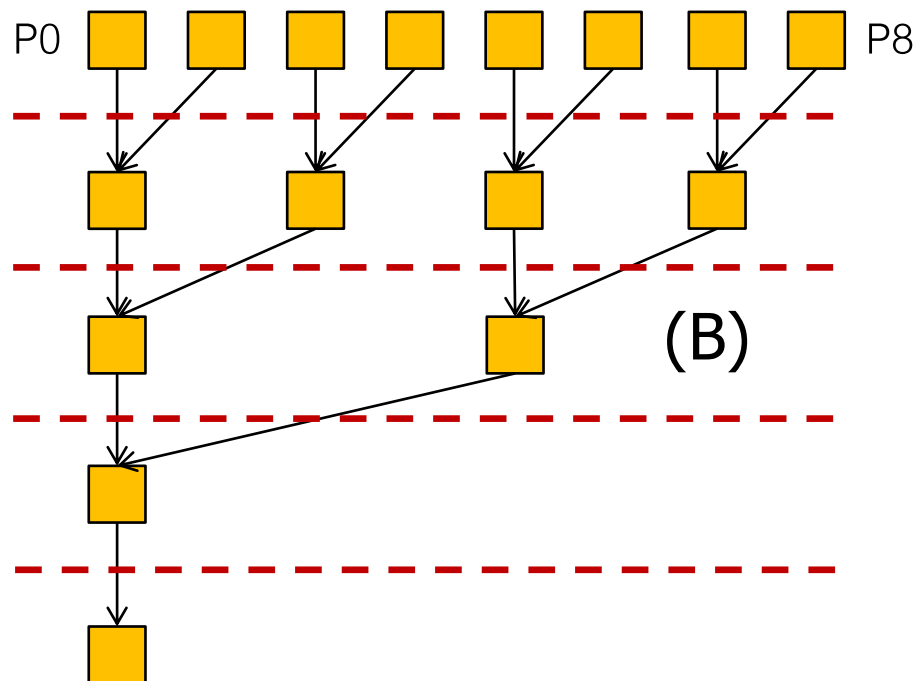
- Processes share data and/or may require data generated by other processes
- Goals:
 - 1. Minimize the volume of interaction overheads (minimize non-local data accesses and maximize local data utilization)
 - Minimize the volume of shared data and maximize cache reuse
 - Use a suitable decomposition and mapping
 - Use local data to store intermediate results (decreases the amount of data exchange)
 - 2. Minimize the frequency of interactions
 - Restructure the algorithm to access and use large chunks of shared data (amortize interaction cost by reducing frequency of interactions)
 - Shared address space: spatial locality in a cache line, etc.

Minimizing contention and hotspots

- Accessing shared data concurrently can generate contention
 - e.g., concurrent accesses to the same memory block, flooding a specific process with messages, etc.
- Solutions:
 - Restructure the program to reorder accesses in a way that does not lead to contention
 - Decentralize the shared data, to eliminate the single point of contention



Reducing the elements of an array:
the threads all need to access the
red block in (A) while (B)
decentralizes this single point of
contention.



Overlap computations with interactions

- Process may idle waiting for shared data => do useful computations while waiting
- Strategies:
 - Initiate an interaction earlier than necessary, so it's ready when needed
 - Grab more tasks in advance, before current task is completed
- May be supported in software (compiler, OS), or hardware (e.g., prefetching hardware, etc.).
- Harder to implement with shared memory models (Pthreads, OpenMP), applies more to distributed and GPU architectures (more on it later in class!)

Replicate data or computations

- To reduce contention for shared data, may be useful to replicate it on each process => no interaction overheads
- Beneficial if the shared data is accessed in read-only fashion
 - Shared-address space paradigm: cache local copies of the data
 - Message-passing paradigm (MPI), more on this later: replicate data to eliminate data transfer overheads
- Disadvantages:
 - Increases overall memory usage by keeping replicas => use sparingly
 - If shared data is read-write, must keep the copies coherent => overheads might dwarf the benefits of local accesses via replication

Parallel Algorithm Design: Outline

- Tasks: Decomposition, Task Dependency, Granularity, Interaction, Mapping, Balance
- Decomposition techniques
- Mapping techniques to reduce parallelism overhead
- Parallel algorithm models
- Parallel program performance model

Parallel algorithm models

Model = 1 decomposition type + 1 mapping type + strategies to minimize interactions

Commonly used parallel algorithm models:

Data parallel model

Work pool model

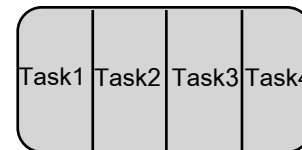
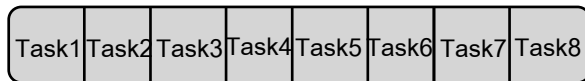
Master slave model

Data parallel model

- **Decomposition:** typically static and uniform data partitioning
- **Mapping:** static (mostly)
- Same operations on different data items, aka **data parallelism**
- Possible optimization **strategies** (depending on the problem and paradigm):
 - Choose a **locality-preserving decomposition**
 - Overlap computation with interaction
- This model scales really well with problem size (by adding more processes)

Work pool model

- Tasks are taken by processes from a common pool of work
- **Decomposition:** highly depends on the problem (data, recursive, etc.)
 - Can be statically available at start, or dynamically create more tasks during execution
- **Mapping:** dynamic
 - Any task can be performed by any process
- Possible **strategies** for reducing interactions:
 - **Adjust granularity:** tradeoff between load imbalance and overhead of accessing work pool



Master slave model

- Commonly used in distributed parallel architectures (more on this later)
- A master process generates work and allocates to worker (slave) processes
 - Could involve several masters, or a hierarchy of master processes
- **Decomposition:** highly depends on the problem (data, recursive)
 - Might be static if tasks are easy to break down a priori, or dynamic
- **Mapping:** Often dynamic
 - Any worker can be assigned any of the tasks
- Possible **strategies** for reducing interactions:
 - Choose granularity carefully so that master does not become a bottleneck
 - Overlap computation on workers with master generating further tasks

Parallel Algorithm Design: Outline

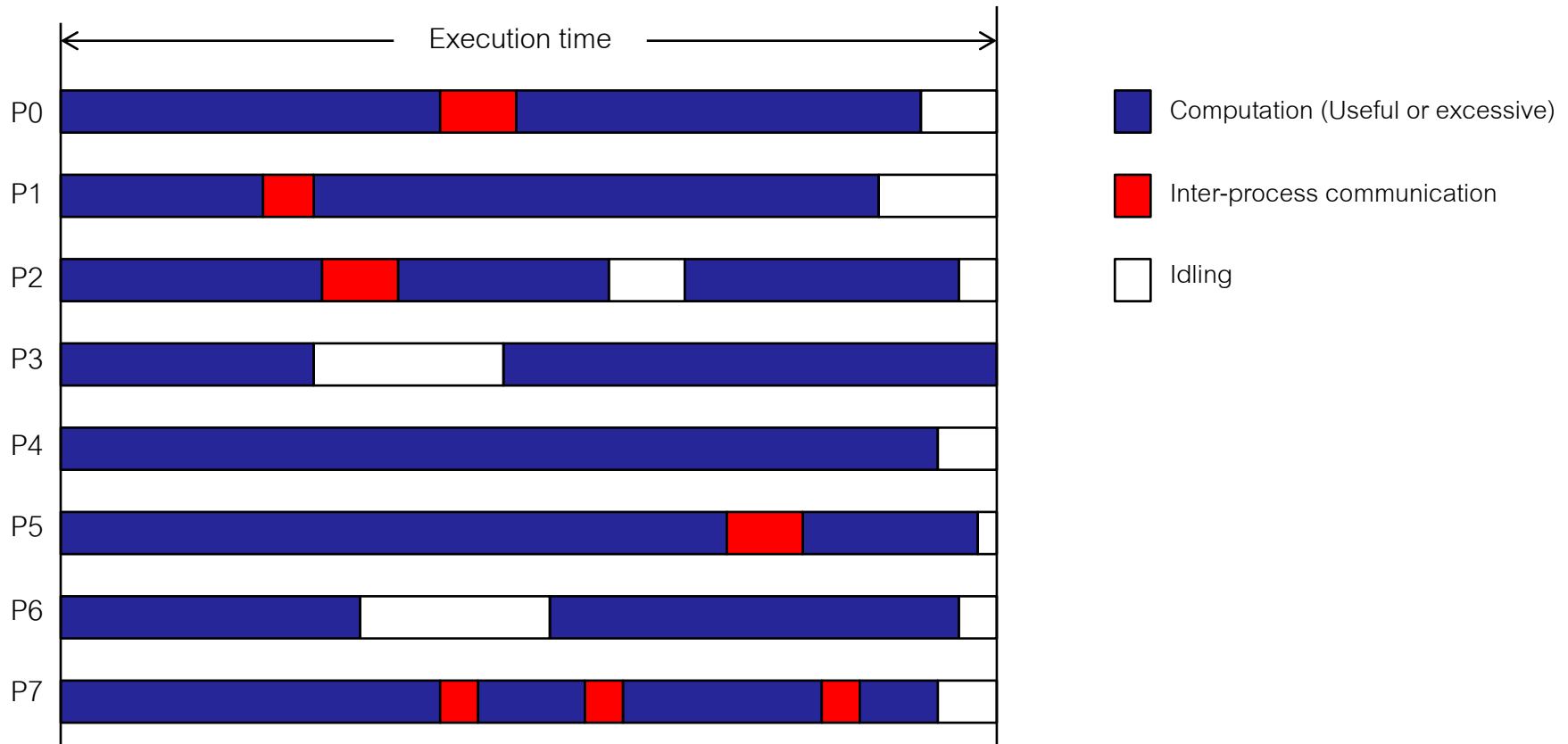
- Tasks: Decomposition, Task Dependency, Granularity, Interaction, Mapping, Balance
- Decomposition techniques
- Mapping techniques to reduce parallelism overhead
- Parallel algorithm models
- Parallel program performance model

Sources of overhead in parallel programs

- Parallel execution over N processes rarely results in N -fold performance gain
- Overheads of **inter-process communication**, **idling**, and/or **excess computation***

*sometimes used to reduce communication

- e.g., some random application profile:

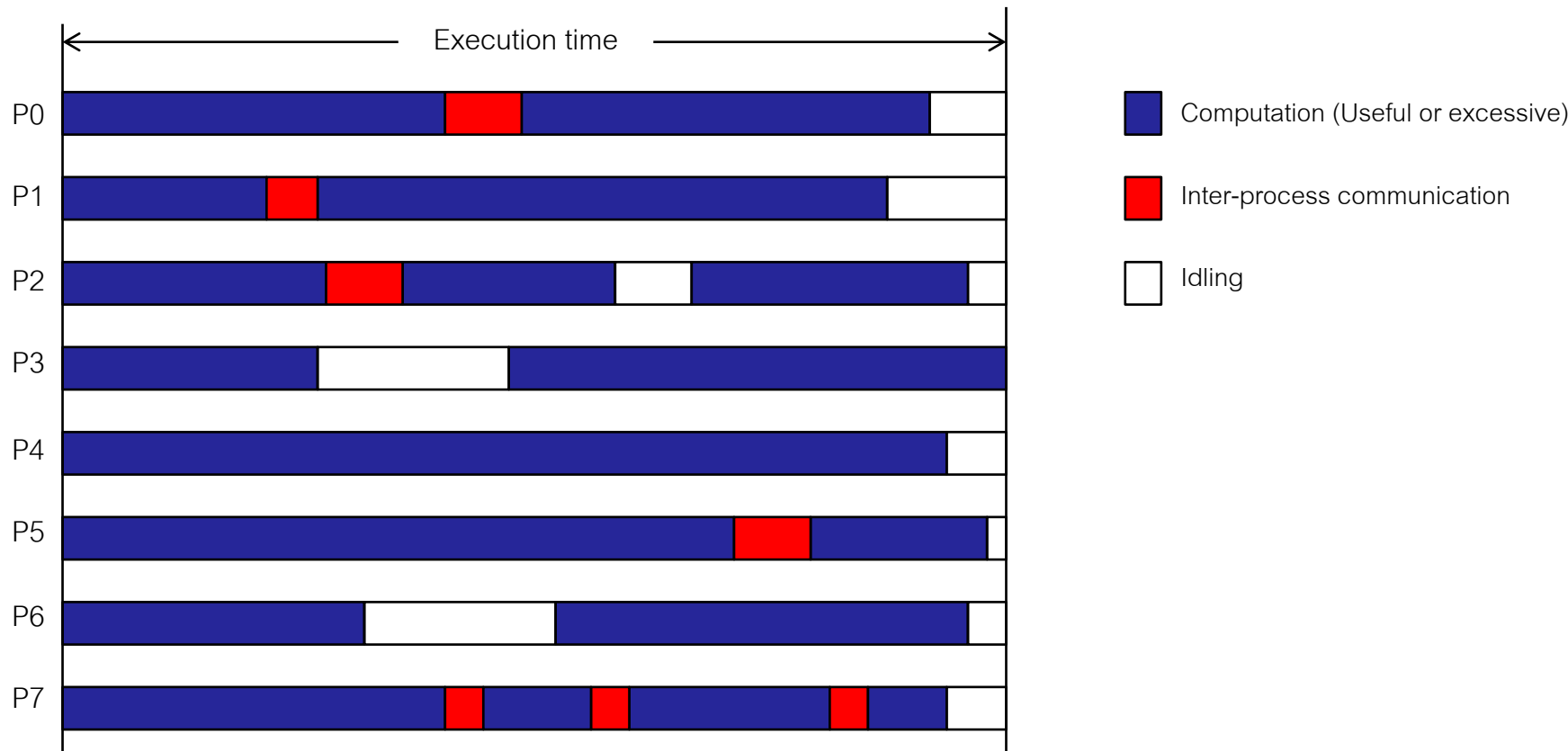


Performance metrics for parallel programs

- Execution time
- Speedup
- Efficiency
- Example performance plots
- Ways to fool the masses with performance results!

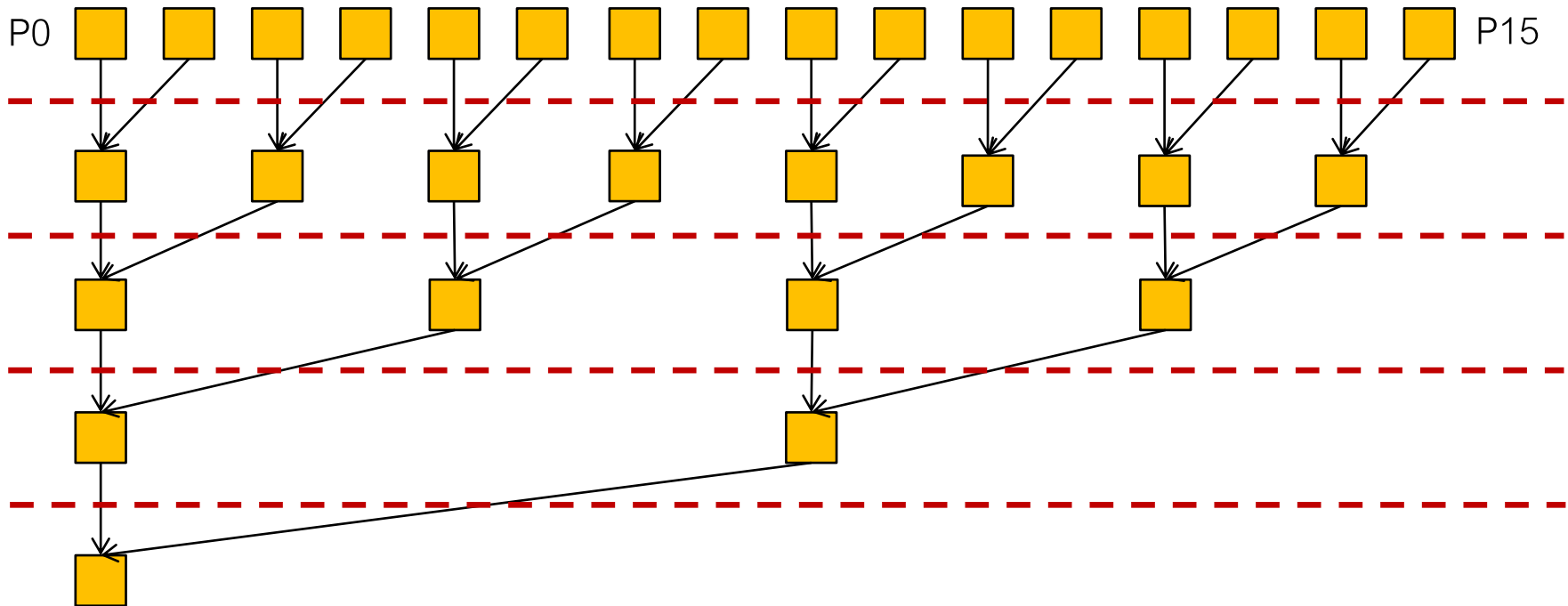
Execution time

- Execution time
 - Time elapsed from start of parallel computation and time when last process finishes
 - Determined by slowest process



Speedup

- Performance gain of the parallel algorithm over the sequential version
 - $S = T_s$ (time of the serial execution) / T_p (time to execute on p processors)
- Example: sum of elements of an array, each process handles one element



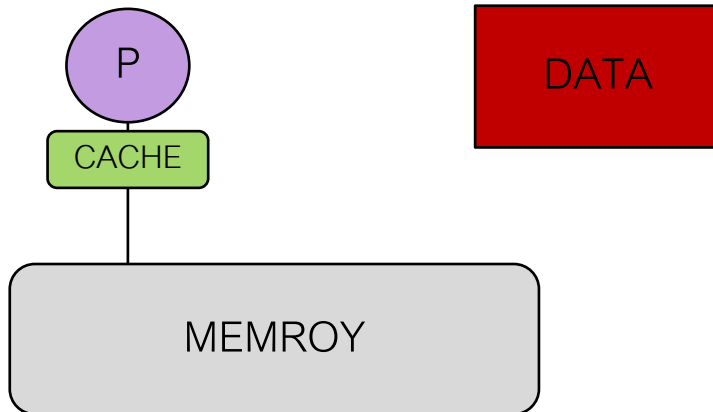
- $T_p = \Theta(\log n) \Rightarrow S = \Theta(n / \log n)$

Speedup considerations

- Speedup is calculated relative to the best serial implementation of the algorithm. First improve your serial implementation, use that as the baseline, and optimize the improved serial implementation on p processors.
- Maximum achievable speedup is called **linear speedup**
 - Each process takes p times less than $T_s \Rightarrow S = p$

Speedup considerations

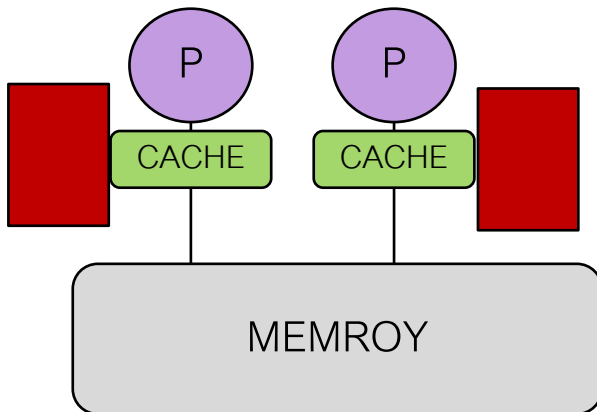
- If $S > p \Rightarrow$ **superlinear speedup**
 - Wait, if all processes spend less than T_s/p , why not use 1 process to run the whole thing then?
- Superlinear speedup can only happen if sequential algorithm is at a disadvantage compared to parallel version
- Data too large to fit into 1 processor's cache \Rightarrow data accesses are slower for serial algorithm



If the program needs to stream the data n times, the data does not fit in the cache so the data has to be moved between the memory and caches n times!

Speedup considerations

- If $S > p \Rightarrow$ **superlinear speedup**
 - Wait, if all processes spend less than T_s/p , why not use 1 process to run the whole thing then?
- Superlinear speedup can only happen if sequential algorithm is at a disadvantage compared to parallel version
 - Data too large to fit into 1 processor's cache \Rightarrow data accesses are slower for serial algorithm



If the half of the data fits in one of the L1 caches and work can be divided between the processors, then the data only gets loaded once into the caches from memory!

Speedup

Old program (unenhanced)



Old time: $T = T_1 + T_2$

T_1 = time that can NOT be enhanced.

T_2 = time that can be enhanced.

New program (enhanced)



New time: $T' = T_1' + T_2'$

T_2' = time after the enhancement.

Speedup: $S_{\text{overall}} = T / T'$

Amdahl's law

- Suppose only part of an application is parallel



- Amdahl's law

$$\text{If } T = T_1 + T_2 = 1$$

- T_1 = fraction of work done sequentially (Amdahl fraction), so $(T_2 = 1 - T_1)$ is fraction parallelizable
- p = number of processors

$$\begin{aligned}\text{Speedup}(P) &= T / T' \\ &\leq 1 / (T_1 + (1 - T_1) / p) \\ &\leq 1 / T_1\end{aligned}$$

- Even if the parallel part speeds up perfectly performance is limited by the sequential part

Efficiency

- The fraction of time when processes doing useful work

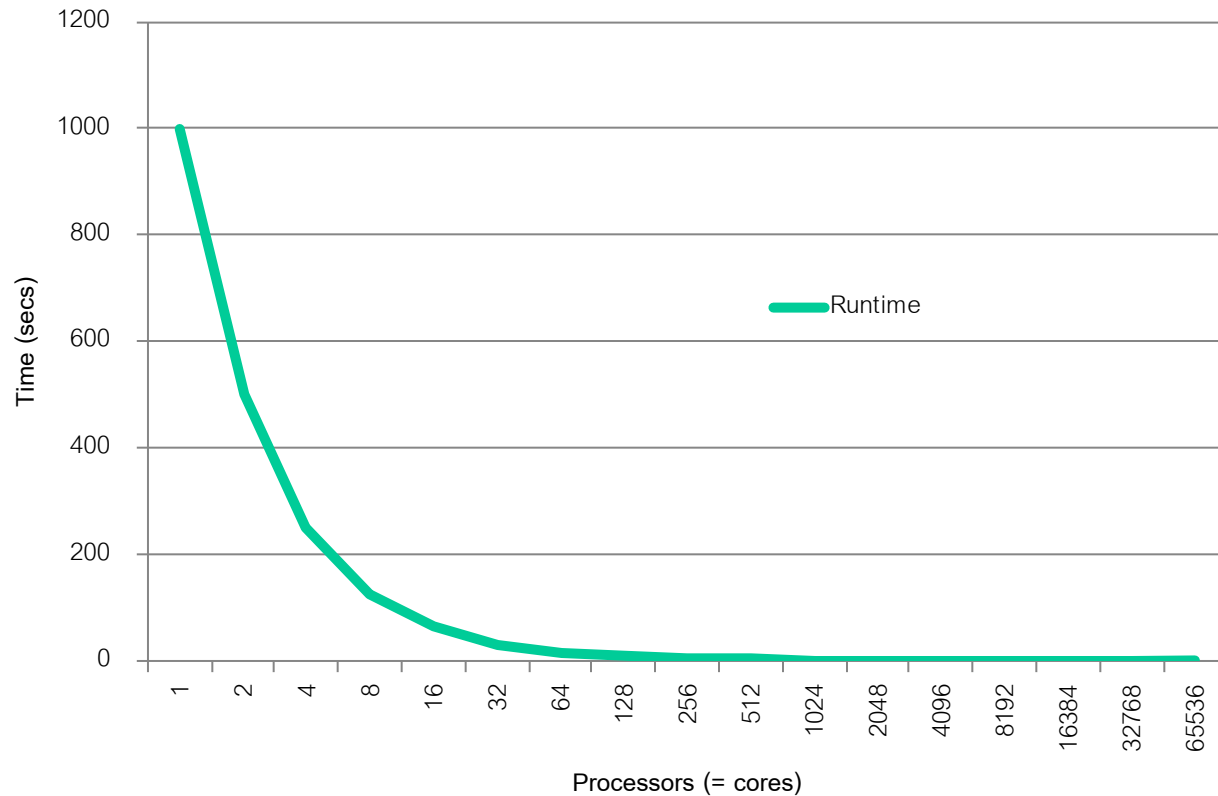
$$E = S / p$$

- What is the ideal efficiency? 1 (sometimes shown as 100%)
- What are the range of values for E? 0 to 1
- What is the efficiency of calculating the sum of n array elements on n processes?

$$E = \Theta(n / \log n) / n$$

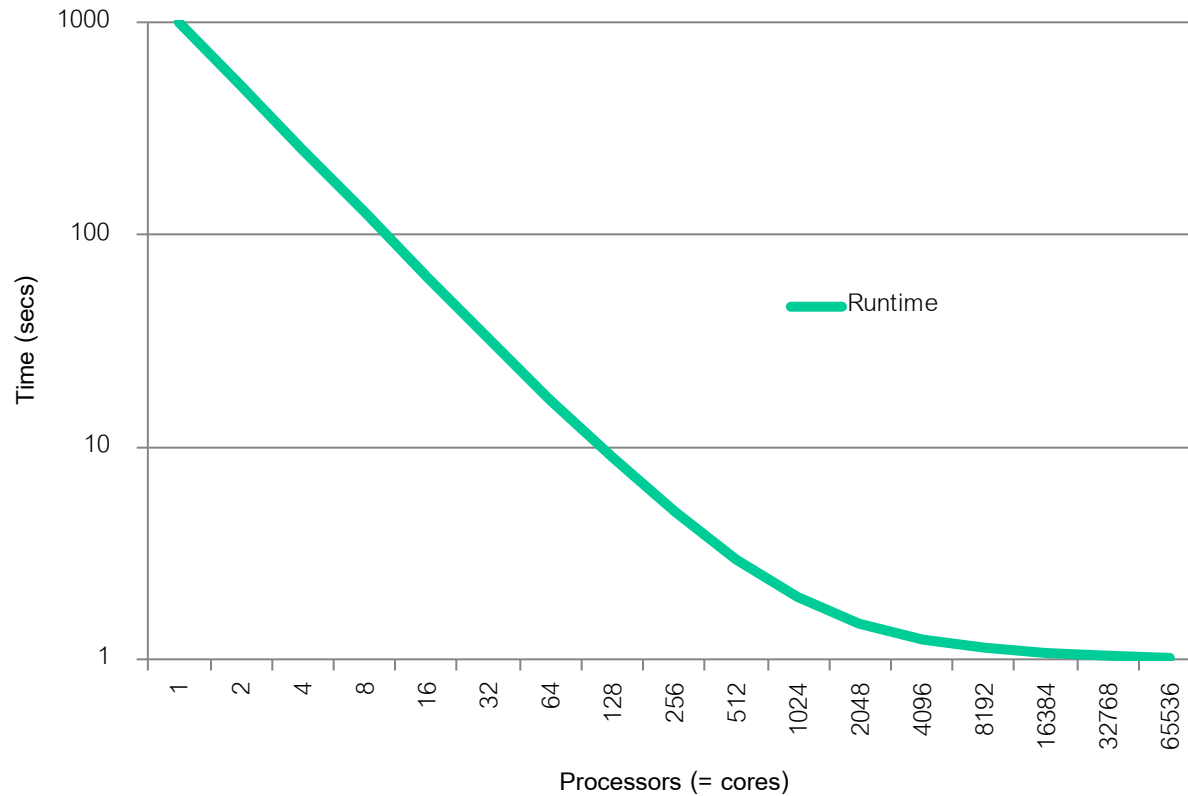
$$E = \Theta(1 / \log n)$$

Reporting running time



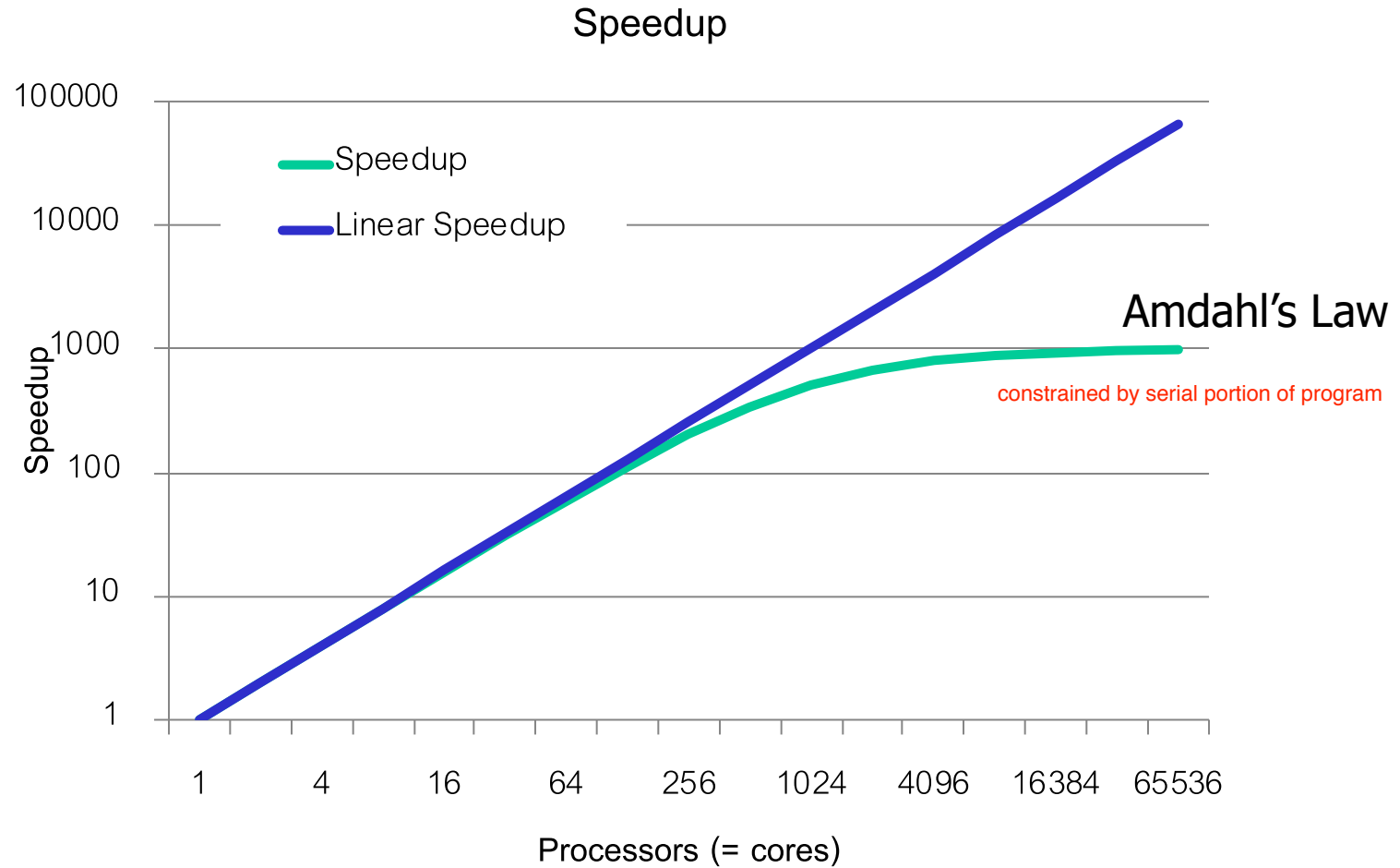
Hard to see performance gains from parallelism after 32 processors!

Reporting running time

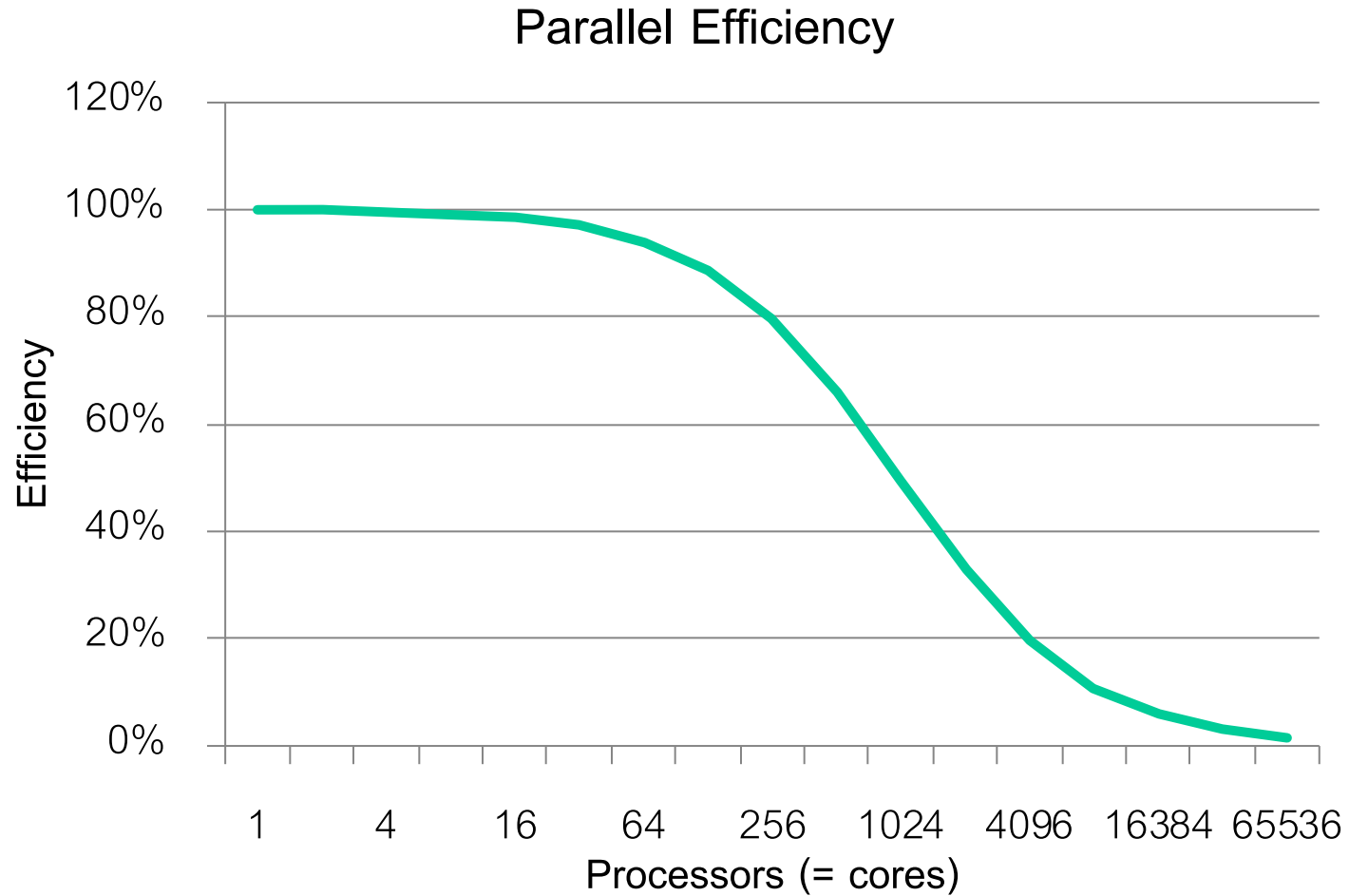


Lets take the y axis (running time) to log scale: A bit better!

Example Speedup Plot



Example Efficiency Plot



Carefully choose and report your serial/baseline

See *David Bailey's Twelve Ways to Fool the Masses*. Below are examples of how to fool the masses when reporting results from your parallel program:

32bit precision is always faster for GPU

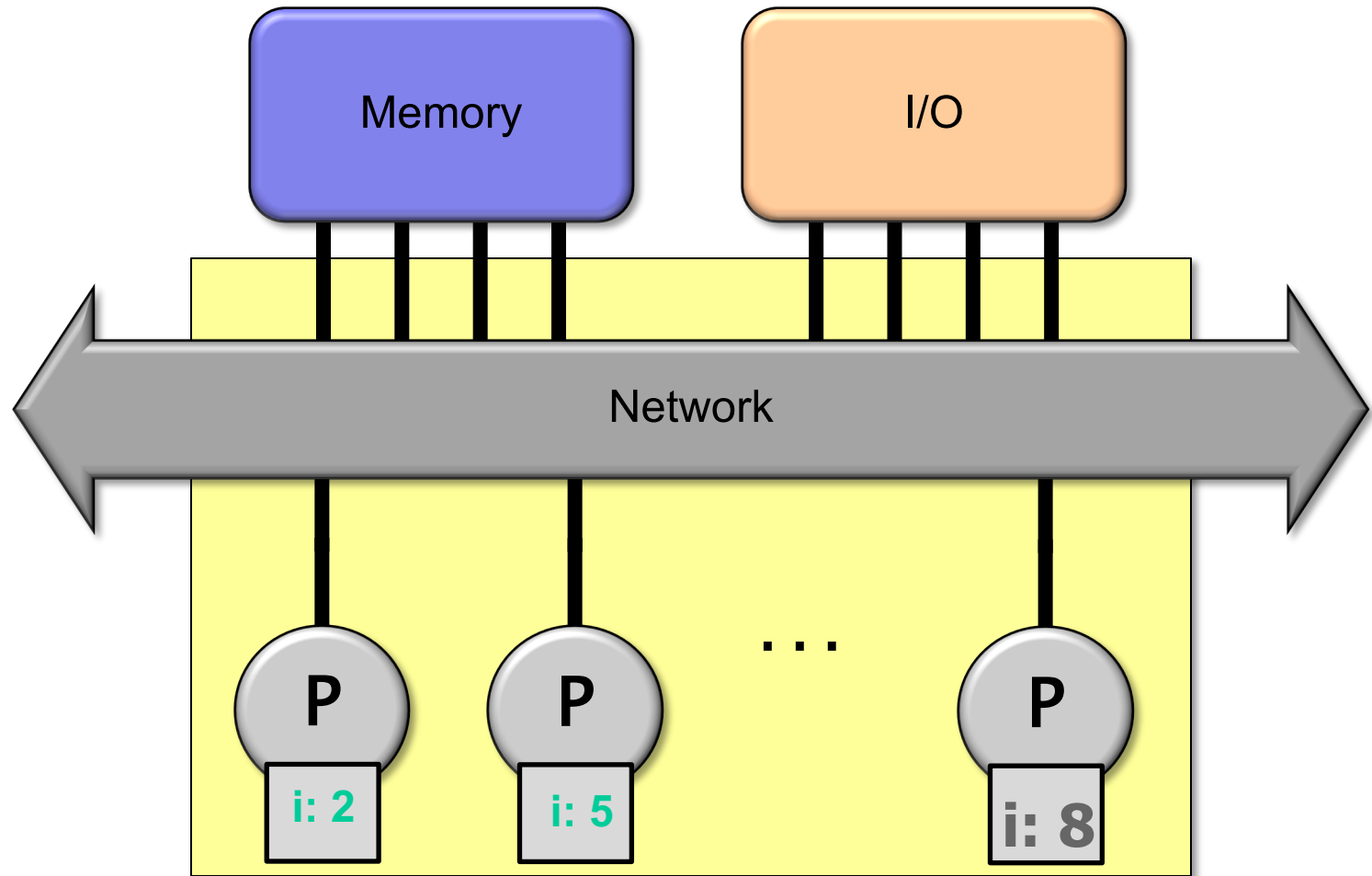
1. Use 64-bit for baseline/serial and 32-bit for parallel numbers:
 - Correct approach: Use the same precision for both the parallel implementation and the serial/baseline: This type of “cheating” in speedup reports often happens in GPU parallel programming, where single-precision is faster than double-precision computing.
2. Use a bad algorithm for the baseline:
 - Correct approach: Always optimize the serial algorithm first and use it as the baseline for speedups.
3. Use a bad implementation for the baseline:
 - Correct approach: While optimizing the parallel code if you realize you could have optimized the serial version better, go back and optimize the serial code and use that as baseline.
4. Don't report running times at all:
 - Correct approach: Report running times as well as speedup.

Shared Memory Architectures and Their Parallel Programming Models!

Next up ...

- Shared memory architecture
- Parallel programming models: shared memory
- Pthreads: Synchronization, Races, Locks
- OpenMP
- Cache coherency

Shared Memory Architecture



Chip Multiprocessor (CMP)

Next up ...

- Shared memory architecture
- Parallel programming models: shared memory
- Pthreads: Synchronization, Races, Locks
- OpenMP
- Cache coherency

Parallel Programming Models

- **Programming model** is made up of the languages and libraries that create an abstract view of the machine: Pthreads!

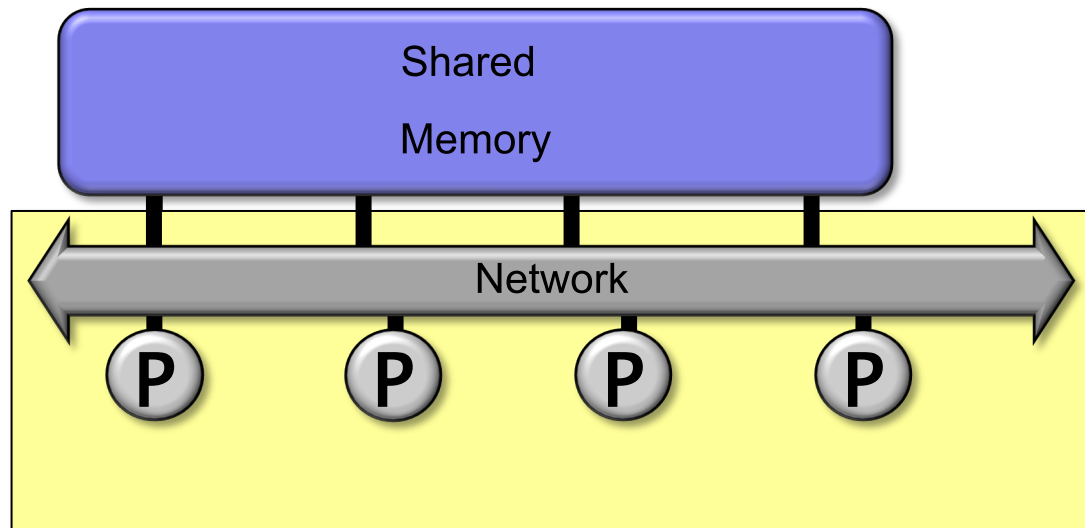
The programming model enables us to identify

- Control
 - How is parallelism **created**?
 - What **orderings** exist between operations?
- Data:
 - What data is **private** vs. **shared**?
 - How is logically shared data accessed or **communicated**?
- Synchronization
 - What operations can be used to coordinate parallelism?
 - What are the **atomic** (indivisible) operations?

Programming Model: Shared Memory

Program is a collection of threads of control, can be created mid-execution.

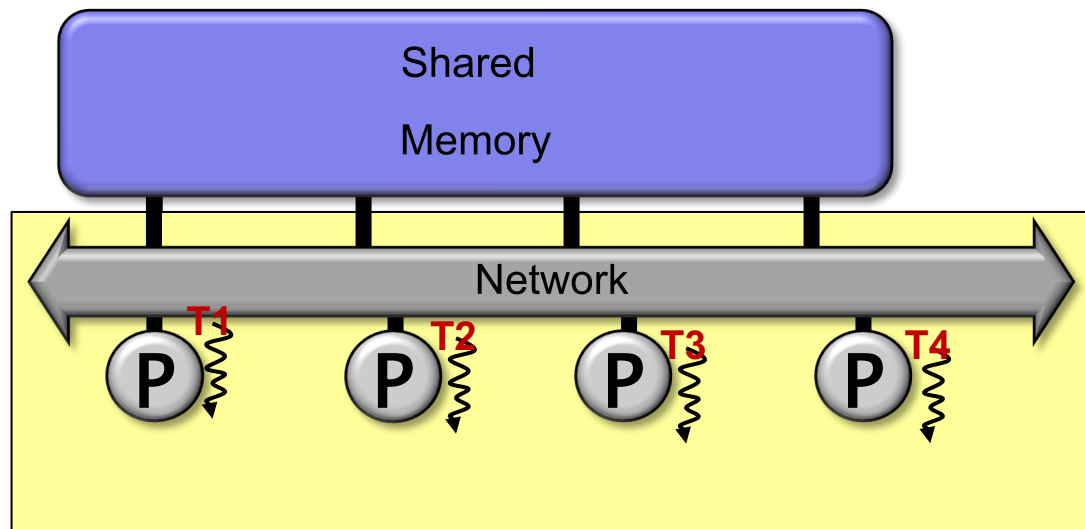
Thread



Programming Model: Shared Memory

Program is a collection of threads of control, can be created mid-execution.

Thread

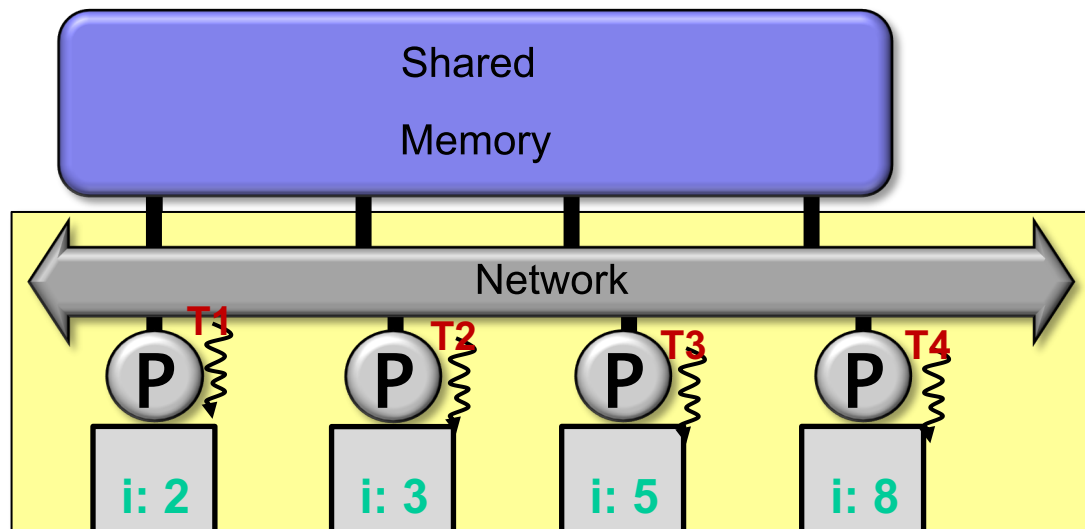


Programming Model: Shared Memory

Program is a collection of threads of control, can be created mid-execution.

Each thread has a set of **private variables**, e.g., local stack variables.

Thread



Programming Model: Shared Memory

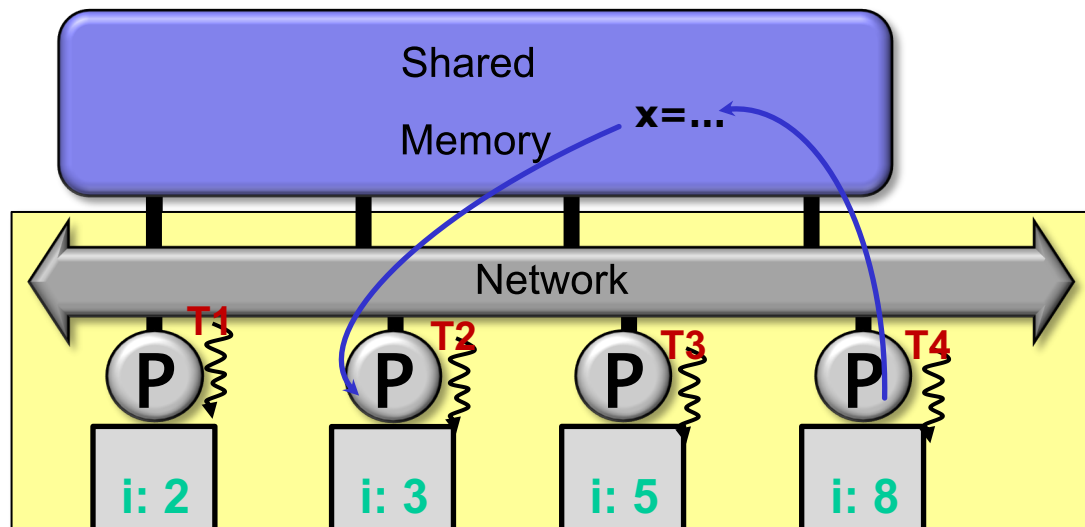
Program is a collection of threads of control, can be created mid-execution.

Each thread has a set of **private variables**, e.g., local stack variables.

Also a set of **shared variables**, e.g., static variables.

- Threads communicate **implicitly** by writing and reading shared variables.

Thread



Slide Source: Demmel