**CSCI-GA.3033-012**
**Multicore Processors:**
**Architecture & Programming**

# Lecture 2: Concurrency and Parallelism

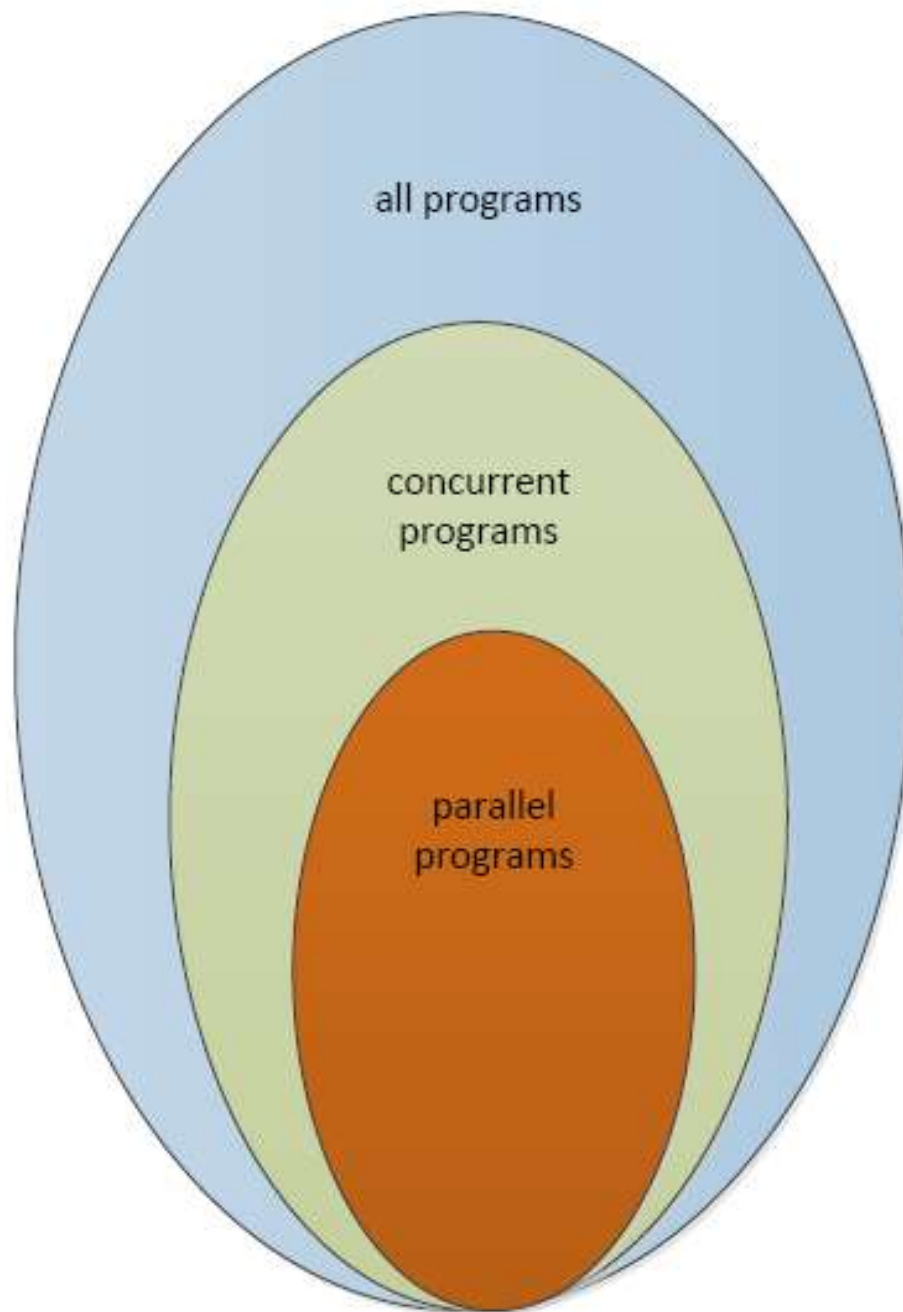Mohamed Zahran (aka Z)

mzahran@cs.nyu.edu

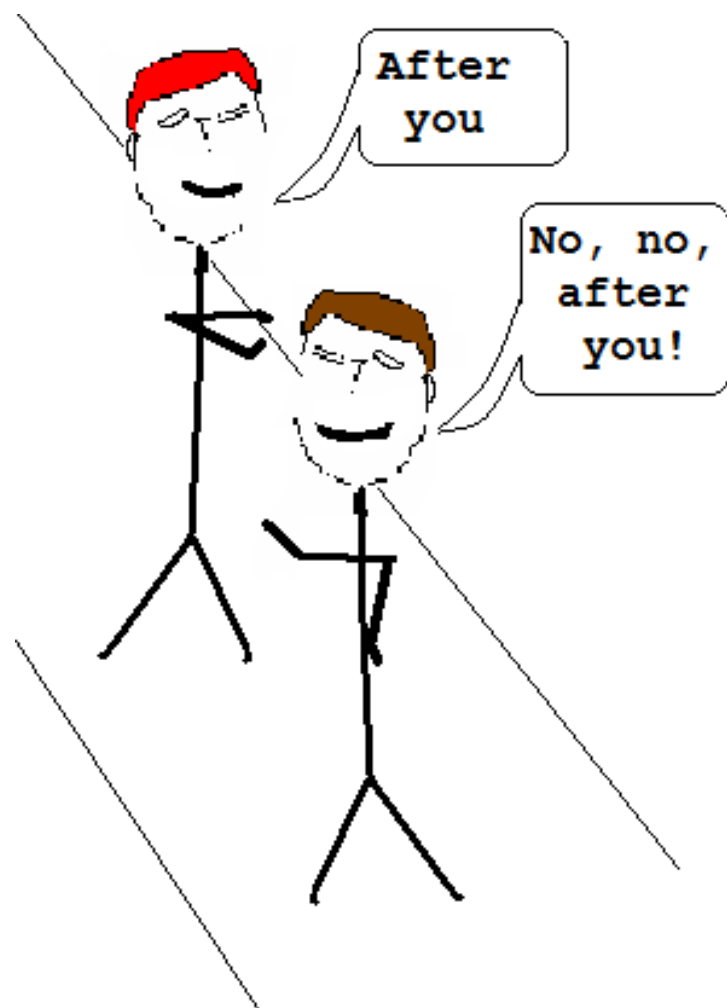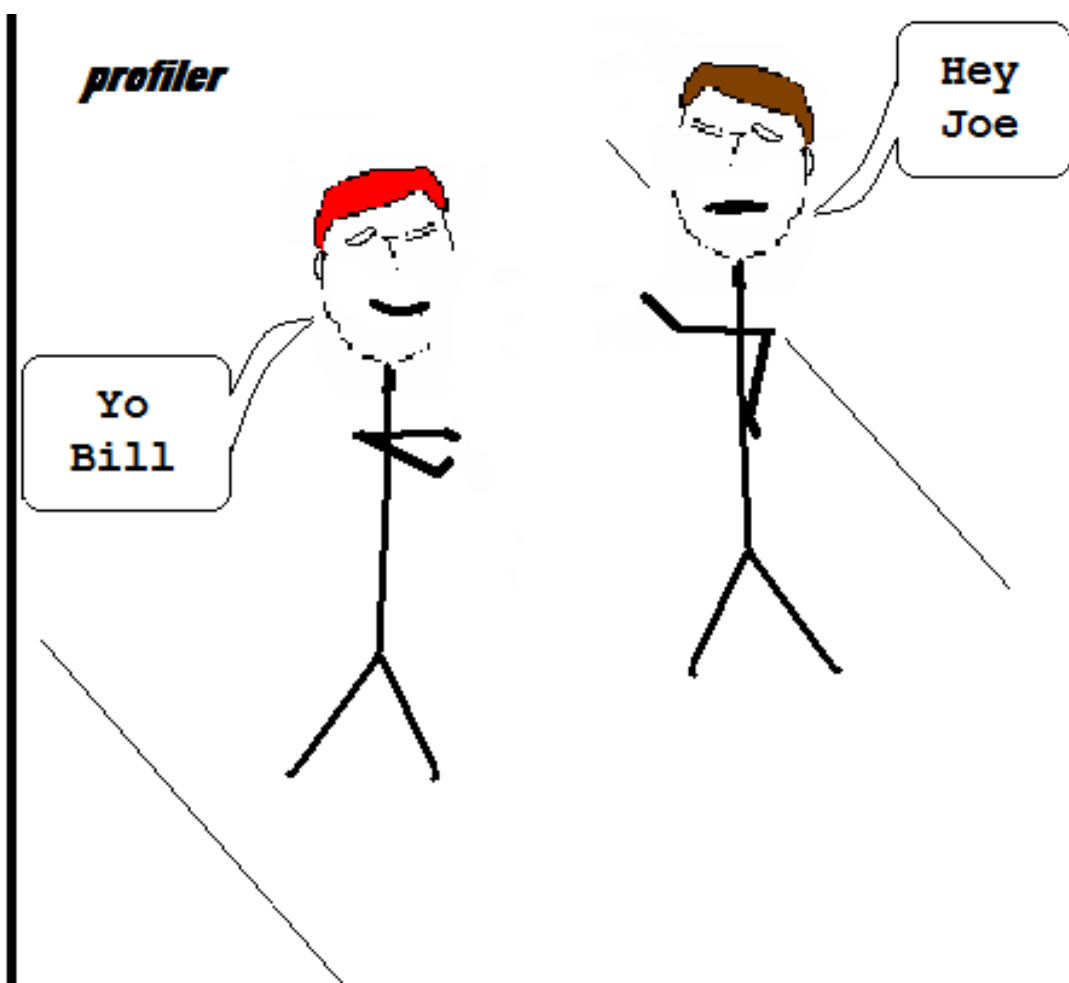http://www.mzahran.com

# Same Meaning?

- <span style="color:red">Concurrency</span>: At least two tasks are making progress at the same time frame.
    - Not necessarily at the same time
    - Include techniques like time-slicing
    - Can be implemented on a single processing unit
    - Concept more general than parallelism
- <span style="color:red">Parallelism</span>: At least two tasks execute literally at the same time.
    - Requires hardware with multiple processing units

Performance tuning technique number 106: Concurrency vs. Parallelism

Copyright © Fasterj.com Limited

# Questions!

If we have as much hardware as we want, do we get as much parallelism as we wish?

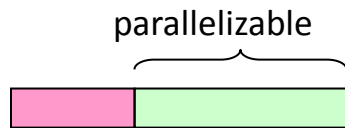If we have 2 cores, do we get 2x speedup?
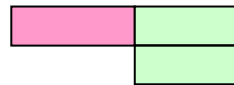
# Amdahl's Law

**Gene M. Amdahl**

- How much of a speedup one could get for a given parallelized task?
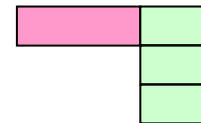
If F is the fraction of a calculation that is sequential then the maximum speed-up that can be achieved by using P processors is **1/(F+(1-F)/P)**

parallelizable

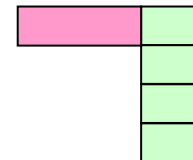1CPU 2CPUs 3CPUs 4CPUs

# What Was Amdahl Trying to Say?

- Don't invest blindly on large number of processors.

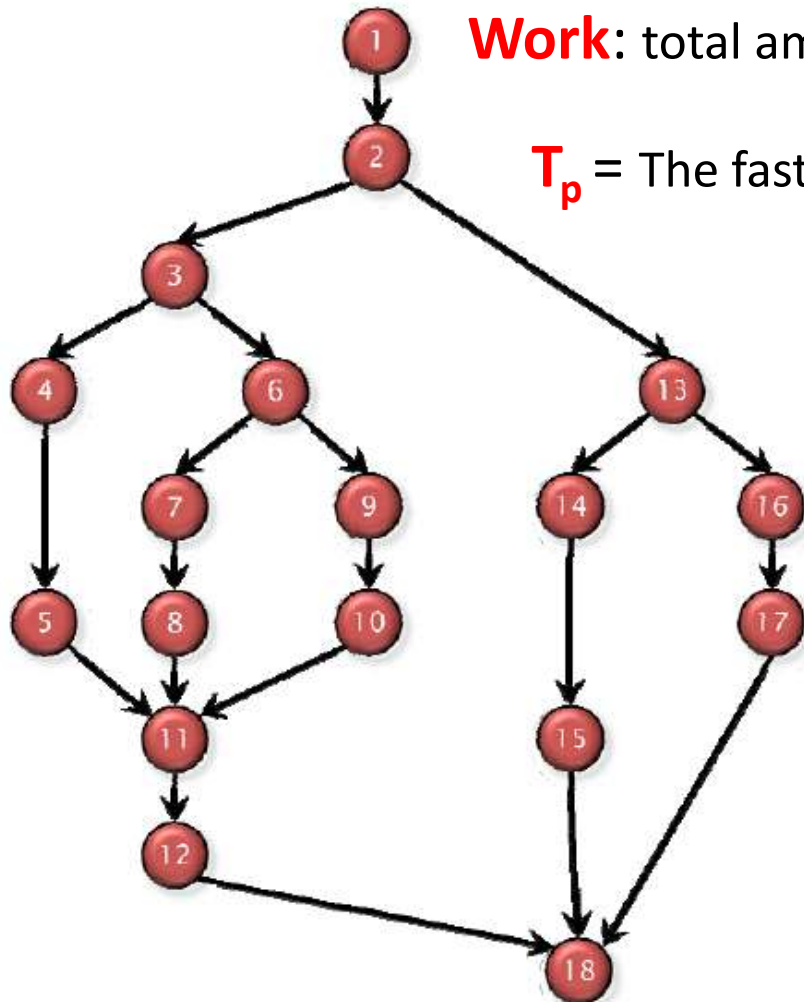- Having faster core (or processor at his time) makes more sense than having many cores.

Was he right?

- At his days (the law appeared 1967) many programs have long sequential parts.

- This is not necessarily the case nowadays.

- It is not very easy to find F (sequential portion)

# So …

- Decreasing the serialized portion is of greater importance than adding more cores
- Only when a program is mostly parallelized, does adding more processors help more than parallelizing the remaining rest
- **Gustafson's law**: computations involving arbitrarily large data sets can be efficiently parallelized
- Both Amdahl and Gustafson do not take into account:
  - The overhead of synchronization, communication, OS, etc.
  - Load may not be balanced among cores
- So you have to use these laws as guideline and theoretical bounds only.

# DAG Model for Multithreading



**Work**: total amount of time spent on all instructions

$T_p$ = The fastest possible execution time on P processors

**Work Law:**   $T_P \geq T_1/P$

# DAG Model for Multithreading



**Span**: The longest path of dependence in the DAG = $T_\infty$

Span Law:   $T_P \geq T_\infty$

# Can We Define Parallelism Now?

How about? $T_1 / T_\infty$

Ratio of work to span

# Can We Define Parallelism Now?



**Work:** $T_1 = 50$

**Span:** $T_\infty = 8$

**Parallelism:** $T_1/T_\infty = 6.25$

# At What Level Can We Reason About Parallelism (algorithm, high-level language, assembly)?
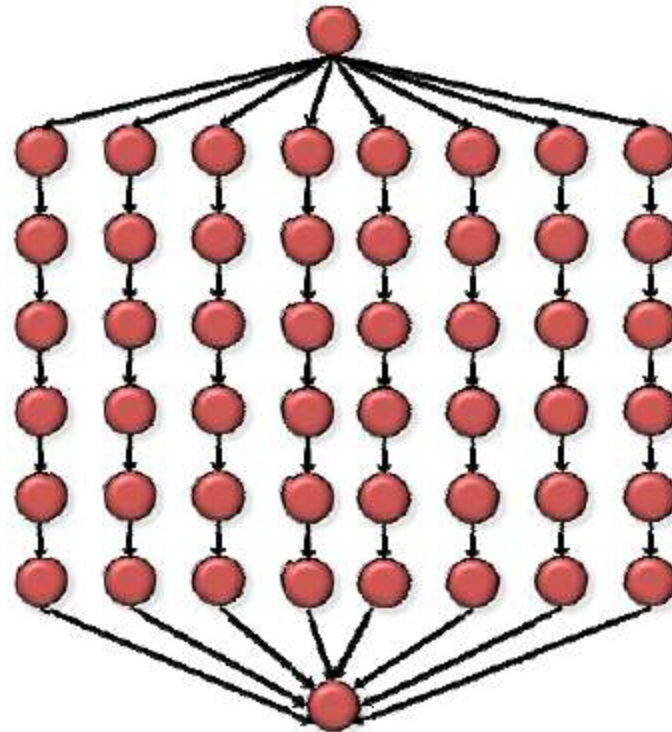
# Is Thread The Only Parallelism Granularity?

- ## Instruction level parallelism (ILP)
  - Superscalar
  - Out-of-order execution
  - Speculative execution

- ## Thread level parallelism
  - Hyperthereading technology (aka SMT)
  - Multicore

- ## Process level parallelism
  - Multiprocessor system
  - Hyperthereading technology (aka SMT)
  - Multicore

# That Was The Software How about the Hardware?

# A Quick Glimpse on: Flynn Classification

- A taxonomy of computer architecture
- Proposed by Michael Flynn in 1966
- It is based on two things:
  - Instructions
  - Data

|  | Single instruction | Multiple instruction |
|---|---|---|
| Single data | SISD | MISD |
| Multiple data | SIMD | MIMD |

**PU = Processing Unit**

# More About MIMD



**Shared-Memory**          **Distributed-Memory**

**Or hybrid**

# Shared Memory

| PE | PE | PE | PE |
|----|----|----|----|
| L1 | L1 | L1 | L1 |
| L2 | L2 | L2 | L2 |
| L3 | L3 | L3 | L3 |

Shared Memory

# Distributed Memory

| MEMORY | CPU | | CPU | MEMORY |
|--------|-----|--|-----|--------|

NETWORK

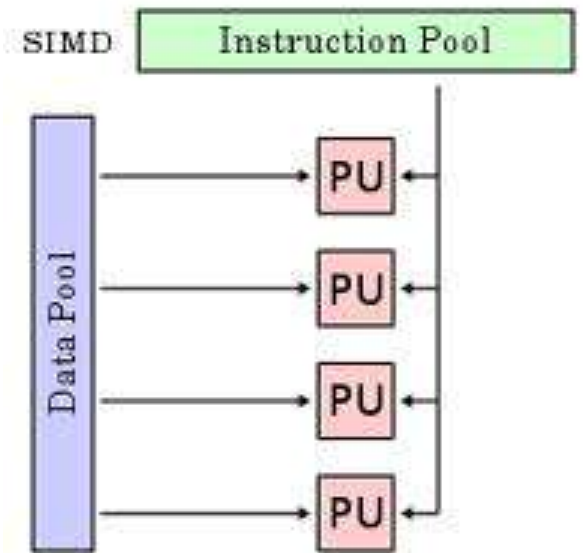| MEMORY | CPU | | CPU | MEMORY |
|--------|-----|--|-----|--------|

# Hybrid

| Core | Core | | Core | Core |
|------|------|--|------|------|
| Shared Memory | | | Shared Memory | |
| Core | Core | | Core | Core |

| Core | Core | | Core | Core |
|------|------|--|------|------|
| Shared Memory | | | Shared Memory | |
| Core | Core | | Core | Core |

# Multicore and Manycore

> We have arrived at many-core solutions *not* because of the success of our parallel software but because of *our failure* to keep increasing CPU frequency*.

**Tim Mattson**

Dilemma:
- Parallel hardware is ubiquitous
- Parallel software is not!
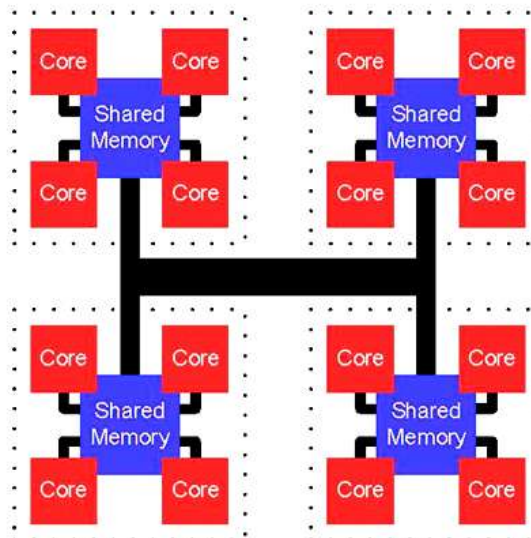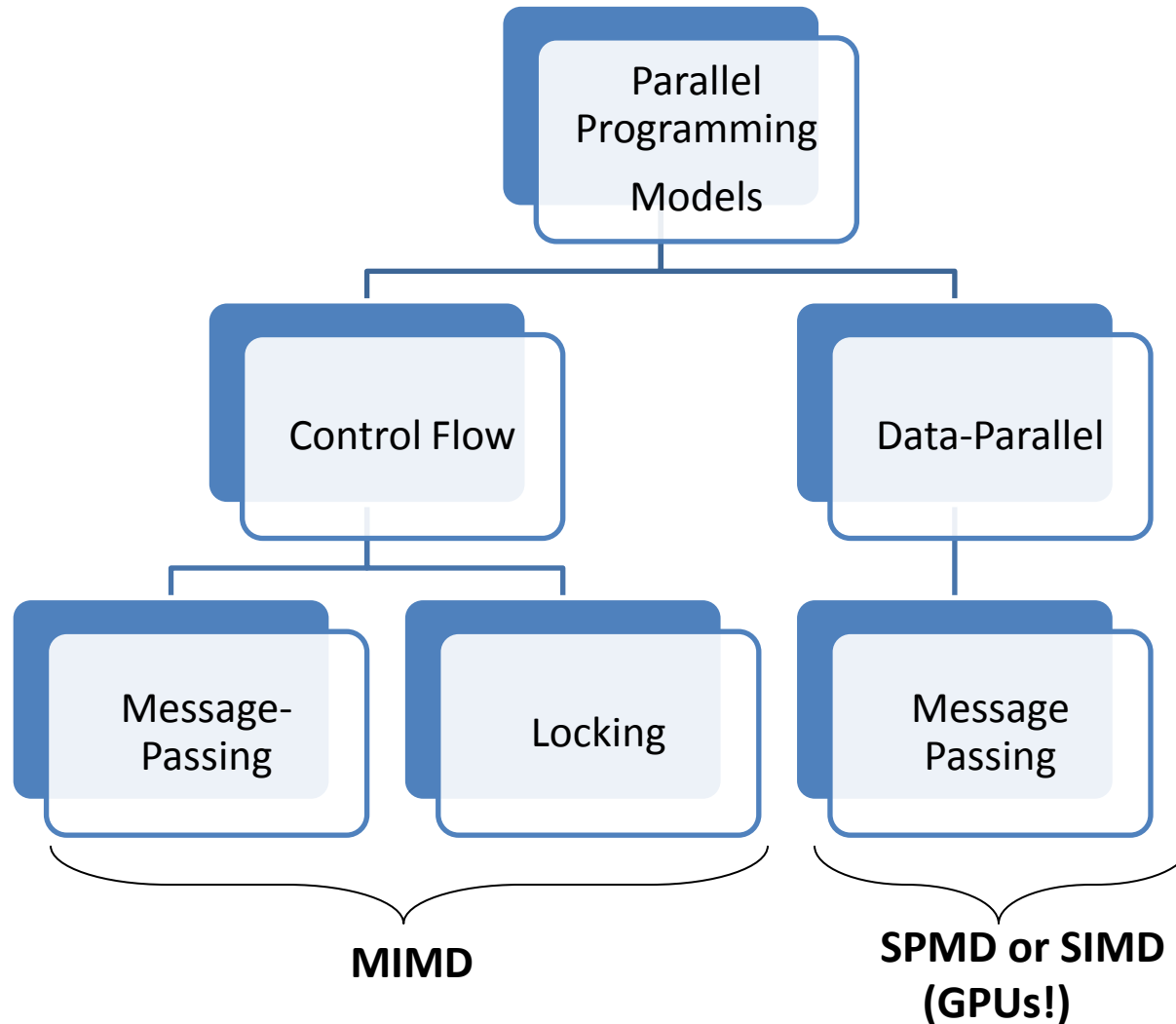
After more than 25 years of research, we are not closer to solving the parallel programming model!

# The Mentality of Yet Another Programming Language … Doesn't work!

| | | | | | |
|---|---|---|---|---|---|
| ABCPL | CORRELATE | GLU | Mentat | Parafrase2 | pC++ |
| ACE | CPS | GUARD | Legion | Paralation | SCHEDULE |
| ACT++ | CRL | HAsL | Meta Chaos | Parallel-C++ | SciTL |
| Active messages | CSP | Haskell | Midway | Parallaxis | SDDA |
| Adl | Cthreads | HPC++ | Millipede | ParC | SHMEM |
| Adsmith | CUMULVS | JAVAR | CparPar | ParLib++ | SIMPLE |
| ADDAP | DAGGER | HORUS | Mirage | ParLin | Sina |
| AFAPI | DAPPLE | HPC | MpC | Parmacs | SISAL |
| ALWAN | Data Parallel C | IMPACT | MOSIX | Parti | distributed smalltalk |
| AM | DC++ | ISIS | Modula-P | pC | SMI |
| AMDC | DCE++ | JAVAR | Modula-2* | PCN | SONiC |
| AppLeS | DDD | JADE | Multipol | PCP | Split-C |
| Amoeba | DICE | Java RMI | MPI | PH | SR |
| ARTS | DIPC | javaPG | MPC++ | PEACE | Sthreads |
| Athapascan-0b | DOLIB | JavaSpace | Munin | PCU | Strand |
| Aurora | DOME | JIDL | Nano-Threads | PET | SUIF |
| Automap | DOSMOS | Joyce | NESL | PENNY | Synergy |
| bb_threads | DRL | Khoros | NetClasses++ | Phosphorus | Telegrphos |
| Blaze | DSM-Threads | Karma | Nexus | POET | SuperPascal |
| BSP | Ease | KOAN/Fortran-S | Nimrod | Polaris | TCGMSG |
| BlockComm | ECO | LAM | NOW | POOMA | Threads.h++ |
| C* | Eiffel | Lilac | Objective Linda | POOL-T | TreadMarks |
| "C* in C | Eilean | Linda | Occam | PRESTO | TRAPPER |
| C** | Emerald | JADA | Omega | P-RIO | uC++ |
| CarlOS | EPL | WWWinda | OpenMP | Prospero | UNITY |
| Cashmere | Excalibur | ISETL-Linda | Orca | Proteus | UC |
| C4 | Express | ParLin | OOF90 | QPC++ | V |
| CC++ | Falcon | Eilean | P++ | PVM | ViC* |
| Chu | Filaments | P4-Linda | P3L | PSI | Visifold V-NUS |
| Charlotte | FM | POSYBL | Pablo | PSDM | VPE |
| Charm | FLASH | Objective-Linda | PADE | Quake | Win32 threads |
| Charm++ | The FORCE | LiPS | PADRE | Quark | WinPar |
| Cid | Fork | Locust | Panda | Quick Threads | XENOOPS |
| Cilk | Fortran-M | Lparx | Papers | Sage++ | XPC |
| CM-Fortran | FX | Lucid | AFAPI | SCANDAL | Zounds |
| Converse | GA | Maisie | Para++ | SAM | ZPL |
| Code | GAMMA | Manifold | Paradigm | | |
| COOL | Glenda | | | | |

# Parallel Programming Models

# Programming Model

- **Definition:** the languages and libraries that create an abstract view of the machine
- Control
  - How is parallelism created?
  - How are dependencies enforced?
- Data
  - Shared or private?
  - How is shared data accessed or private data communicated?
- Synchronization
  - What operations can be used to coordinate parallelism
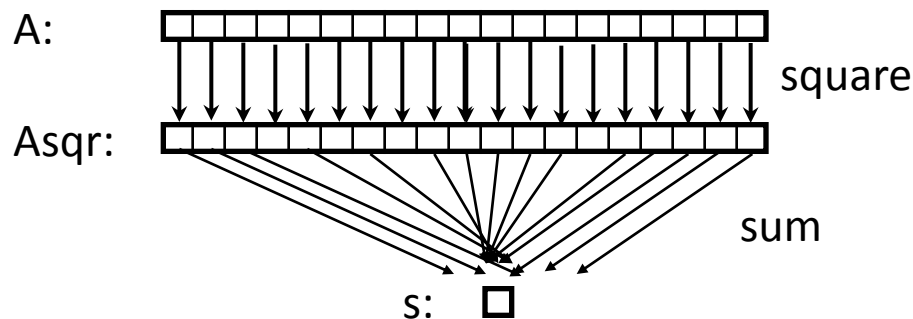  - What are the atomic (indivisible) operations?

# It Is Important to Note

- You can run any paradigm on any hardware (e.g. an MPI on shared-memory)
- The same program can have different type of parallel paradigms
- The hardware itself can be heterogeneous

The whole challenge of parallel programming is to make the best use of the underlying hardware to exploit the different type of parallelisms

# Example

We have a matrix A. We need to form another matrix Asqr that contains the square of each element of A. Then we need to calculate S, which is the sum of the elements in Asqr.
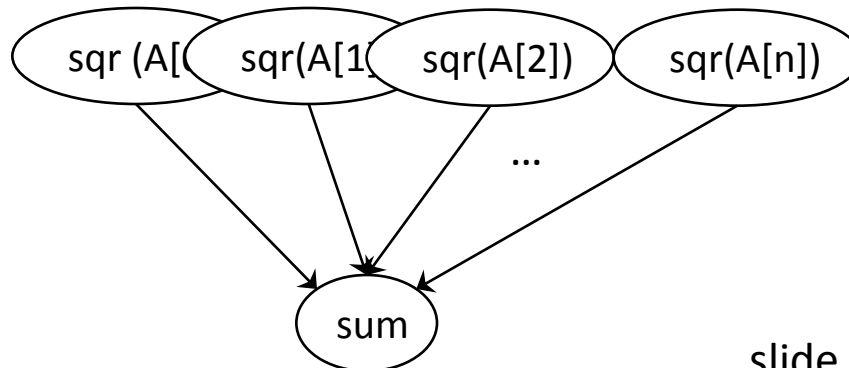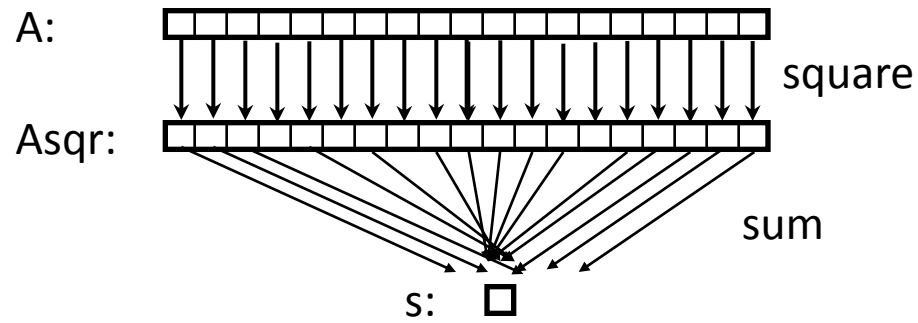
A:

square

Asqr:

sum

s:

- **How can we parallelize this?**
- **How long will it take if we have unlimited number of processors?**

slide derived from Katherine Yelick

# Example

- First, decompose your problem into a set of tasks
  - There are many ways of doing it.
  - Tasks can be of the same, different, or undetermined sizes.
- Draw a task-dependency graph (do you remember the DAG we saw earlier?)
  - A directed graph with Nodes corresponding to tasks
  - Edges indicating dependencies, that the result of one task is required for processing the next.

slide derived from Katherine Yelick

# Example



slide derived from Katherine Yelick

Does your knowledge of the underlying hardware change your task dependency graph? If yes, how?

# Conclusions

- Concurrency and parallelism are not exactly the same thing.
- There is parallelism at different granularities, with methods to exploit each parallelism granularity.
- You need to know the difference among: threads/processors/tasks.
- Knowing the hardware will help you generating a better task dependency graph.