# Moore's Law

*Forecast*: #transistors/chip when cost-per-transistor is lowest



"ICs will lead to such wonders as home computers"



**1965** *forecast: double every year for minimum component cost*



*Revised*

*Original*

**1975** *forecast: double every two years*
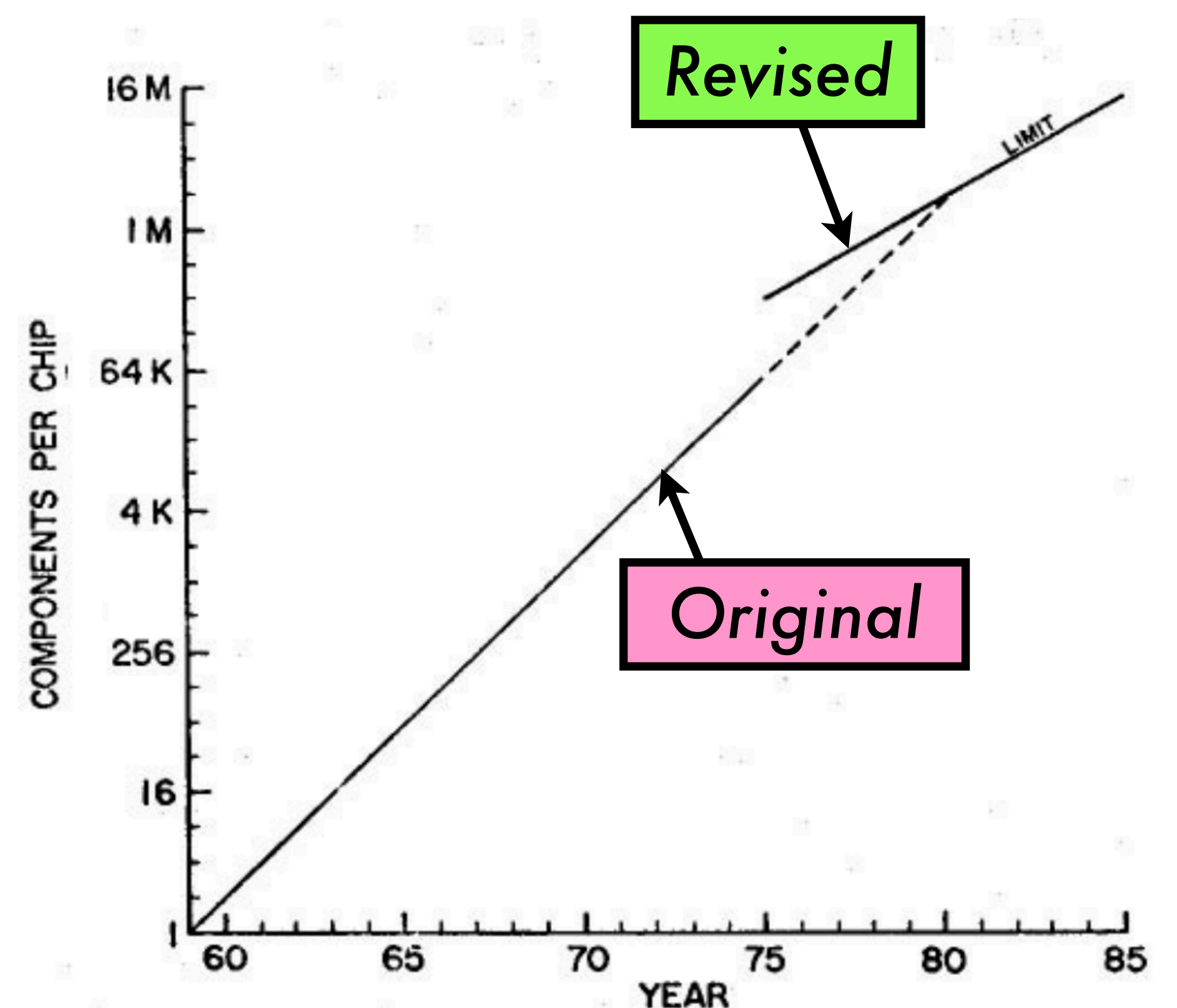
*Progress in Digital Integrated Electronics* (1975)

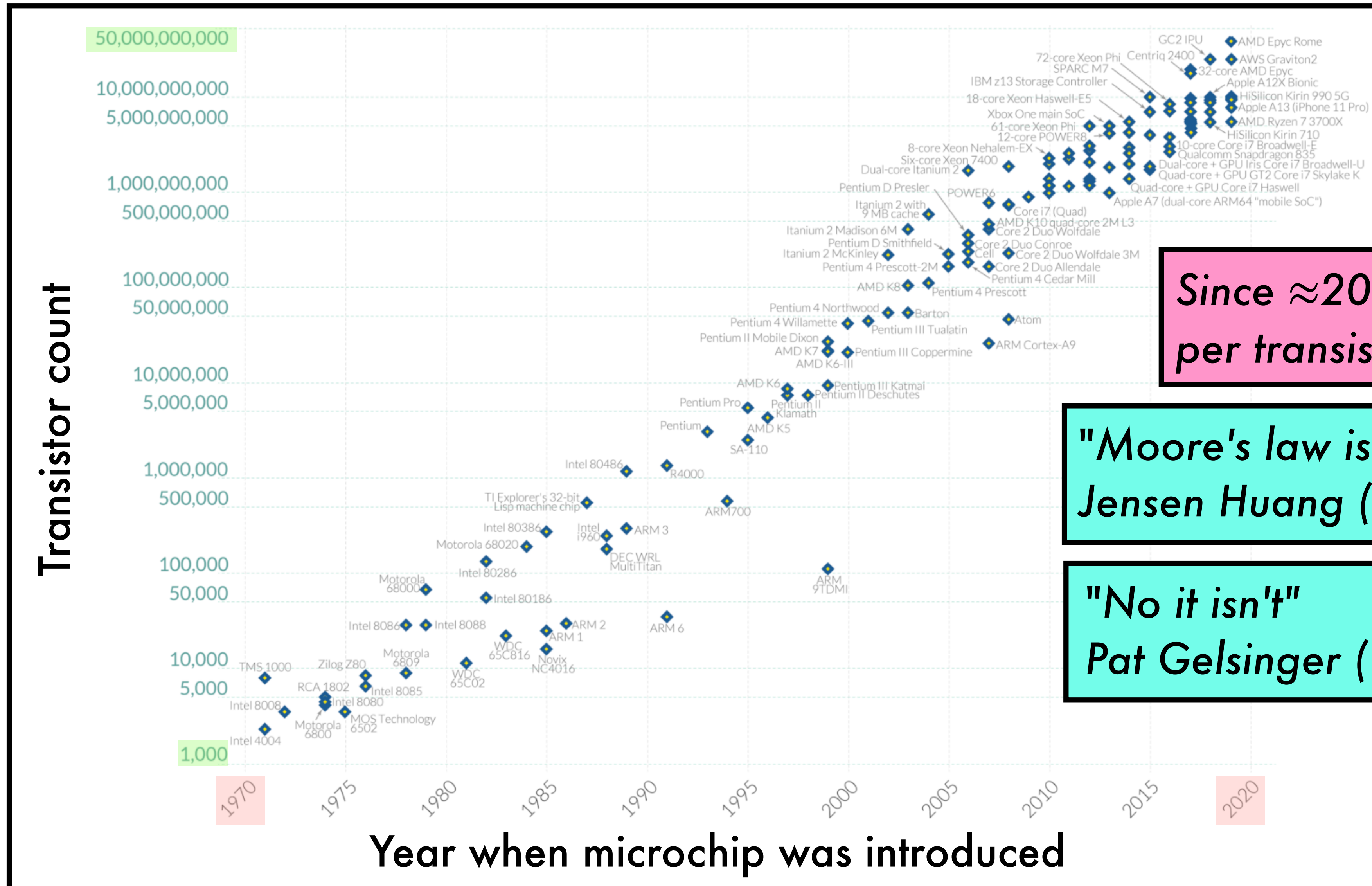References/Notes/Image credits:
(G. Moore image) https://www.businesswire.com/news/home/20161102005463/en/Gordon-and-Betty-Moore-Foundation-Announces-Inaugural-Moore-Inventor-Fellows
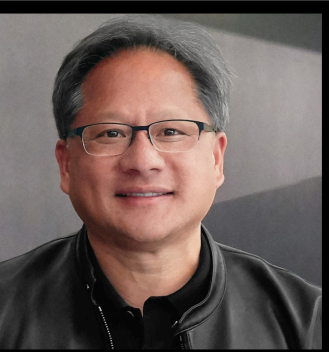(Projection figure from 1965) G. E. Moore, "Cramming more components onto integrated circuits" (1965)
(Projection figure from 1975) G. E. Moore, "Progress in digital integrated electronics", Electron devices meeting (1975)

# Moore's Law - Historical Data



Scatter plot of transistor count vs. year when microchip was introduced (log scale on y-axis).

**Y-axis (Transistor count):** 50,000,000,000 — 10,000,000,000 — 5,000,000,000 — 1,000,000,000 — 500,000,000 — 100,000,000 — 50,000,000 — 10,000,000 — 5,000,000 — 1,000,000 — 500,000 — 100,000 — 50,000 — 10,000 — 5,000 — 1,000

**X-axis (Year when microchip was introduced):** 1970, 1975, 1980, 1985, 1990, 1995, 2000, 2005, 2010, 2015, 2020

Selected labeled data points: GC2 IPU, AMD Epyc Rome, 72-core Xeon Phi, Centriq 2400, AWS Graviton2, SPARC M7, 32-core AMD Epyc, IBM z13 Storage Controller, Apple A12X Bionic, 18-core Xeon Haswell-E5, HiSilicon Kirin 990 5G, Xbox One main SoC, Apple A13 (iPhone 11 Pro), 61-core Xeon Phi, AMD Ryzen 7 3700X, 12-core POWER8, HiSilicon Kirin 710, 8-core Xeon Nehalem-EX, 10-core Core i7 Broadwell-E, Six-core Xeon 7400, Qualcomm Snapdragon 835, Dual-core Itanium 2, Dual-core + GPU Iris Core i7 Broadwell-U, Quad-core + GPU GT2 Core i7 Skylake K, Pentium D Presler, POWER6, Quad-core + GPU Core i7 Haswell, Itanium 2 with 9 MB cache, Core i7 (Quad), Apple A7 (dual-core ARM64 "mobile SoC"), Itanium 2 Madison 6M, AMD K10 quad-core 2M L3, Pentium D Smithfield, Core 2 Duo Wolfdale, Itanium 2 McKinley, Core 2 Duo Conroe, Cell, Core 2 Duo Wolfdale 3M, Pentium 4 Prescott-2M, Core 2 Duo Allendale, AMD K8, Pentium 4 Cedar Mill, Pentium 4 Prescott, Pentium 4 Northwood, Barton, Atom, Pentium 4 Willamette, Pentium III Tualatin, Pentium II Mobile Dixon, ARM Cortex-A9, AMD K7, Pentium III Coppermine, AMD K6-III, AMD K6, Pentium III Katmai, Pentium II Deschutes, Pentium Pro, Pentium II Klamath, Pentium, AMD K5, SA-110, Intel 80486, R4000, TI Explorer's 32-bit Lisp machine chip, ARM700, Intel 80386, Intel i960, ARM 3, Motorola 68020, DEC WRL MultiTitan, Intel 80286, Intel 80186, ARM 9TDMI, Motorola 68000, Intel 8086, Intel 8088, ARM 2, ARM 6, WDC 65C816, ARM 1, Motorola 6809, WDC 65C02, Novix NC4016, TMS 1000, Zilog Z80, Intel 8085, RCA 1802, Intel 8080, Intel 8008, Motorola 6800, MOS Technology 6502, Intel 4004

Annotations:
- (pink box) *Since ≈2010, cost decline per transistor is slowing*
- (cyan box) *"Moore's law is dead"* Jensen Huang ('17, '22)
- (cyan box) *"No it isn't"* Pat Gelsinger ('22)

# Dennard Scaling



*Inventor of DRAM*

As transistors shrink, power density remains constant

## Dennard's model of MOSFET scaling

With each generation, transistor dimensions shrink by 30%

| | |
|---|---|
| Device area shrinks by 50% | area $=$ length $\times$ width |
| Capacitance shrinks by 30% | capacitance $\propto$ area/distance |
| Voltage is reduced by 30% | electric field $\propto$ voltage/distance |
| Circuit delay reduces by 30% | due to reduced gate delays |
| Frequency increases by 40% | frequency $=$ $1$/time period |
| Active power reduces by 50% | active power $\propto CV^2f$ |

Each generation: double # transistors | same power | 40% faster

References/Notes/Image credits:
(image source) https://www.ibm.com/blogs/think/2019/11/ibms-robert-h-dennard-and-the-chip-that-changed-the-world/
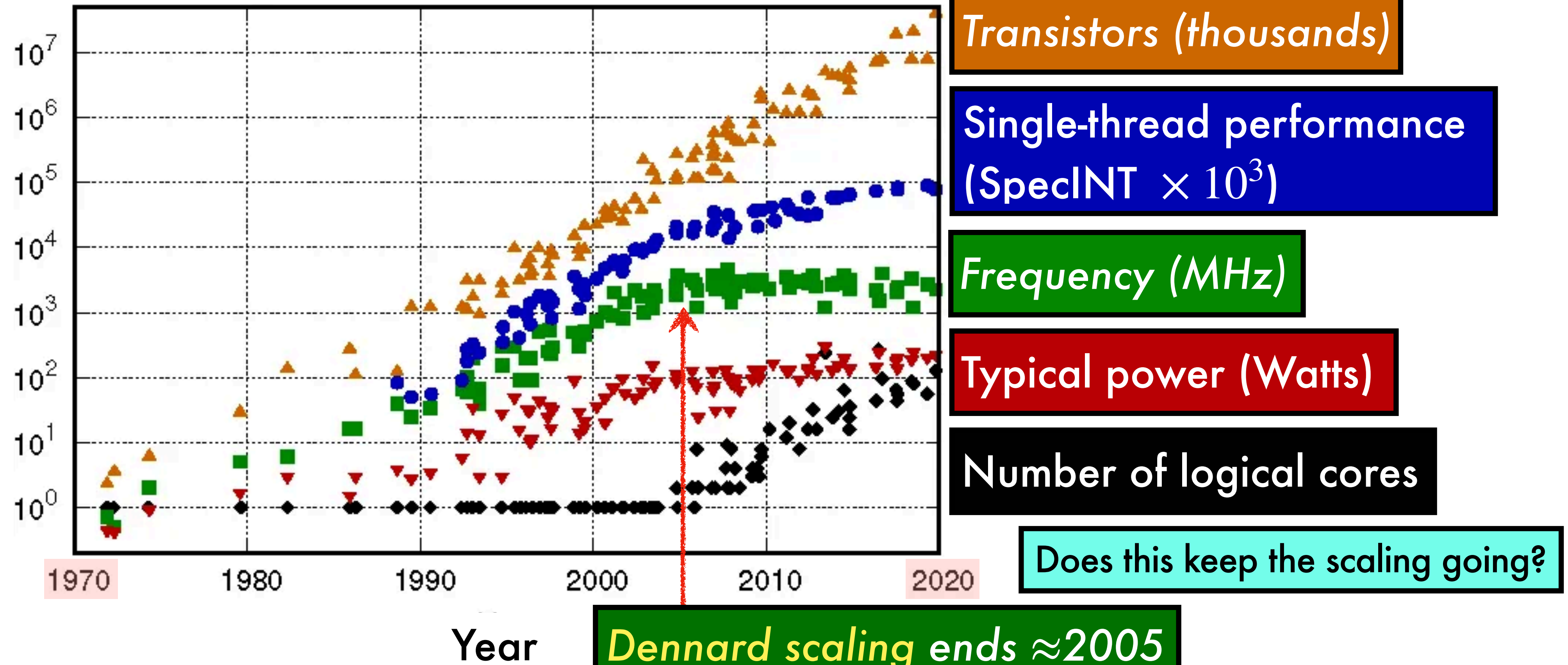R. H. Dennard et al., "Design of ion-implanted MOSFET's with very small physical dimensions", *IEEE Journal of solid-state circuits* (1974)
https://en.wikipedia.org/wiki/Dennard_scaling
S. Borkar et al., "The future of microprocessors", *Communications of the ACM* (2011)

# The End Of Dennard Scaling



Microprocessor trends

**Transistors (thousands)**

**Single-thread performance (SpecINT $\times 10^3$)**

*Frequency (MHz)*

Typical power (Watts)

Number of logical cores

Does this keep the scaling going?

Year

*Dennard scaling ends $\approx 2005$*

*Dennard's model ignored "leakage current"*

# Amdahl's Law 💔

Many programs contain code that **cannot be parallelised**   **Notation**   **Processors, $P$**   **Time, $T$**

**Serial code**

*Time to execute program*

**Parallelisable code**

**Amdahl's law**

$P = 1$   $T = 8$

$P = 2$   $T = 6$

$P = 4$   $T = 5$

Let $\alpha \in [0,1]$ be the **fraction** of code that can be parallelised

$$T_{new} = T_{orig} \times \left( (1 - \alpha) + \frac{\alpha}{P} \right)$$

$P = 4$   $T_{new} = 8$   $\alpha = 0.5$

$$T_{new} = 8 \cdot (0.5 + 0.125) = 5$$

**Amdahl's law is an (often loose) <u>upper bound</u> on parallelism**

Reference:
G. M. Amdahl, "Validity of the single processor approach to achieving large scale computing capabilities", Proceedings of the Spring Joint Computer Conference (1967)
J. Hennessy and D. Patterson, "Computer Architecture: A Quantitative Approach", Chap. 1 (2017)

# Gustafson's Law

There is a key idea underpinning the interpretation of Amdahl's law:

Problem size stays fixed even as more processors become available

Given more processing, the problem expands to use it

Better: assume runtime (not problem size) is constant

This is virtually never the case!

In practice, the problem scales with # processors

**Parkinson's law (1955)**

*"Work expands so as to fill the time available for its completion."* - Cyril Parkinson

Often, the parallel part of the program scales with problem size

$P$ Processors    $speedup = 1 + \alpha \cdot (P - 1)$    *A linear scaling law!*

Key difference to Amdahl - we assume that we will scale up the parallel part of the problem

*Which "law" is a better fit depends on the domain*    Booting your Operating System    Amdahl

Reference:
https://en.wikipedia.org/wiki/Gustafson's_law
J. Gustafson, "Reevaluating Amdahl's law", *Comms. of the ACM* (1988)
https://en.wikipedia.org/wiki/Parkinson's_law

Train model on 1000s of GPUs    Gustafson    Boot stronger machine

# Memory Models For Parallel Computing

Since the end of Dennard scaling (≈2005), multicore computers have become pervasive

Important design choice with multiple cores/multiple multiprocessors: how to organise memory

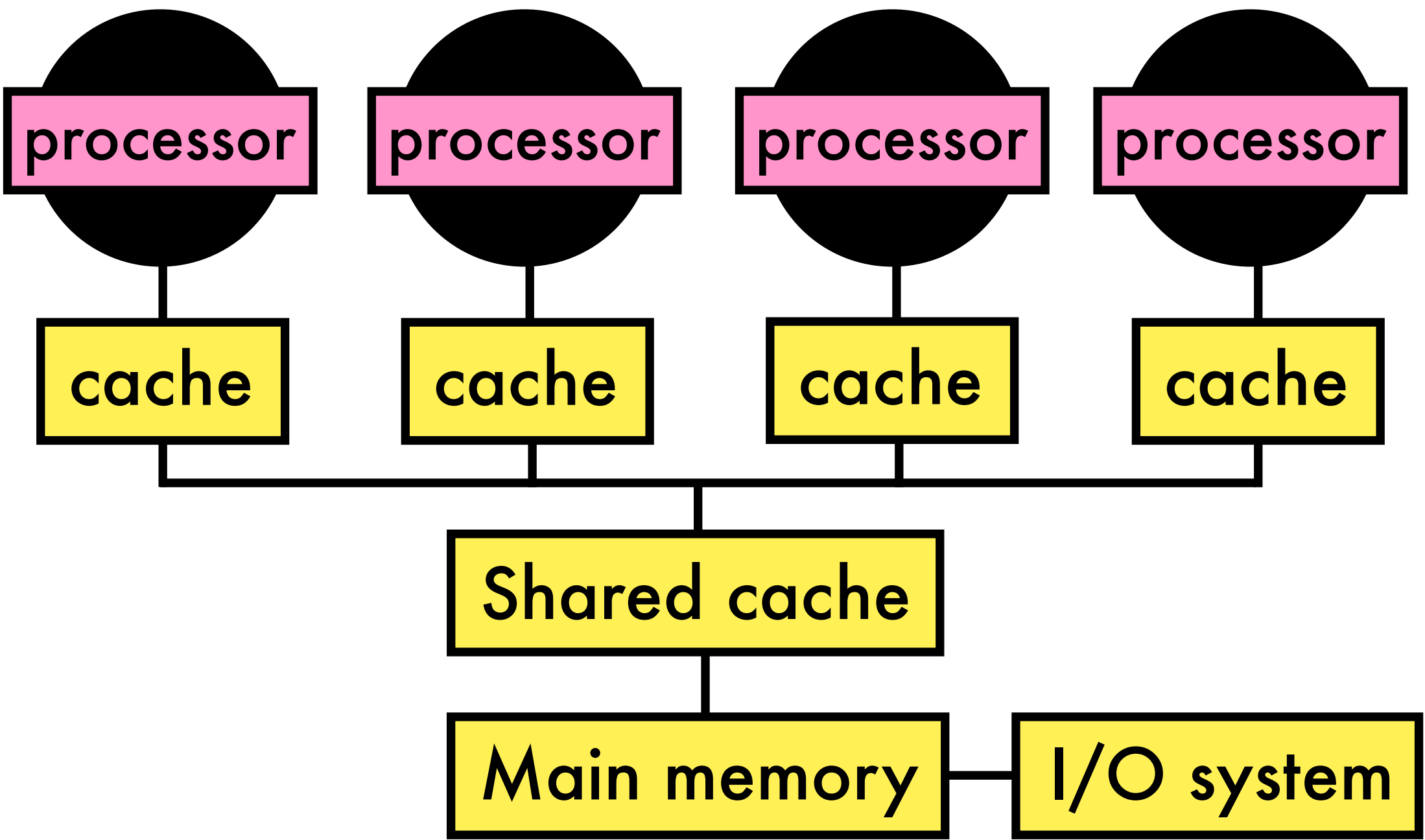| Shared Memory | Distributed Memory |
|---|---|
| **Key idea:** any core can directly access any location in a *shared address space* | **Key idea:** each core sends **message** (over network) to access memory belonging to another core |
| Typical hardware: phones laptops | Typical hardware: compute clusters |
| *"Hard to build, easy to program"* | *"Easy to build, hard to program"* |

References:
J. Hennessy and D. Patterson. "Computer Architecture: A Quantitative Approach", Chap. 5 (2017)

# Shared Memory Variants

## Symmetric Multiprocessing (SMP)
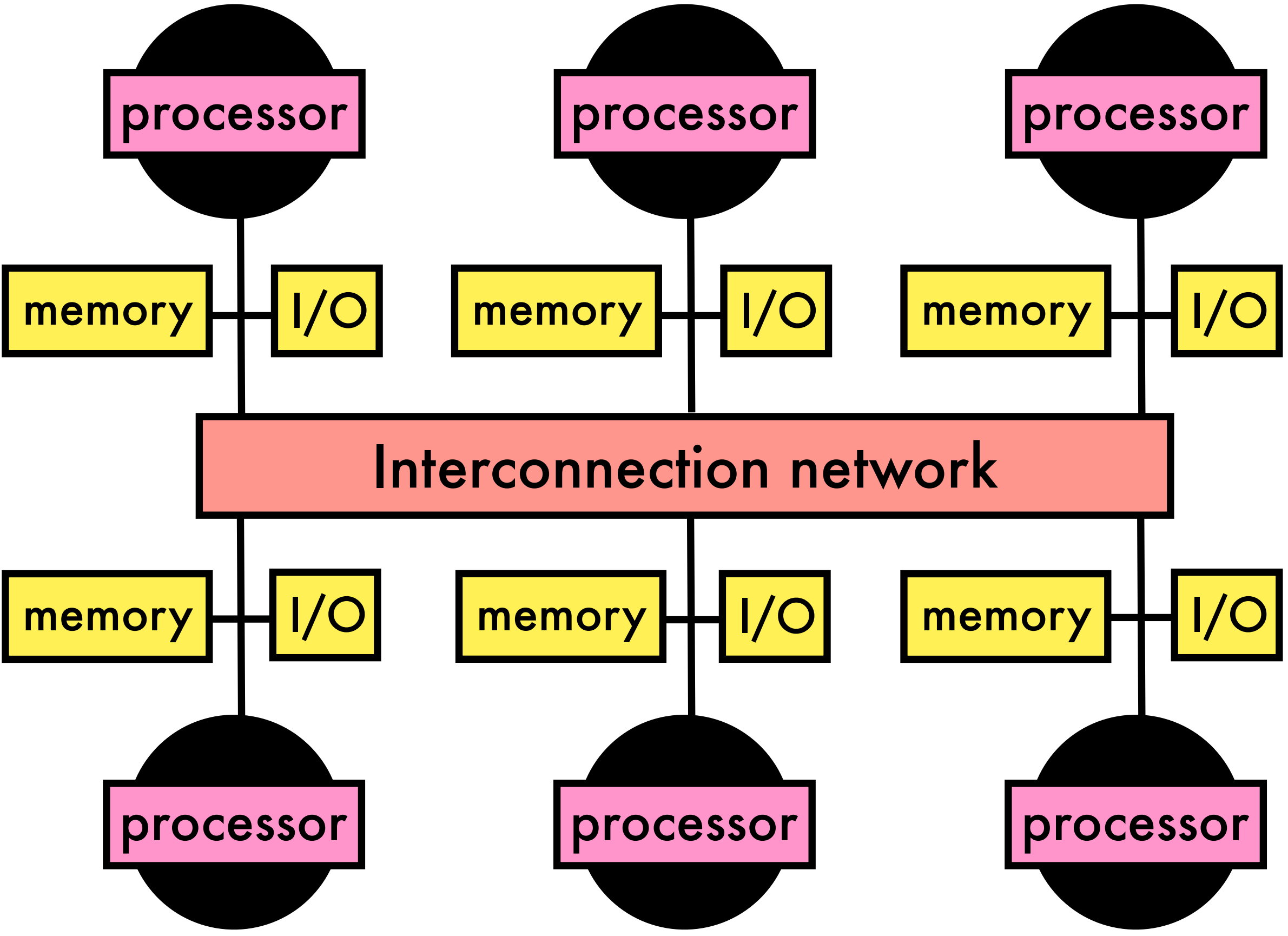
Key idea: uniform access time to all memory

Uniform Memory Access (UMA)

| processor | processor | processor | processor |
| --- | --- | --- | --- |
| cache | cache | cache | cache |

Shared cache

Main memory — I/O system

## Distributed Shared Memory (DSM)

Key idea: access time depends on location of data

Non-Uniform Memory Access (NUMA)

| processor | processor | processor |
| --- | --- | --- |
| memory — I/O | memory — I/O | memory — I/O |

Interconnection network

| memory — I/O | memory — I/O | memory — I/O |
| --- | --- | --- |
| processor | processor | processor |

References:
J. Hennessy and D. Patterson. "Computer Architecture: A Quantitative Approach", Chap. 5 (2017)

# Forms Of Parallelism

## Data-level parallelism

Distribute data across processors

Processors perform the same task on different subsets of data in parallel

## Task-level parallelism

Distribute tasks across processors

Different tasks may be run on the same data (as well as across different subsets of data)

more general    but also more complex

We will focus on task-level parallelism with a shared memory model

References:
https://en.wikipedia.org/wiki/Data_parallelism
https://en.wikipedia.org/wiki/Task_parallelism

# Task-Parallel Platforms

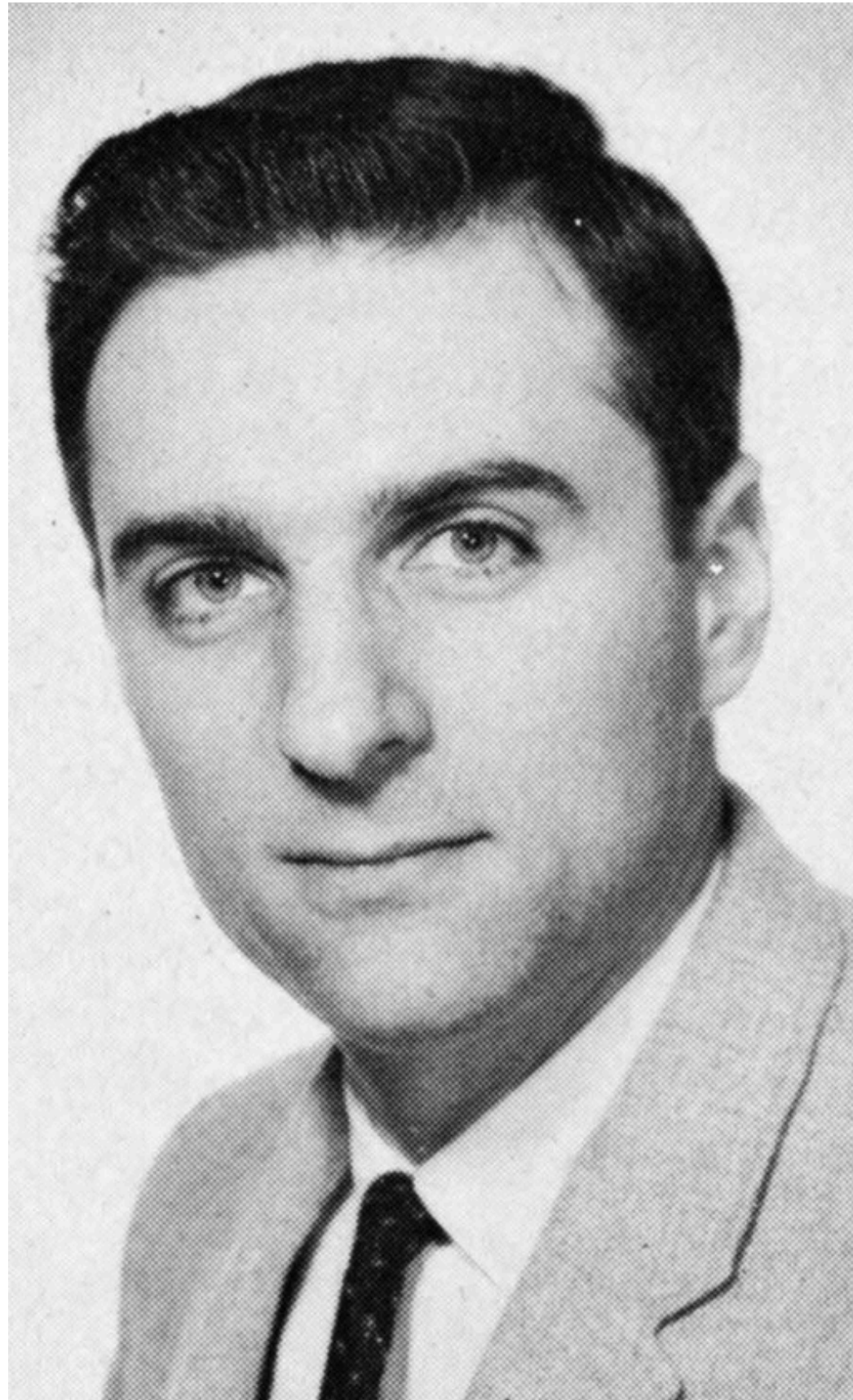| Implementing task-level parallelism with threads |
|---|
| Task parallelism can be implemented with threads ("virtual processors") that share memory However, this has proven difficult to program: Scheduling/load-balancing is a challenging job |

| Task-parallel platforms |
|---|
| Add abstraction layer on top of threads Programmer specifies which tasks can run in parallel (but not where they run) Platform manages scheduling, balancing etc. |

References:
(CLRS) T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)

# Fork-Join Parallelism

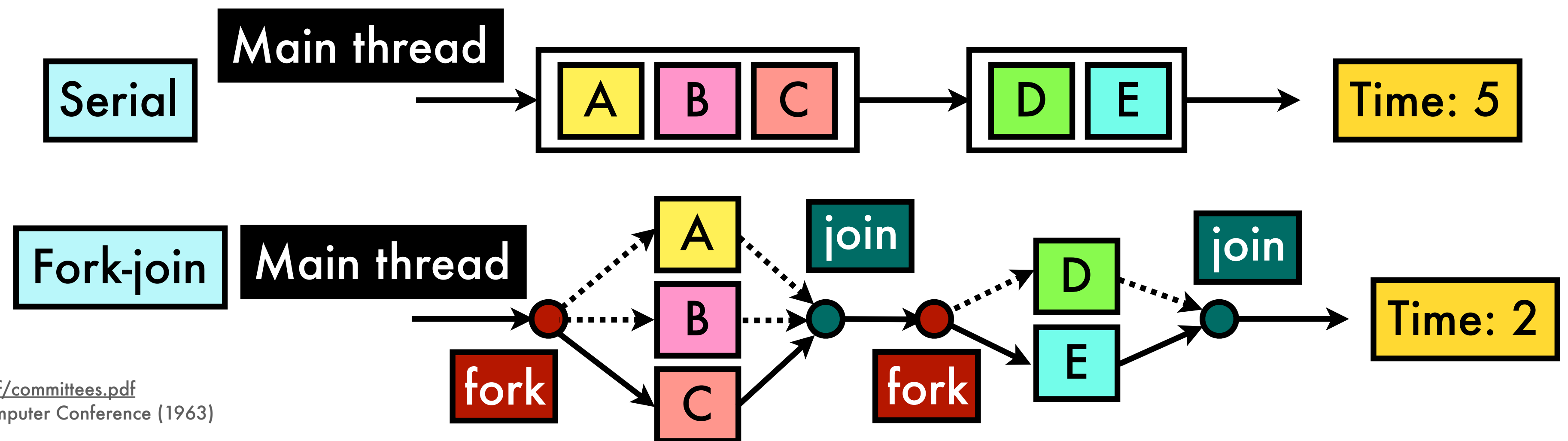Most task-parallel platforms support fork-join   | Cilk | OneTBB | OpenMP |

Spawn: "forks" - executes function while caller continues to run in parallel

Sync: "joins" - waits for spawned threads to finish before proceeding

Key concept: programmer only specifies which tasks can run in parallel, not which tasks must run in parallel

Parallel sections can fork recursively until reaching a given task granularity

**Fork-join parallelism (1963)**

References:
(M. Conway image) http://www.melconway.com/Home/pdf/committees.pdf
M. Conway, "A multiprocessor system design", Fall Joint Computer Conference (1963)
https://en.wikipedia.org/wiki/Fork-join_model
(CLRS) T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)
R. Blumofe et al., "Cilk: An efficient multithreaded runtime system", ACM SigPlan (1995)
(OneTBB) https://github.com/oneapi-src/oneTBB

Serial — Main thread → A B C → D E → Time: 5

Fork-join — Main thread → fork → A B C → join → fork → D E → join → Time: 2

# An Example: Fibonacci

## Fibonacci Numbers

```python
def fib(n):
    if n < 2:
        return n
    x = fib(n - 1)
    y = fib(n - 2)
    return x + y
```
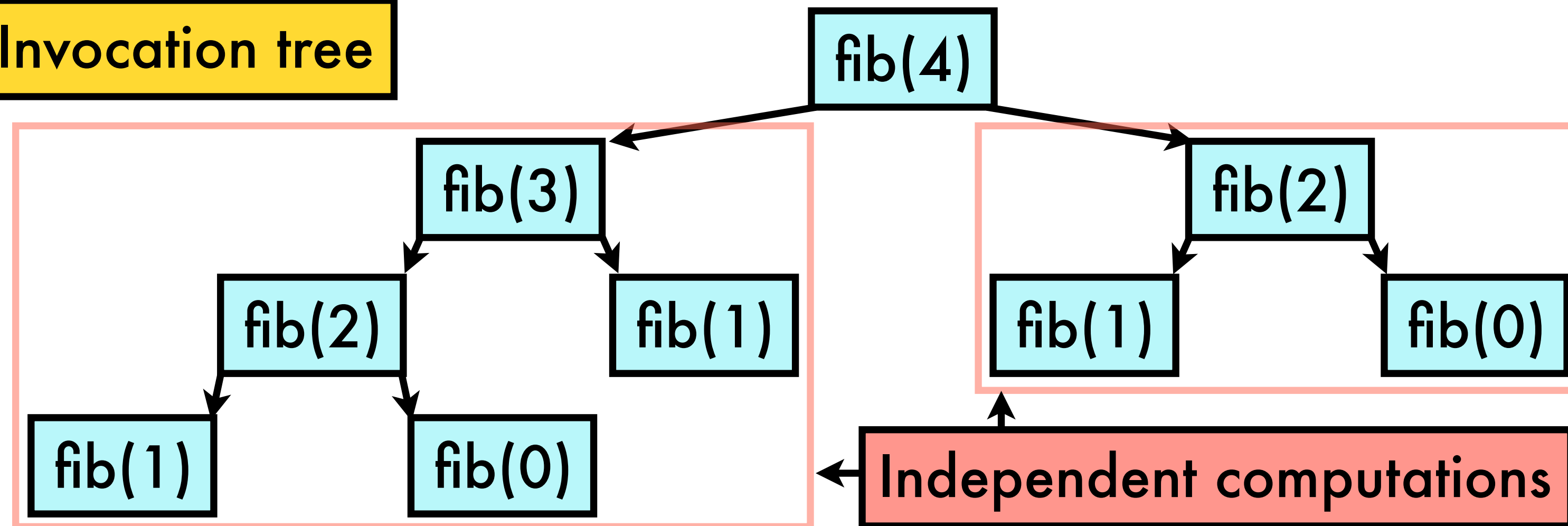
Not very efficient (no memoization)

$$T(n) = T(n-1) + T(n-2) + \Theta(1)$$

$$T(n) = \Theta\left(\left(\frac{1+\sqrt{5}}{2}\right)^n\right) \text{ (exponential)}$$

## Invocation tree

fib(4)

fib(3)        fib(2)

fib(2)   fib(1)   fib(1)   fib(0)

fib(1)   fib(0)

Independent computations

References:
(CLRS) T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)

# Parallel Code

**Parallel pseudocode (CLRS/cilk)**

```
def par_fib(n):
    if n < 2:
        return n
    x = spawn par_fib(n - 1)
    y = par_fib(n - 2)
    sync
    return x + y
```

**Parallel version (Python)**   CPython GIL 🐢

Use `nogil` + coarsening for speedup

```
def par_fib(n):
    if n < 2:
        return n
    with ThreadPoolExecutor() as exec:
        x_future = exec.submit(par_fib, n - 1) # spawn
        y = par_fib(n - 2)
        x = x_future.result() # sync
    return x + y
```

`spawn` says main thread **can** execute in parallel with the spawned child, not that it **must**

`sync` says parent must **wait** for all spawned children to finish (join)

`spawn/sync` express the **logical parallelism** of the tasks

It is the responsibility of the **scheduler** to assign the tasks to processors

References:
(CLRS) T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)
R. Blumofe et al., "Cilk: An efficient multithreaded runtime system", ACM SigPlan (1995)
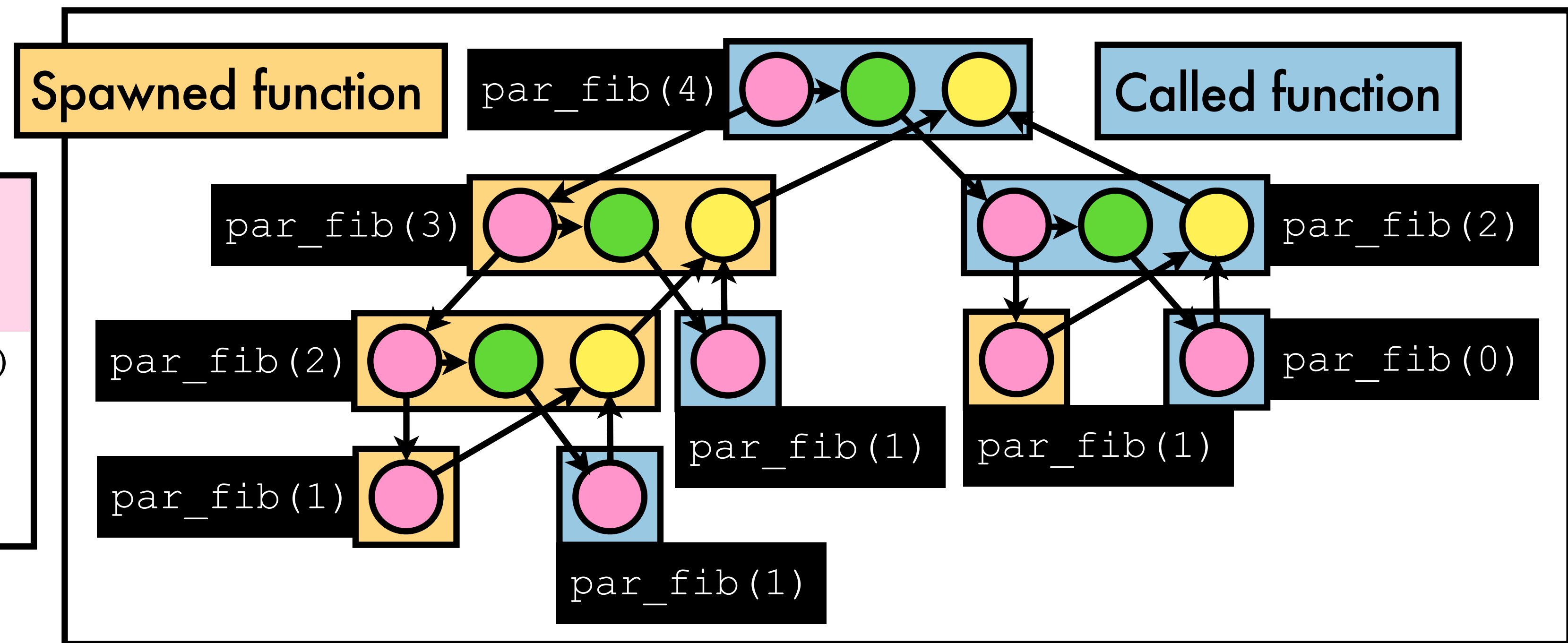(nogil python) https://nogil.dev/

# Computation DAG

We can view execution as a computation DAG (a.k.a. "parallel trace"): $G = (V, E)$

Executed instructions: vertices in $V$ | Dependencies between instructions: edges in $E$

To avoid clutter: group chains of instructions with no parallel/procedural control into "strands"



```
def par_fib(n):
    if n < 2:
        return n
    x = spawn par_fib(n - 1)
    y = par_fib(n - 2)
    sync
    return x + y
```

References:
(CLRS) T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)
https://www.csd.uwo.ca/~mmorenom/cs3101_Winter_2015/Multithreaded_Parallelism_and_Performance_Measures.pdf

# Parallel Computation Analysis: Assumptions

Assumptions for analysis

1. We have an ideal parallel computer:
- multiple processors
- sequentially consistent memory
2. Processors have equal computing power
3. No overhead for scheduling

Sequentially consistent (Lamport, 1979 ):

Instruction execution preserves the partial ordering of DAG

Attain via sequential processors; FIFO memory

(processors communicate through memory)

References/image credits:

T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)

(Sequential consistency) L. Lamport, "How to make a multiprocessor computer that correctly executes multiprocess programs", *IEEE ToC* (1979)

(Image of L. Lamport) https://en.wikipedia.org/wiki/Leslie_Lamport#/media/File:Leslie_Lamport.jpg
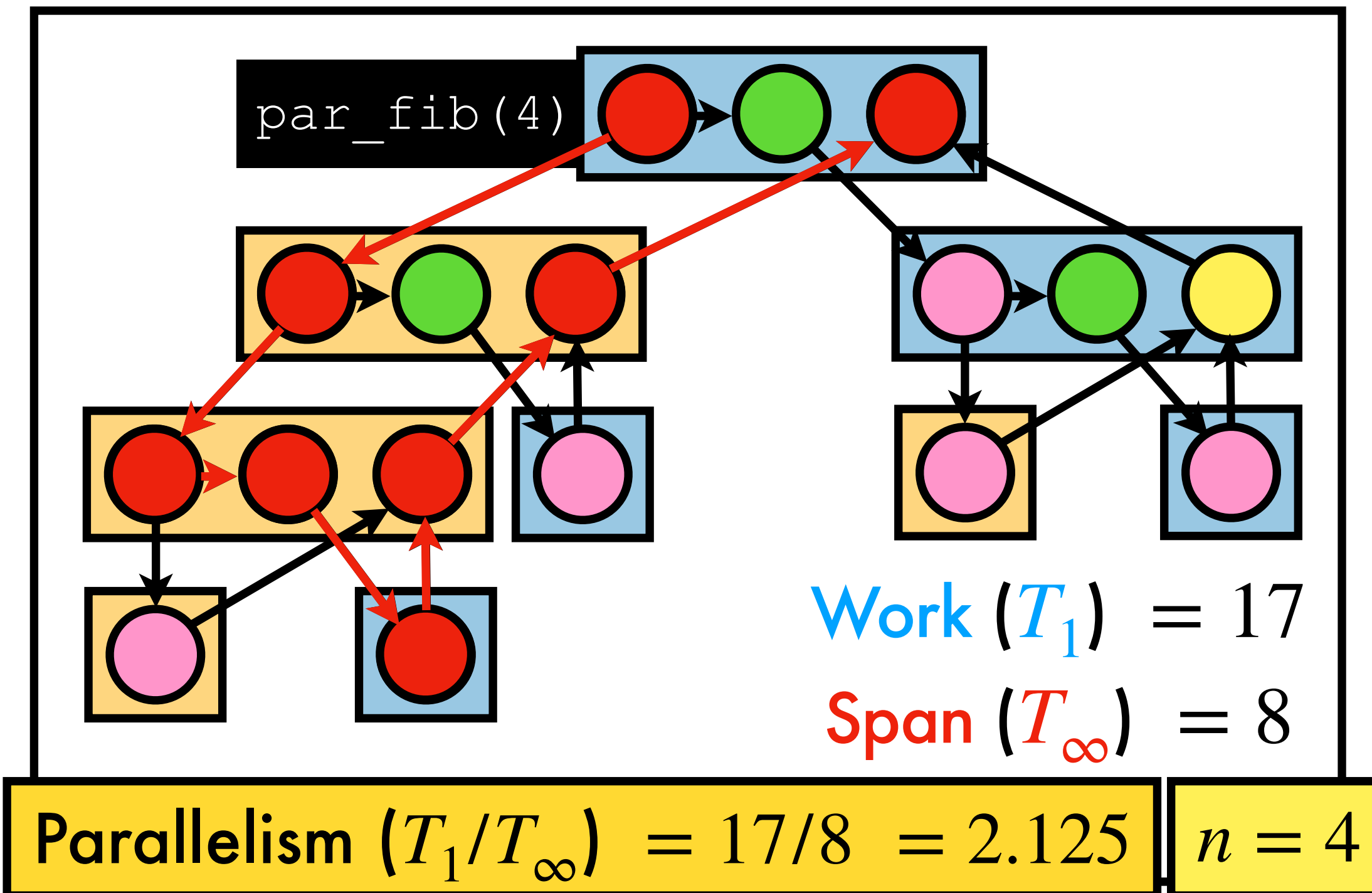
# Work/Span Analysis



Let $T_P$ denote runtime of program on $P$ processors

Work ($T_1$): time to execute program on 1 processor

Span ($T_\infty$): time to execute program on $\infty$ processors

The span is the sum of the runtimes of strands on the

"critical path" - the longest path in the computation DAG

`par_fib(4)`

Work ($T_1$) $= 17$

Span ($T_\infty$) $= 8$

Parallelism ($T_1/T_\infty$) $= 17/8 = 2.125$  |  $n = 4$

Work law: $T_P \geq T_1/P$  ($P$ processors can achieve at most $P$ units of work per time step)

Span law: $T_P \geq T_\infty$  (with $\infty$ processors, we can emulate the $P$ processors and leave rest idle)

Speedup: $T_1/T_P$  (by work law, $T_1/T_P \leq P$ - speedup on our ideal parallel machine is at most $P$)

Linear speedup: $T_1/T_P = \Theta(P)$  Perfect linear speedup: $T_1/T_P = P$  |  super-linear impossible in this model

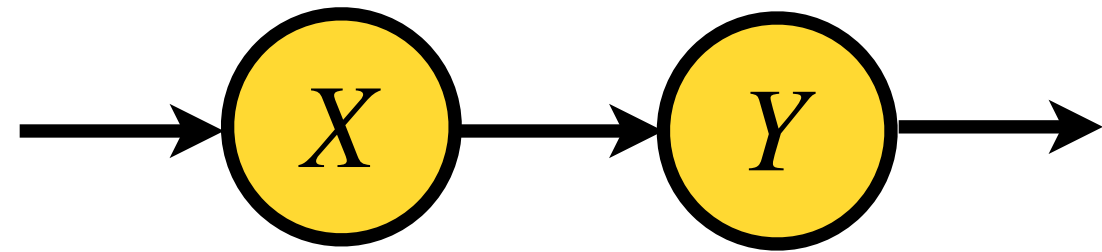Parallelism:  $T_1/T_\infty$  (maximum possible speed up with any number of processors)

References:
T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)
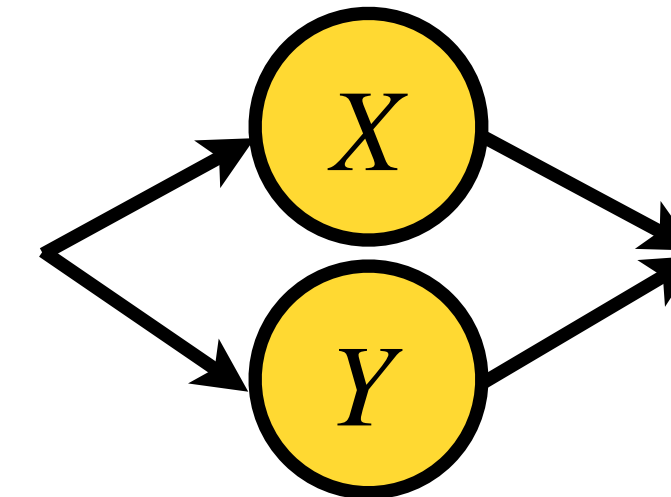https://en.wikipedia.org/wiki/Analysis_of_parallel_algorithms

# Parallel Analysis



**Serial**

Work: $T_1(X \cup Y) = T_1(X) + T_1(Y)$

Span: $T_\infty(X \cup Y) = T_\infty(X) + T_\infty(Y)$

**Parallel**

Work: $T_1(X \cup Y) = T_1(X) + T_1(Y)$

Span: $T_\infty(X \cup Y) = \max(T_\infty(X), T_\infty(Y))$

Work $T_1(n)$ for `par_fib(n)`:

$$T_1(n) = \Theta\left(\left(\frac{1+\sqrt{5}}{2}\right)^n\right)$$

Span $T_\infty(n)$ for `par_fib(n)`:

$$T_\infty(n) = \max(T_\infty(n-1), T_\infty(n-2)) + \Theta(1)$$
$$= T_\infty(n-1) + \Theta(1) \qquad \boxed{T_\infty = \Theta(n)}$$

Parallelism: $\dfrac{T_1(n)}{T_\infty(n)} = \Theta\left(\left(\dfrac{1+\sqrt{5}}{2}\right)^n / n\right)$

Grows fast with $n$

Lots of parallelism

Reference:
T. Cormen et al., "Introduction to algorithms", Chap 26, MIT press (2022)