

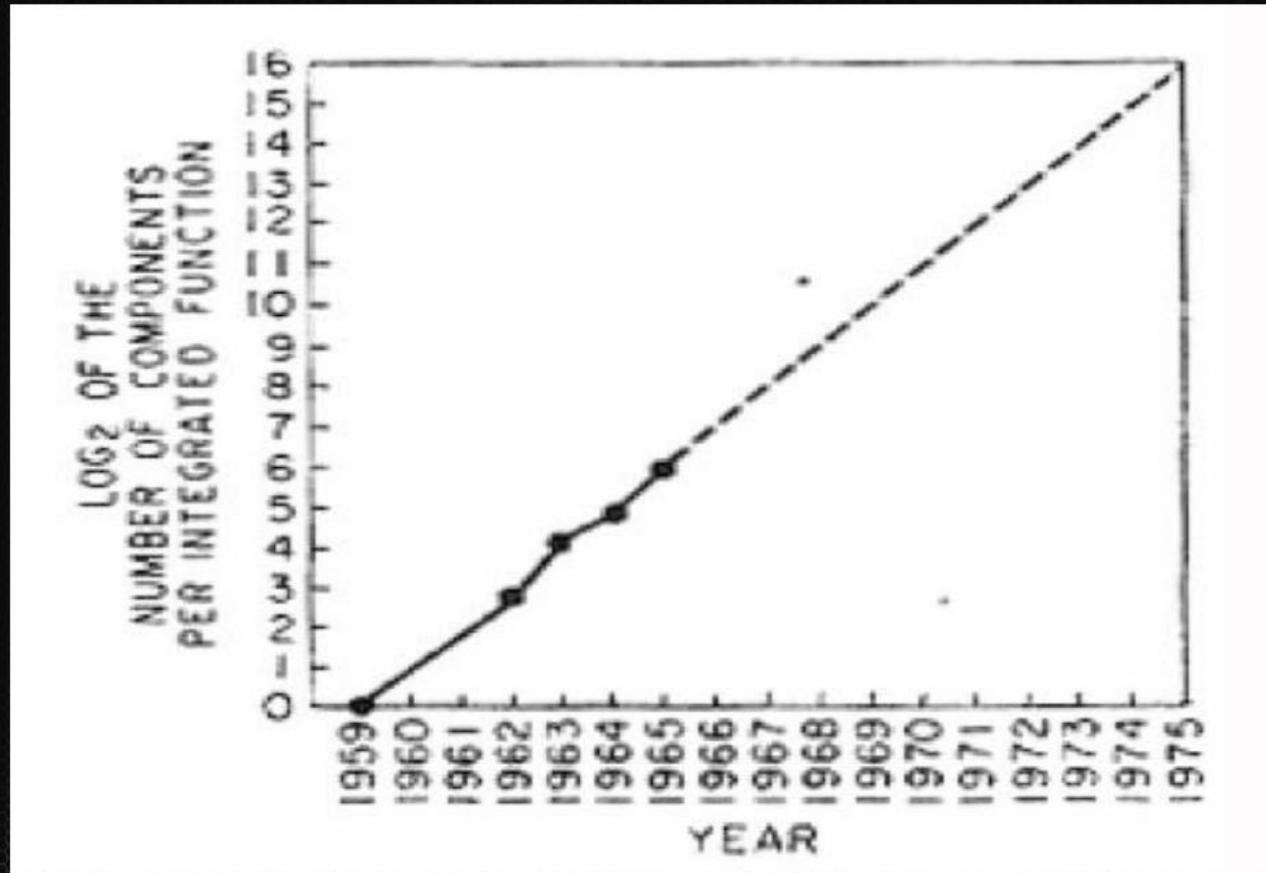
# The Progress in CPU Performance

**Moore's Law doubles processor  
performance every few years, right?**

**Wrong!**



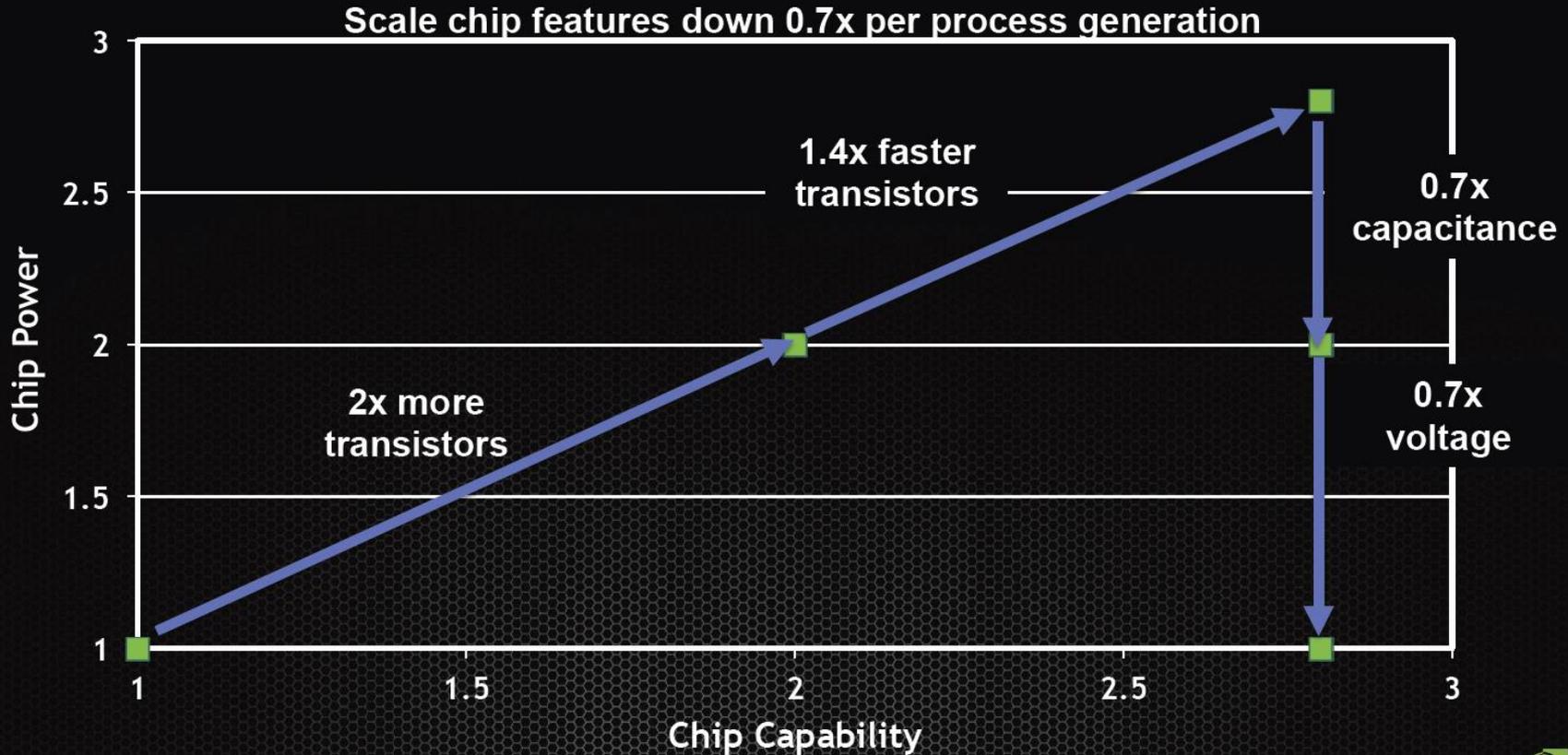
# Moore's Law gives us transistors Which we used to turn into scalar performance

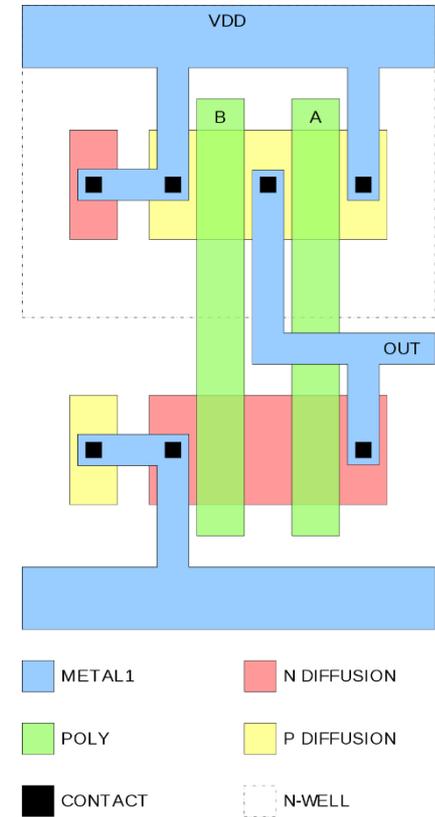
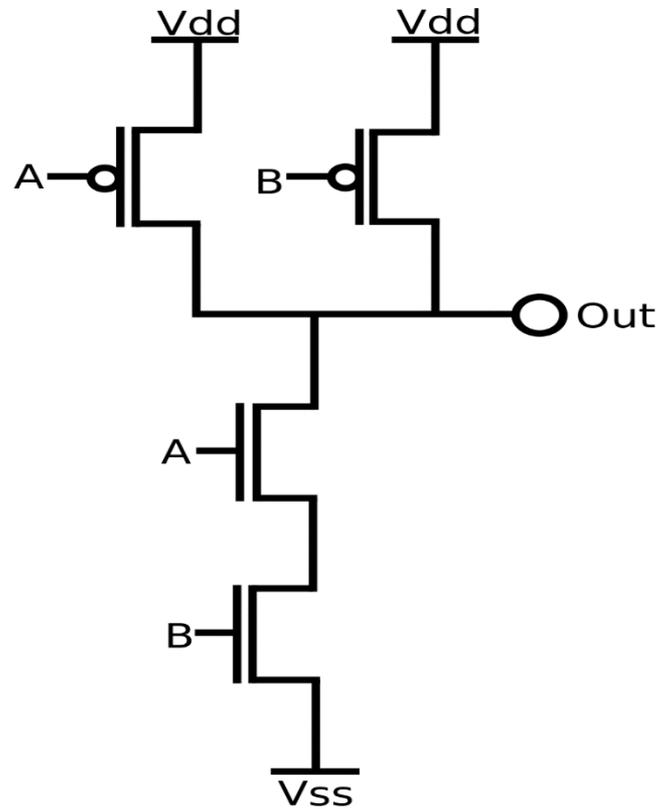
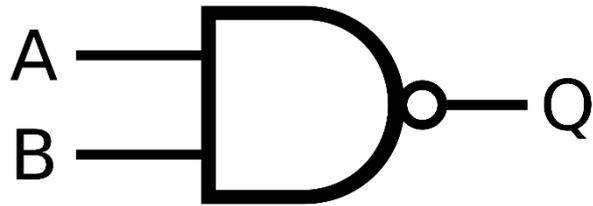


Moore, Electronics 38(8) April 19, 1965

# Classic Dennard Scaling

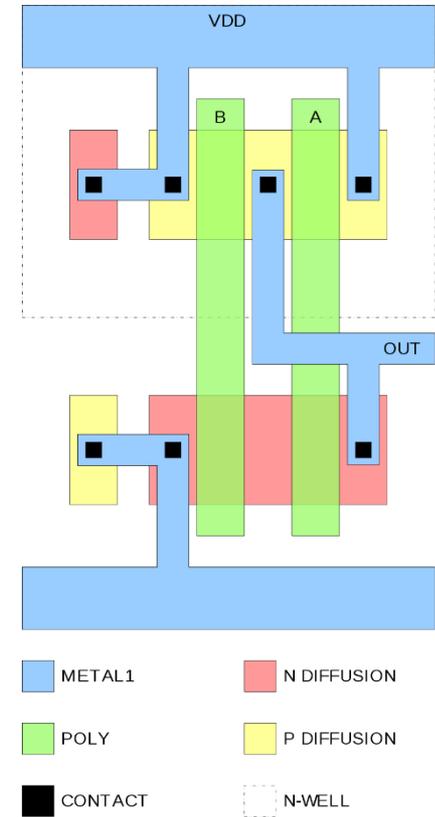
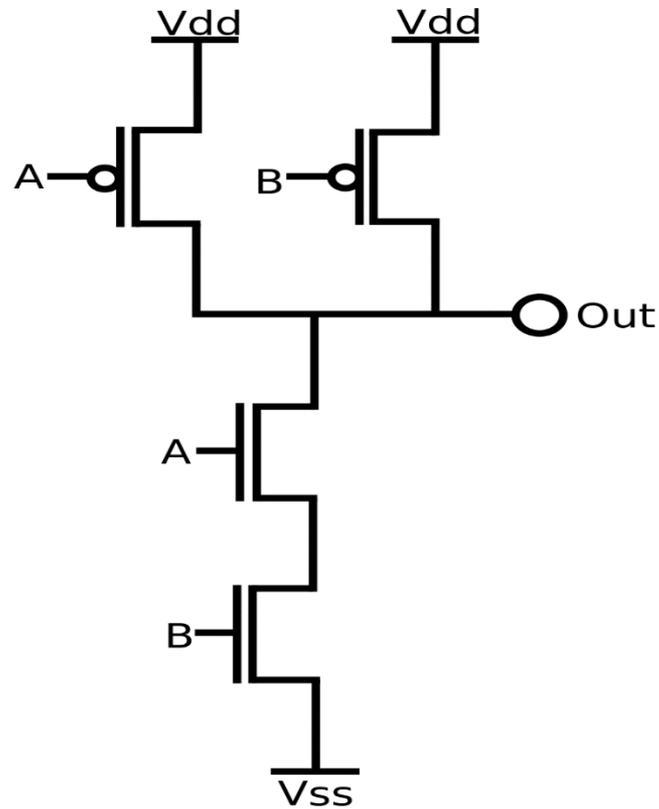
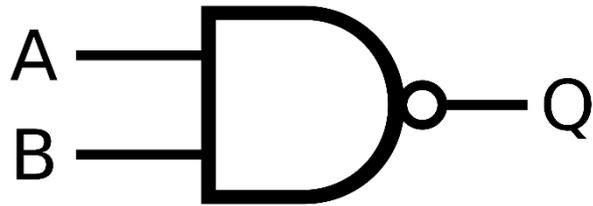
2.8x chip capability in same power





## ► Four transistors constitute a **NAND** gate

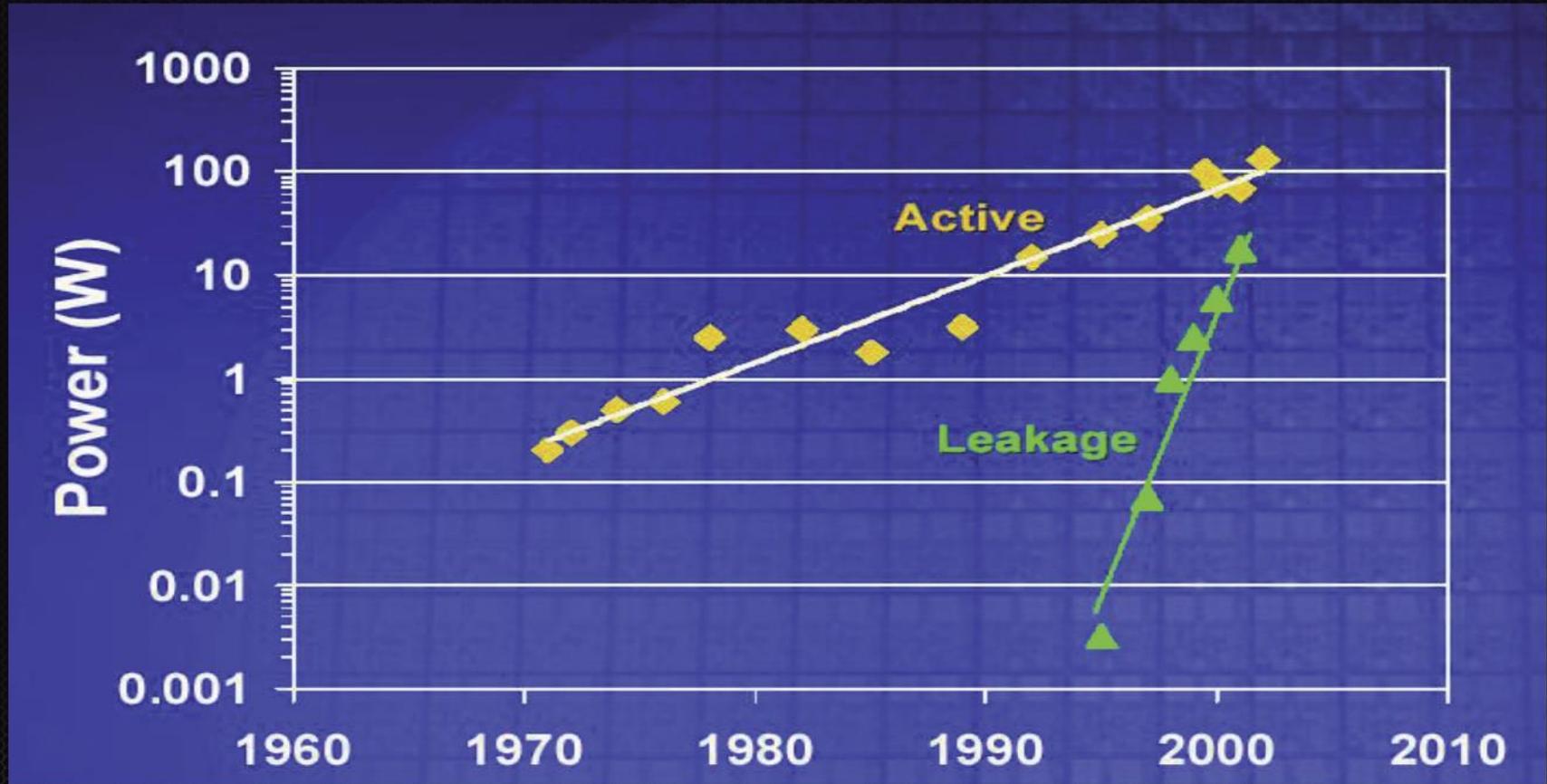
- Two inputs: A, B
- Output =  $\!(A\&\&B)$
- High voltage on both inputs closes the p-type transistors (current sources, top) and opens the n-type transistors (current sink, bottom)



## ► Power consumption

- Gates and wires connected to Out form a capacitor which must be charged and discharged in each cycle – proportional to frequency, reduced by Dennard scaling
- Short-circuit current (all transistors are open during state changes)
- Leakage through closed transistors – increased by geometry scaling down

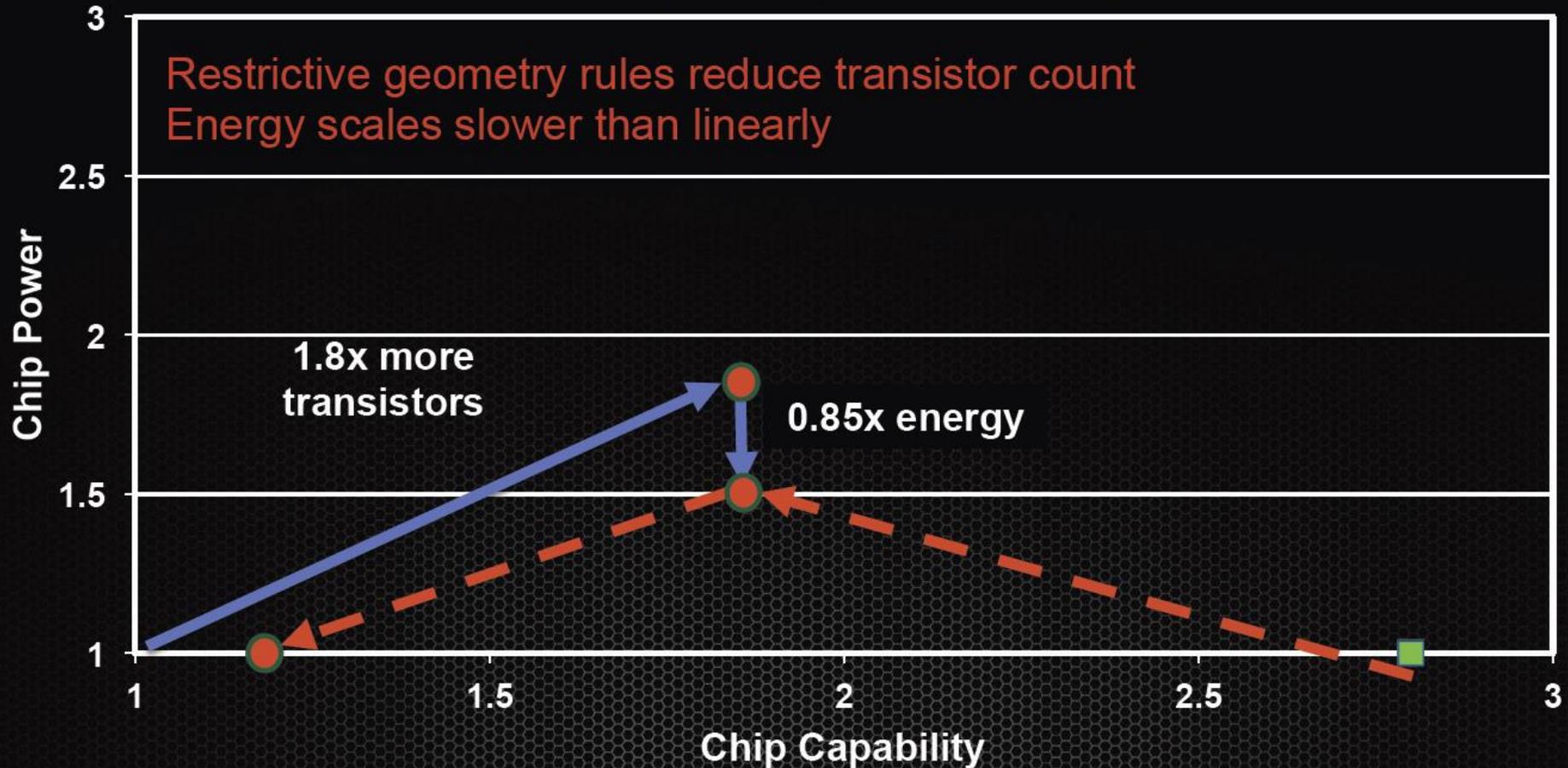
# But $L^3$ energy scaling ended in 2005



Moore, ISSCC Keynote, 2003

# Reality isn't even this good

1.8x chip capability at 1.5x power  
1.2x chip capability at same power



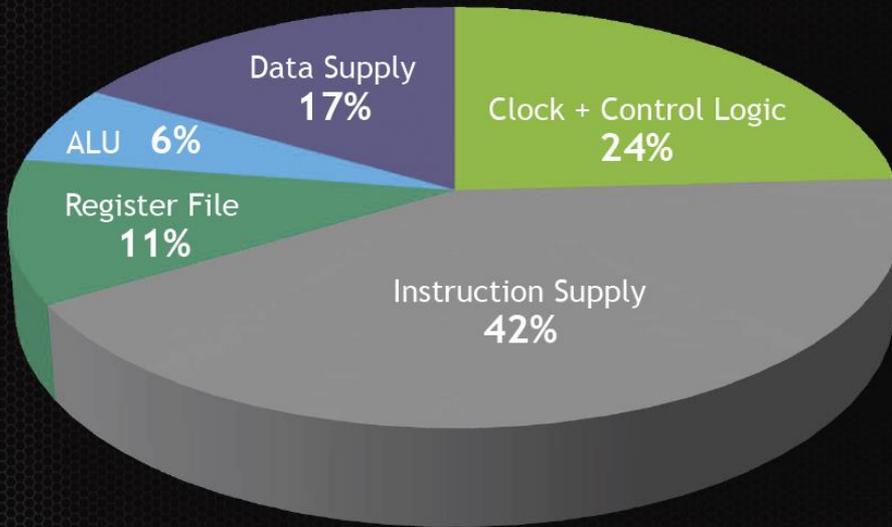
# The End of Dennard Scaling

- Processors aren't getting faster, just wider
  - Future gains in performance are from parallelism
- Future systems are energy limited
  - Efficiency *IS* Performance
- Process matters less
  - One generation is 1.2x, not 2.8x



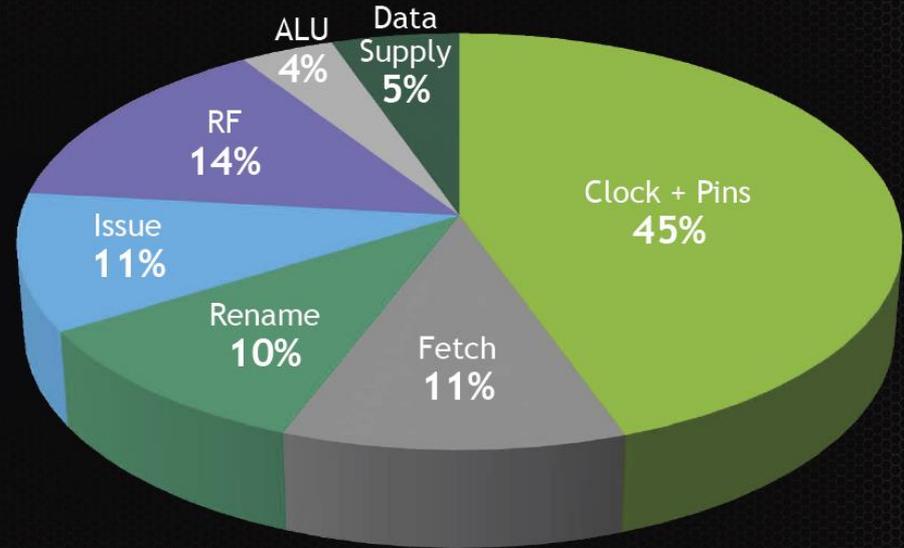
# How is Power Spent in a CPU?

## In-order Embedded



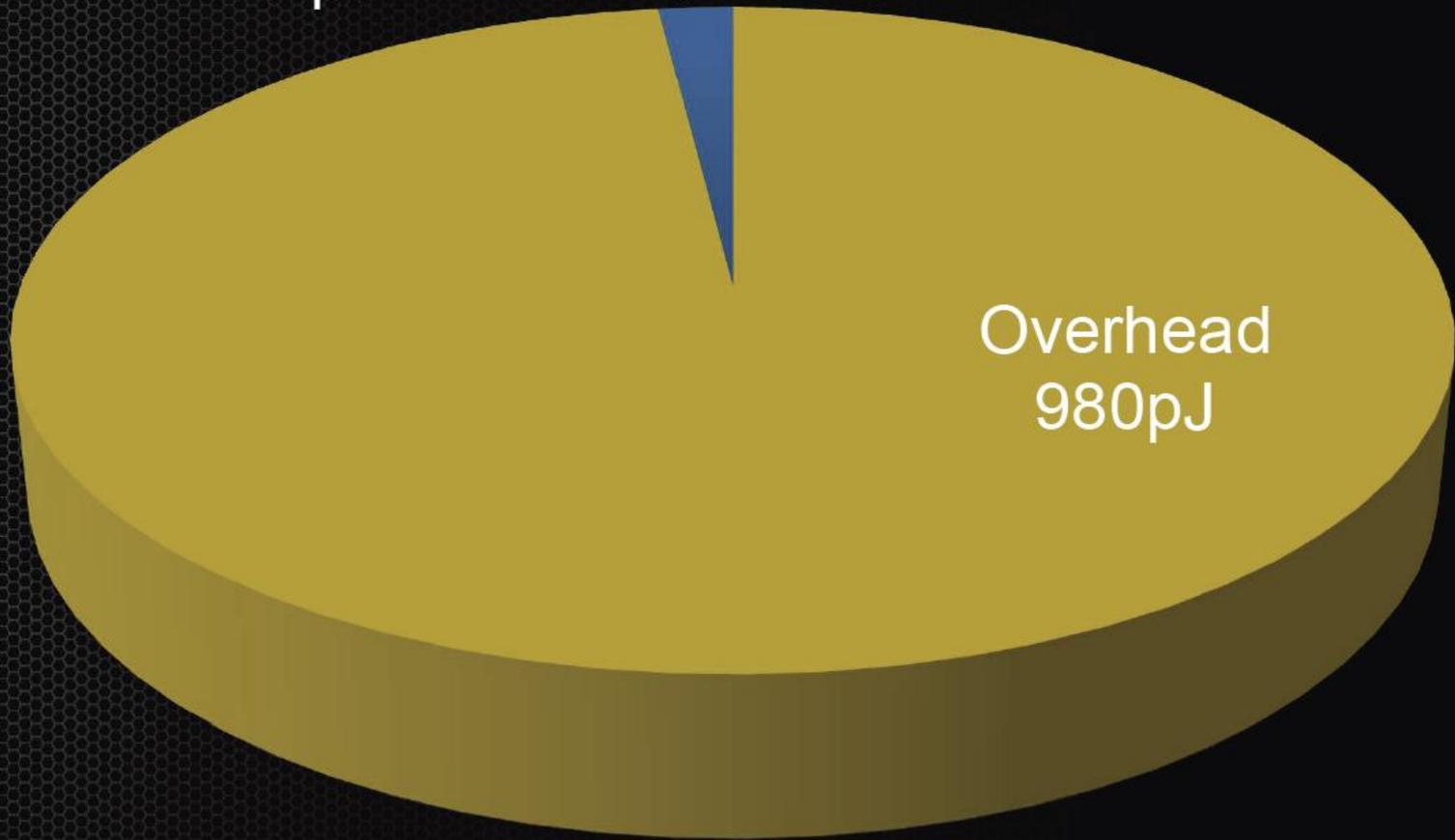
Dally [2008] (Embedded in-order CPU)

## OOO Hi-perf



Natarajan [2003] (Alpha 21264)

Payload  
Arithmetic  
20pJ

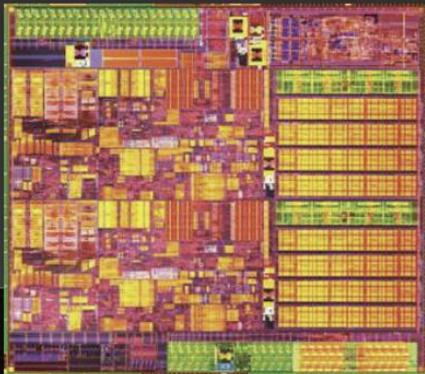


Overhead  
980pJ

# CPU

1690 pJ/flop

Optimized for Latency  
Caches

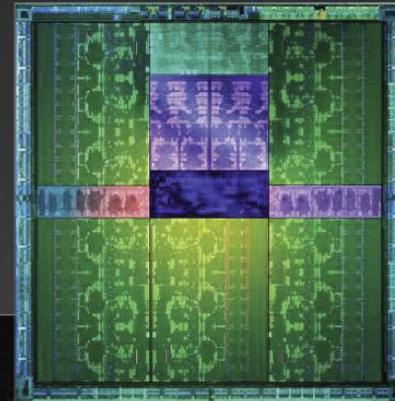


Westmere  
32 nm

# GPU

140 pJ/flop

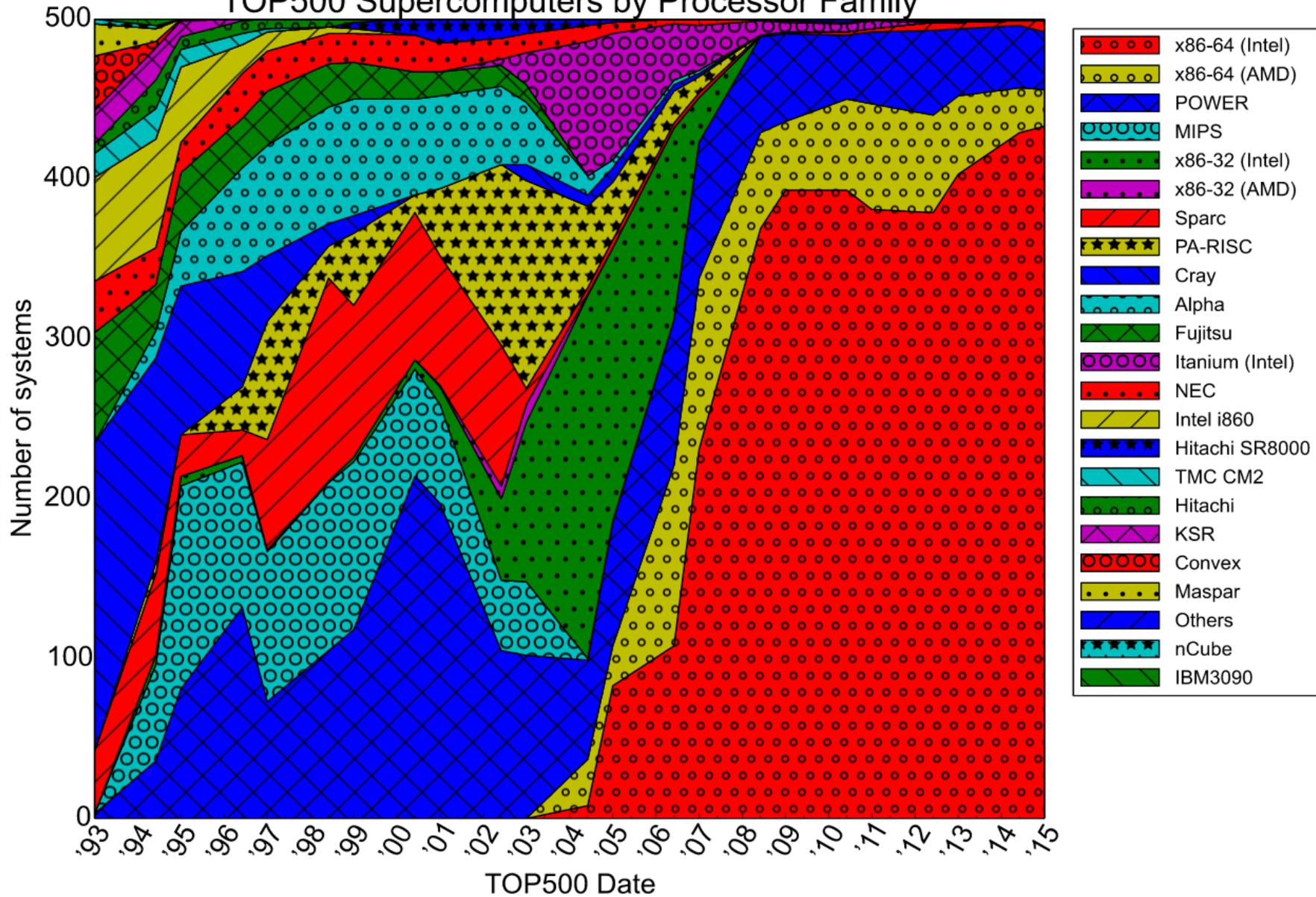
Optimized for Throughput  
Explicit Management  
of On-chip Memory



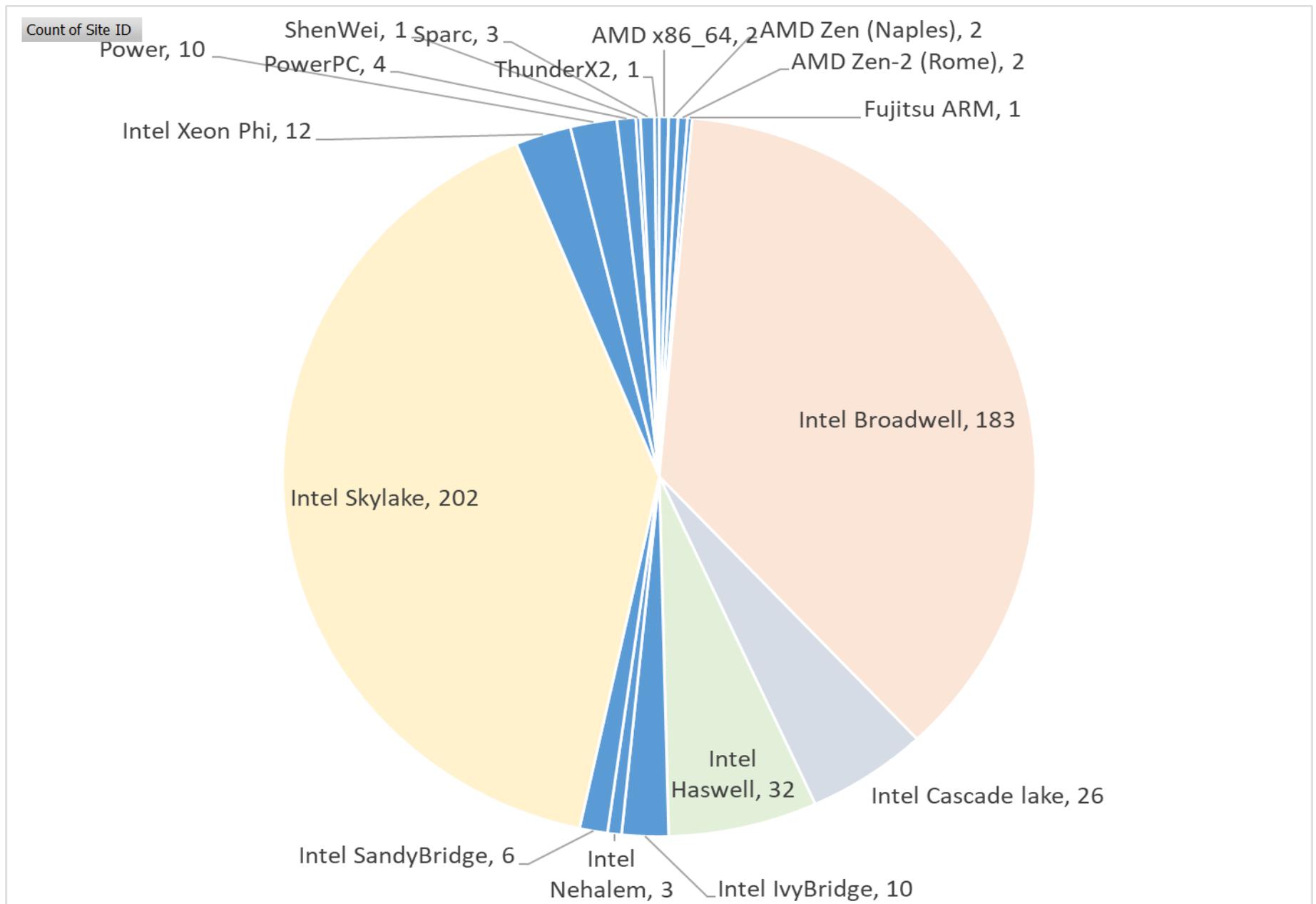
Kepler  
28 nm

# CPU architectures used in high-performance computing

# TOP500 Supercomputers by Processor Family



# Top500 Supercomputers (Nov 2019) – CPU types



# TOP500 list – Nov 2019

| Rank | TFlop/s | Name                                    | Country        | Total Cores | Cores per Socket | Processor Generation              | Accelerator/Co-Processor Cores | Accelerator/Co-Processor | Interconnect                      |
|------|---------|---|----------------|-------------|------------------|-----------------------------------|--------------------------------|--------------------------|-----------------------------------|
| 1    | 148600  | Summit                                  | United States  | 2414592     | 22               | IBM POWER9                        | 2211840                        | NVIDIA Volta GV100       | Dual-rail Mellanox EDR Infiniband |
| 2    | 94640   | Sierra                                  | United States  | 1572480     | 22               | IBM POWER9                        | 1382400                        | NVIDIA Volta GV100       | Dual-rail Mellanox EDR Infiniband |
| 3    | 93015   | Sunway TaihuLight                       | China          | 10649600    | 260              | Sunway                            |                                |                          | Sunway                            |
| 4    | 61445   | Tianhe-2A                               | China          | 4981760     | 12               | Intel Xeon E5 (IvyBridge)         | 4554752                        | Matrix-2000              | TH Express-2                      |
| 5    | 23516   | Frontera                                | United States  | 448448      | 28               | Xeon Platinum 82xx (Cascade Lake) |                                |                          | Mellanox InfiniBand HDR           |
| 6    | 21230   | Piz Daint                               | Switzerland    | 387872      | 12               | Intel Xeon E5 (Haswell)           | 319424                         | NVIDIA Tesla P100        | Aries interconnect                |
| 7    | 20159   | Trinity                                 | United States  | 979072      | 68               | Intel Xeon Phi                    |                                |                          | Aries interconnect                |
| 8    | 19880   | AI Bridging Cloud Infrastructure (ABCI) | Japan          | 391680      | 20               | Xeon Gold                         | 348160                         | NVIDIA Tesla V100 SXM2   | Infiniband EDR                    |
| 9    | 19477   | SuperMUC-NG                             | Germany        | 305856      | 24               | Xeon Platinum                     |                                |                          | Intel Omni-Path                   |
| 10   | 18200   | Lassen                                  | United States  | 288288      | 22               | IBM POWER9                        | 253440                         | NVIDIA Tesla V100        | Dual-rail Mellanox EDR Infiniband |
| 375  | 1458    | Salomon                                 | Czech Republic | 76896       | 12               | Intel Xeon E5 (Haswell)           | 52704                          | Intel Xeon Phi 7120P     | Infiniband FDR                    |