



Performance Optimization

Why should we care?

Servesh Muralidharan

IT-DI-WLCG, UP Team

CERN

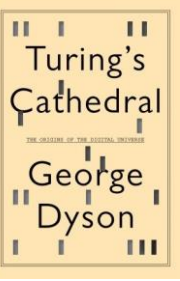
Nov 2018



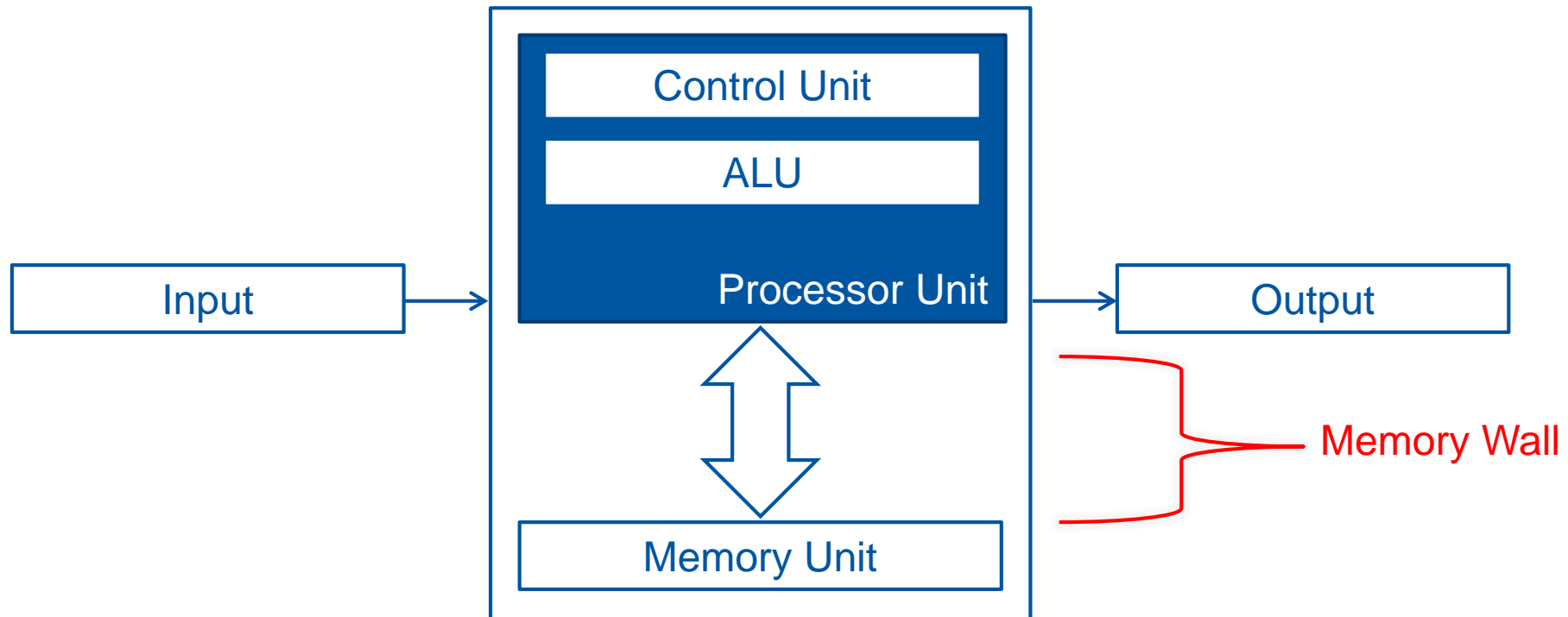
Optimization

The process of making something *as good or effective as possible*.

The Unchanged Computer Architecture

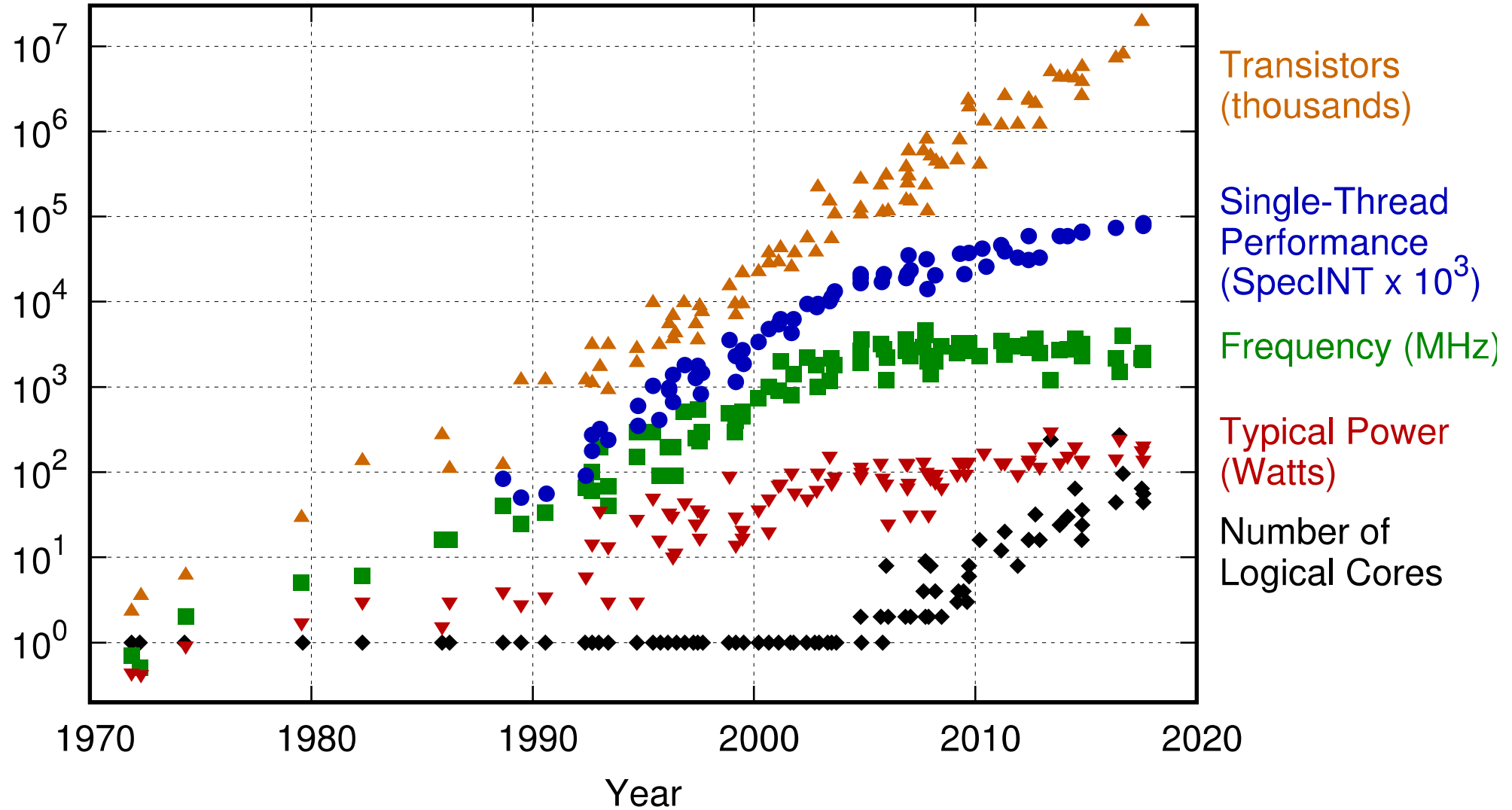


- ❑ Stored-Program computer
 - Von Neumann model
 - Program and data are stored in memory and then processed



Moore's Law

42 Years of Microprocessor Trend Data



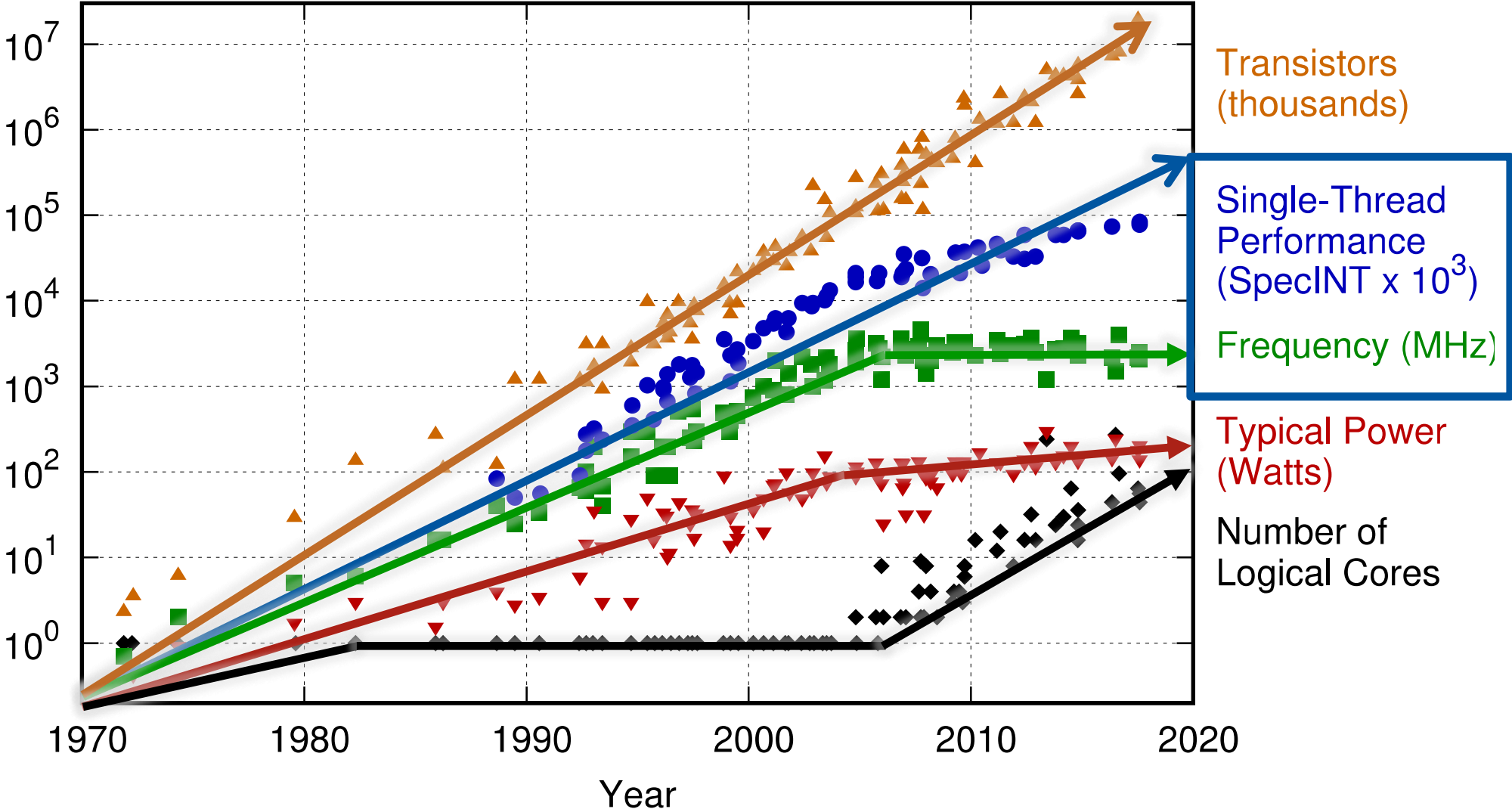
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

<https://github.com/karlrupp/microprocessor-trend-data>



Performance Trend

42 Years of Microprocessor Trend Data



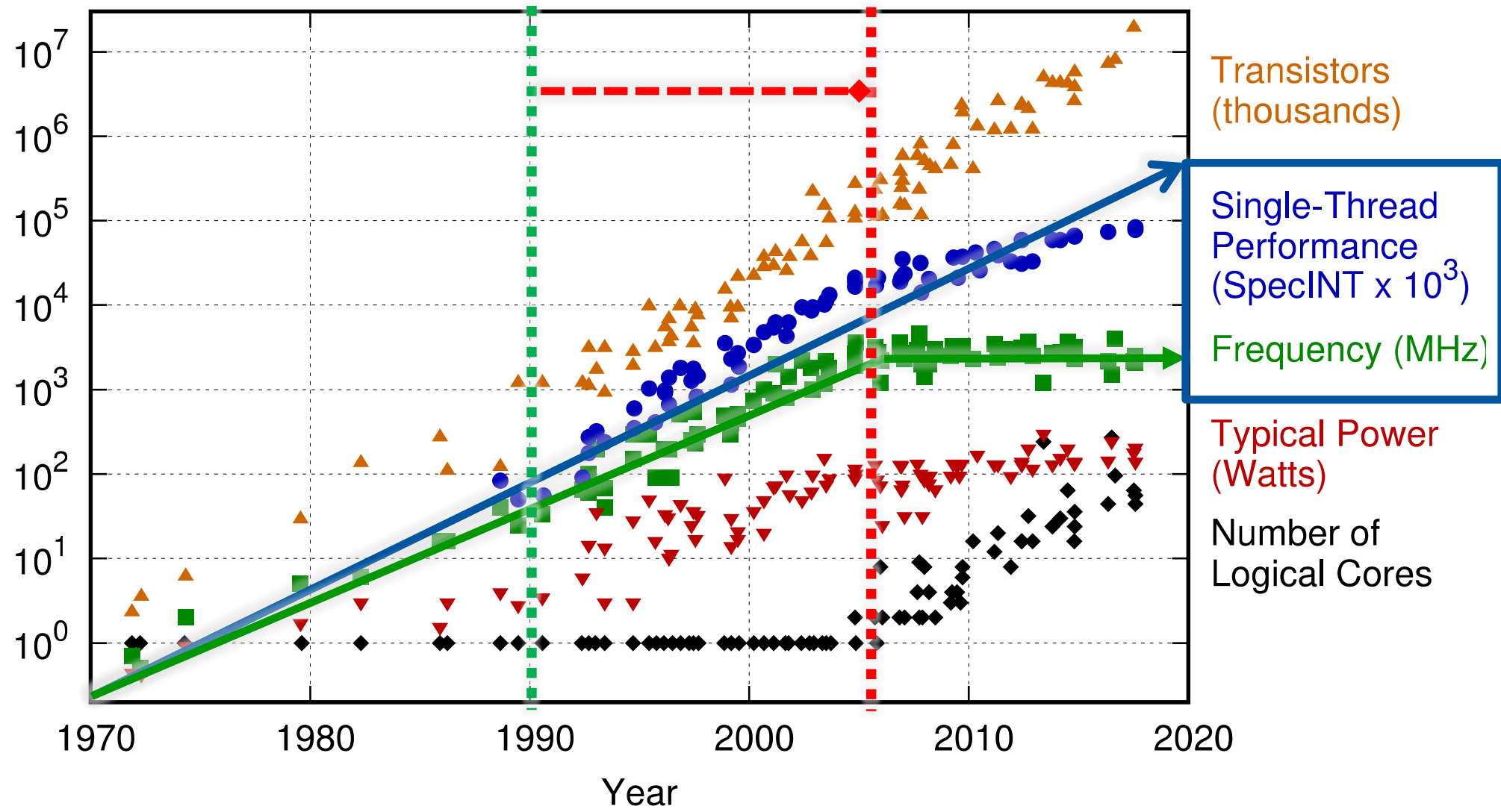
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

<https://github.com/karlrupp/microprocessor-trend-data>



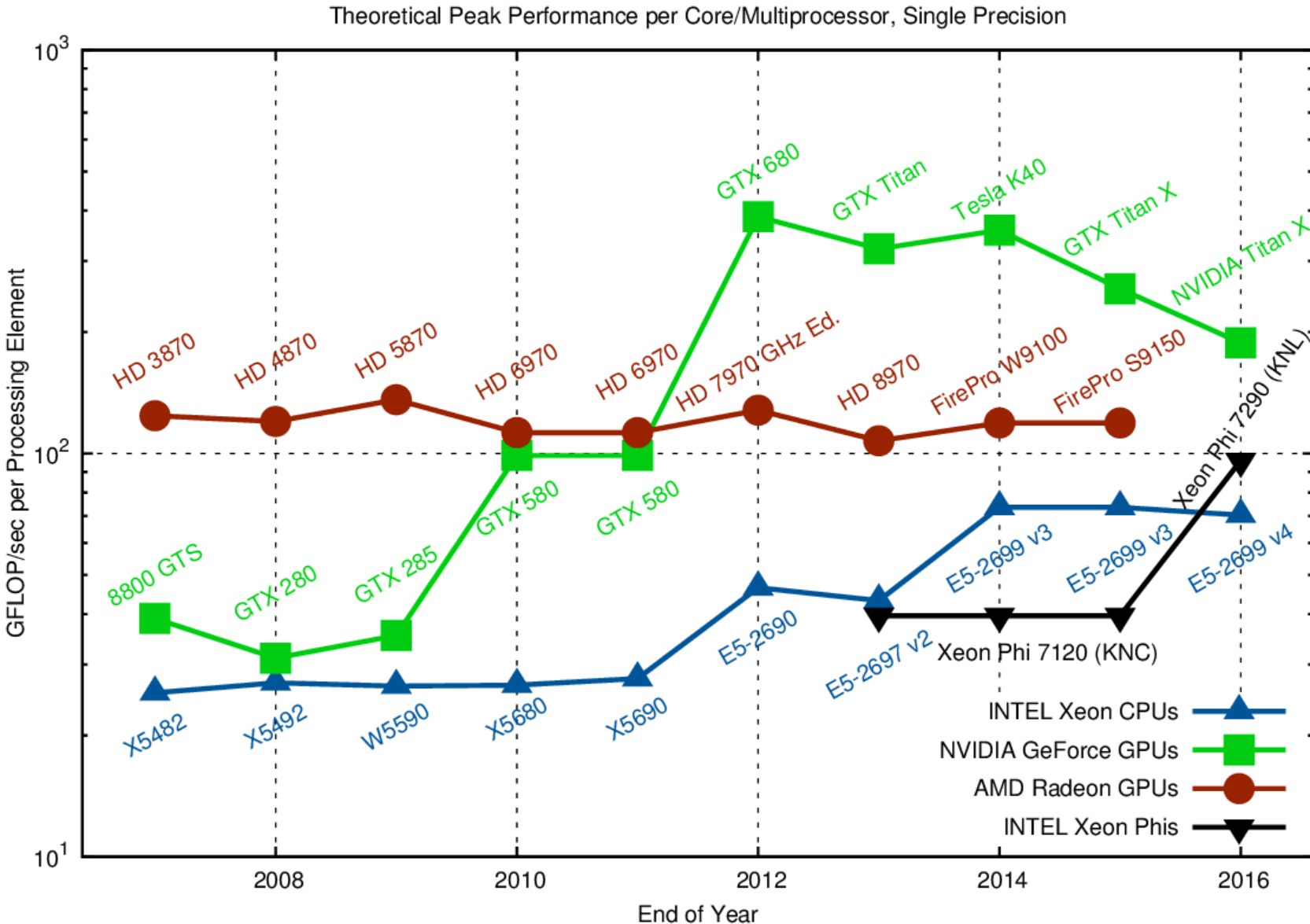
The Free Lunch – When Did It Stop?

42 Years of Microprocessor Trend Data



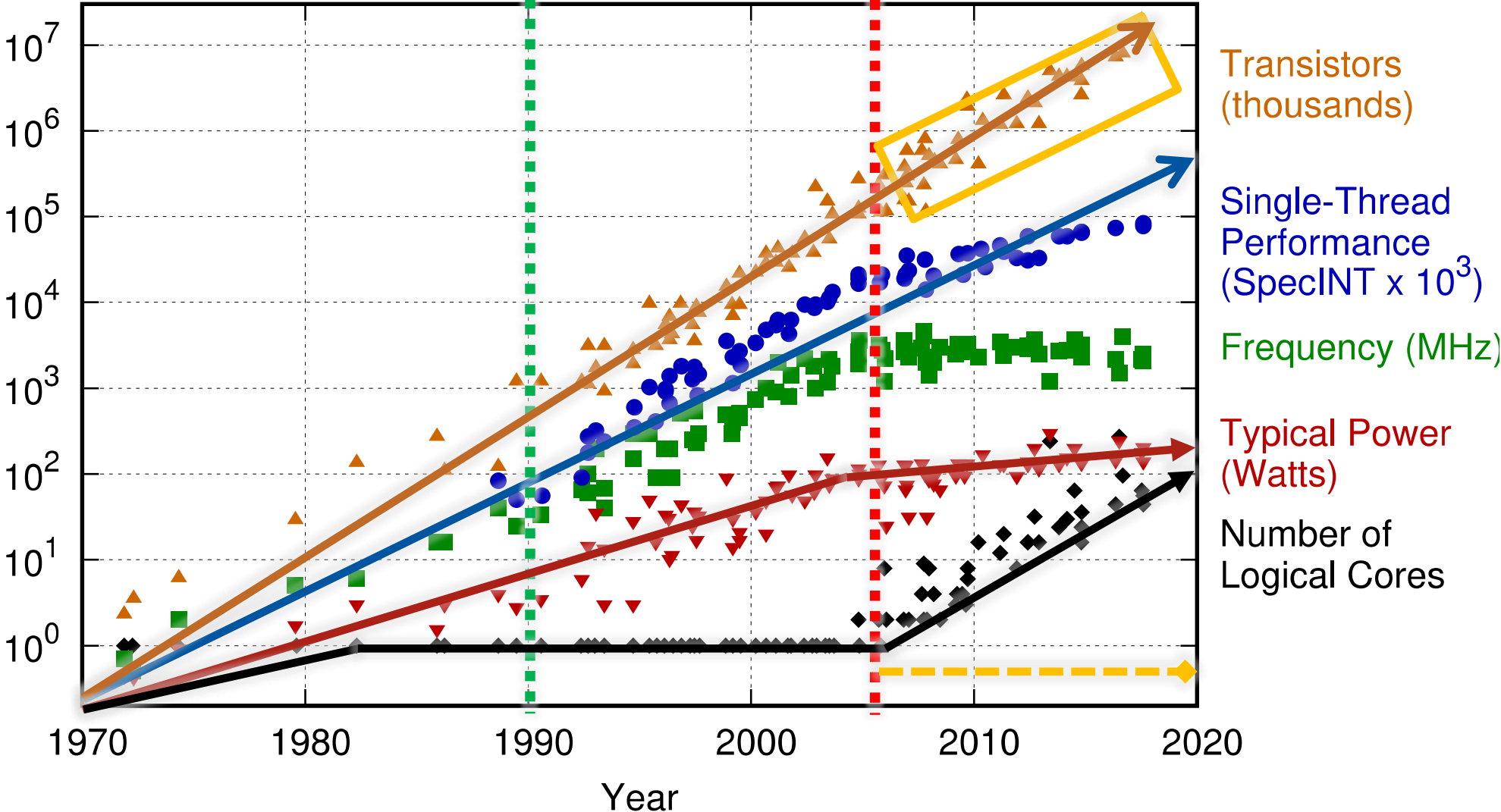
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp
<https://github.com/karlrupp/microprocessor-trend-data>

Single Core Performance



Post Free Lunch

42 Years of Microprocessor Trend Data

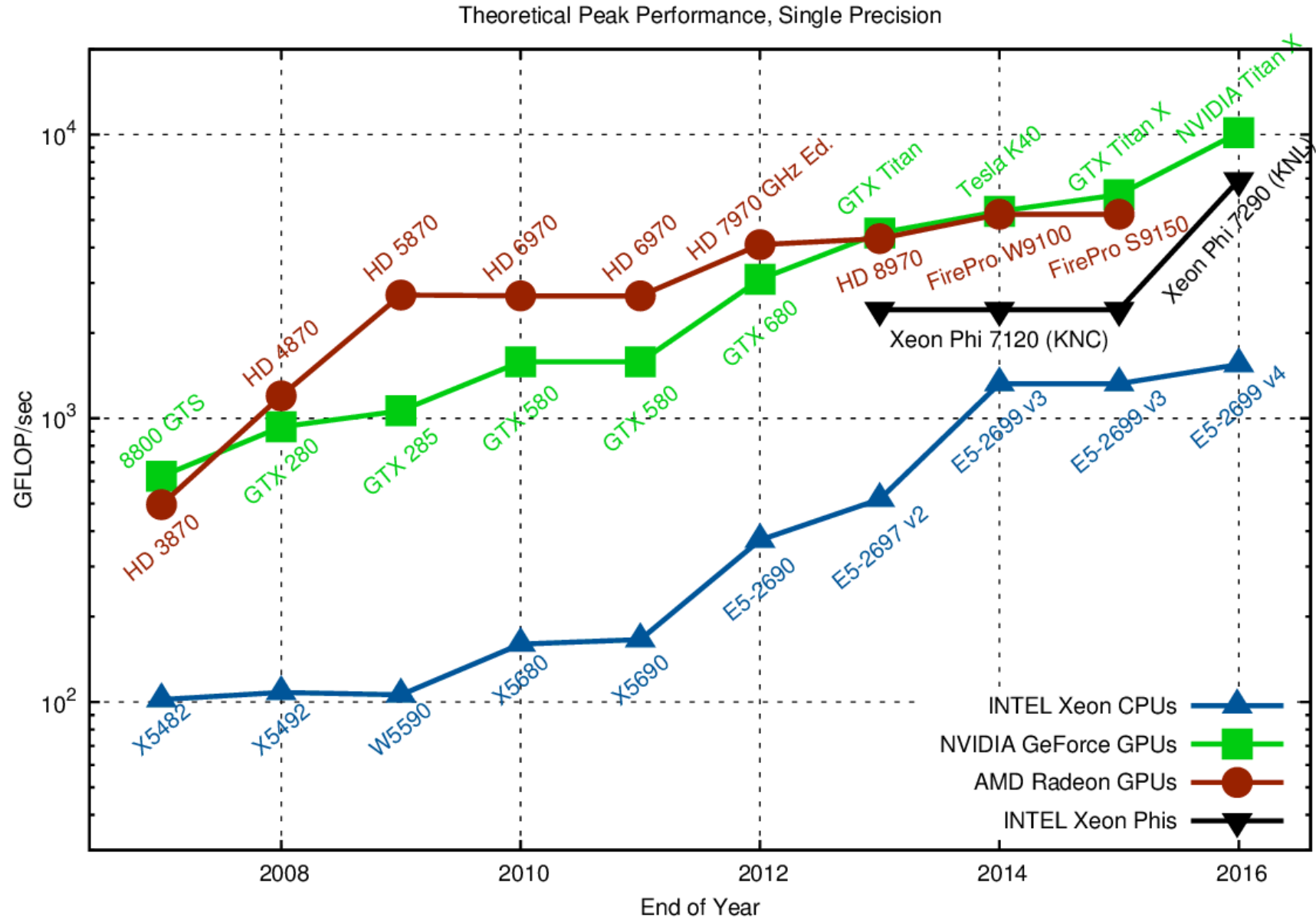


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

<https://github.com/karlrupp/microprocessor-trend-data>



Architectural Performance Trend



The Future : Heterogeneous Architecture

Single-Core Era

Enabled by:

- ✓ Moore's Law
- ✓ Voltage Scaling

Constrained by:

- ✗ Power
- ✗ Complexity

Assembly ➔ C/C++ ➔ Java ...

Single-thread Performance

Time

Multi-Core Era

Enabled by:

- ✓ Moore's Law
- ✓ SMP architecture

Constrained by:

- ✗ Power
- ✗ Parallel SW
- ✗ Scalability

pthread ➔ OpenMP / TBB ...

Throughput Performance

Time (# of processors)

Heterogeneous Systems Era

Enabled by:

- ✓ Abundant data parallelism
- ✓ Power efficient GPUs

Temporarily Constrained by:

- ✗ Programming models
- ✗ Comm.overhead

Shader ➔ CUDA ➔ OpenCL ➔ C++ and Java

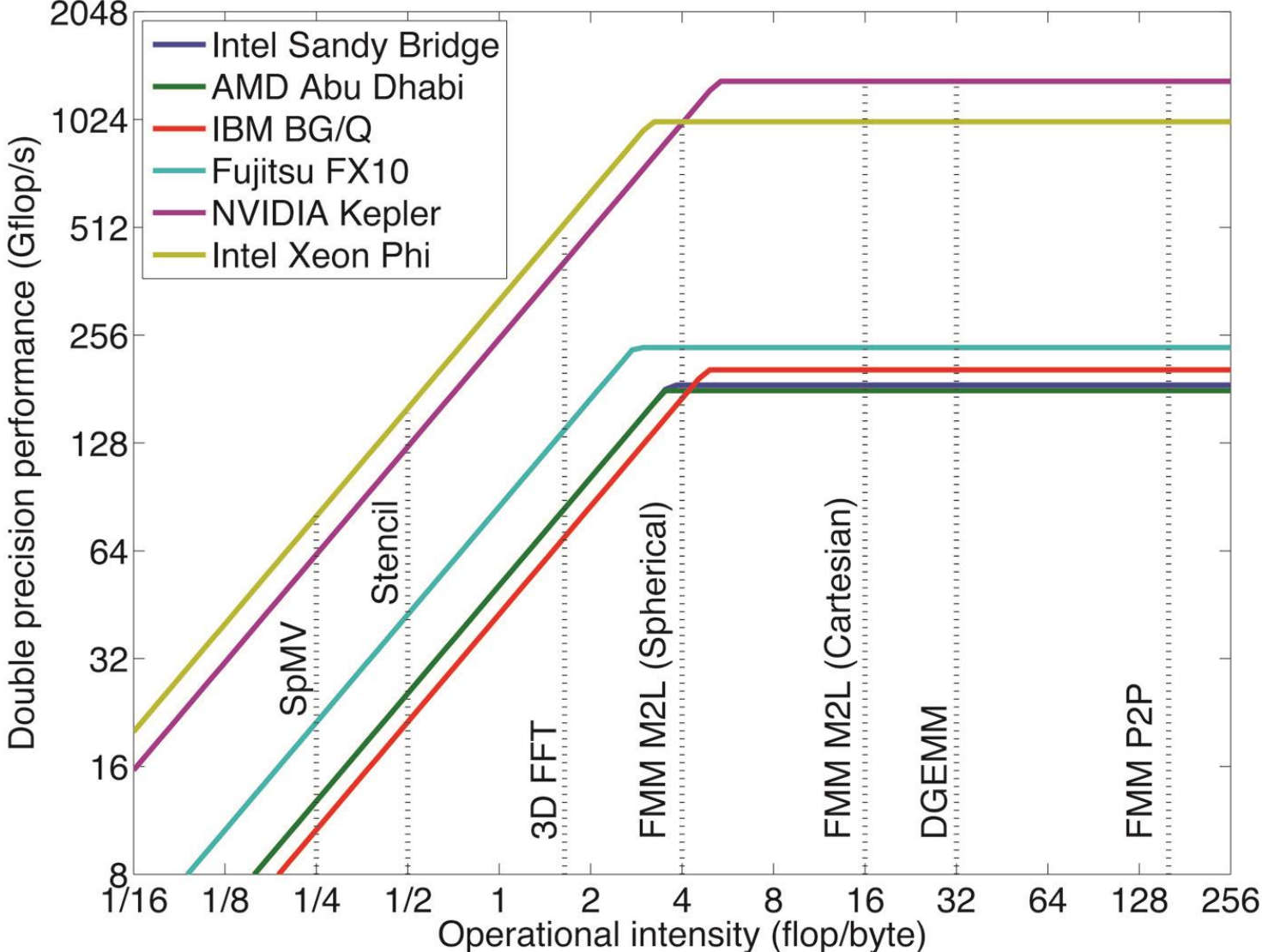
Modern Application Performance

Time (Data-parallel exploitation)

<https://opensourceforu.com/2016/12/how-heterogeneous-systems-evolved-and-the-challenges-going-forward/>



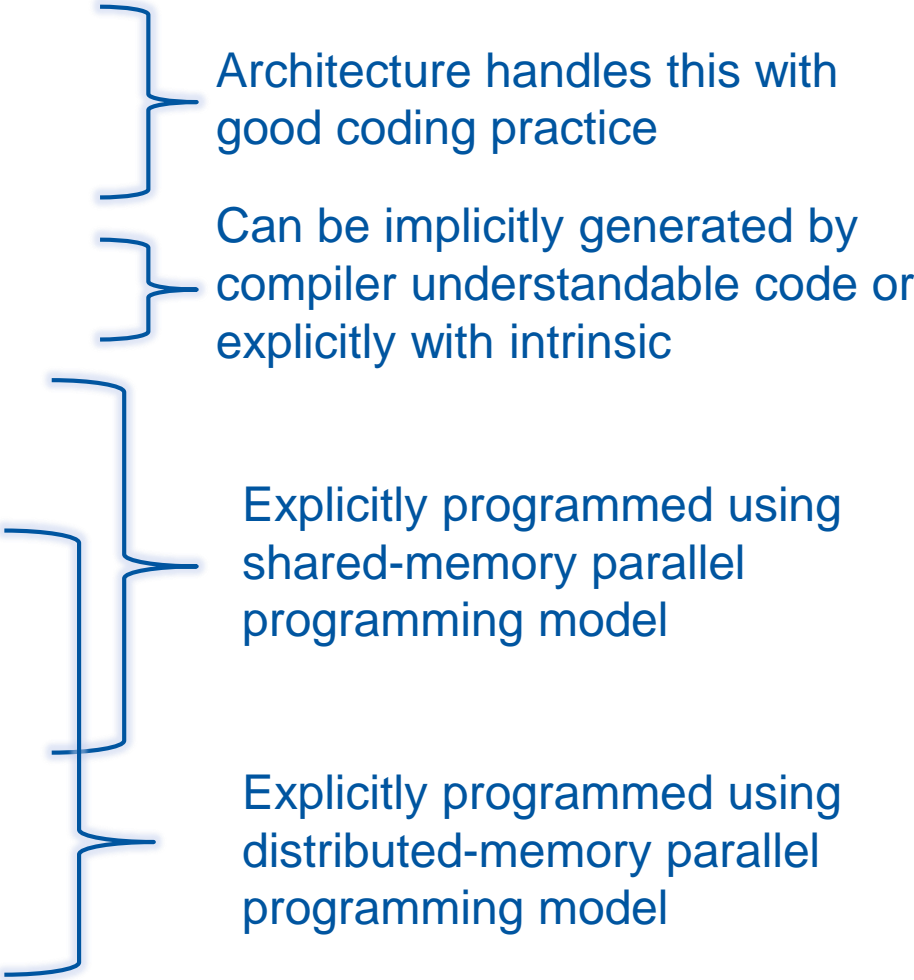
Roofline Performance Comparison



<https://sinews.siam.org/About-the-Author/how-will-the-fast-multipole-method-fare-in-the-exascale-era>

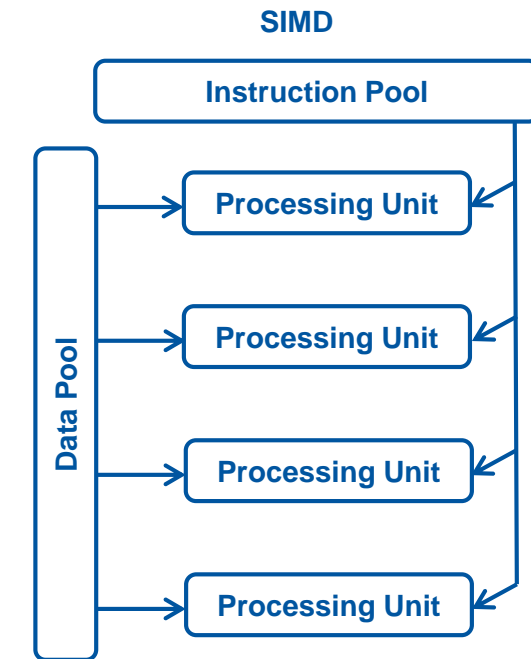
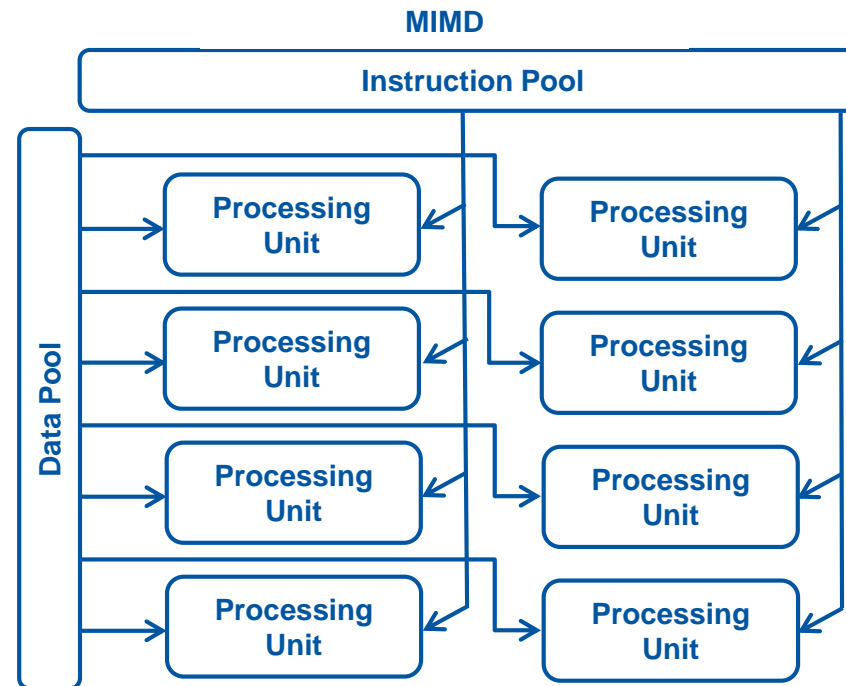
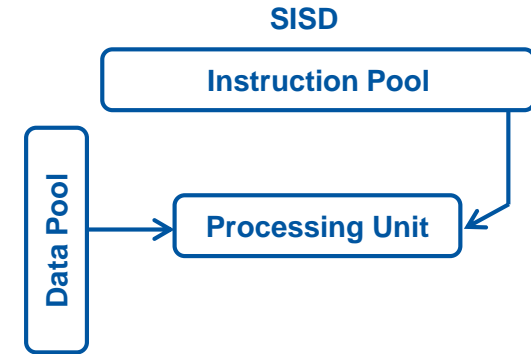
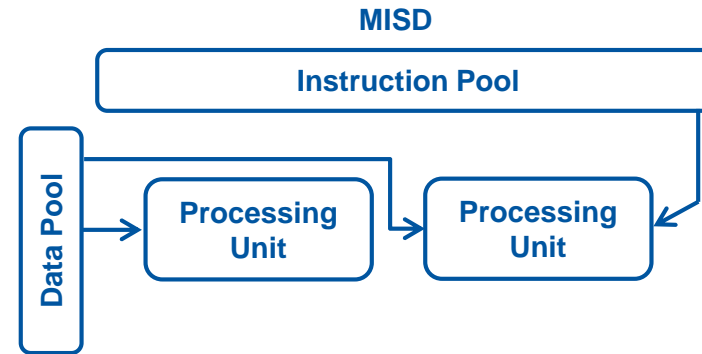
Sources of Parallelism in Modern Architectures

- 1. Instruction level parallelism (ILP)
- 2. Pipelining
- 3. Vector Operations
- 4. Hardware Threads
- 5. Multicore
- 6. Multi Socket
- 7. Cluster
- 8. Grid



Flynn's taxonomy: Can be a programmer's guideline

- ❑ Proposed in 1966
- ❑ Instruction streams
 - single (SI) or multiple (MI)
- ❑ Data streams
 - single (SD) or multiple (MD)
- ❑ Types
 - SISD
 - SIMD
 - MISD
 - MIMD



Concept of an Ideal Program

❑ Readability Vs Performance

❑ Compute Algorithm

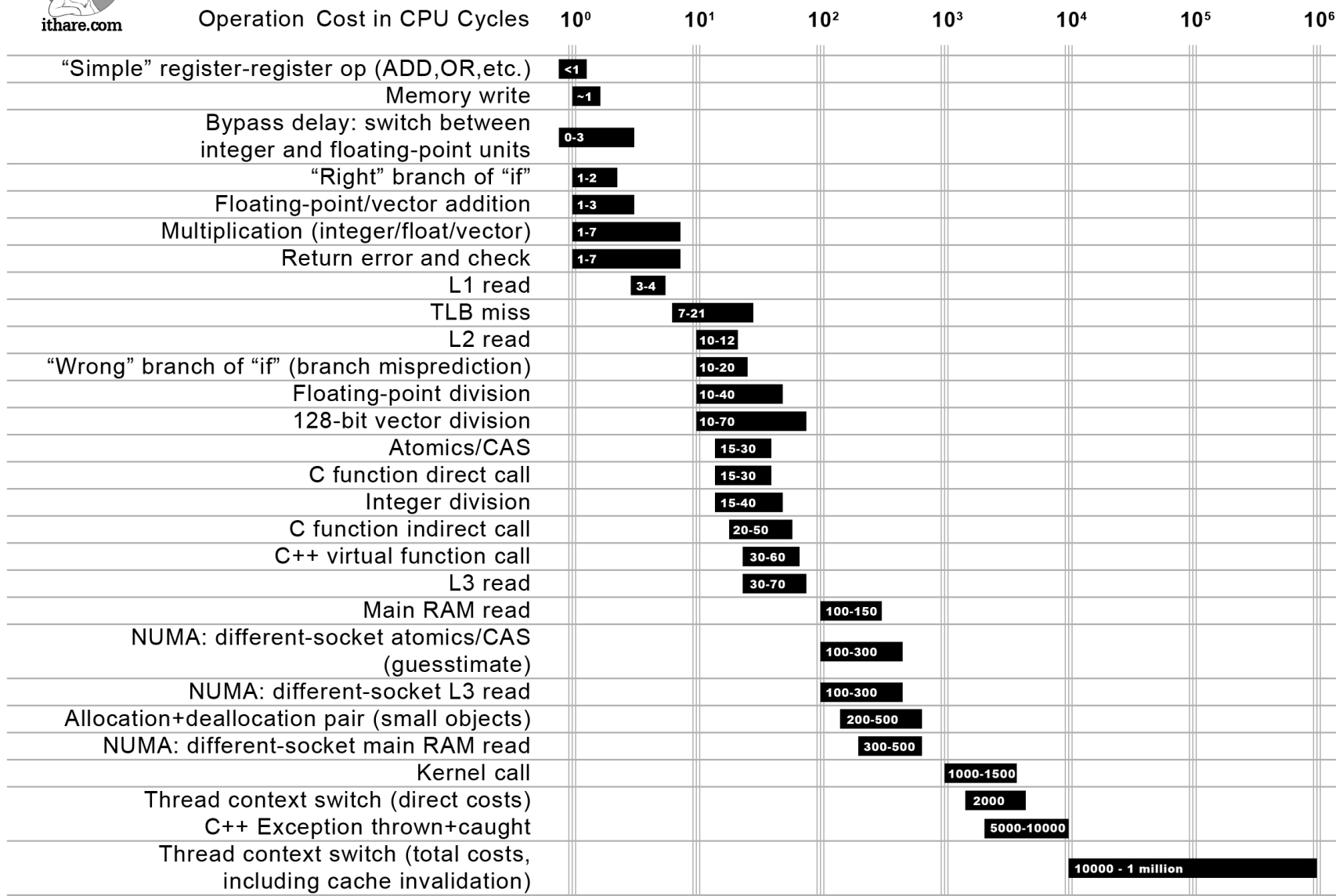
- **20% – 30% of code but has >90% of run time**
- Most optimizations are applied here
- *Big O notation*
 - Describes the worst case performance in terms of input size
 - Can represent time or space
 - Example: $O(n)$ – Linear (Finding an item in an unsorted array)
- **Extremely readable code for **compilers****
 - Elegance and Obscurity
 - Compilers will love it and we only care about performance!!!
 - **Be nice and use explicit comments for fellow human beings to understand**

❑ Glue code

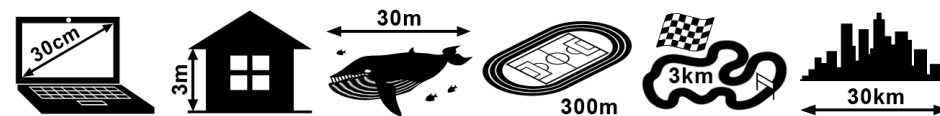
- **70% – 80% of code but has <10% of run time**
- Contains code used for structuring and connecting different blocks
- **Extremely readable code for **humans****
 - Compilers will hate you but that's okay, we don't care about performance here



Not all CPU operations are created equal



Distance which light travels while the operation is performed



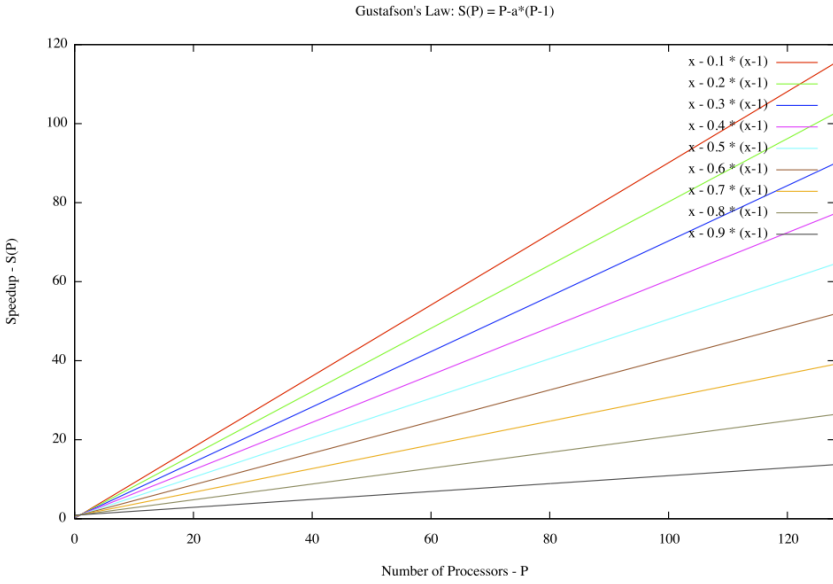
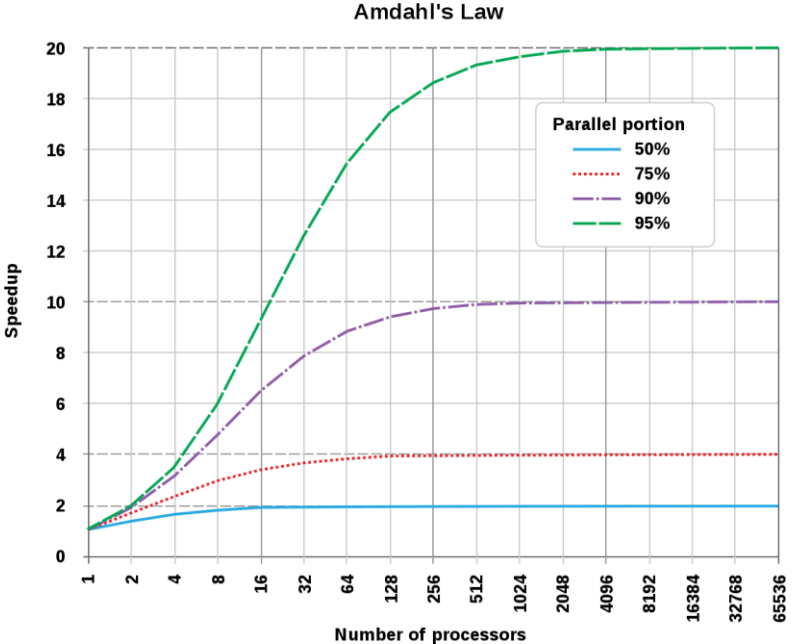
Parallel computing

- ❑ Performing certain computations simultaneously using multiple resources
- ❑ **Amdahl's law & Gustafson's law**
 - Speedup only comes from the parallelizable part of the code
 - i.e. Serial part of the code will impact or limit the achievable performance

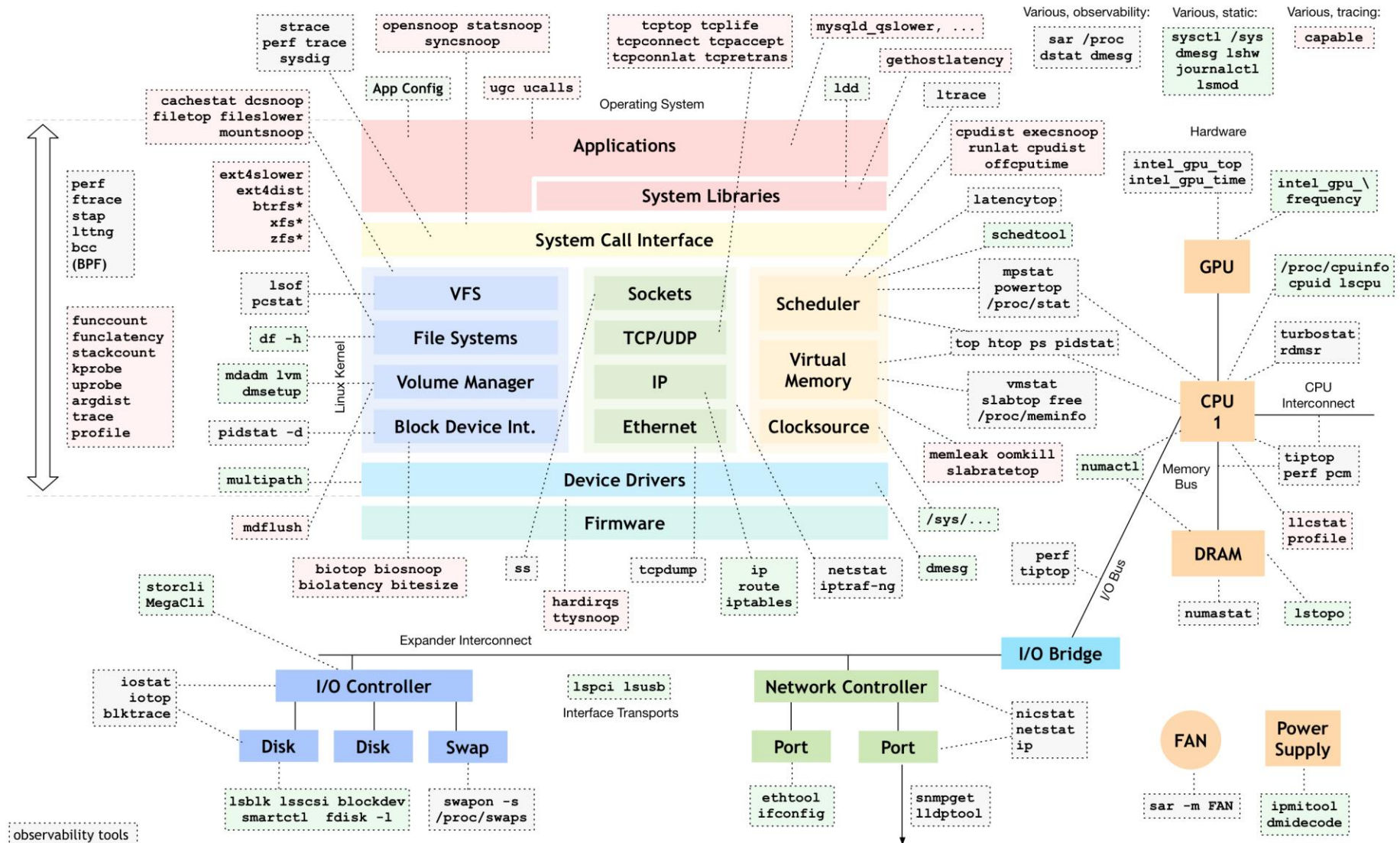
Amdahl's law
$$S_{\text{latency}}(s) = \frac{1}{(1 - p) + \frac{p}{s}}$$

Gustafson's law
$$S_{\text{latency}}(s) = 1 - p + sp,$$

Where,
 S_{latency} is the theoretical overall speedup
 s is the speedup in the parallel part
 p is the percentage of the execution time of serial part



Profiling



- observability tools
- static performance tools
- perf-tools/bcc tracing tools

these can observe the state of the system at rest, without load
<https://github.com/brendangregg/perf-tools> <https://github.com/iovisor/bcc>

style inspired by reddit.com/u/redct
<http://www.brendangregg.com/linuxperf.html> 2017



Conclusions

- ❑ Optimization is not an extra step
- ❑ Identification of a suitable form of parallelism for the algorithm
- ❑ Architecture agnostic implementation
- ❑ Reversible customizations
- ❑ Profile and Optimize

