

High-Performance Computing

What?
(Goal)

How to write ***fast*** code (not how to write code fast)

Why?

Fast code saves *time* and *energy*

How?

Parallelism

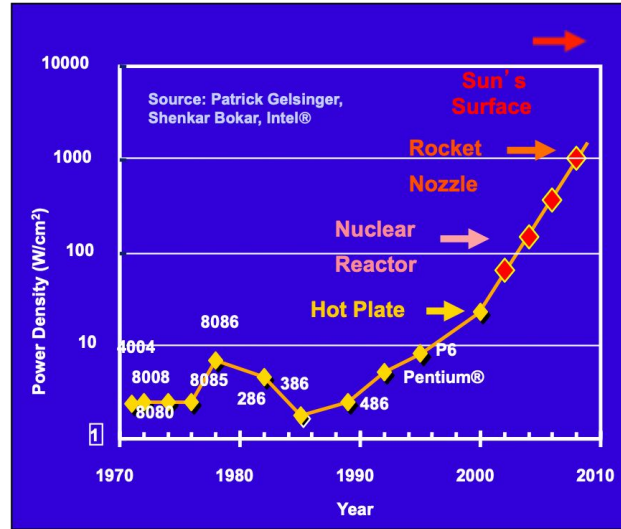
Data Locality

Specialization

Motivation

Currently, parallelism is driven by
power, memory, and physics.

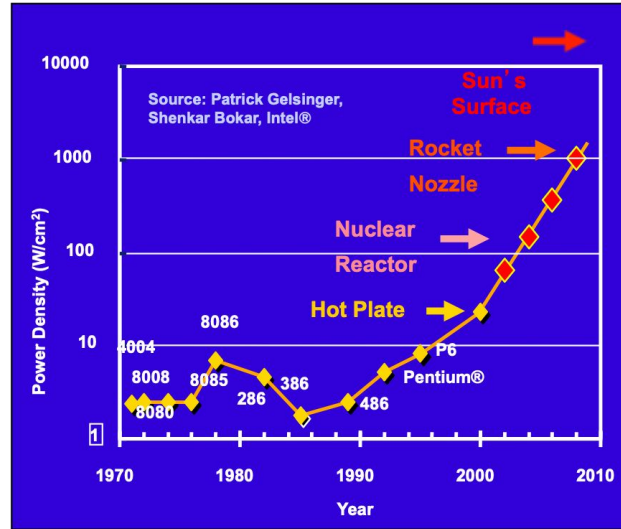
Performance & Power



Performance \propto (cores) \times (freq)

Power \propto (cores) \times (freq^{2.5})

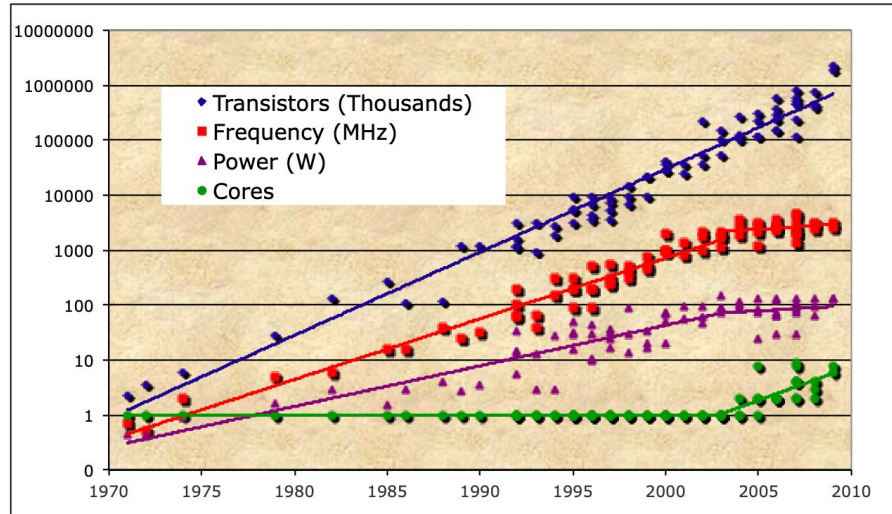
Performance & Power



$$\text{Performance} \propto (\text{cores}) \times (\text{freq})$$
$$\text{Power} \propto (\text{cores}) \times (\text{freq}^{2.5})$$

*Is it better to increase performance by doubling **frequency** or the number **cores**?*

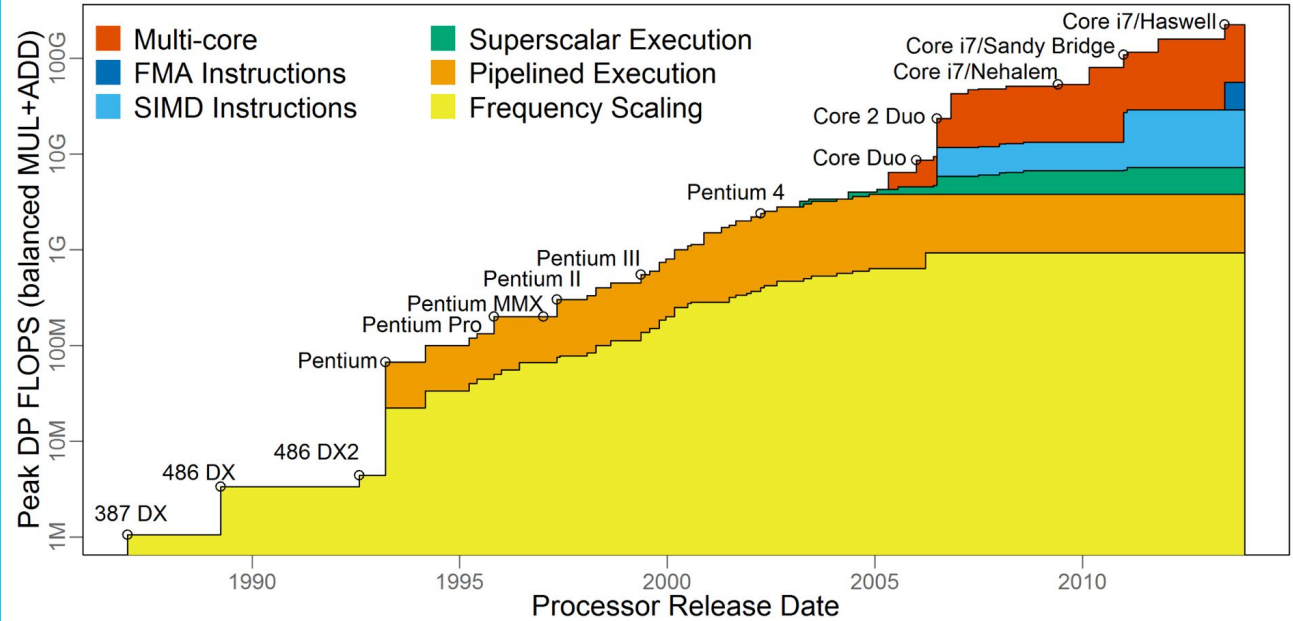
Parallelism & Power



Notice the transition around 2004.

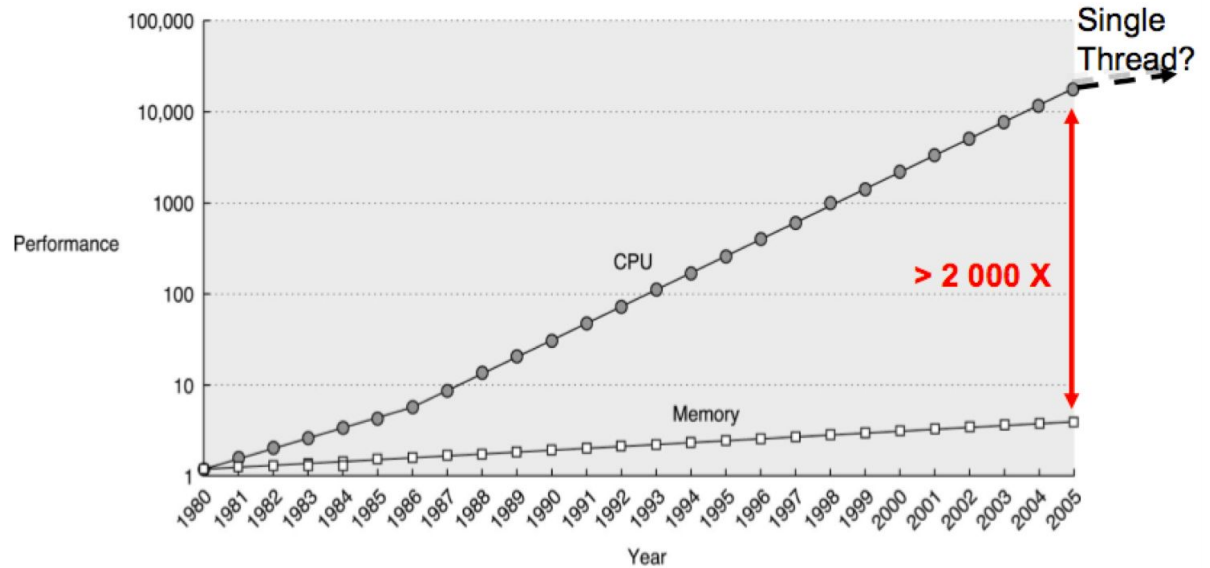
End of *Dennard Scaling*.

Performance Factors



What else has been done to increase performance?

Memory



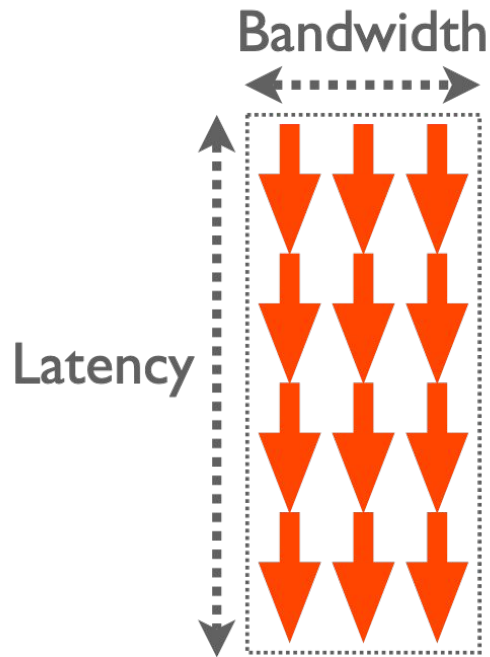
Memory performance has scaled much more slowly than instruction performance (i.e., FLOPS)

Historically, memory

latency halves every ~9 years

bandwidth doubles every ~3 years

Memory



Little's Law (queueing theory)

$$L = \lambda W$$

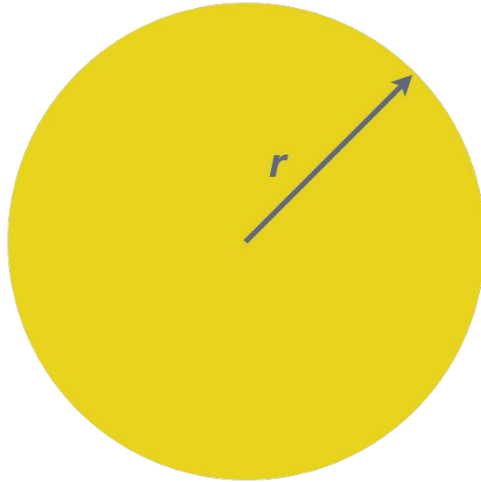
Equivalently,

Concurrency =

Latency \times Bandwidth

That is, larger the latency and bandwidth, more concurrency (i.e., threads) is required to fully utilize the memory bandwidth.

Physical Limits



Speed of light $c = 3 \times 10^8$ m/s

1 "op" needs to travel distance $2r$

$1 \text{ Top/s} \rightarrow 1 \text{ op} = 1 \times 10^{-12} \text{ s}$

$2r / c = 1 \times 10^{-12} \text{ s} \rightarrow$

$r = (1 \times 10^{-12} \times c) / 2 = 0.15 \text{ mm}$

Given this size, what about memory?

$1 \text{ TB} \rightarrow \sim 6 \text{ Å}^2 / \text{byte} \rightarrow O(10) \text{ atoms}$

Physical Limits

1. Quantum mechanics (e.g., tunneling)
2. EUV (extreme ultraviolet) lithography
3. 3D transistors/chips

Conclusion

As a result of these trends, parallelism has become ubiquitous.
No matter what system scale you care about, parallelism affects you.



Qualcomm Snapdragon 855 Processor

1. **4x** Kryo 485 high-efficiency cores @ **1.8 GHz**
2. **3x** Kryo 485 high-performance cores @ **2.42 GHz**
3. **1x** Kryo 485 high-performance cores @ **2.84 GHz**
4. **1x** Adreno 640 GPU @ **600 MHz**
 - a. **384x** ALU
5. **~10 Watts**



Intel Core i7-9700k

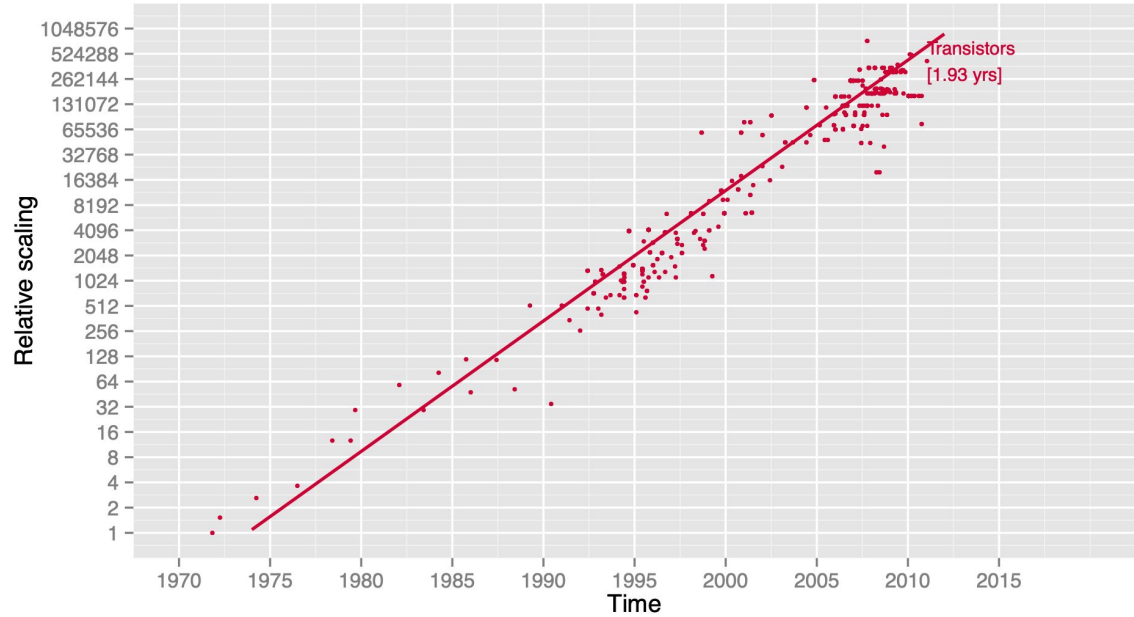
1. **8x Cores @ 4.90 GHz**
2. **1x Intel UHD Graphics 630 @ 1.20 GHz**
 - a. **24x execution units**
3. **95 Watts**



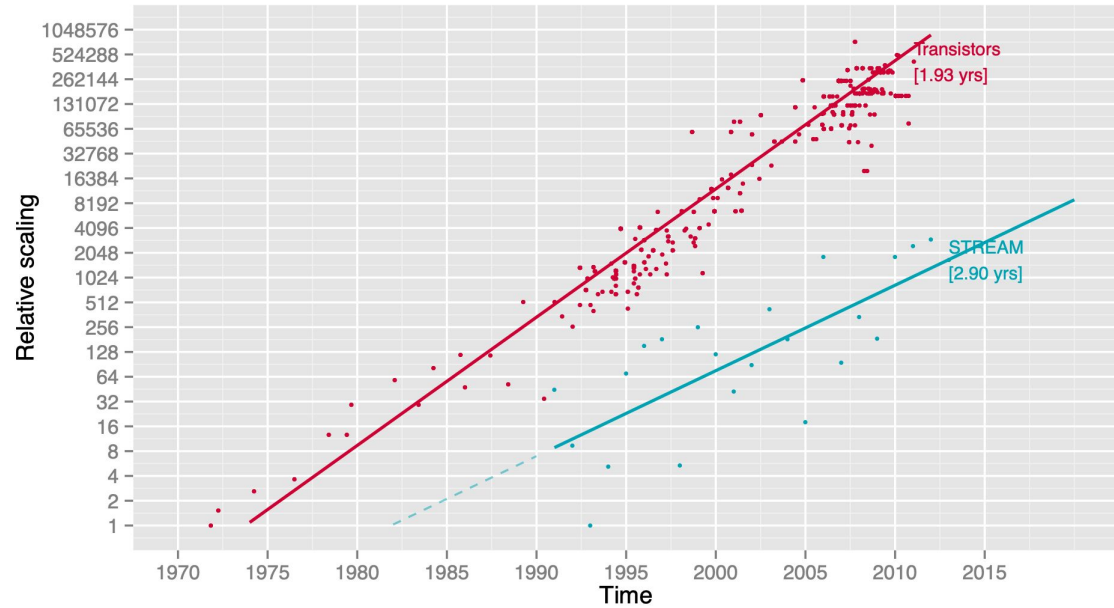
Summit Supercomputer (Department of Energy)

1. **4608x** POWER9 nodes
 - a. **2x** POWER9 CPUs @ **4 GHz**
 - i. **20x** Cores per CPU
 - b. **6x** Nvidia Tesla V100 GPUs @ **1.53 GHz**
 - i. **5120x** CUDA cores per GPU
2. **13 MWatts**

Data Movement



Data Movement



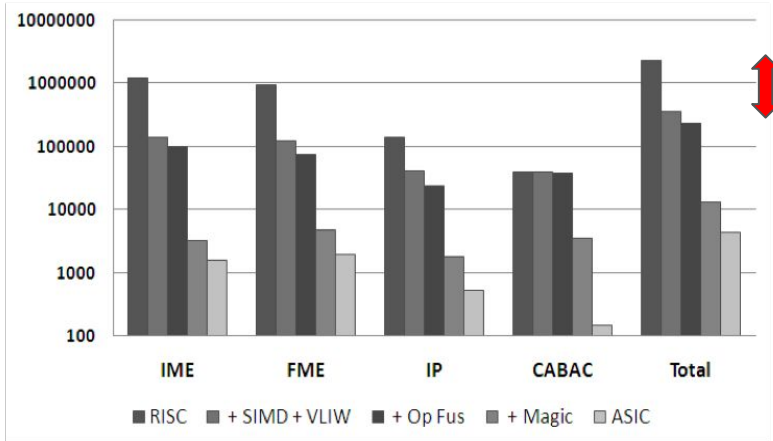


Figure 2. Each set of bar graphs represents energy consumption (μJ) at each stage of optimization for IME, FME, IP and CABAC respectively. Each optimization builds on the ones in the previous stage with the first bar in each set representing RISC energy dissipation followed by generic optimizations such as SIMD and VLIW, operation fusion and ending with “magic” instructions

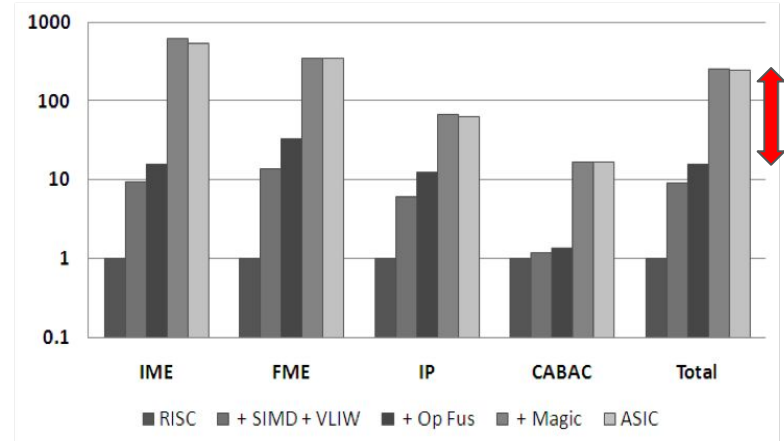


Figure 3. Each set of bar graphs represents speedup at each stage of optimization. Each optimization builds on those of the previous stage with the first bar in each set representing RISC speedup, followed by generic optimizations such as SIMD and VLIW, then operation fusion and finally “magic” instructions

~ 10x in energy (left) and time (right) from generic optimizations (hardware & software)

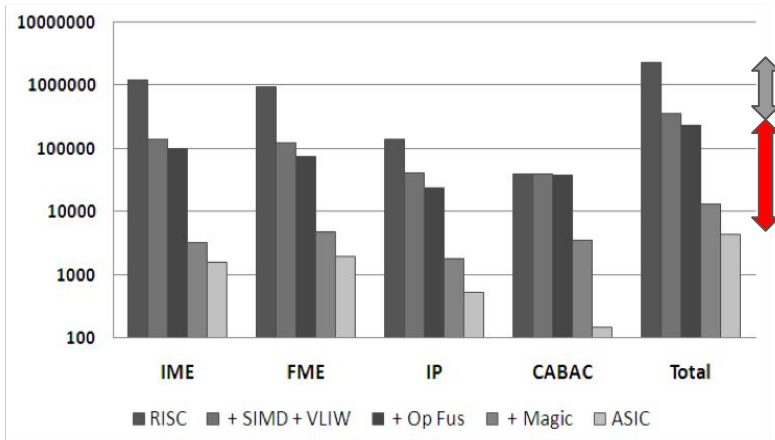


Figure 2. Each set of bar graphs represents energy consumption (μJ) at each stage of optimization for IME, FME, IP and CABAC respectively. Each optimization builds on the ones in the previous stage with the first bar in each set representing RISC energy dissipation followed by generic optimizations such as SIMD and VLIW, operation fusion and ending with “magic” instructions

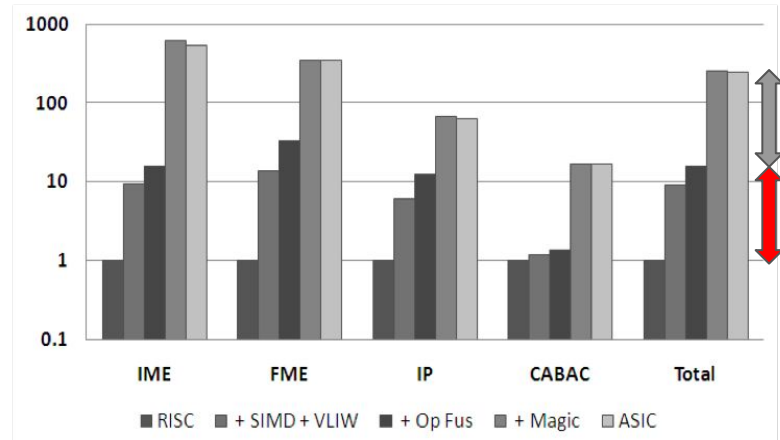


Figure 3. Each set of bar graphs represents speedup at each stage of optimization. Each optimization builds on those of the previous stage with the first bar in each set representing RISC speedup, followed by generic optimizations such as SIMD and VLIW, then operation fusion and finally “magic” instructions

~ 10x in energy (left) and time (right) from generic optimizations (hardware & software)

~ 10x+ in energy (left) and time (right) from application-specific customization (of the hardware)

Specialization

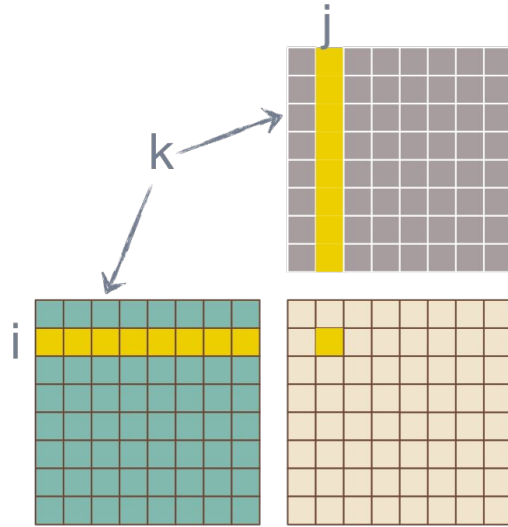
1. Application-specific integrated circuit (ASIC) is **500x** more energy efficient than general-purpose chip multi-processors (CMP)
2. Over **90%** of energy is “**overhead**” (e.g., instruction fetch/decode, etc.) because the cost of actual FLOP is very cheap.

Conclusion

These observations imply we need to pay attention to data movement and exploit custom units (e.g., GPUs) to improve energy efficiency and performance (in addition to parallelization)

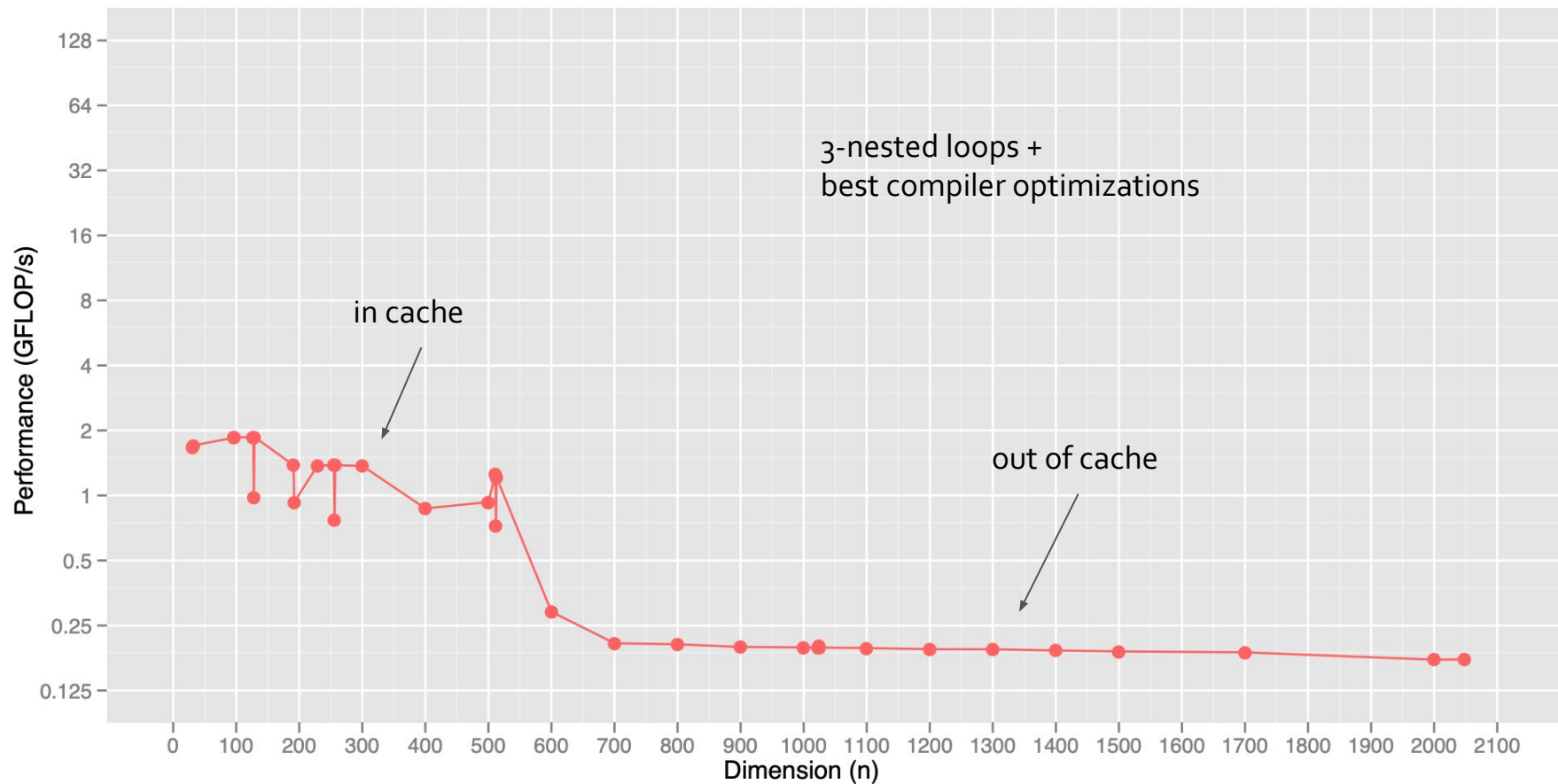
Theory vs. Practice

Example: **Matrix multiply** (non-Strassen)



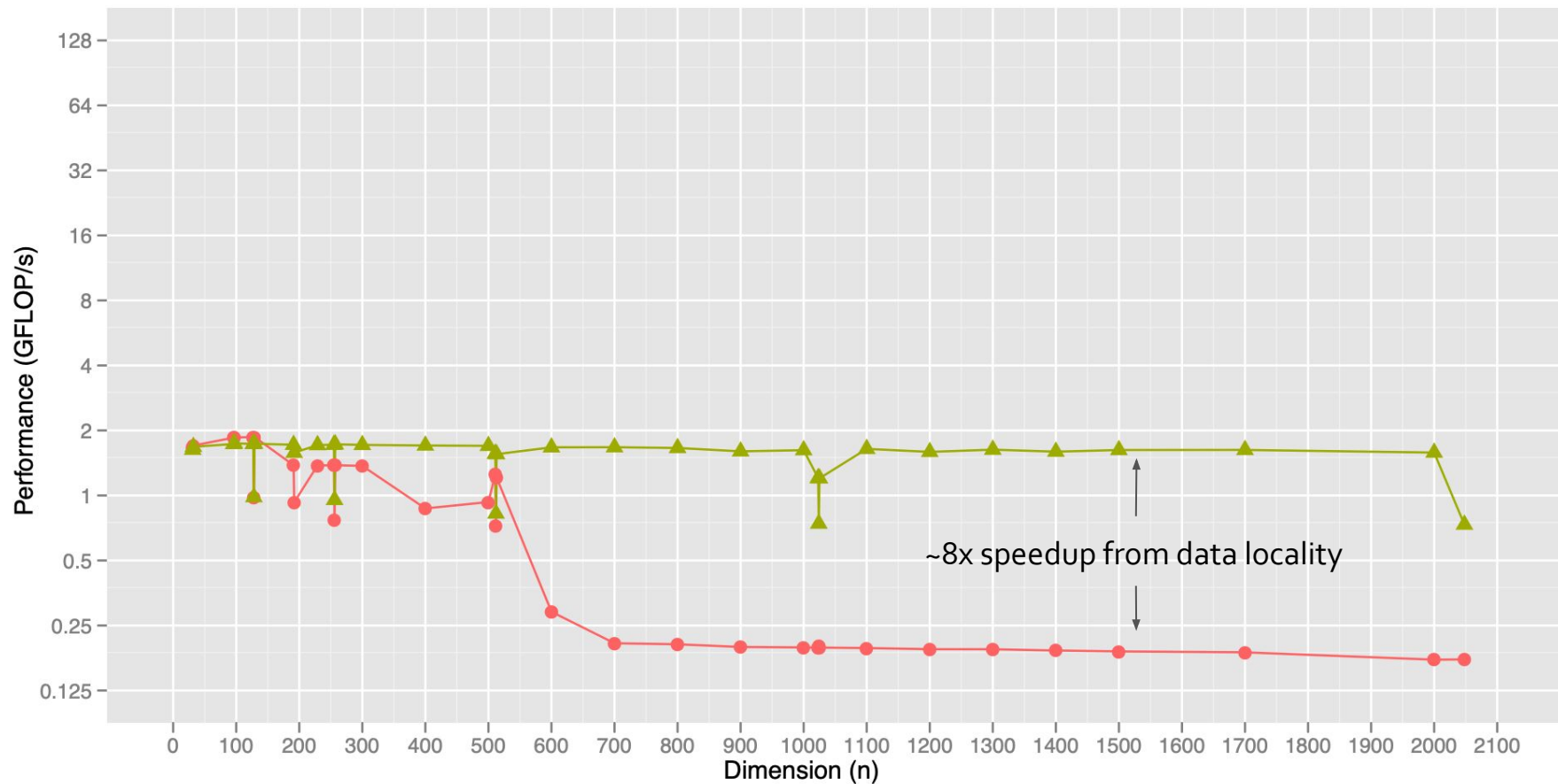
$$C \leftarrow C + A * B$$

```
for i = 1 to n do  
  for j = 1 to n do  
    for k = 1 to n do  
       $C[i,j] \leftarrow C[i,j] + A[i,k] \cdot B[k,j]$ 
```

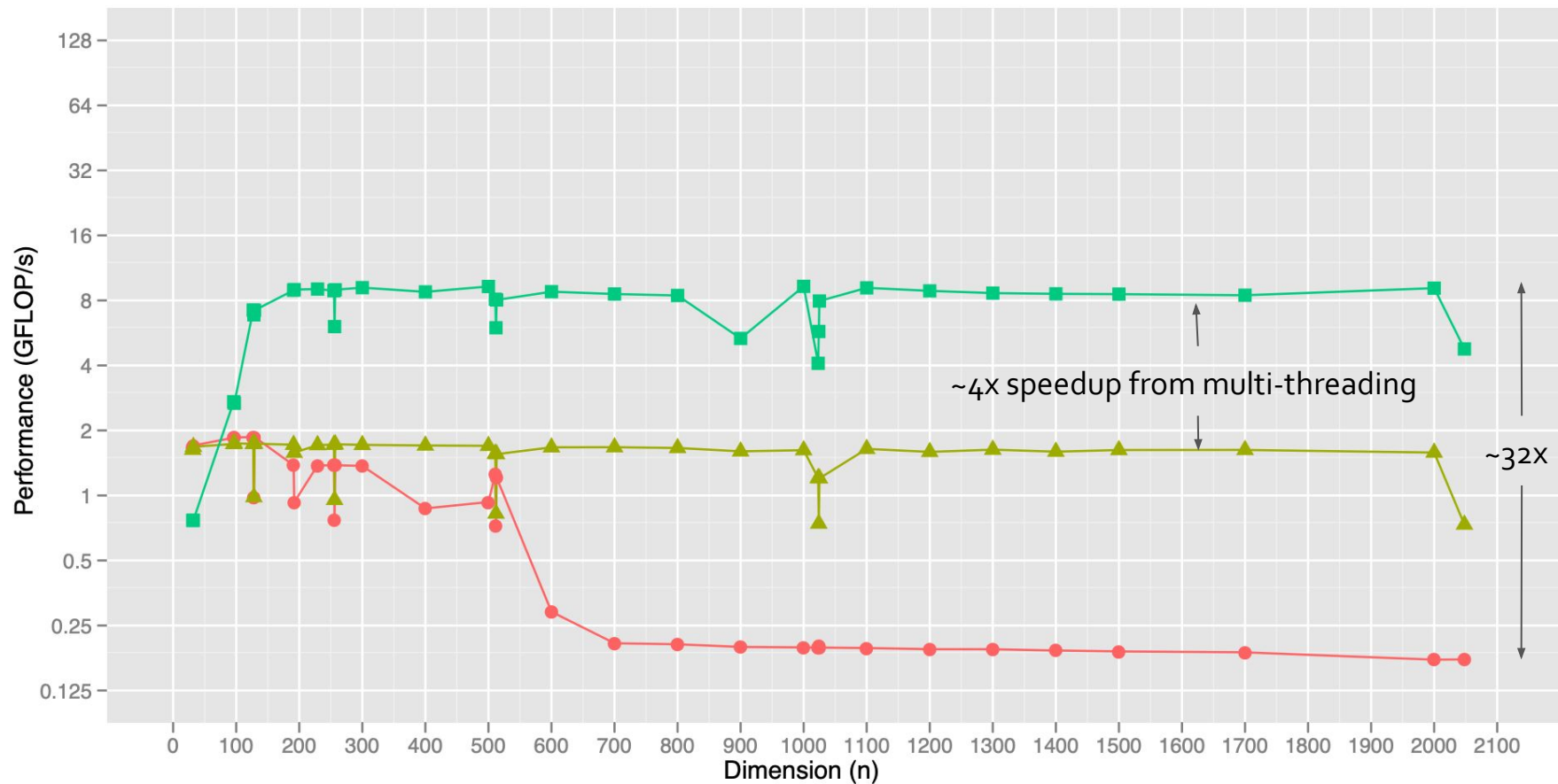
Code

Baseline Blocked Cilk++ (rec): p=16 Intel MKL Intel MKL: p=16



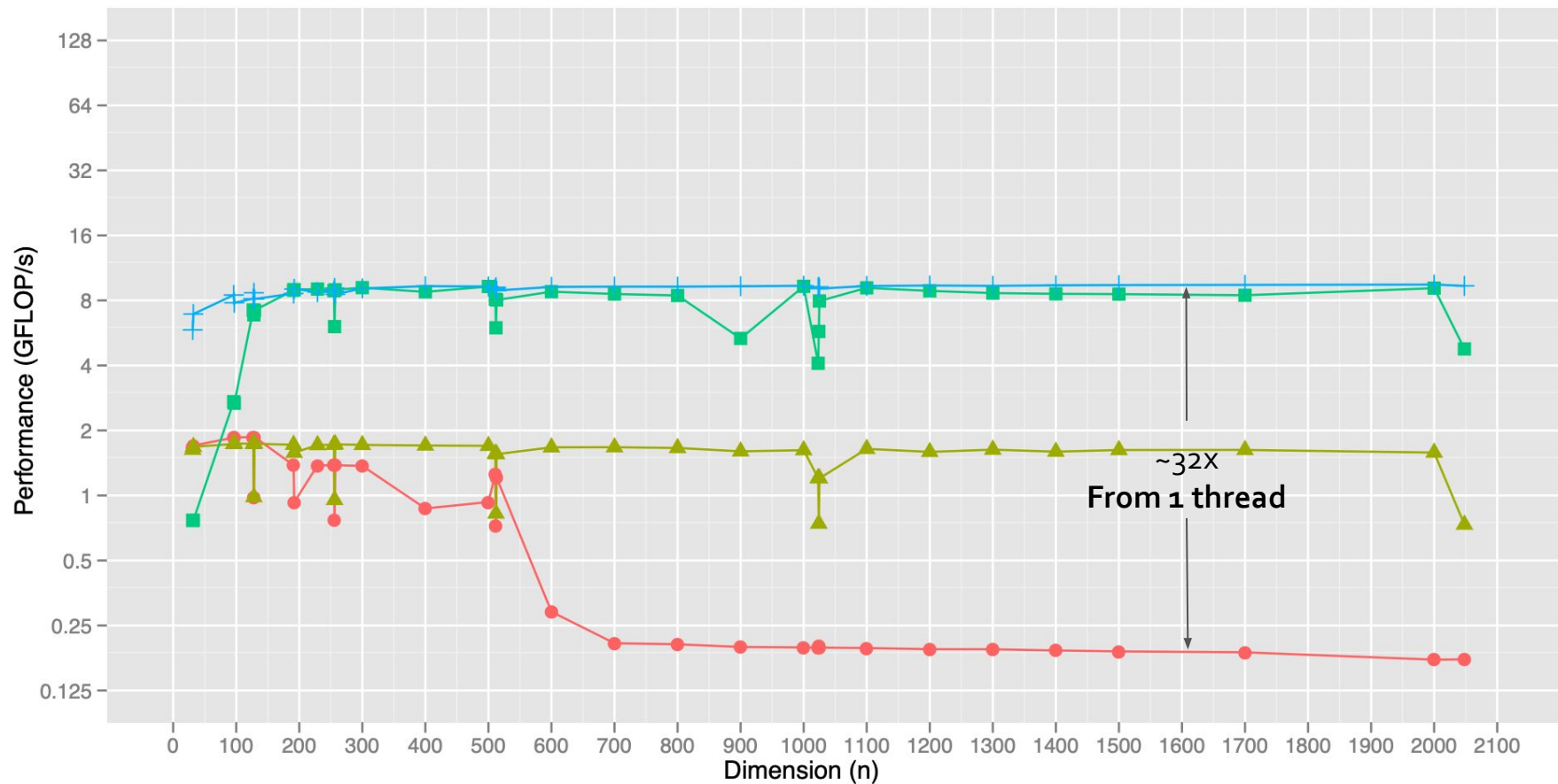
Code

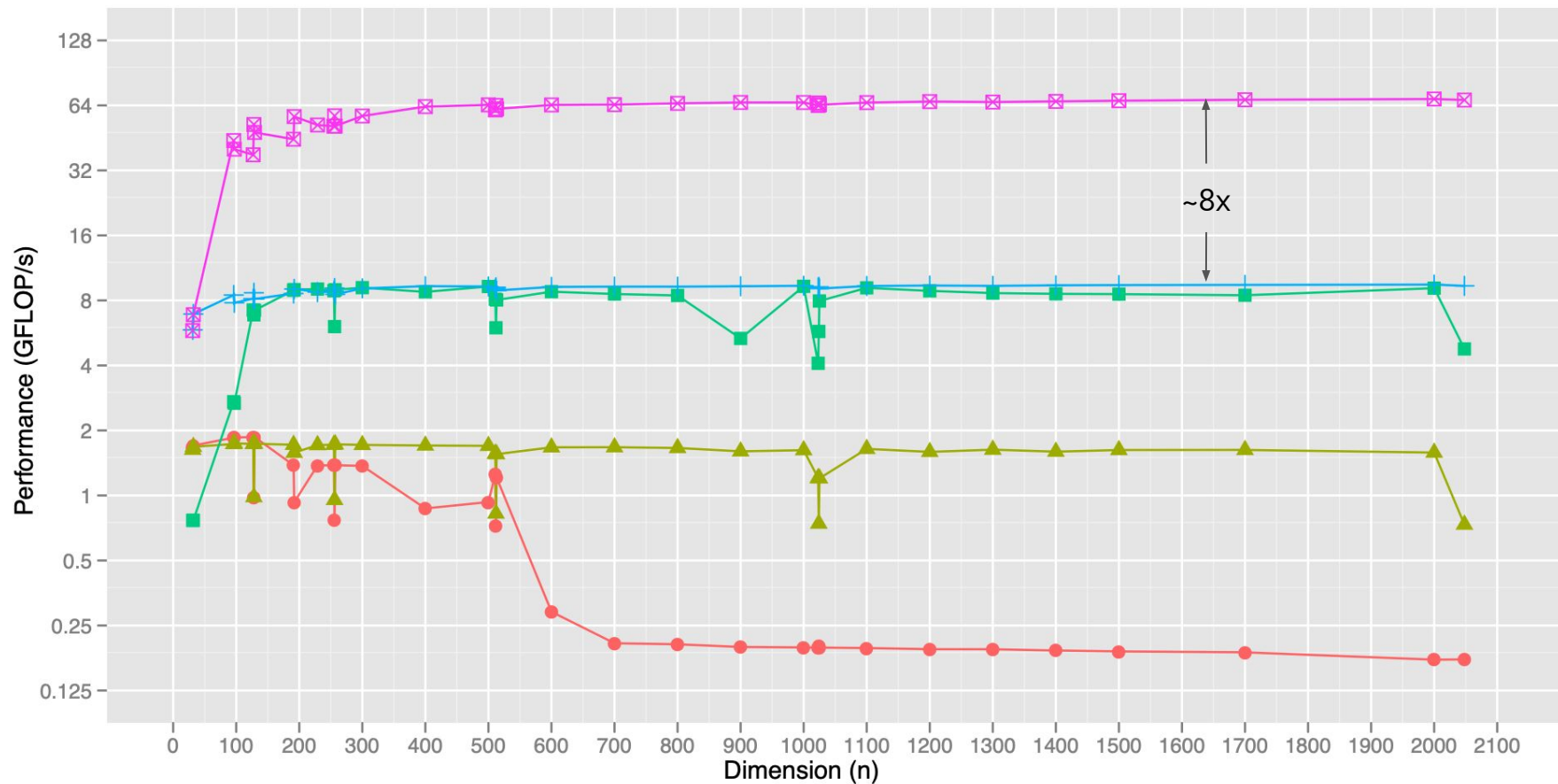
Baseline Blocked Cilk++ (rec): p=16 Intel MKL Intel MKL: p=16



Code

Baseline Blocked Cilk++ (rec): p=16 Intel MKL Intel MKL: p=16





Can Compilers Do This?

Proebsting's Law

Compilers double code performance every **18 years**.

vs. **2 years** for Moore's Law and **3 years** for memory

OpenMP

OpenMP - Open Multi-Processing

API that supports multi-platform parallel shared-memory multiprocessing programming in C, C++, etc.

Uses `#pragma` to denote places that can be parallelized

OpenMP Syntax

Most OpenMP constructs are compiler directives

```
#pragma omp <directive> [clause ...]
```

```
#pragma omp parallel default(shared) private(a, b)
```

Library Functions

Thread queries (number of threads, thread ID, etc.)

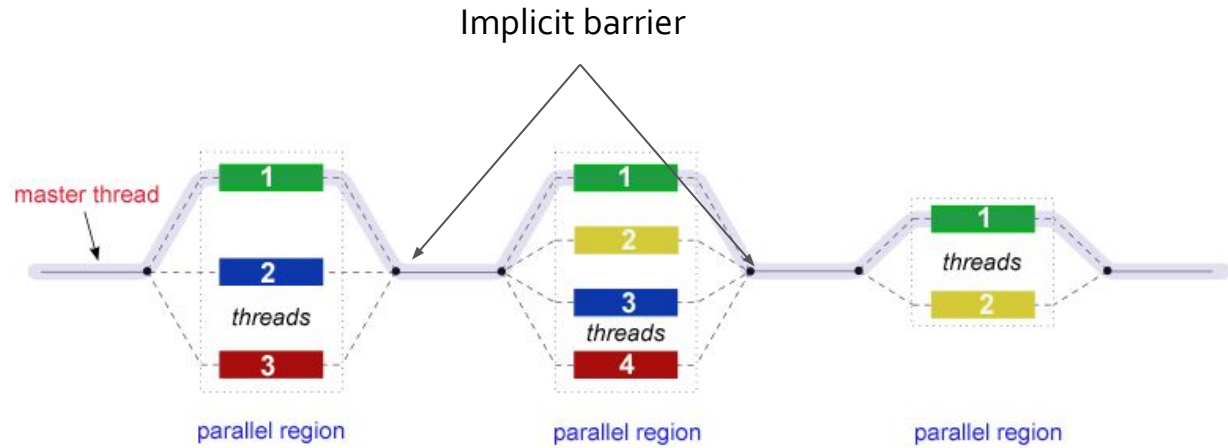
```
int omp_get_num_threads(void)
```

Environment Variables

Setting number of threads, affinity, etc.

```
export OMP_NUM_THREADS=8
```

OpenMP Model

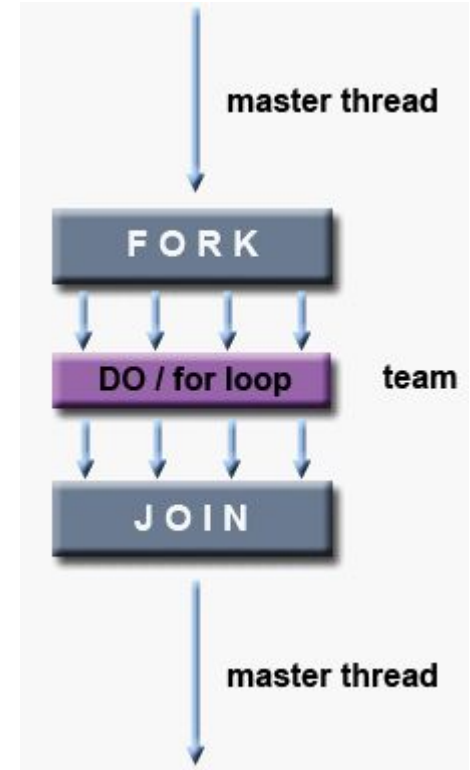


Work-Sharing - Do/For

Share iterations of the loop across the threads (i.e., data parallelism)

```
#pragma omp parallel
{
    #pragma omp for
    for(int i = 0; i < 100; i++) {
        x[i] = 1;
    }
}
OR
#pragma omp parallel for
for(int i = 0; i < ARR_SIZE; i++) {
    x[i]++;
}
```

Also an implicit barrier at the end of the loop



OpenMP Directives

`schedule (type, chunk size)`

Specify how the threads are assigned to the loop iterations

`type`

`static`

`dynamic`

`guided`

`runtime`

`Auto`

`chunk size` – denotes how many loop iterations are assigned to each thread at a time

Kind	Description
static	Divide the loop into equal-sized chunks or as equal as possible in the case where the number of loop iterations is not evenly divisible by the number of threads multiplied by the chunk size. By default, chunk size is <code>loop_count/number_of_threads</code> . Set chunk to 1 to interleave the iterations.
dynamic	Use the internal work queue to give a chunk-sized block of loop iterations to each thread. When a thread is finished, it retrieves the next block of loop iterations from the top of the work queue. By default, the chunk size is 1. Be careful when using this scheduling type because of the extra overhead involved.
guided	Similar to dynamic scheduling, but the chunk size starts off large and decreases to better handle load imbalance between iterations. The optional chunk parameter specifies the minimum size chunk to use. By default the chunk size is approximately <code>loop_count/number_of_threads</code> .
auto	When <code>schedule (auto)</code> is specified, the decision regarding scheduling is delegated to the compiler. The programmer gives the compiler the freedom to choose any possible mapping of iterations to threads in the team.
runtime	Uses the <code>OMP_schedule</code> environment variable to specify which one of the three loop-scheduling types should be used. <code>OMP_SCHEDULE</code> is a string formatted exactly the same as would appear on the parallel construct.

Next Homework

Quicksort

- You should all be familiar with this
- Recursively...
 - **Partition** your list around a pivot - smaller numbers go on the left, larger numbers go on the right
 - Call quicksort on your left and right sides of the array

Is there a different way of partitioning the data?

Partitioning

A:

9	5	7	11	1	3	8	14	4	21
---	---	---	----	---	---	---	----	---	----

 $x = 8$

Partitioning

Is g less than 8?

Yes = 1

No = 0

A:

9	5	7	11	1	3	8	14	4	21
---	---	---	----	---	---	---	----	---	----

 $x = 8$

B:

0	1	2	3	4	5	6	7	8	9
9	5	7	11	1	3	8	14	4	21

lt:

0	1	2	3	4	5	6	7	8	9
0	1	1	0	1	1	0	0	1	0

gt:

0	1	2	3	4	5	6	7	8	9
1	0	0	1	0	0	0	1	0	1

Partitioning

A:

9	5	7	11	1	3	8	14	4	21
---	---	---	----	---	---	---	----	---	----

x = 8

B:

0	1	2	3	4	5	6	7	8	9
9	5	7	11	1	3	8	14	4	21

lt:

0	1	2	3	4	5	6	7	8	9
0	1	1	0	1	1	0	0	1	0

0	1	2	2	3	4	4	4	5	5
---	---	---	---	---	---	---	---	---	---

prefix sum

gt:

0	1	2	3	4	5	6	7	8	9
1	0	0	1	0	0	0	1	0	1

gt:

1	1	1	2	2	2	2	3	3	4
---	---	---	---	---	---	---	---	---	---

prefix sum

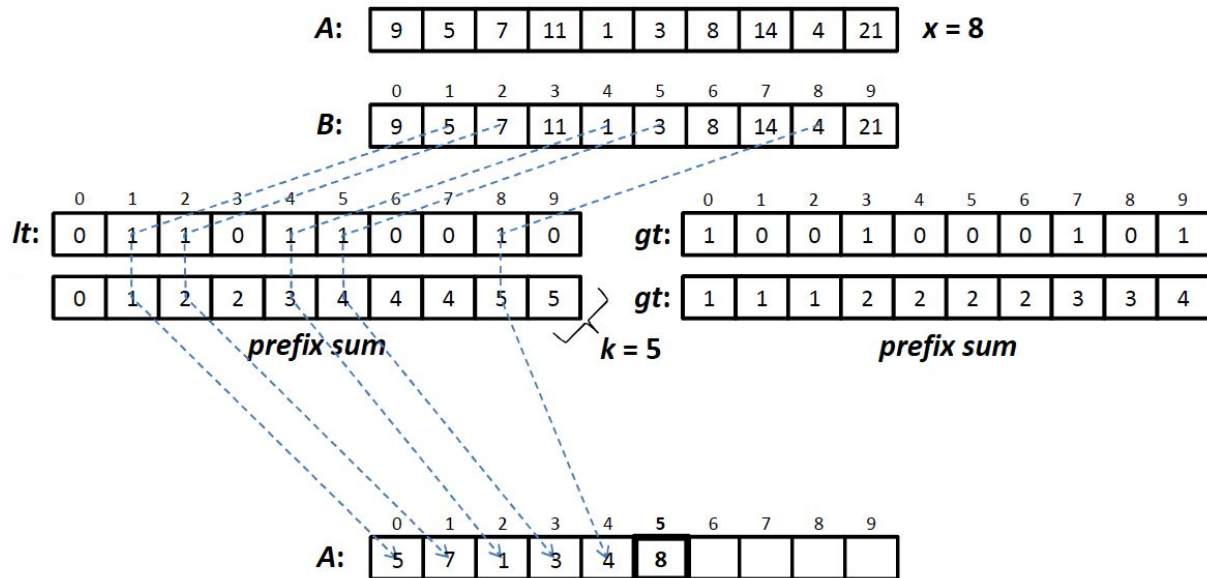
$$y_0 = x_0$$

$$y_1 = x_0 + x_1$$

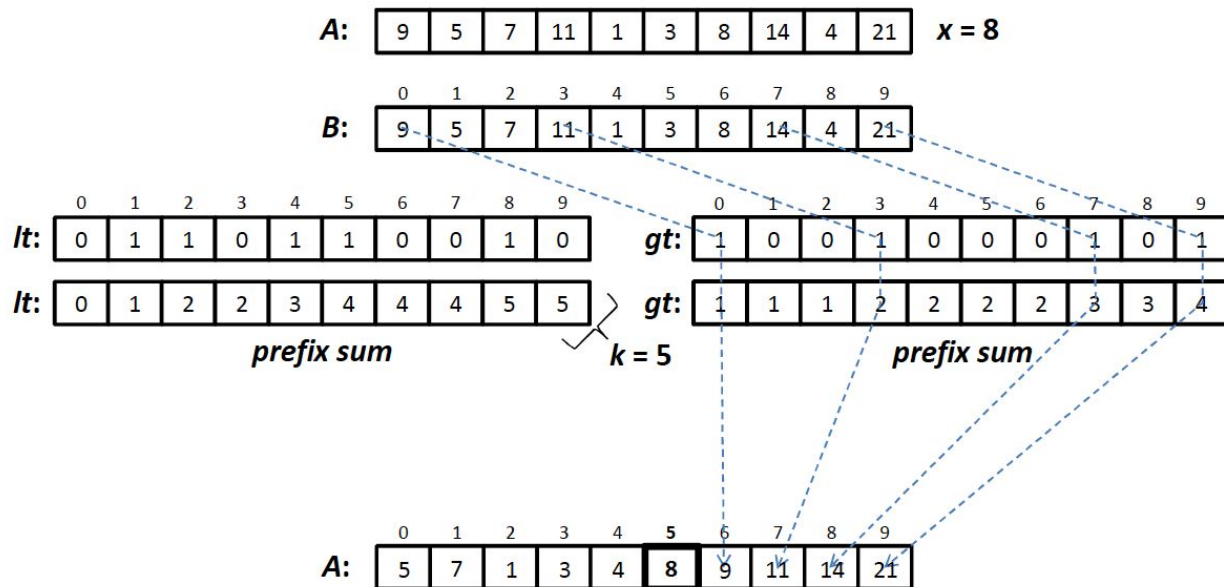
$$y_2 = x_0 + x_1 + x_2$$

...

Partitioning



Partitioning



Partitioning

Homework 8 (last one)

- Implement naive quicksort
- Implement quick sort with this new partitioning method
 - Parallelize it with OpenMP
 - The code WILL NOT be faster on ix-dev, but it should work correctly regardless of the # of threads (i.e., OMP_NUM_THREADS)
 - Try up to 1M elements to make sure it works correctly