

# Multi-GPU Short-Read Alignment

Radhakrishna “RK” Bettadapura

Strand Life Sciences



# Next Generation Sequencing (NGS)



Parallel sequencing

Read 1: **AGGT**    Read 2: **GAAATTG**    ...    Read n: **GCGGTTCCG**

Alignment

Reference Genome: **AGGCCTTACCGGAAATTGCGGTTCCGAA**....

Reassembly

Sequenced Genome: **AGGCTCTTAGCCAGATATGCGGTTCCGA**....



# Next Generation Sequencing (NGS)

Parallel sequencing

Read 1: **AGGT**    Read 2: **GAAATTG**    ...    Read n: **GCGGTTCCG**

Alignment

Reference Genome: **AGGCCTTACCGGAAATTGCGGTTCCGAA...**

Reassembly

Sequenced Genome: **AGGCTCTTAGCCAGATATGCGGTTCCGA....**

BIG Data  
10<sup>9</sup> reads  
1 Read ~=  
100-1000bp  
NA 12878:  
1 GB  
reference,  
100GB  
reads



# Short-read alignment

Read: ACTAGAATGGCT

Search (Burrows-Wheeler  
Index)

Genome: AGGCCTTACCGGAAATTGCGGTTCCGAA....

Matching locations: 2, 7, 10...  
Corresponding genomic substrings:  
Reference substring: CCATACTGAACTGACTAAC

Smith-Waterman local  
alignment

Reference: CCATACT GAACTGACTAAC  
Read: ACTAGAA TGGCT

“The core algorithms haven’t changed.”



# Short-read alignment

Read: **ACTAGAA**TGGCT

Search (Burrows-Wheeler Index)

Genome: **AGGCCTTACCGGAAATTGCGGTTCCGAA**....

Reference substring: **CCATACTGAACTGACTAAC**

Local alignment Smith-Waterman dynamic program (DP)

Reference : **CCATACT GAACTGACTAAC**

Read : **ACTAGAA TGGCT**

←→ ←→ ←→ ←→ ←→ ←→ ←→ ←→  
Start 3M 1I 3M 1D 2M 1m 2M



# Short-read alignment

Read 1: **ACTAGAA**TGGCT

Read 2: **ACTATGTAT**GGCT

Read n: **ACTATGTTAG**GCCT

Search (Burrows-Wheeler Index)

Genome: **AGGCCTTACCGGAAATTGCGGTTCCGAA**....

Reference 1: **CCATACTGAACTGACTAAC** 2, 3,...m

Smith-Waterman local alignment

Reference 1: **CCATACT** **GAACTGACTAAC**  
Read 1: **ACTAGAA** **TGGCT**

$m \sim 20 \times 10^9$

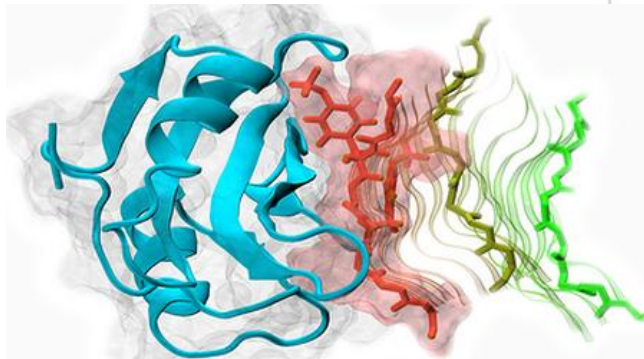
Takes 40-50 hours  
on single-thread  
CPU!

With 30 threads, 3-  
5 hours

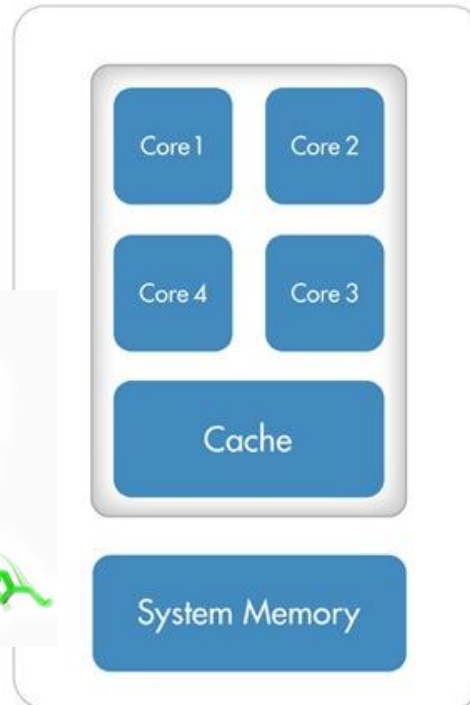
Can we use the  
GPU?



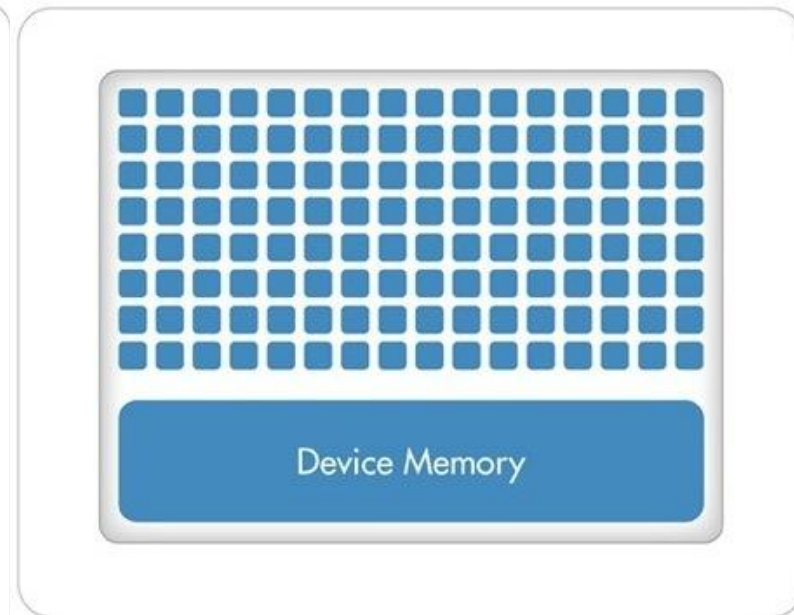
# Graphics Processing Unit (GPU)



**CPU (Multiple Cores)**



**GPU (Hundreds of Cores)**



- Scalable, massively parallel computing
- Initially used to speed up graphics
- General purpose took off with CUDA
- Speedups of upto 500x, depending on the problem

K10 Tesla:  
2.8GB global mem.  
1536 cores  
745MhZ clock



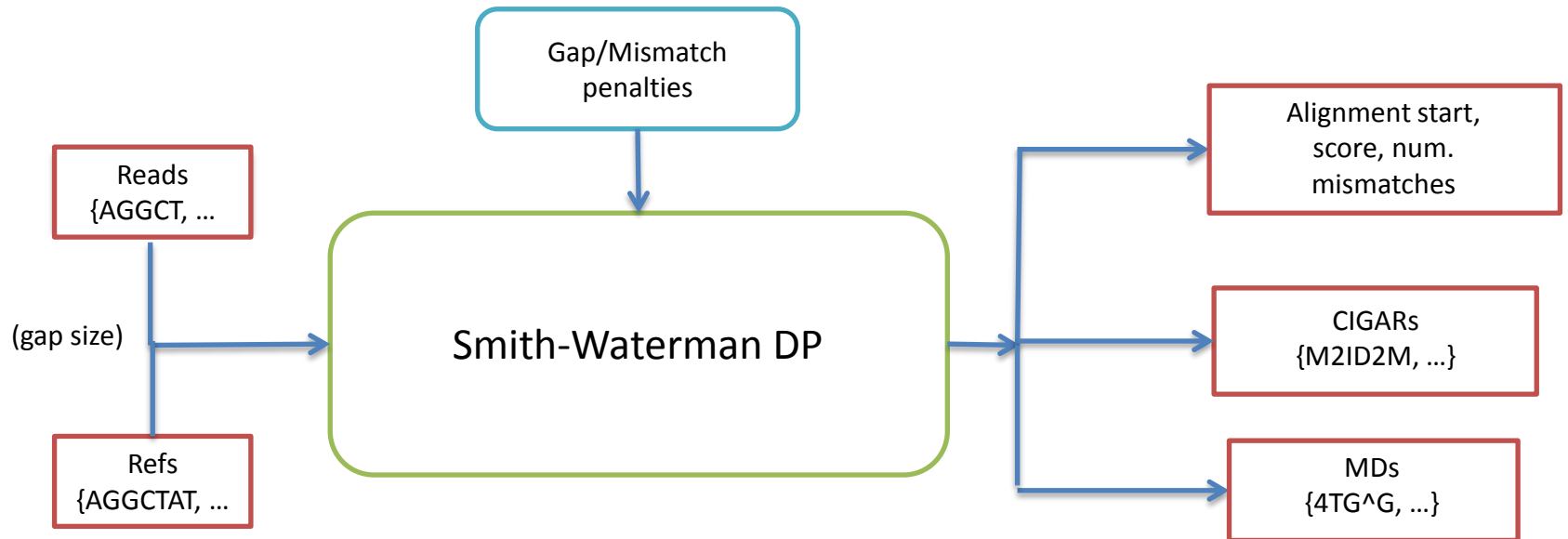
# Goal: Implement Smith Waterman on multiple GPUs

Secondary Goal: Prove that it's  
(much) faster!

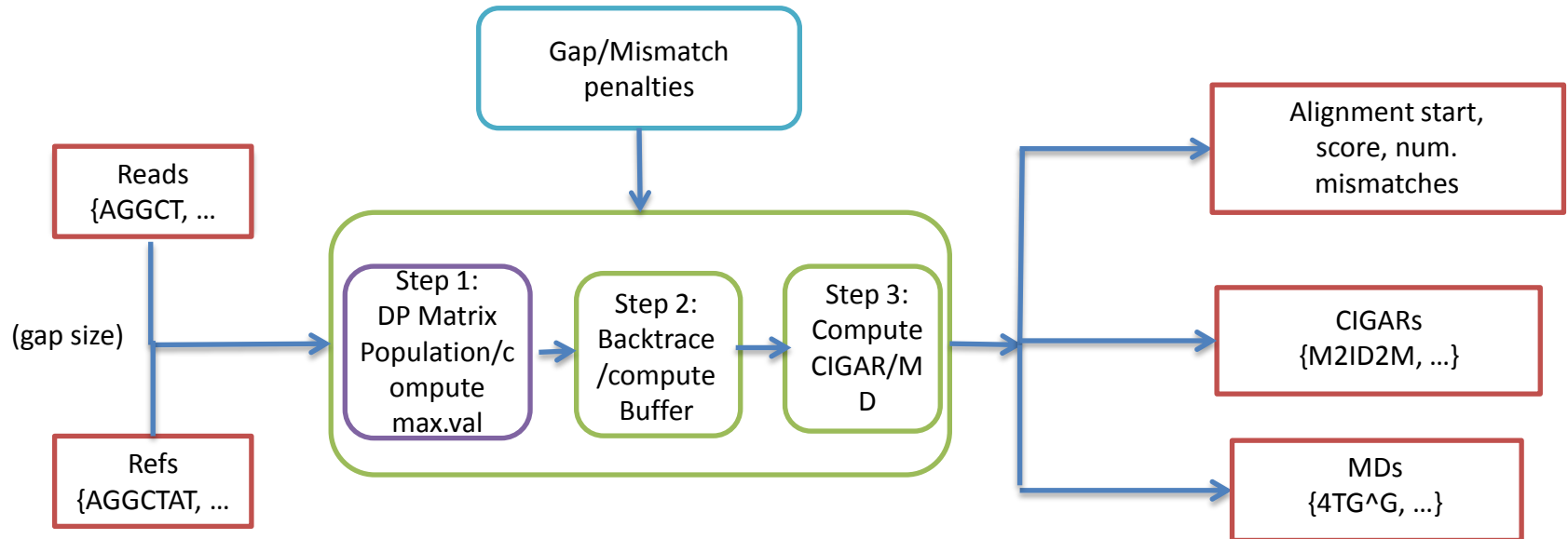




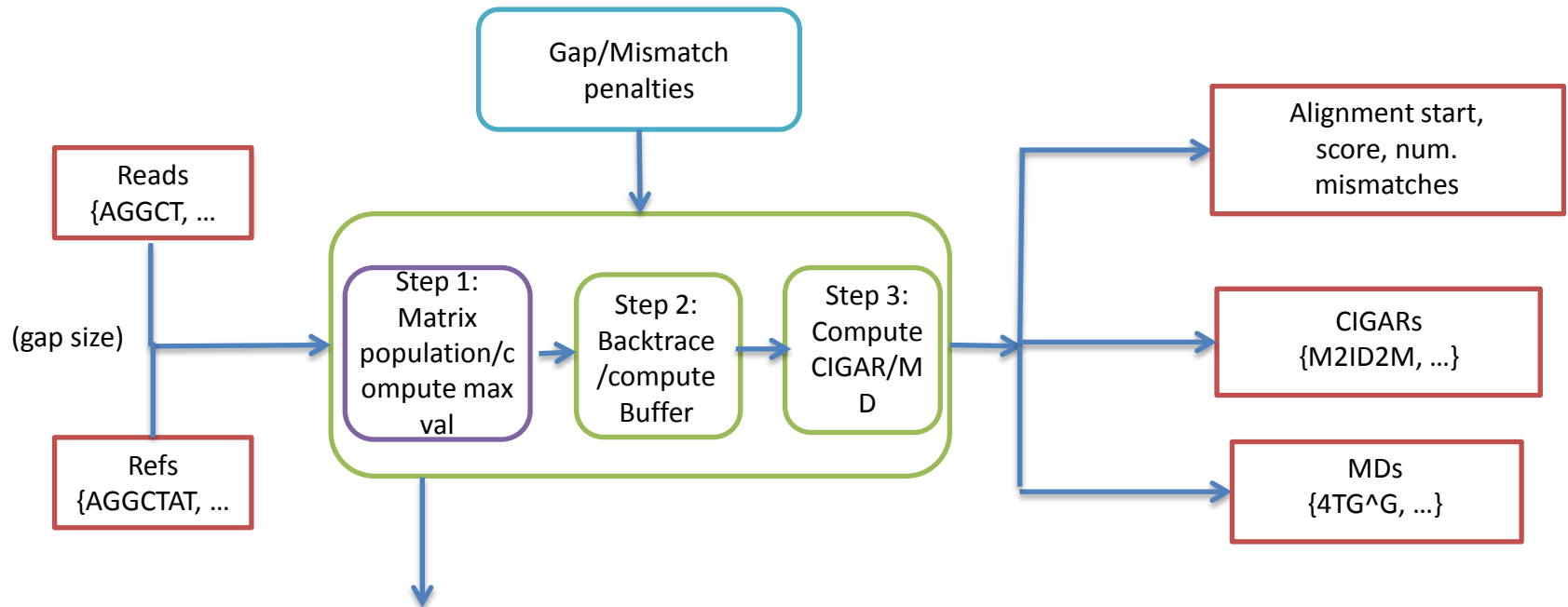
# Algorithm



# Algorithm



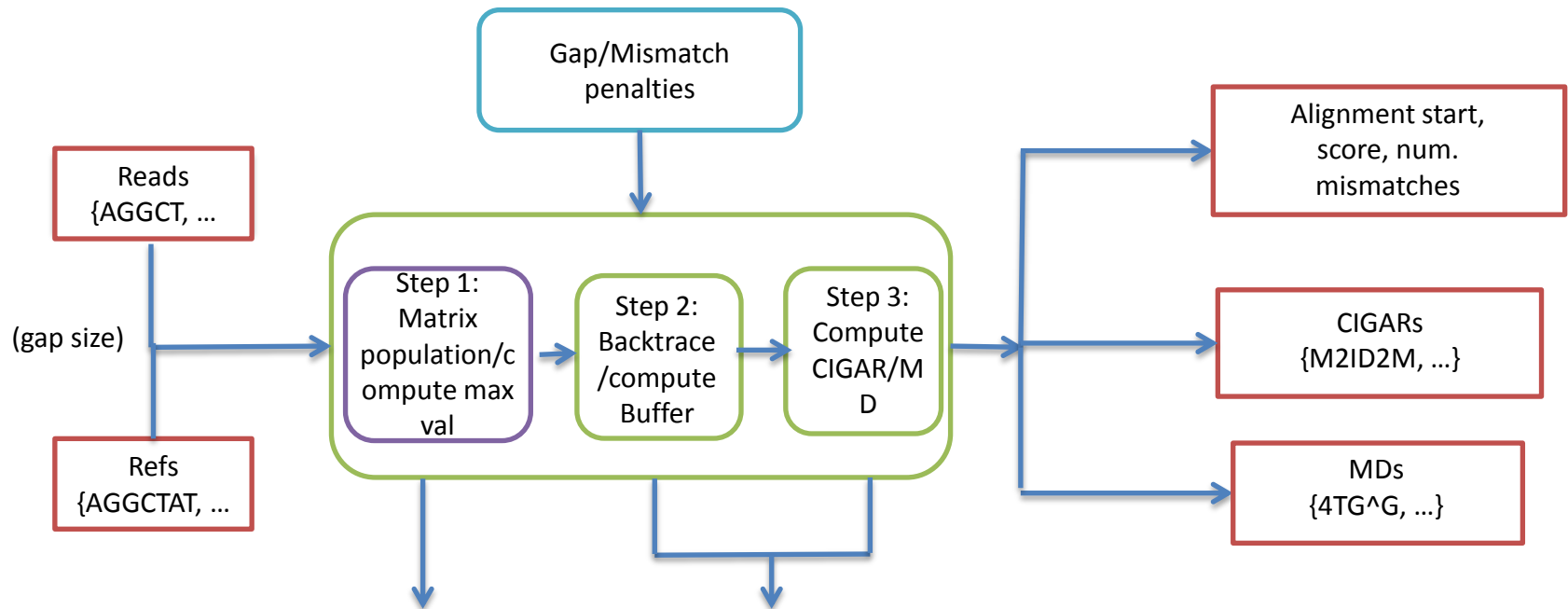
# Algorithm



- Prior work (protein alignment)
- Quadratic computation but linear storage



# Algorithm

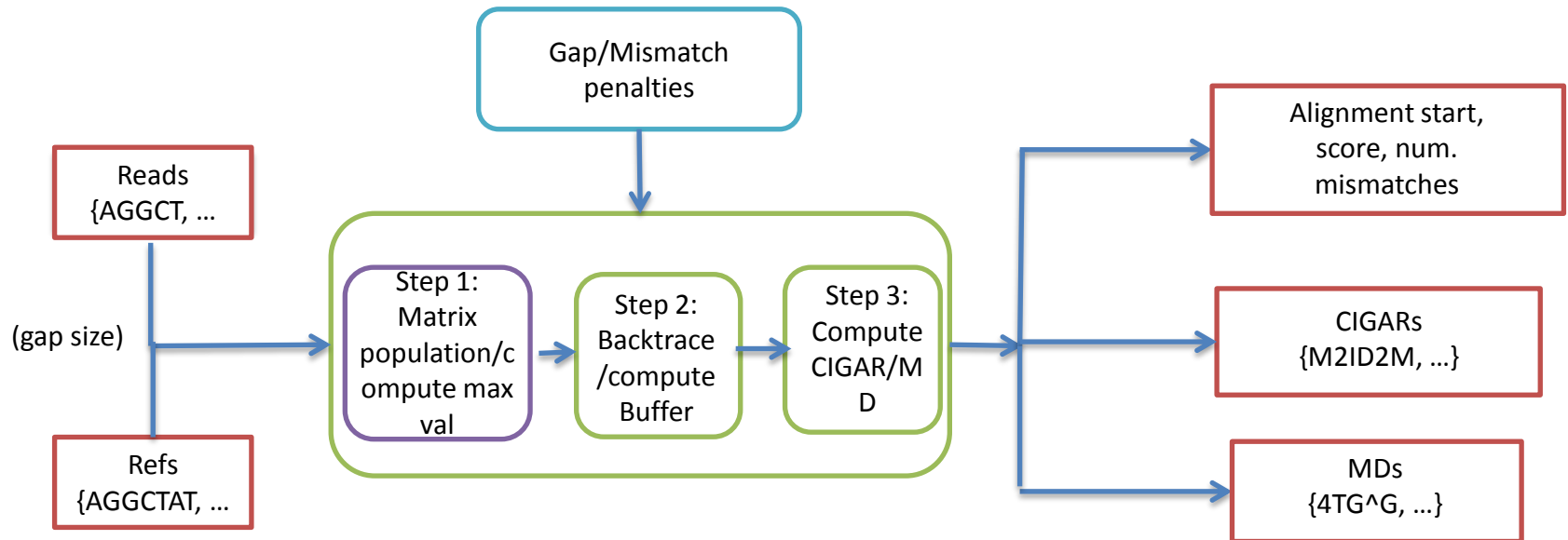


- Prior work (protein alignment)
- Quadratic computation but linear storage

- Quadratic storage
- Cannot be done in parallel with step 1 (though see final slide)
- Very high mem/math ratio (no floating point math + loads/stores + string manipulation)



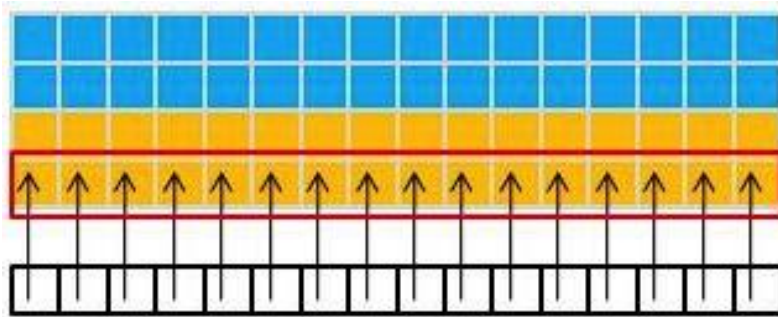
# Algorithm



Parallelisation scheme: 1 thread per DP  
(Store DP matrix/CIGAR/MD in thread local storage)

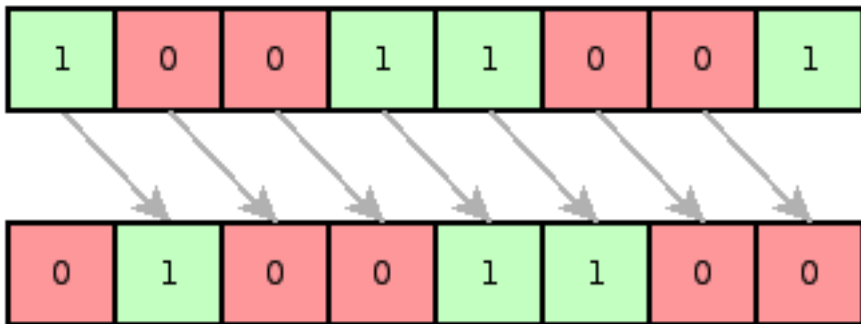
# Implementation

Coalesced memory access:

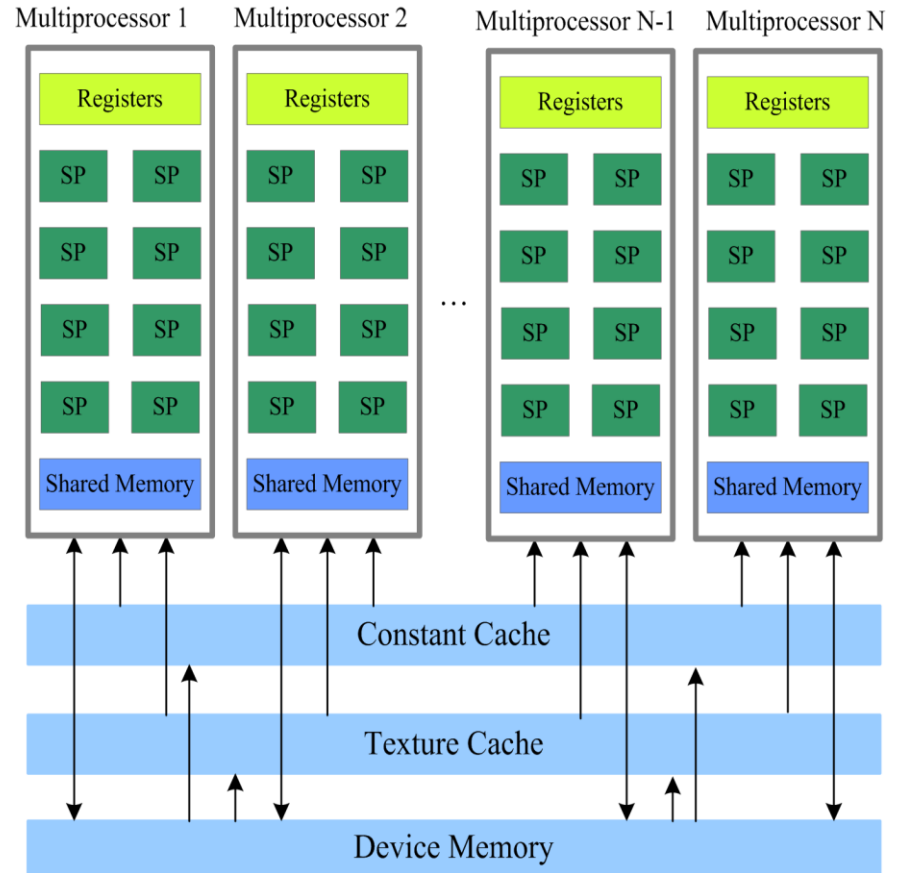


Thread: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

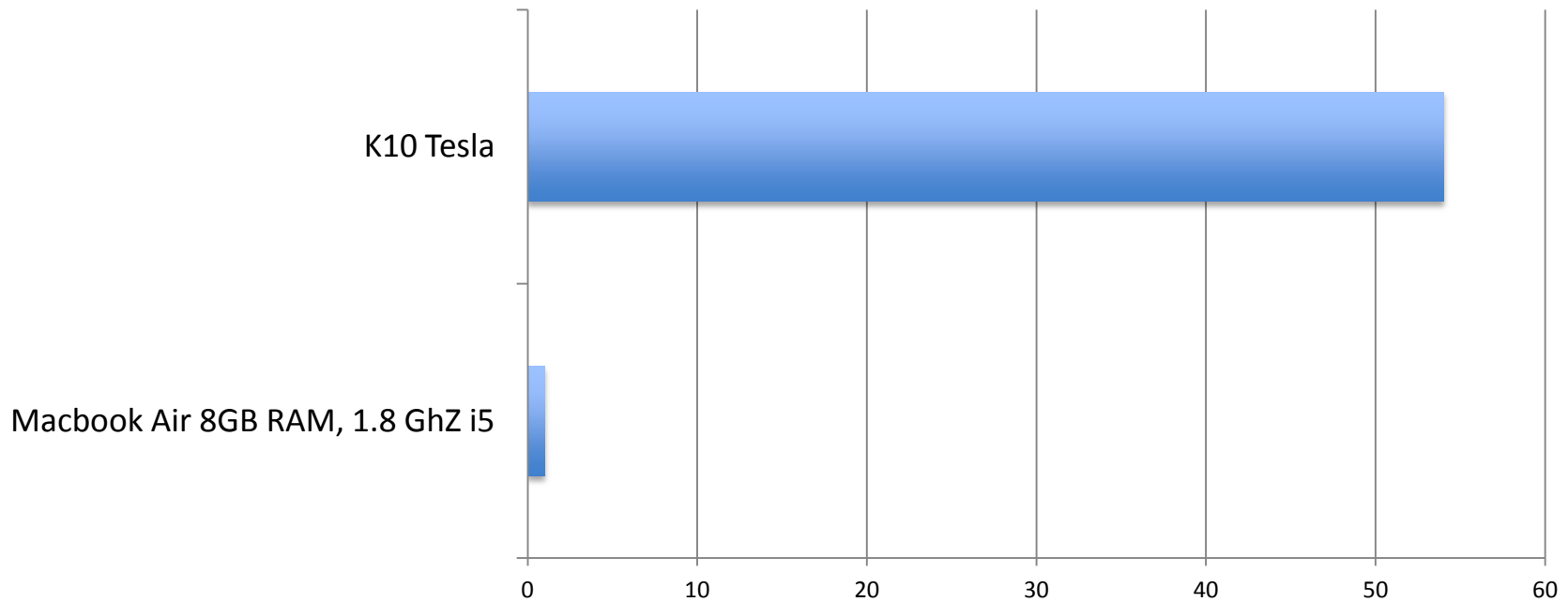
Judicious use of bitwise operations (prove them optimal first!)



Judicious use of constant/register mem.



# Implementation



Bottom line: SW GPU is 54x CPU.



# Implementation

## Performance metrics

- 7.5 billion DPs/hr for read length = 100 and gap size = 5
- Some more metrics:
  - L1 cache miss ratio = 2%
  - L2 cache miss ratio = 0%
  - No register spills
  - No. of instructions executed/total number of instructions issued = 95% (indicates low warp serialisation)

Bottom line: Novel implementation of SW (1, 2, 3) on the GPU results in extremely fast execution times

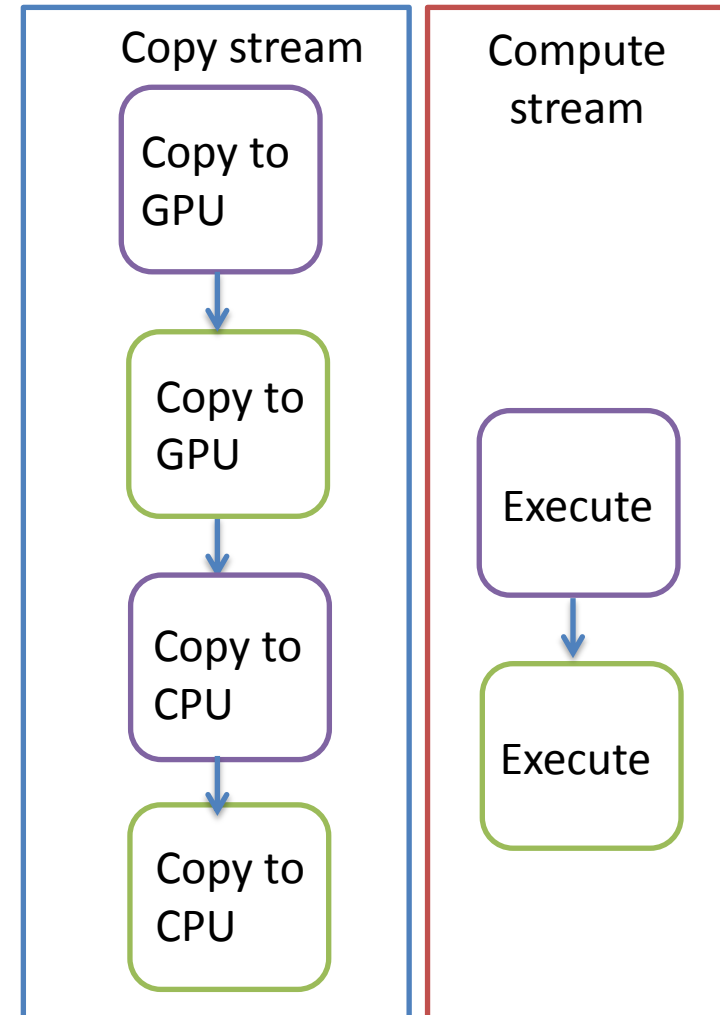
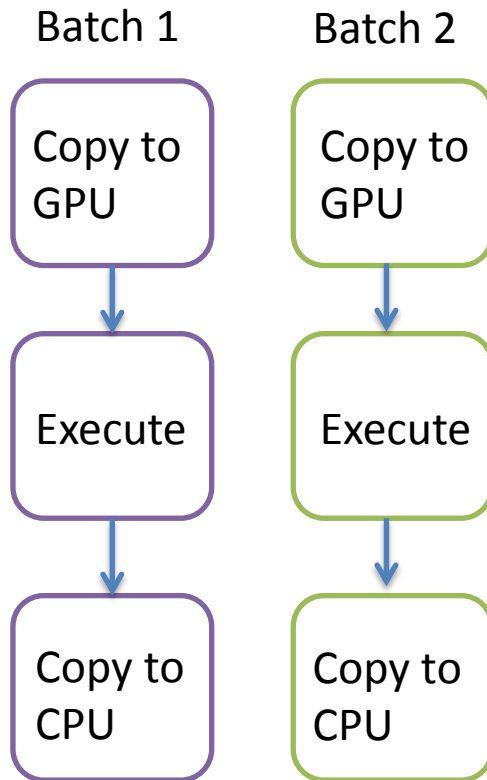
“That’s all right in practice, but does it work in theory?”





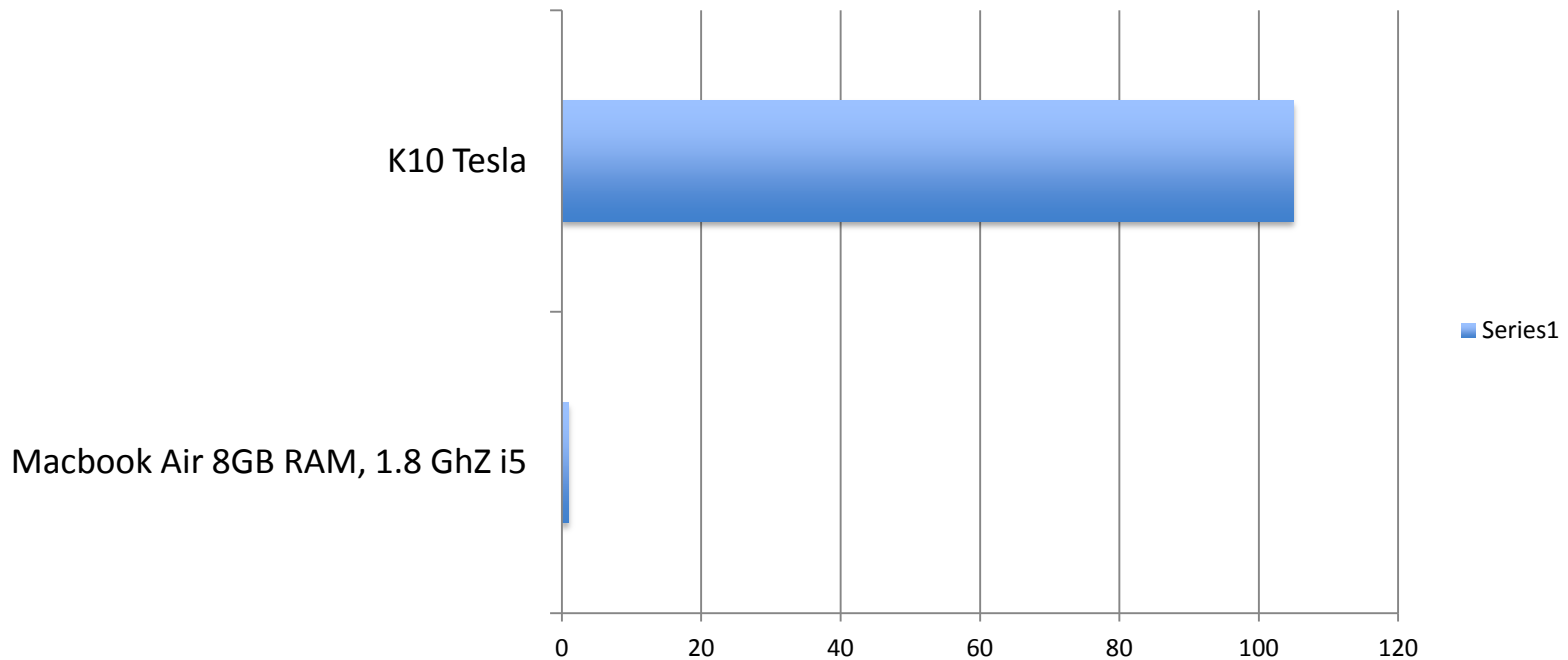
# Multiple GPUs

Use streams to schedule work and reduce copy overhead!



On multiple GPUs, streams mitigate copy overhead via GPU-wise asynchronicity

# Implementation

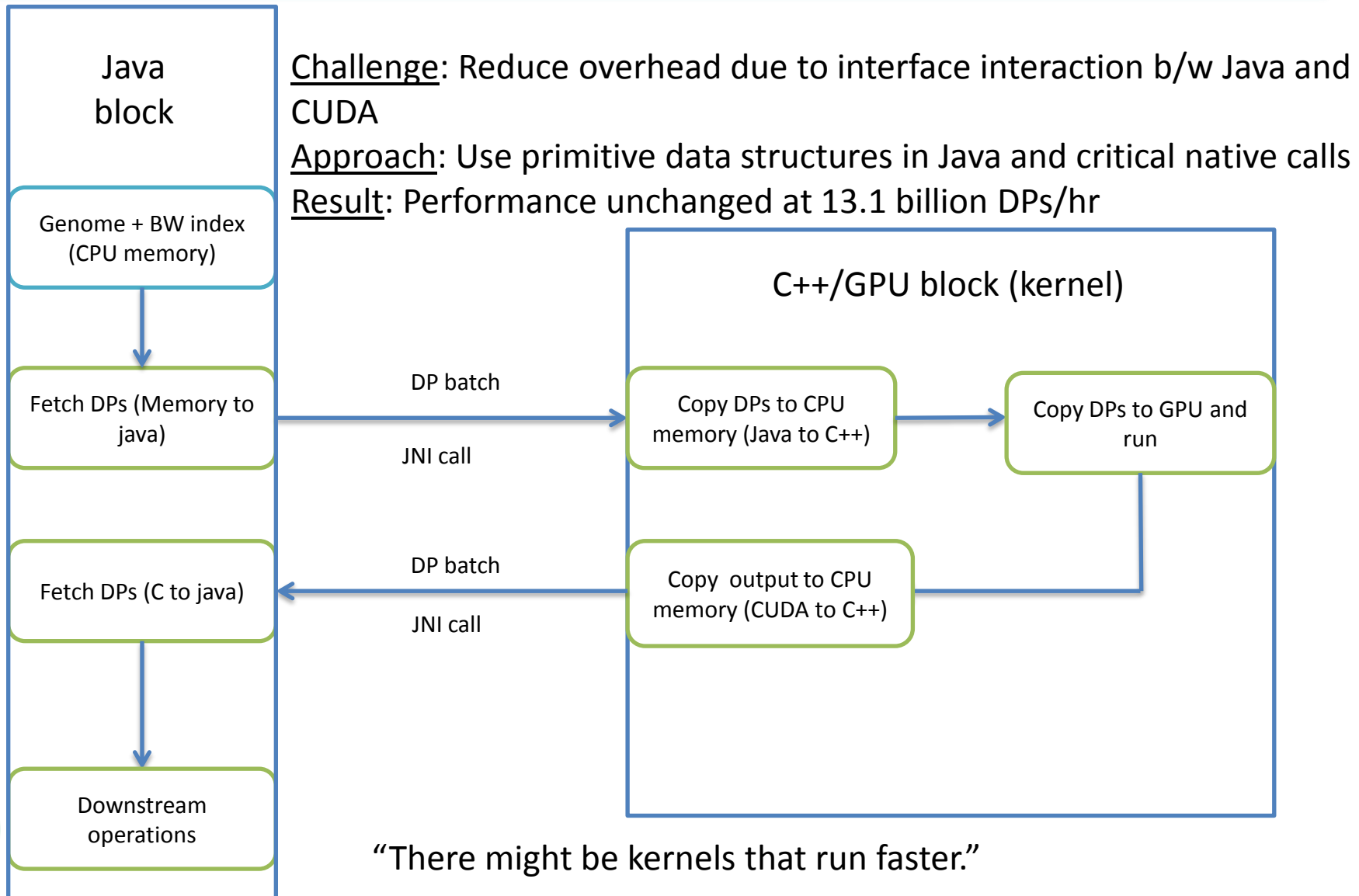


Peak rate: 13.1 billion DPs/hr for read length = 100 and gap size = 5

Bottom line: SW 2xGPU is 105x CPU.

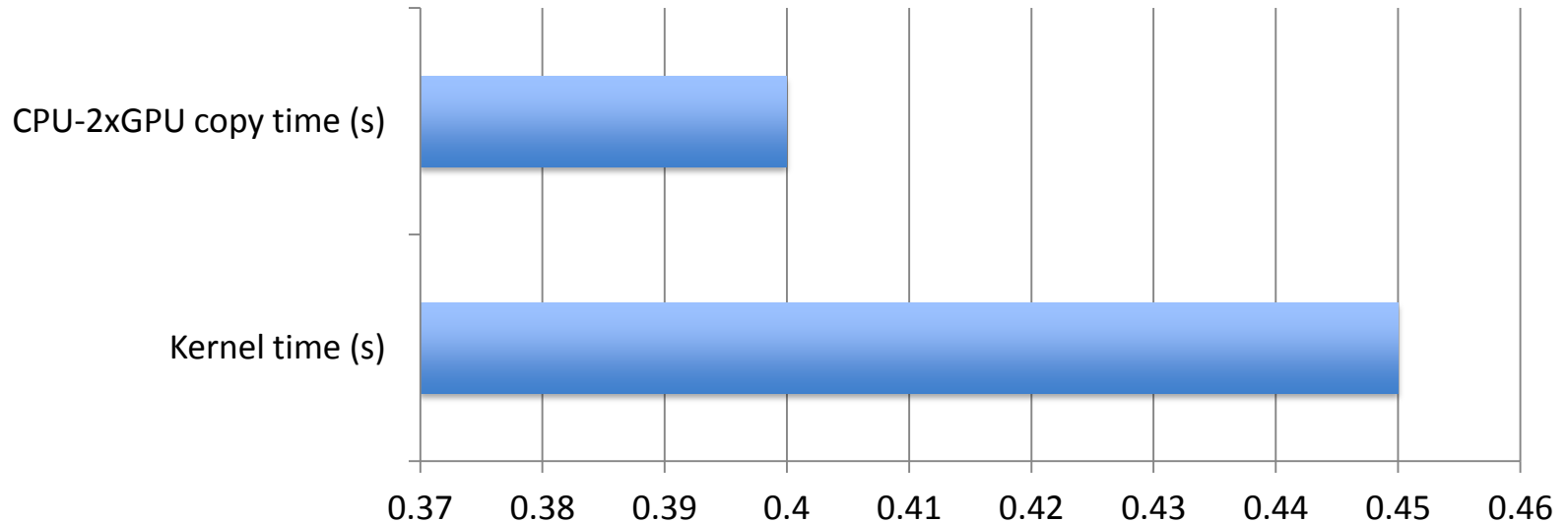


# Implementation



“There might be kernels that run faster.”

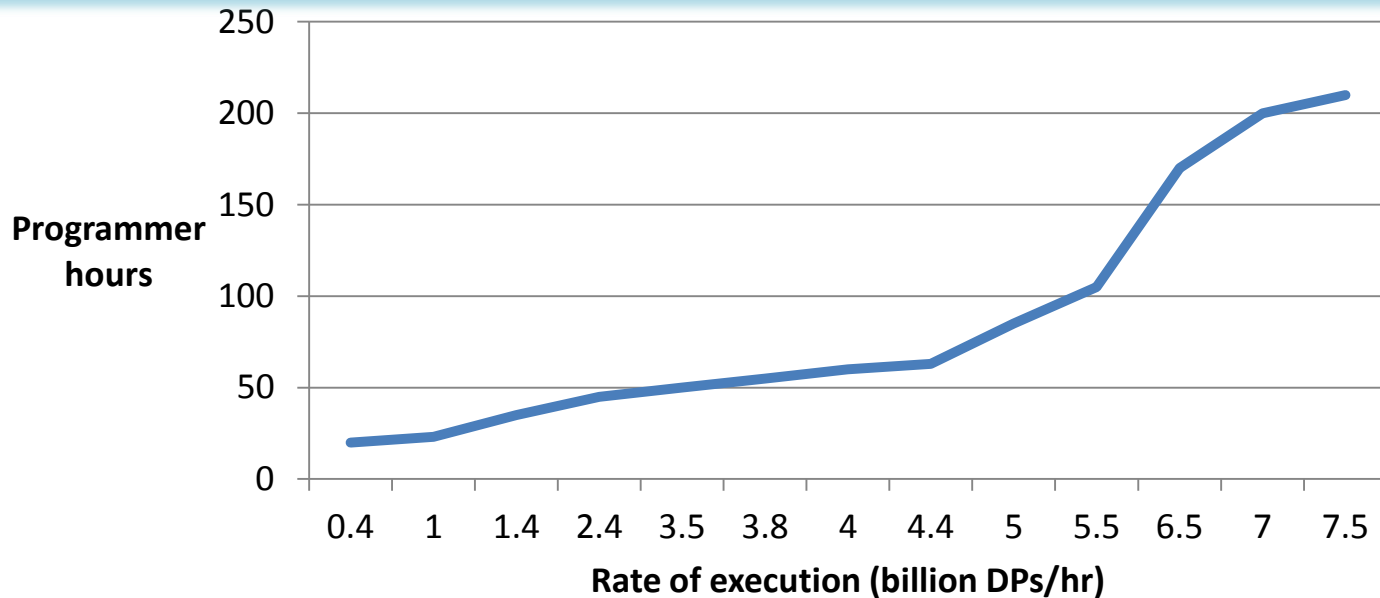
# Lessons and Provisos



- Can't get much better (Limited by GPU-CPU bandwidth)



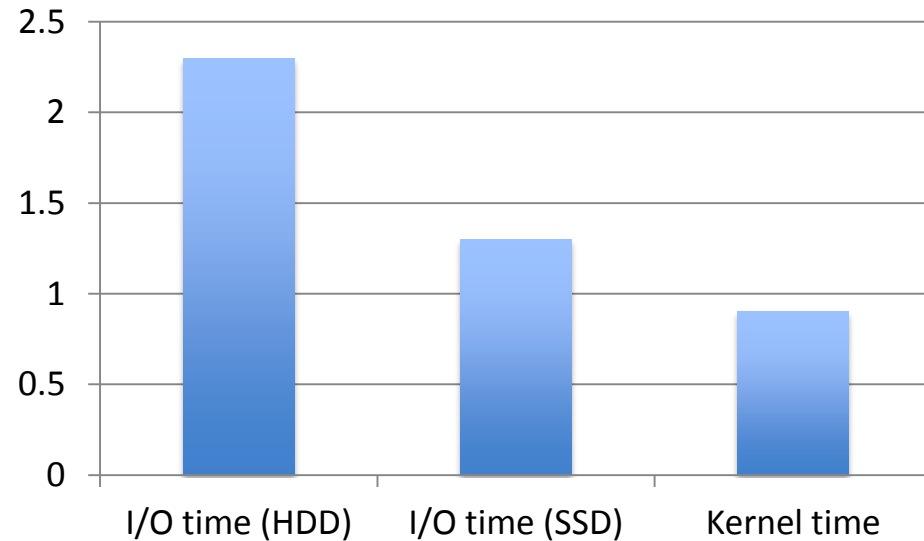
# Lessons and Provisos



- GPU isn't cure-all
  - Since core algorithms haven't changed, tweaking implementation on GPU requires many "clever" changes)



# Lessons and Provisos



- I/O limitations are dire

“Que cela manque souvent d’architecture.”



Thanks!