

GPU Microarchitecture Note Set 1c—Program Needs and System Cap

Our Goal:

Determine whether we should be **satisfied**...

... with **the performance** of our program...

... on our system.

We are not satisfied when the program runs **more slowly** than we **expect**.

Knowing when to be satisfied is a **key skill**.

The Idea:

- Estimate the *computation needs* of our problem, algorithm, or program.

For example, 10^{20} floating-point operations.

- Determine the *computation capabilities* of our system.

For example, 10^{17} FLOPS.

- Use these to estimate execution time.

For example, $10^{20}/10^{17} = 1000$ s.

If the estimate does not match the measured value, find the cause.

Computation Needs

These refer to a program, algorithm, or problem.

Determined by analyzing a problem, algorithm, or program.

Common Computation Needs:

- Floating-Point Operations
- Instruction Count (Number of executed instructions)
- Data Transfer (bytes read from and written to memory).

Computation Capabilities

These refer to the capabilities of a system (CPU, GPU, etc.)

Common Computation Capabilities

Floating Point Execution Rate - FLOPS

Number of floating-point operations divided by amount of time.

Instruction Execution Rate - IPC

Number of executed instructions divided by amount of time.

Data Transfer Bandwidth - B/s

Number of bytes crossing some line, divided by amount of time.

Computation Need: Number of Floating Point Operations (FLOPS)

Commonly used for scientific programs...

... because such operations are considered essential (can't be avoided) and...

... because of the historic high cost of FP hardware.

Example: fragment of some scientific code.

```
double *x, *a, *b, *c;  
for ( i=0; i<1e9; i++ ) x[i] = a[i] + b[i] * c[i];
```

Uses $10^9 \times 2 = 2 \times 10^9$ FP operations.

Consider a system capable of 100 GFLOPS (10^{11} FLOPS).

Execution time bound is 20 ms.

If measured execution time was 100 ms we would **not be satisfied**.

Computational Needs \gg Data Transfer

Computation Need: Data Transfer

Used for all kinds of programs.

Easy to measure when each data item read exactly once.

Otherwise, need to account for cache and scratchpad performance.

Data Transfer Example:

```
double *x, *a, *b;  
for ( i=1; i<1e9-1; i++ ) x[i] = a[i-1] * l + a[i] + a[i+1] * r;
```

The size of a double is 8 bytes.

In an iteration, three elements of **a** are accessed...

... two were accessed in a prior iteration...

... and so shouldn't need to be read from memory...

... leaving one new element of **a** accessed per iteration.

Total data transfer needs: $10^9 \times 8 \times 2 = 16$ GB.

Suppose system can transfer 30 GB/s.

Performance bound: $16/30 = 533$ ms.

Data Transfer Example, Continued.

Suppose measured execution time were just 1.07 s. We are not satisfied.

Maybe we were wrong about each element of **a** being read once.

Maybe we should look at number of FP operations.

Maybe we should look at number of instructions.

Computation Needs, Matrix Multiplication

Code for multiplying two $d \times d$ matrices.

```
for ( int r=0; r<d; r++ )  
  for ( int c=0; c<d; c++ )  
    for ( int k=0; k<d; k++ )  
      c[r,c] += a[r,k] * b[k,c];    // This line executed  $d^3$  times.
```

Data Transfer

Size of each matrix: d^2 elements.

Naïve data bound: $3d^3$.

Minimum data bound: $3d^2$.

Computation

Straightforward: d^3 multiply/adds.

But, $O(d^{\log_2 7}) \approx O(d^{2.807})$ practically possible. (Strassen's Algorithm.)

Computational Needs \gg Example: Matrix Multiplication

Instruction Count and Execution Rate

Considered a less fundamental bound than FLOPS and data transfer.

Number of instructions determined by many factors:

How program is written.

Quality of compiler.

Instruction set of processor.

Computational Needs \gg Example: Matrix Multiplication

Performance Estimation

With above needs and capabilities can compute performance bounds.

Actual performance will be no greater than the smallest of these.

There's another measure that prevents this peak performance from being achieved.

Latency:

The time needed to do something from start to finish.

In discussion of GPUs, latency often refers...

... to the **execution time** of a **single instruction** or of a **single thread**.

A processor is called *latency oriented* ...

... if it has **low latency** (runs a single thread quickly).