

GPU Computing

E. Carlinet, J. Chazalon (firstname.lastname@rde.epita.fr)

Oct 21

EPITA Research & Development Laboratory (LRDE)



Slides generated on October 6, 2020

Fifty shades of Parallelism

1

How to get things done quicker

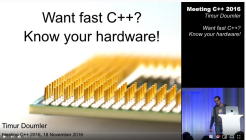
- 1. Do less work
- 2. Do some work better (i.e. the one being the more time-consuming)
- 3. Do some work at the same time
- 4. Distribute work between different workers

Fifty shades of Parallelism

How to get things done quicker

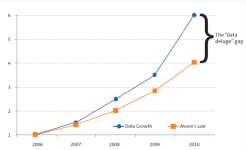
- 1. Do less work
- 2. Do some work better (i.e. the one being the more time-consuming)
- 3. Do some work at the same time
- 4. Distribute work between different workers

- (1) Choose the most adapted algorithms, and avoid re-computing things
- (2) Choose the most adapted data structures
- (3,4) Parallelism



2

Toward data-oriented programming

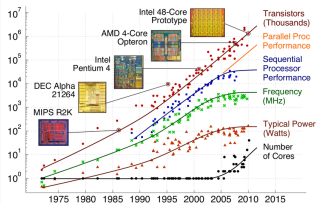


- while the CPU clock rate got bounded...
- ... the quantity data to process has shot up!

We need another way of thinking "speed"

Why parallelism ?

- Moore's law: processors are **not** getting twice as powerful every 2 years anymore



- So the processor is getting smarter:
  - Out-of-order execution / dynamic register renaming
  - Speculative execution with branch prediction
- And the processor is getting super-scalar:

3

The burger factory assembly line



How to make several sandwiches as fast as possible ?

3

5

The burger factory assembly line



How to make several sandwiches as fast as possible ?

- Optimize for **latency**: time to get 1 sandwich done.
- Optimize for **throughput**: number of sandwiches done during a given duration

Data-oriented programming parallelism

Flynn's Taxonomy

	Single Instruction	Multiple Instruction
Single Data	SISD	MISD
Multiple Data	SIMD	MIMD

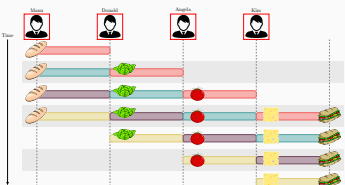
- SISD: no parallelism
- SIMD: same instruction on data group (vector)
- MISD: rare, mostly used for fault tolerant code
- MIMD: usual parallel mode

6

Optimize for throughput (MIMD Vertical Pipelining)



- Manu cuts the bread
- Donald slices the salads
- Angela slices the tomatoes
- ...



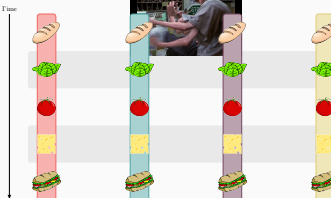
Time to make 4 sandwiches: s (400% speed-up)

7

Optimize for throughput (SIMD DLP)



A worker has many arms and make 4 sandwiches at a time



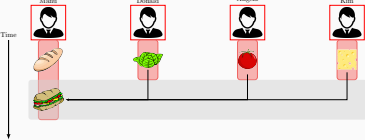
Time to make 4 sandwiches: s (400% speed-up)

10

Optimize for latency (MIMD with collaborative workers)



- 4 **super-workers** (4 CPU cores) collaborate to make 1 sandwich.
- Manu gets the bread and cuts and waits for the others
- Donald slices the salad
- Angela slices the tomatoes
- Kim slices the cheeses



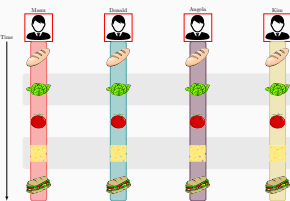
Time to make 1 sandwich:  $\frac{s}{4}$  (400% speed-up)

This is optimized for **latency** (CPU are good for that).

Optimize for throughput (MIMD Horizontal with multiple jobs)



- Manu makes sandwich 1
- Donald makes sandwich 2
- ...



Time to make 4 sandwiches: s (400% speed-up)

This is optimized for **throughput** (GPU are good for that).

8

More cores is trendy

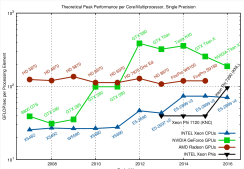
Data-oriented design have changed the way we make processors (even CPUs):

- Lower clock-rate
- Larger vector-size, more vector-oriented ISA
- More cores (processing units)

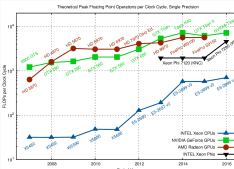
	64bits Intel Xeon	Xeon 5100 series	Xeon 5500 series	Xeon 5600 series	Xeon ES 2600 series	Xeon Phi 7120P
Freq	3.6 Ghz	3.0 Ghz	3.2 Ghz	3.3 Ghz	2.7 Ghz	1.24 Ghz
Cores	1	2	4	6	12	61
Threads	2	2	8	12	24	244
SIMD	128 bits	128 bits	128 bits	128 bits	256 bits	512 bits
Width (2 clocks)		(1 clock)	(1 clock)	(1 clock)	(1 clock)	(1 clock)

9

More cores is trendy



Peak performance / core is getting lower



Global peak performance is getting higher (with more cores!)

12

13

The diagram illustrates the OMAP5430 system architecture. At the top, interfaces for LPDDR2, LPDDR3, NAND Flash, SD, and USB 2.0 HS are shown. The core components include EMIF 1, EMIF 2, GPU, MMIO/SD, and Mail 2.0. The central processing unit is the OMAP5430, which contains a Dynamic memory manager, L2 cache, ARM Cortex-M4, ARM Cortex-M3, ARM Cortex-A15 MPCore, and ARM Cortex-A9. It also features a POWERVR SGX540 GPU, DRA750 DSP, and a Video Input Processor (VIP). The diagram shows connections to external components like 3G/4G modems, WiLink™, 30/45 modems, and various sensors (accelerometer, gyroscope, compass, etc.). It also depicts the system's power management and clock distribution, including a PLL, DPLL, and various clock sources. The bottom section shows the system's connectivity to external devices like UART, GPIO, I2C, and USB.

## Toward Heterogeneous Architectures (1/2)

The diagram illustrates a neural network architecture. On the left, a blue curved shape represents the input, which feeds into a grid of orange squares representing the hidden layer. From the hidden layer, red lines connect to a single blue square representing the output layer. To the right of the output layer, a horizontal blue arrow points to the right, indicating the output of the network.