

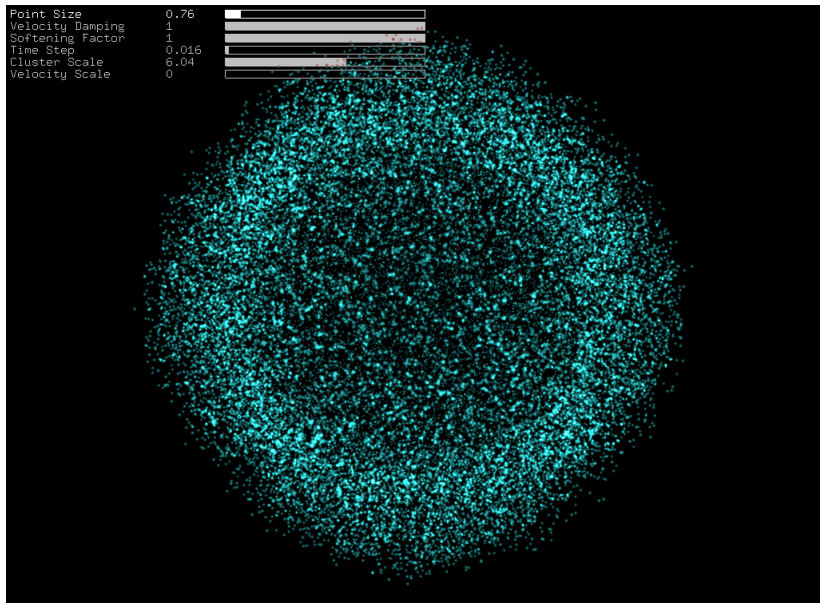
Introduction to GPU programming

Martin Dybdal
dybber@dybber.dk

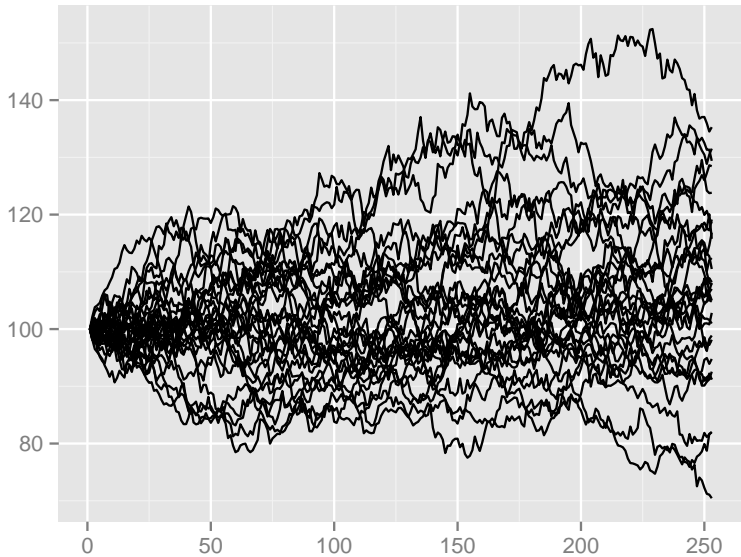
HIPERFIT research center
DIKU
University of Copenhagen

4 March 2016

Physics simulation

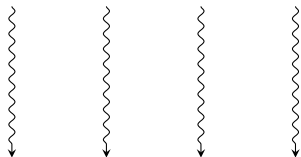


Financial simulation

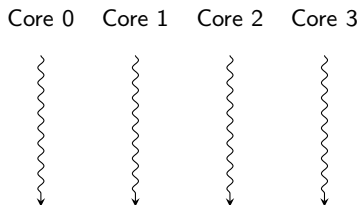


Standard CPU

Core 0 Core 1 Core 2 Core 3

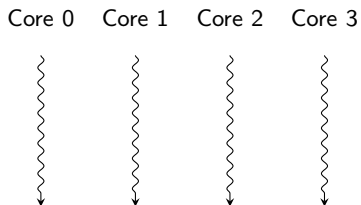


Standard CPU



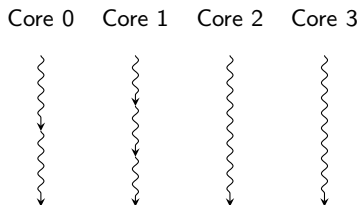
- One operation at a time

Standard CPU



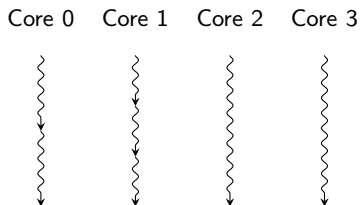
- ▶ One operation at a time
- ▶ Few compute units (cores)

Standard CPU

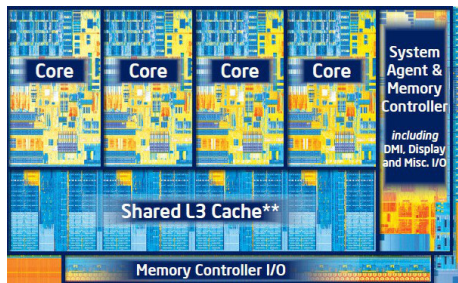


- ▶ One operation at a time
- ▶ Few compute units (cores)
- ▶ Fast at switching between tasks

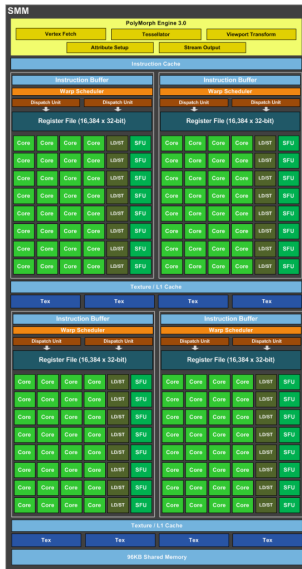
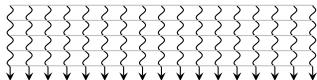
Standard CPU



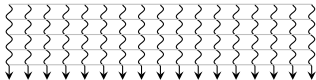
- ▶ One operation at a time
- ▶ Few compute units (cores)
- ▶ Fast at switching between tasks
- ▶ Most transistors used for “recalling”



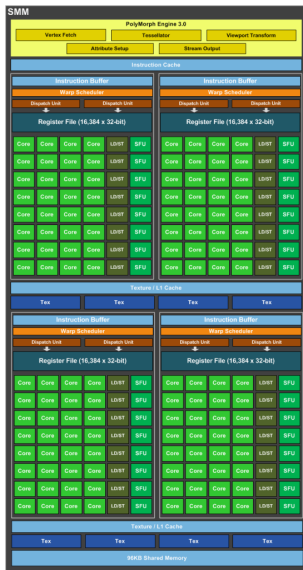
GPU



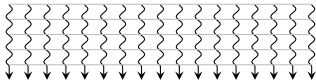
GPU



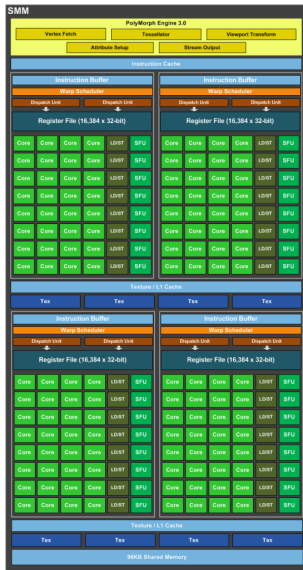
- Identical operations on diff. data



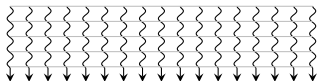
GPU



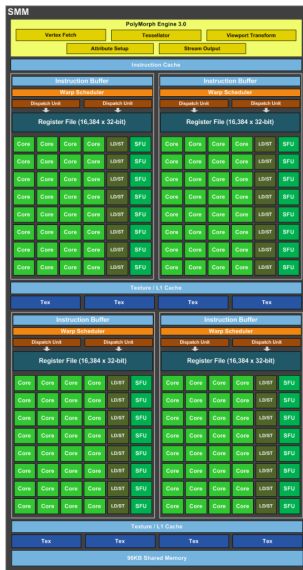
- ▶ Identical operations on diff. data
- ▶ Thousands of compute units (cores)



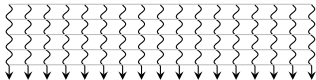
GPU



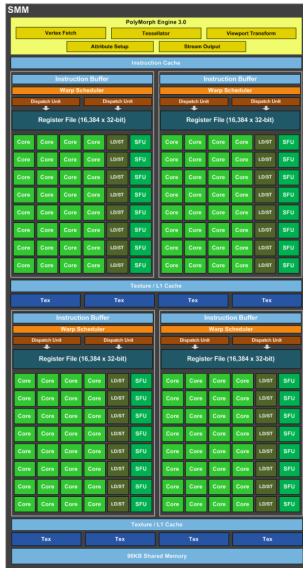
- ▶ Identical operations on diff. data
- ▶ Thousands of compute units (cores)
- ▶ Tasks executed in order (queued)



GPU



- ▶ Identical operations on diff. data
- ▶ Thousands of compute units (cores)
- ▶ Tasks executed in order (queued)
- ▶ Most transistors used for computing



CPU vs. GPU programming

CPU programming

5+9

14

14+3

17

17+22

39

GPU programming

(2 4 6 8 10) + 100

102 104 106 108 110

(102 104 106 108 110) * 2

204 208 212 216 220

GPU programming

Problem:

- ▶ GPU cores are bad at “recalling”
- ▶ manual control of “scratch pad”

Fusion

```
((2 4 6 8 10) + 100) * 2  
204 208 212 216 220
```


Summary

- ▶ GPUs require many similar computations on different data
- ▶ GPUs require attention to memory transactions (fusion)
- ▶ GPU programming: as hard as programming CPUs in the 60s/70s