



CSCI-GA.3033-012

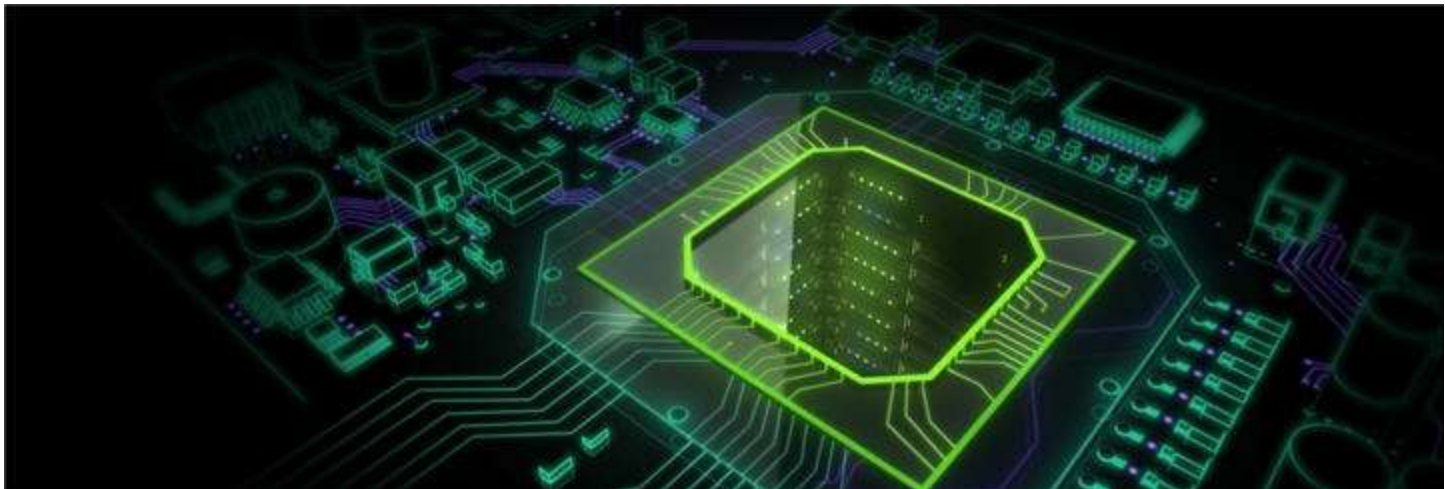
Graphics Processing Units (GPUs): Architecture and Programming

Lecture 10: Heterogeneous Systems

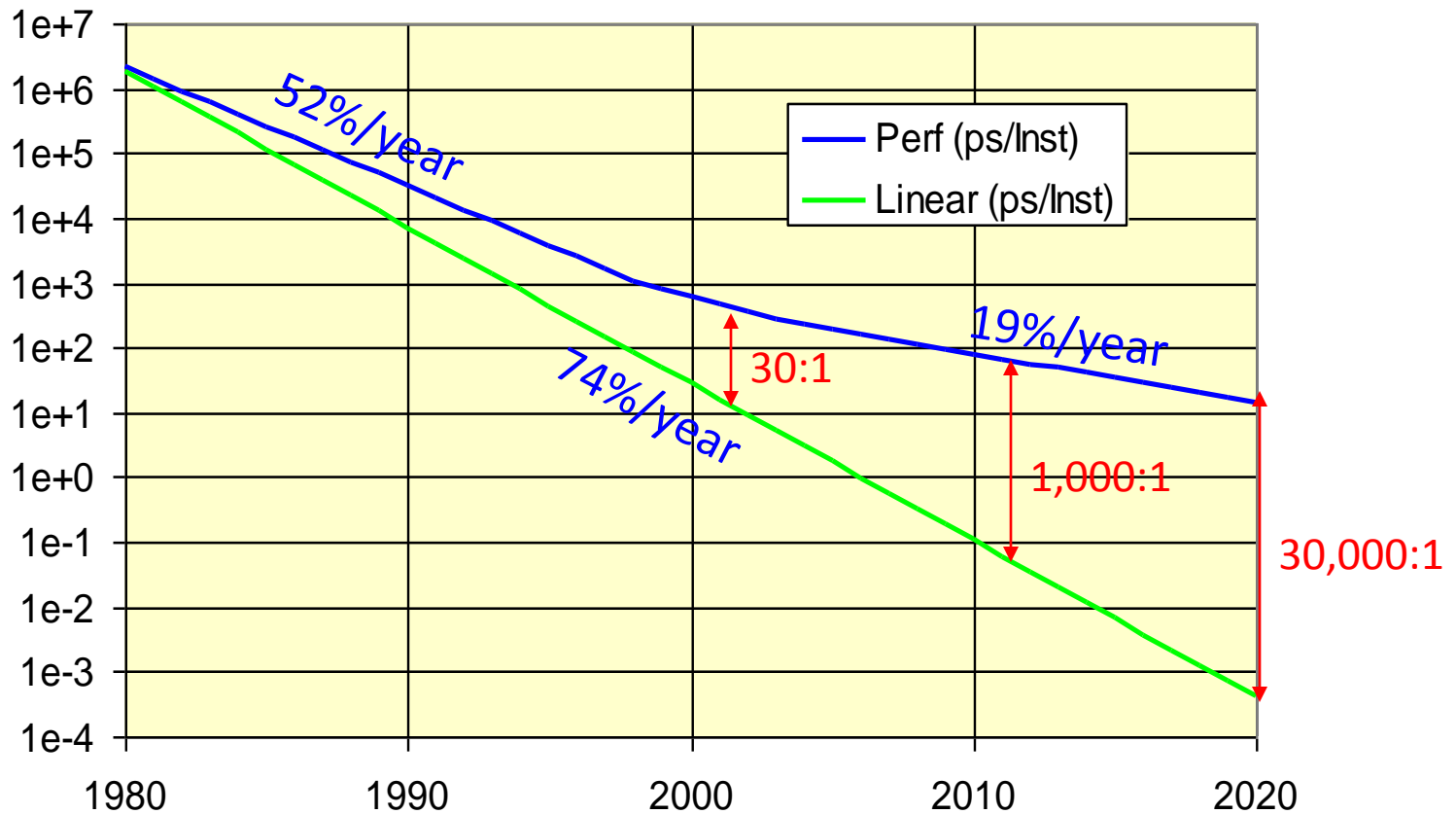
Mohamed Zahran (aka Z)

mzahran@cs.nyu.edu

<http://www.mzahran.com>

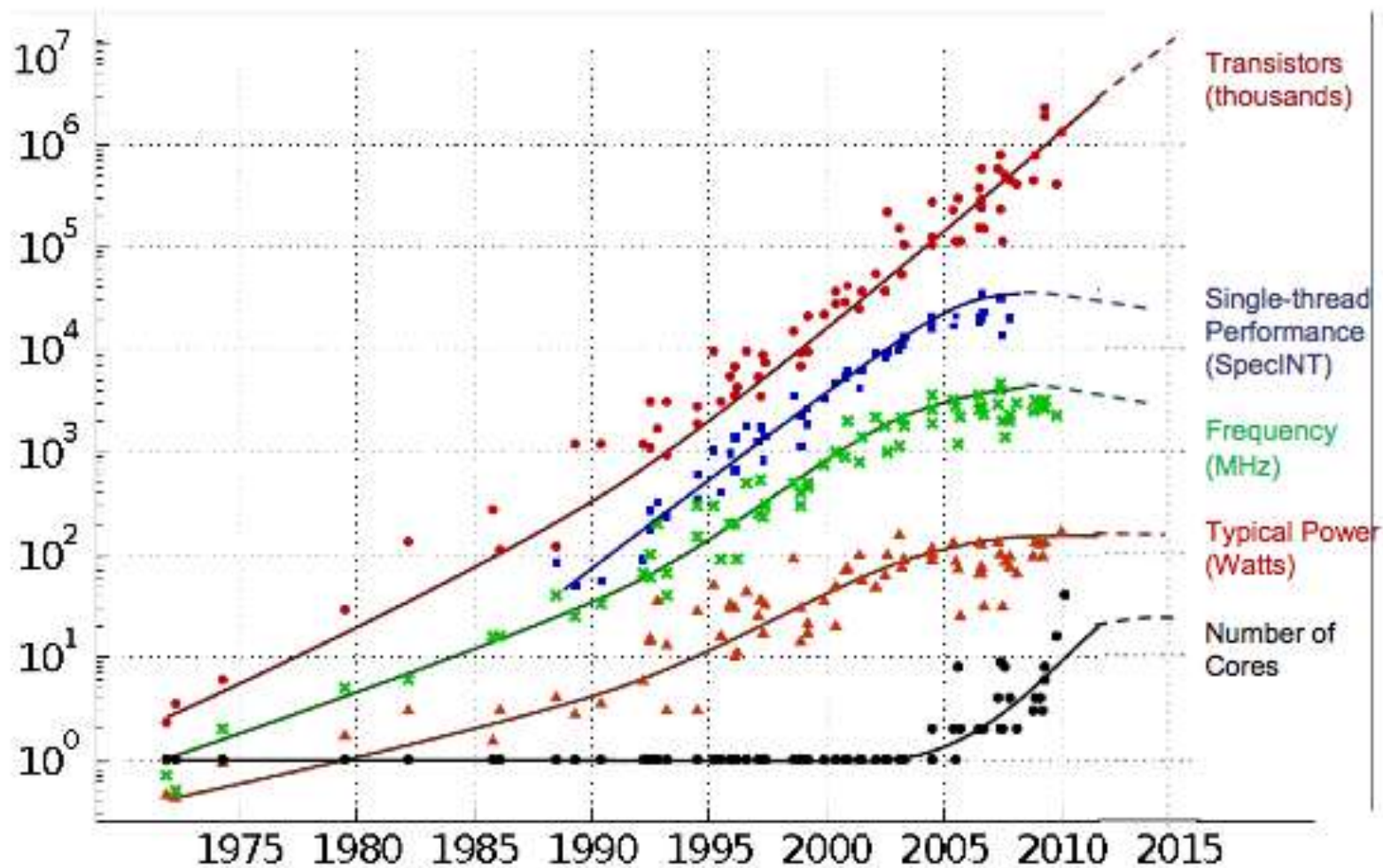


Not Enough ILP

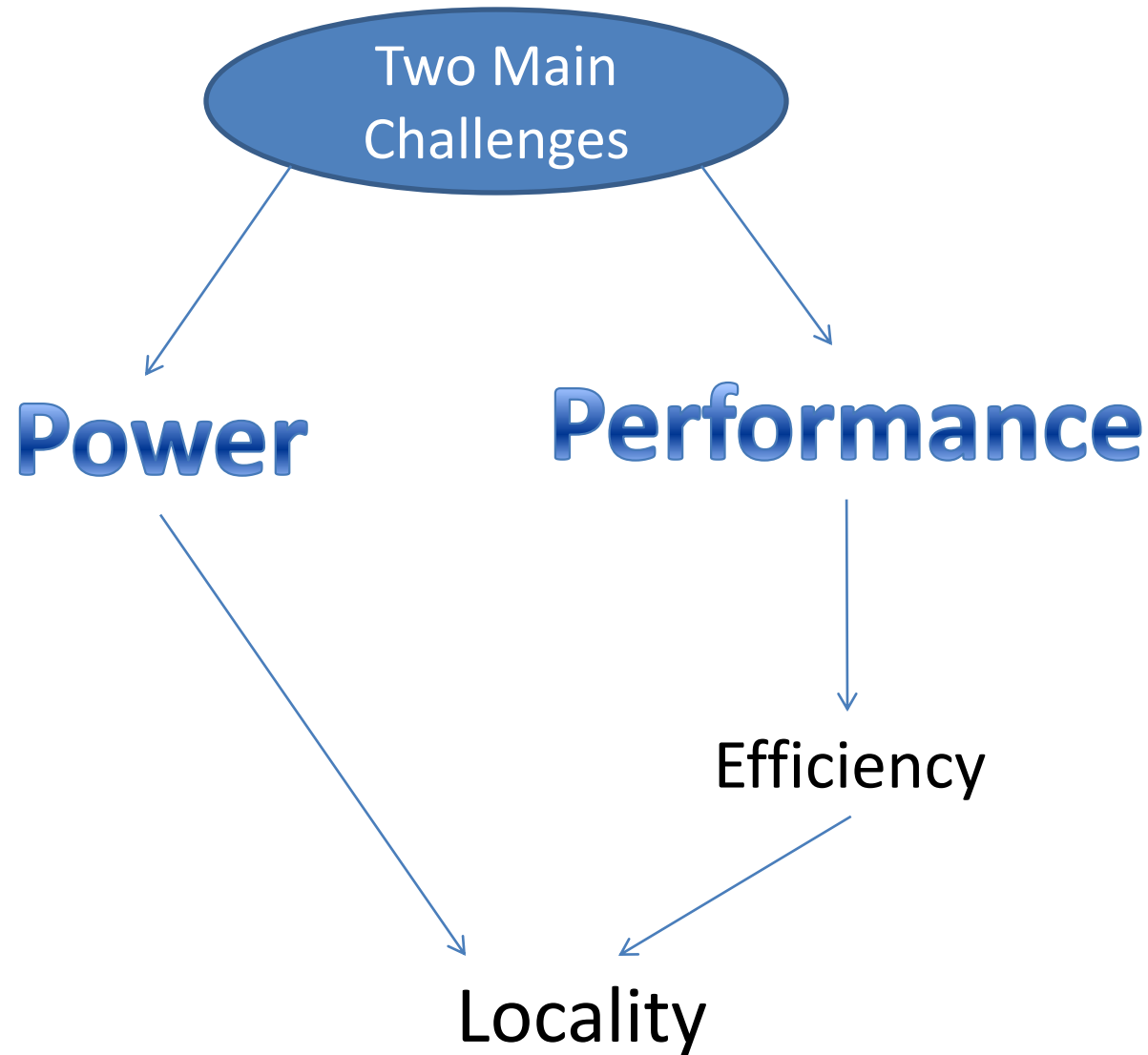


Dally et al. "The Last Classical Computer", ISAT Study, 2001

Historical scaling has ended in 2010



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore



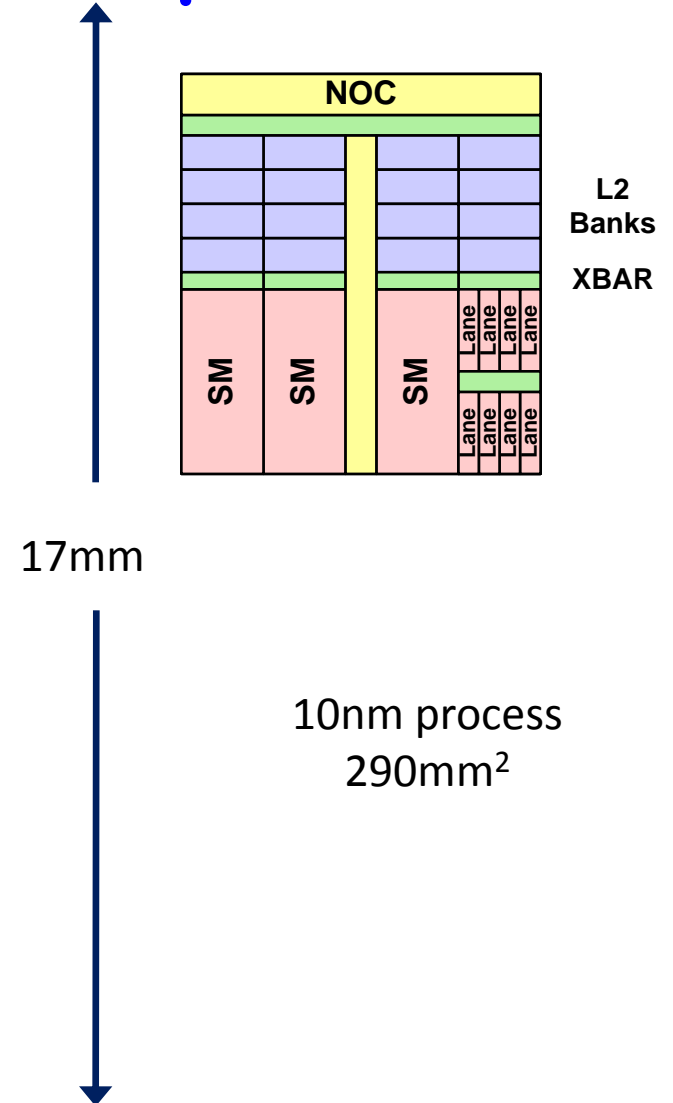
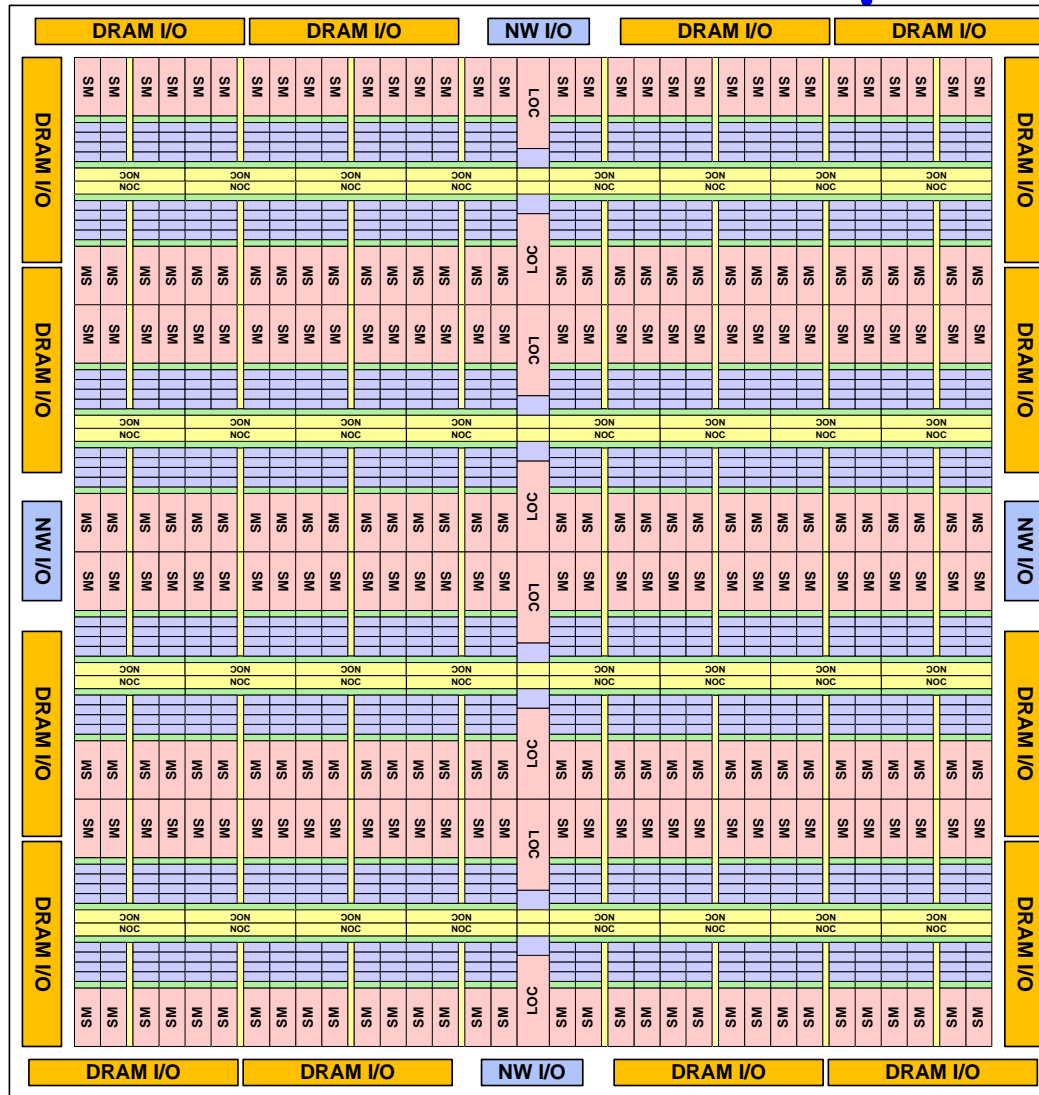
Data movement costs more than computation.

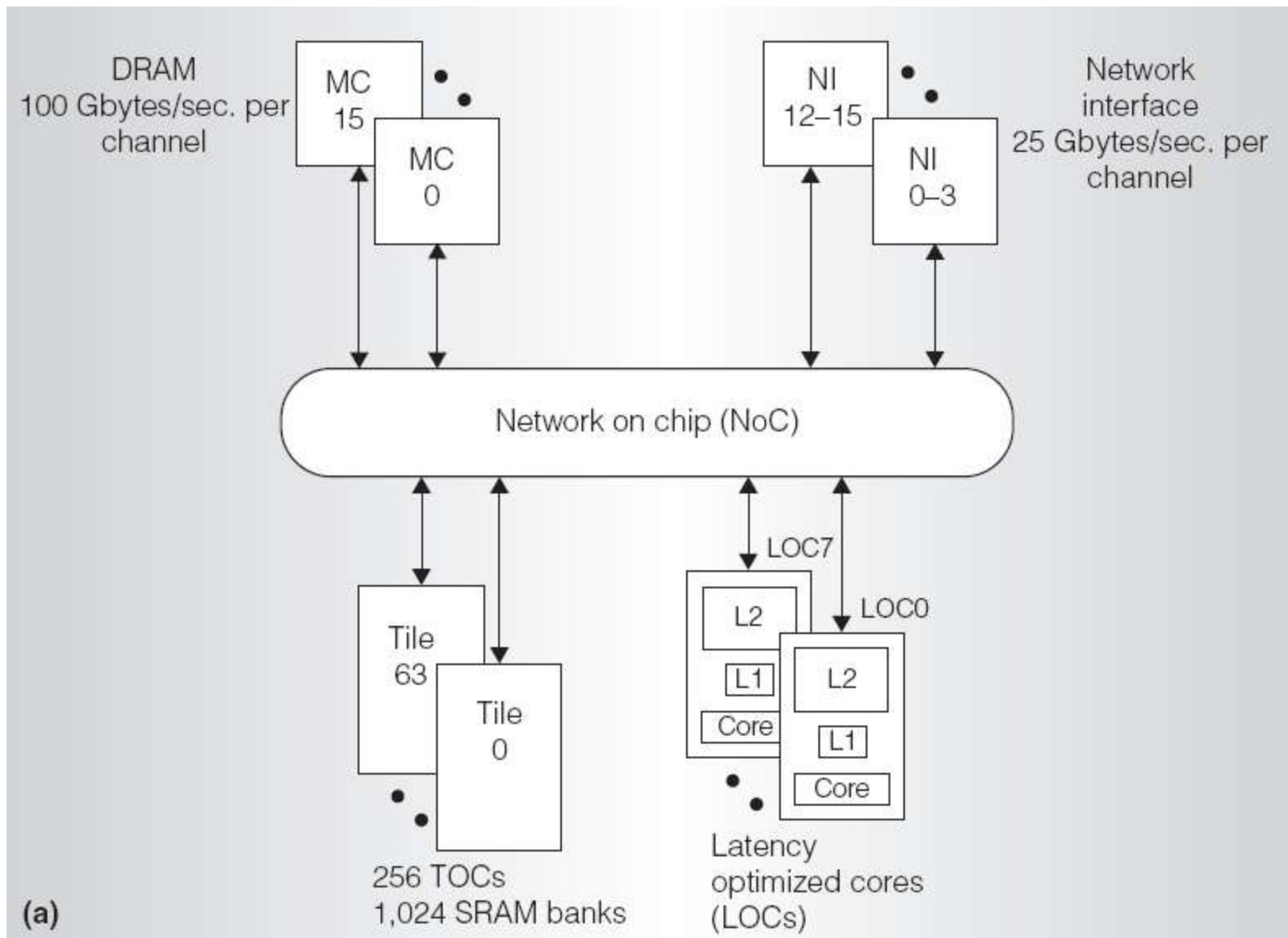
Why Heterogeneous Multicore/Manycore

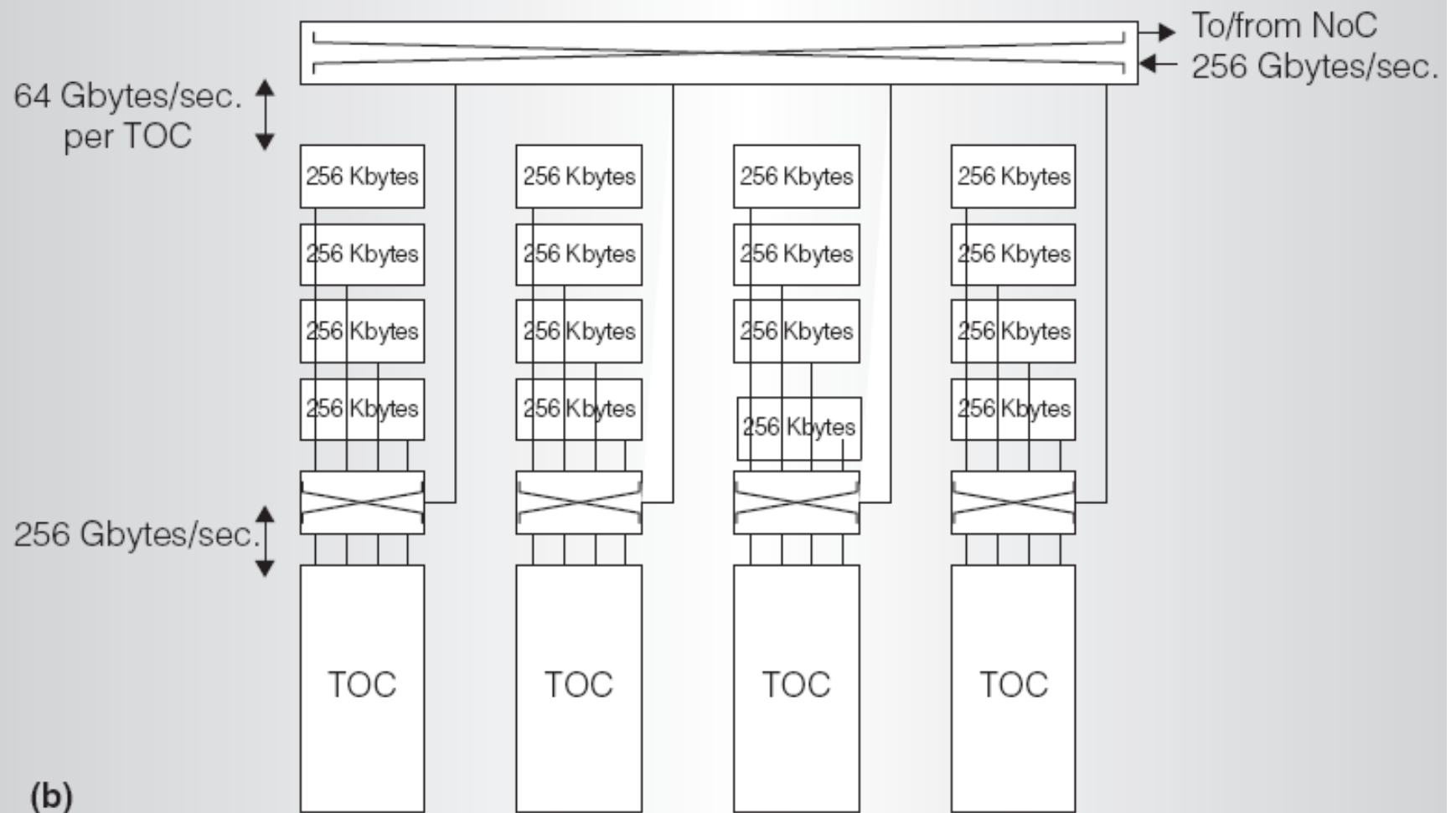
- Wide variety of software applications
- Application with fewer and sophisticated threads -> Traditional multicore with **latency optimize cores**
- Application with high concurrency -> large number of **throughput optimized cores**
- Energy efficiency
- Heterogeneous computing is needed to reach ExaScale computing

NVIDIA: ECHELON

Echelon Chip Floorplan







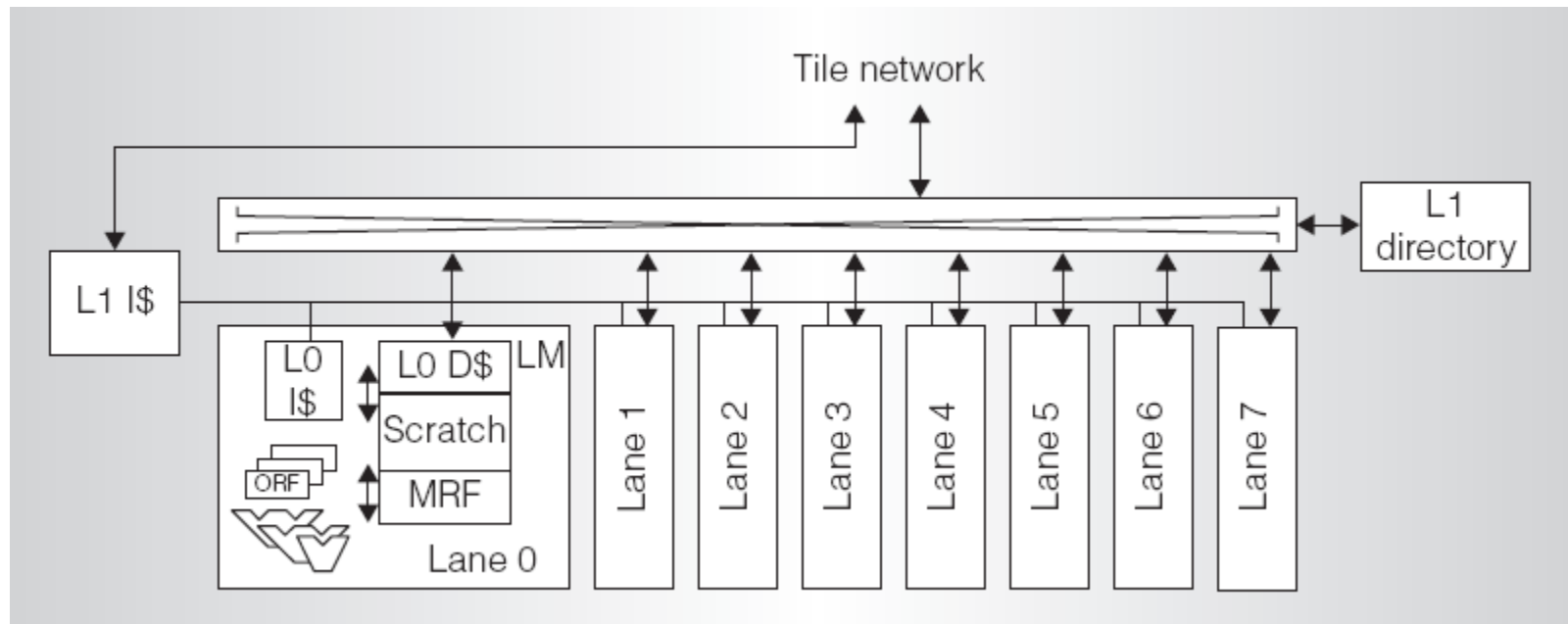
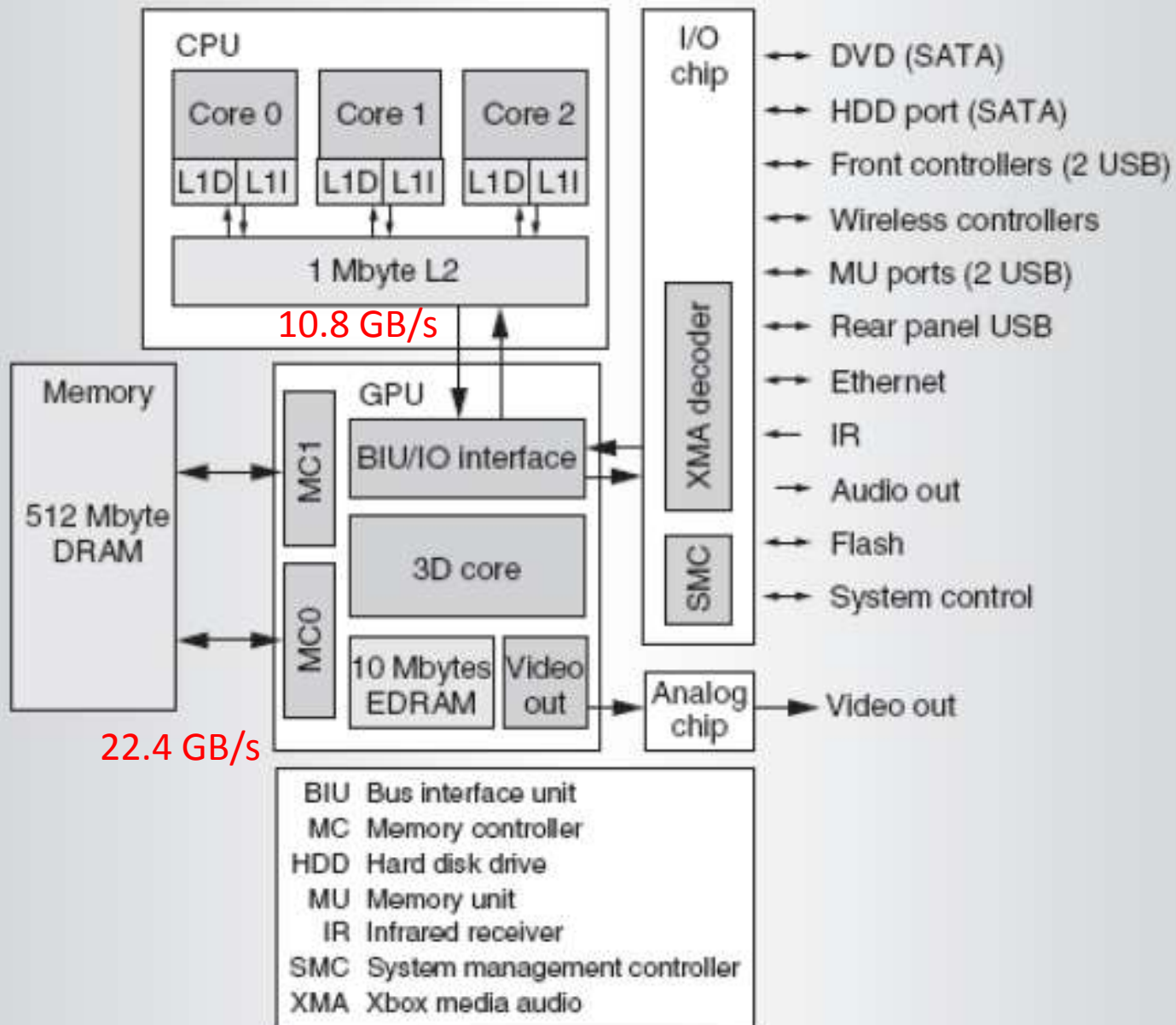


Table 3. Echelon's parameters.

Capability	Per lane	Per TOC	Per tile	Per chip
Instructions per clock	3	24	96	6,144
Instructions per second at 2 GHz	6	48	192	12,288
Double-precision floating-point operations per clock	4	32	128	8,192
Gflops at 2 GHz	8	64	256	16,384
Total threads	64	512	2,048	131,072
Active threads	4	32	128	8,192
Data SRAM (Kbytes)	32	256	5,120	327,680
L0 instruction cache (Kbytes)	1	8	32	2,048
L1 instruction cache (Kbytes)	—	32	128	8,192
DRAM bandwidth (Gbytes/second)	—	—	—	1,600
Network interface bandwidth (Gbytes/second)	—	—	—	400

XBOX 360



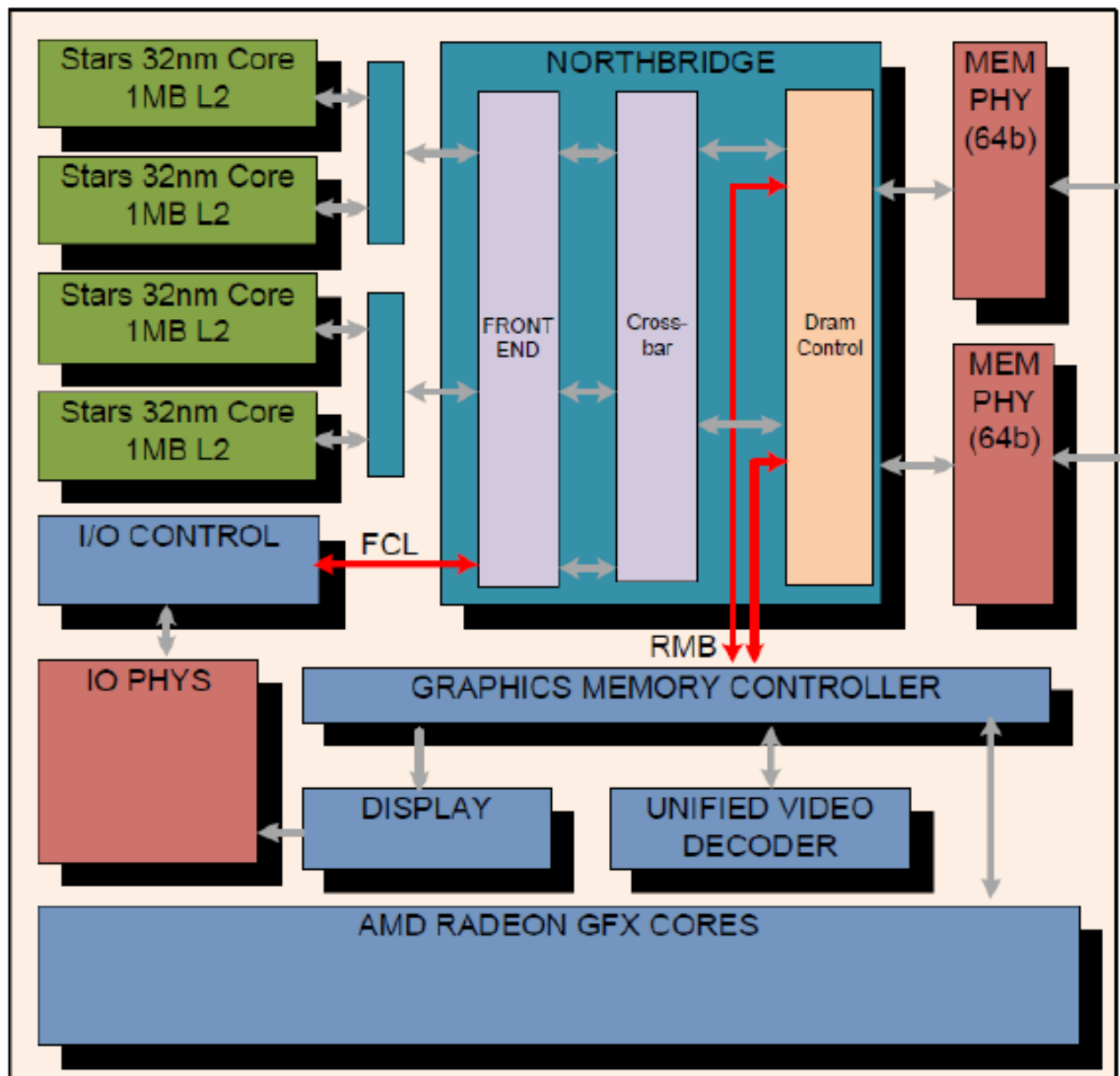
Specifications

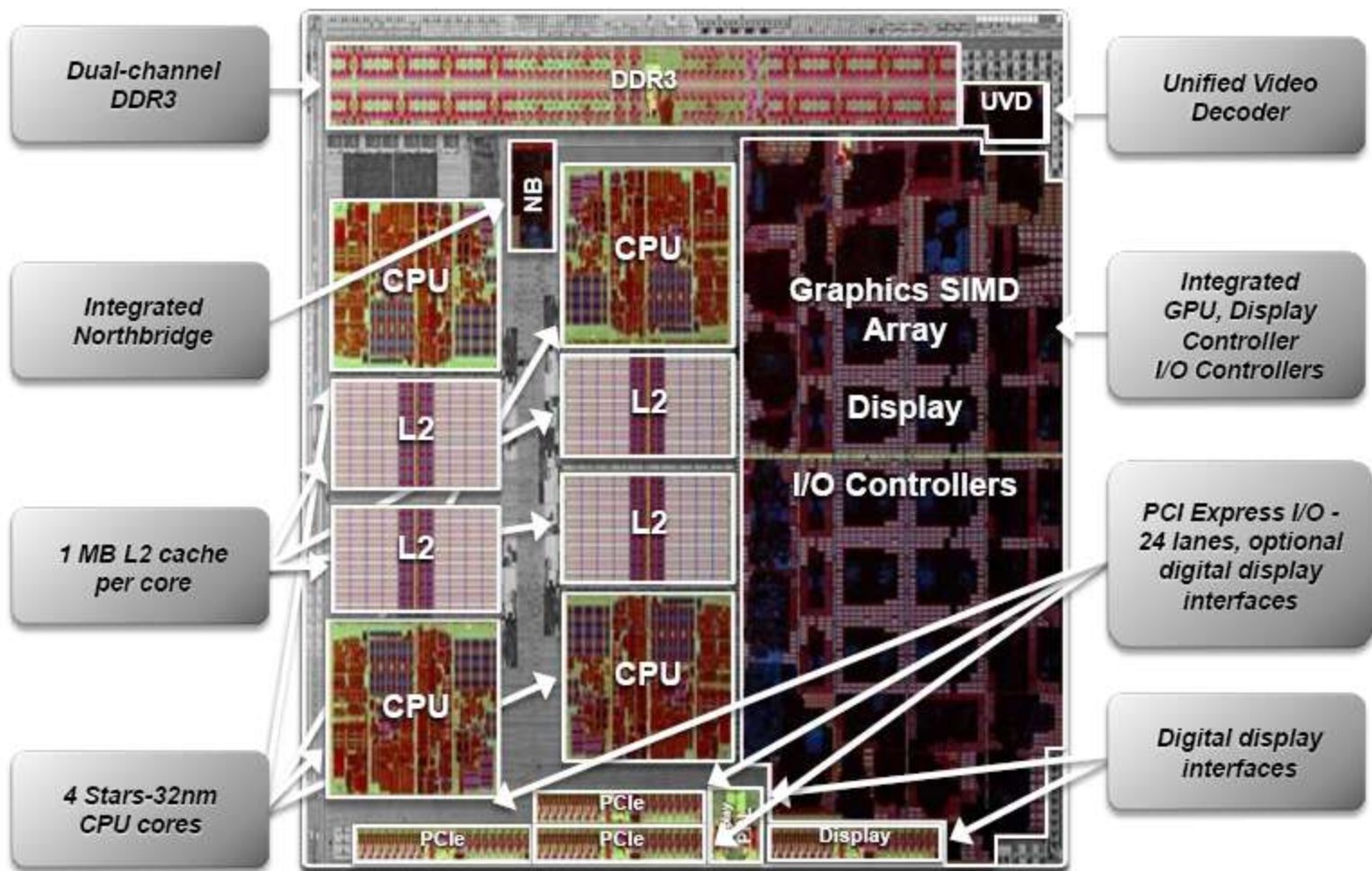
- 3 CPU cores
 - 4-way SIMD vector units
 - 8-way 1MB L2 cache (3.2 GHz)
 - 2 way SMT
 - In-order
 - 2 Instructions/cycle
- ATI GPU with embedded EDRAM
- 3D graphics units
- 512-Mbyte DRAM main memory

Philosophy

- Value for 5-7 years
- Big performance increase over last generation
- Support high-definition video
- extremely high pixel fill rate (goal: 100+ million pixels/s)
- Flexible to suit dynamic range of games
- balance hardware, homogenous resources
- Programmability (easy to program)

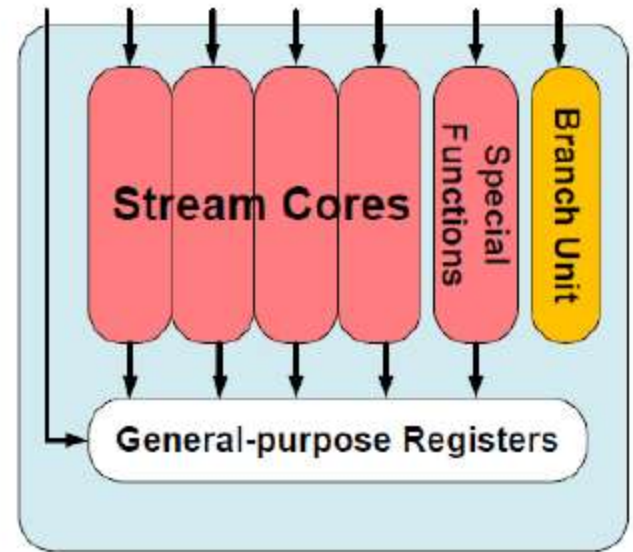
AMD'S "LLANO" FUSION APU





AMD - Radeon GPU VLIW-5

- Includes 4 Stream Cores
 - 1 Special Functions Stream Core
 - Branch Unit
 - General Purpose Registers
- 4 Stream Cores are capable of
 - 4 32-bit FP MULADD per clock
 - 4 24-bit IntMUL or ADD per clock
 - 2 64-bit FP MUL or ADD per clock
 - 1 64-bit FP MULADD per clock
- Additional special function core 32b-FP MULADD per clock

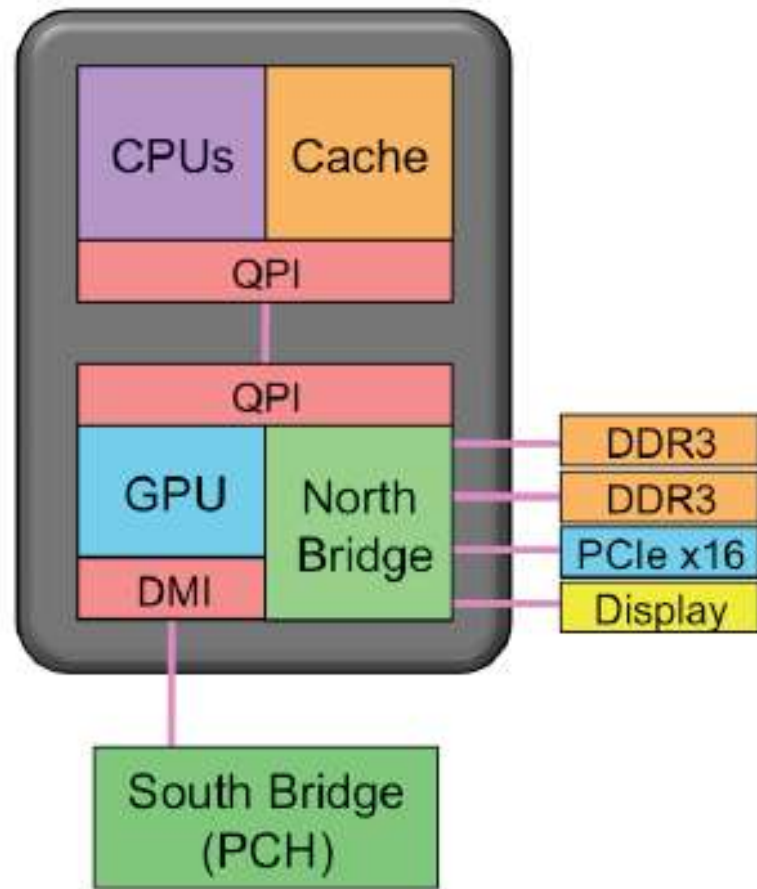


INTEL: SANDY BRIDGE

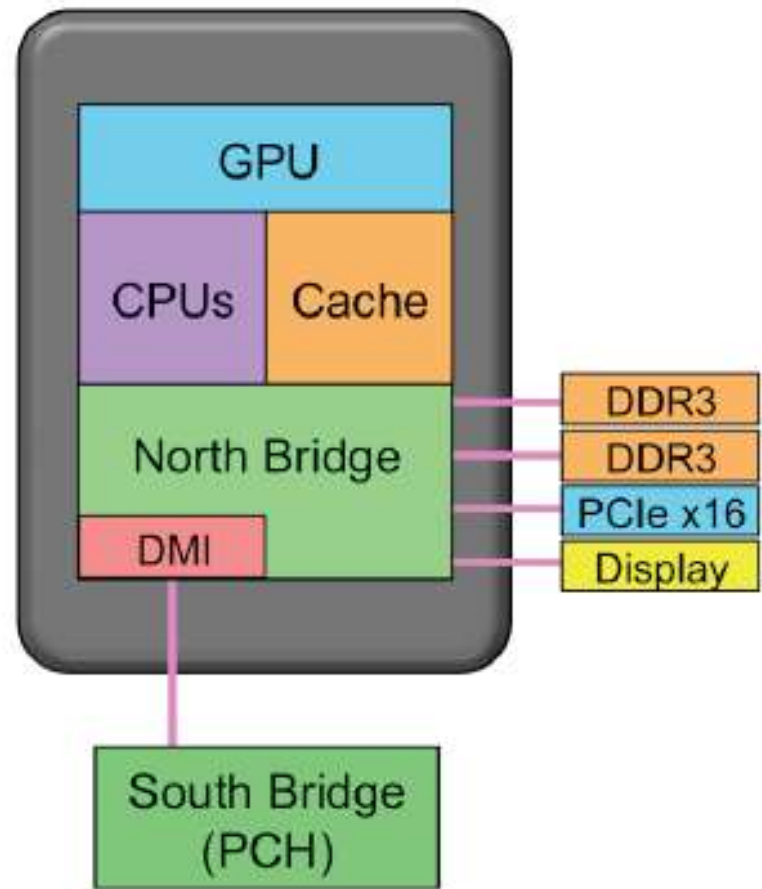
Sandy Bridge

- Intel Tock
 - Intel's first to GPU on the processor itself.
- Improvement over its predecessor Nehalem
- Targeting multimedia applications
 - Introduced Advanced Vector Extensions (AVX)
- More power-efficient than Westmere

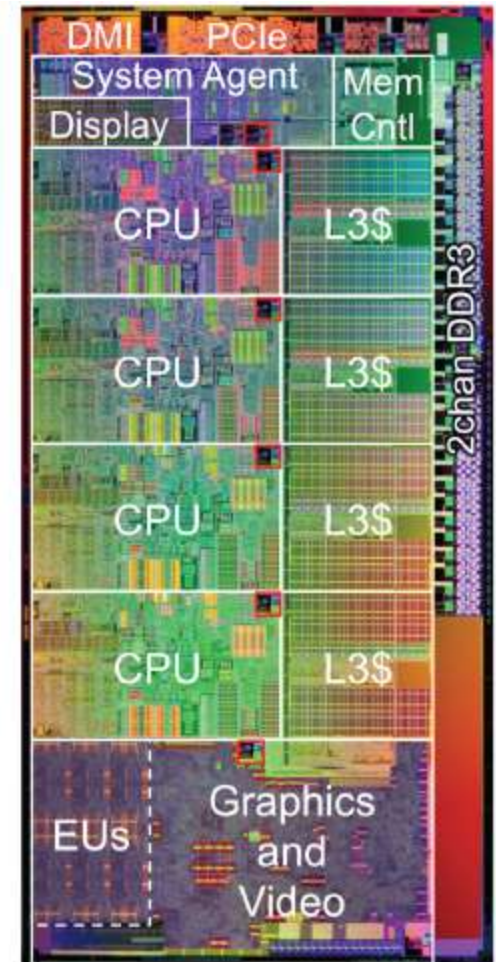
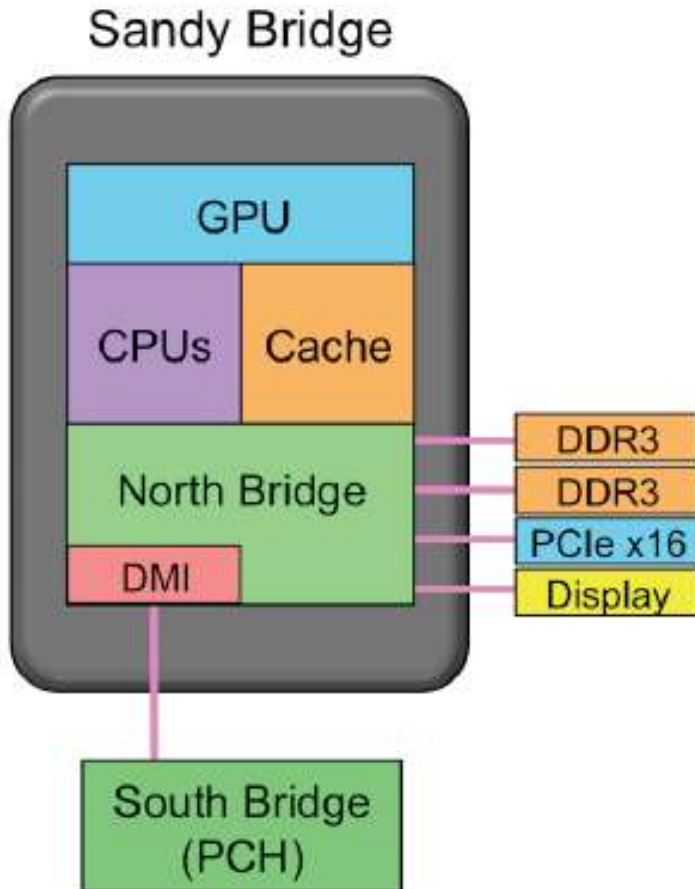
Arrandale



Sandy Bridge



- The GPU can access the large L3 cache
- Intel's team totally re-designed the GPU



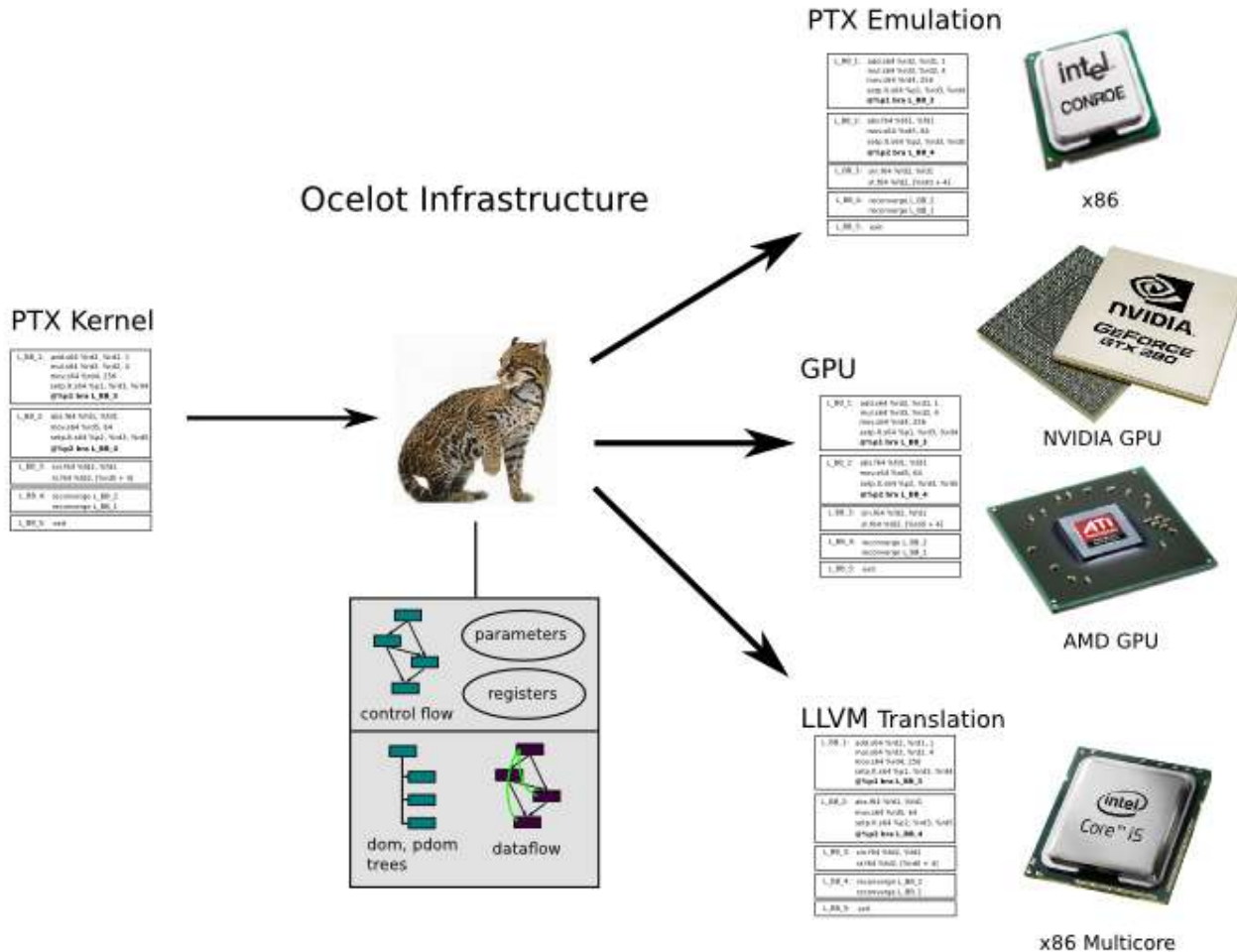
What is/are the main design philosophy/philosophies of all those heterogeneous systems?

Which programs benefit the most from them?

Programming Heterogeneous Systems

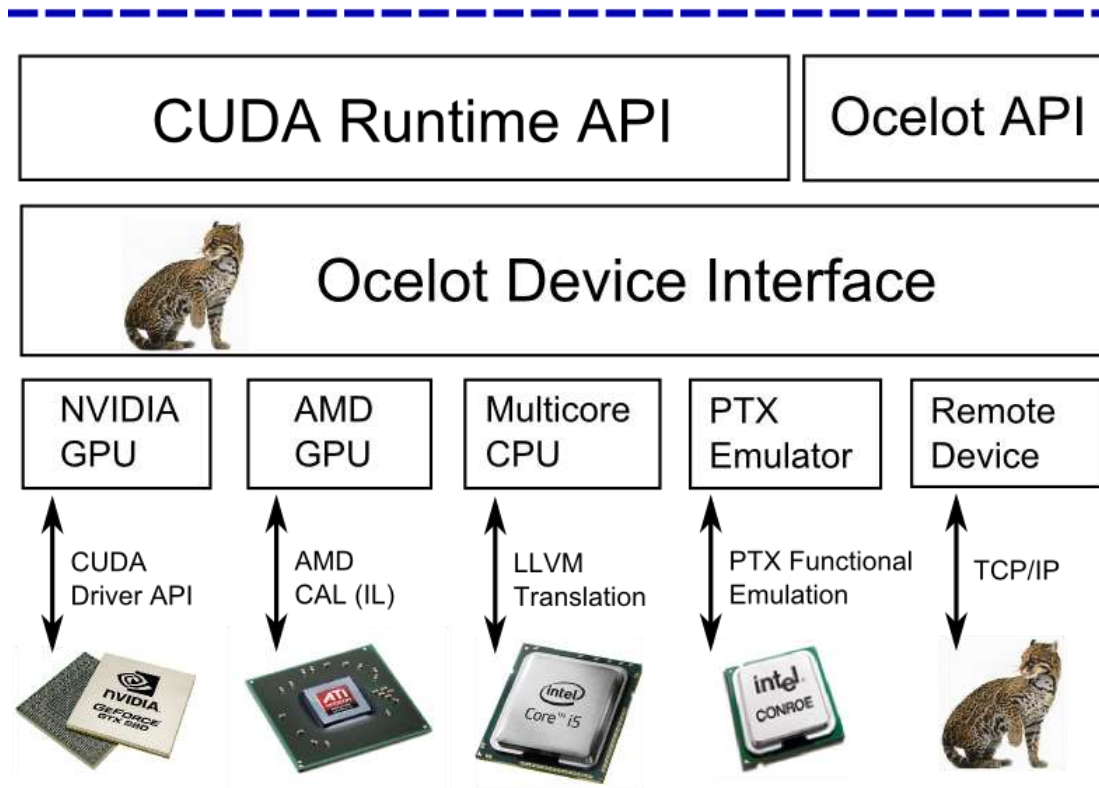
- Who's doing what?
 - programmer
 - OS
 - Compiler
- SSE/AVX should be appropriate if the process is suitable for SIMD parallelization and not for GPU
- Several choices:
 - OpenMP
 - pthreads
 - CUDA
 - OpenCL
 - Intel threading Building Block (TBB)

Useful Tool: gpuocelot



Useful Tool: gpuocelot

CUDA Application



Conclusions

- It is all about performance and power
- Algorithms should be designed to perform more work per unit data movement.
 - Data movement dominates power.
- Issues:
 - How do GPU and CPU share memory?
 - Bandwidth (on-chip and off-chip)
 - Coherence?
 - Portability of code?