

GPU ARCHITECTURE

FOR SLEEP DEPRIVED COLLEGE STUDENTS...

AND ALSO DUMMIES

JEFFREY GRANGE

A GP WHAT? WHAZZAT?

A GPU (GRAPHICS PROCESSING UNIT) IS USED TO ACCELERATE GRAPHICS PROCESSES IN A COMPUTER

THE EARLIEST PREDECESSORS APPEARED IN THE 1970s, THEY HELPED, BUT STILL RELIED ON THE CPU FOR CALCULATIONS

IN THE 1990s, THE FIRST ACTUAL GPUS APPEARED, WHICH HAD THEIR OWN HARDWARE FOR PROCESSING



HOW GPU BETTER?

GPUS ACHIEVE HIGHER THROUGHPUT BY SACRIFICING CONTROL AND SINGLE INSTRUCTION LATENCY

THERE'S A LOT OF SIMILAR CALCULATIONS IN GRAPHICS, SO YOU CAN GROUP THEM TOGETHER AND PERFORM ONE CALCULATION ON THE WHOLE GROUP

THIS IS CALLED SIMD (SINGLE INSTRUCTION MULTIPLE DATA)



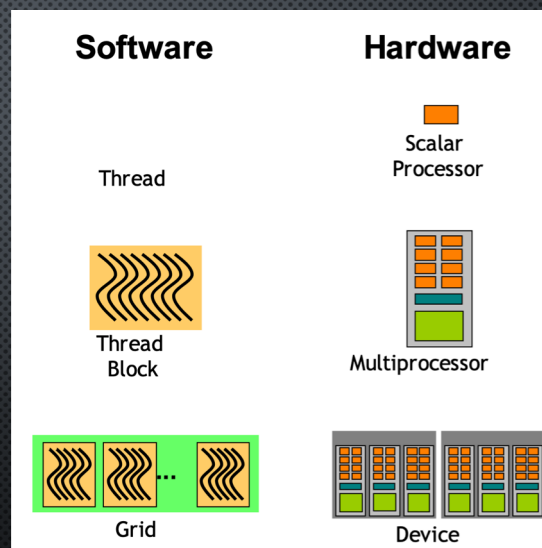
DATA (CUDA 2.0)

THREADS: DATA AND A SERIES OF OPERATIONS DONE ON IT

WARPS: SETS OF 32 SIMILAR THREADS RUN ON A STREAMING MULTIPROCESSOR

BLOCKS: A WHOLE BUNCH OF THREADS, MULTIPLE OF 32

GRID: WHOLE BUNCH OF BLOCKS



STREAMING MULTIPROCESSOR

HAS 32 CUDA CORES, WHICH CAN DO FLOATING POINT OR INTEGER CALCULATIONS

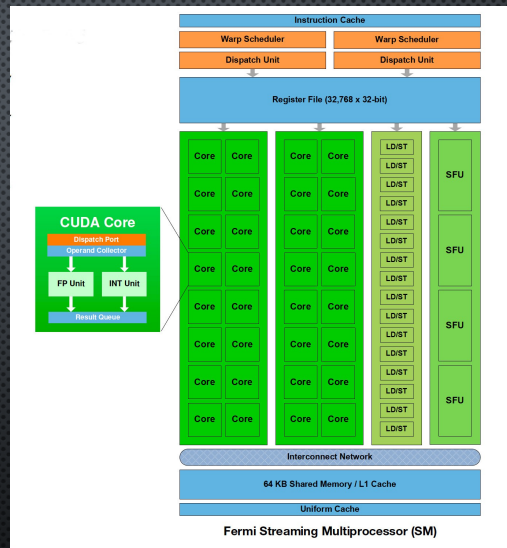
A WHOLE BUNCH OF REGISTERS TO HOLD BLOCKS AND WARPS

2 WARP SCHEDULERS, EACH FOR HALF THE CORES

4 SPECIAL FUNCTION UNITS, DO SINE, COSINE...ETC

L1 CACHE

UNIFORM CACHE (HOLDS DUPLICATE OF REQUIRED ROM CONTENTS)



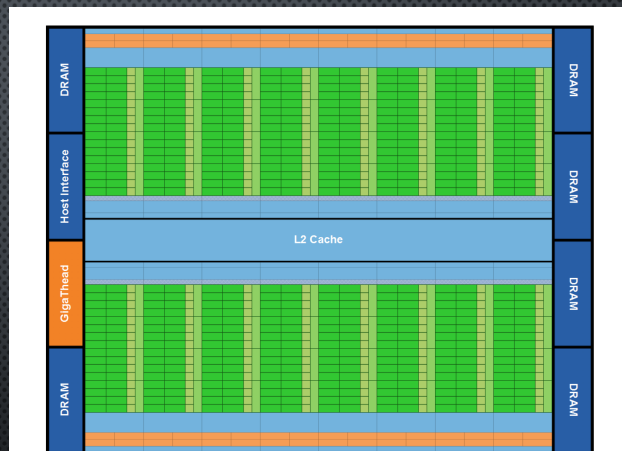
GPU - ALMOST

BUNCH OF STREAMING MULTIPROCESSORS

L2 CACHE

ACCESS TO VRAM

ACCESS TO OTHER STUFF



Fermi's 16 SM are positioned around a common L2 cache. Each SM is a vertical rectangular strip that contain an orange portion (scheduler and dispatch), a green portion (execution units), and light blue portions (register file and L1 cache).

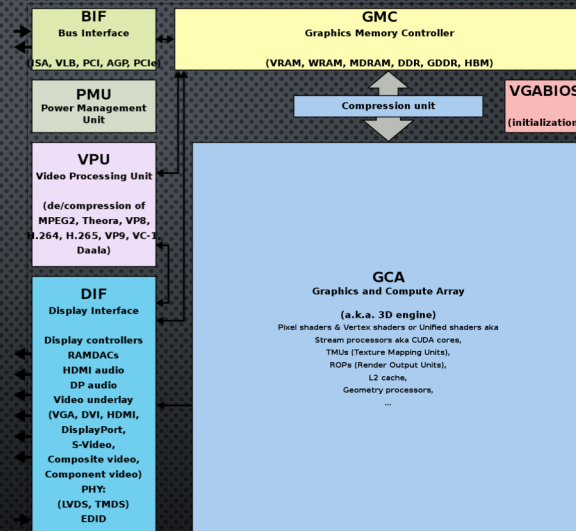
THE OTHER STUFF

GCA: EVERYTHING WE JUST COVERED

VPU: USED FOR COMPRESSING AND DECOMPRESSING DIFFERENT VIDEO STORAGE FORMATS

DISPLAY CONTROLLER: CONVERTS VRAM DATA TO SOMETHING INTERPRETABLE BY SCREEN AND OUTPUTS

MEMORY CONTROLLER: CONTROLS MEMORY...



REFERENCES:

- [HTTP://WWW.ICL.UTK.EDU/~LUSZCZEK/TEACHING/COURSES/FALL2016/COSC462/PDF/GPU_FUNDAMENTALS.PDF](http://www.icl.utk.edu/~LUSZCZEK/TEACHING/COURSES/FALL2016/COSC462/PDF/GPU_FUNDAMENTALS.PDF)
- [HTTP://DEVELOPER.DOWNLOAD.NVIDIA.COM/COMPUTE/CUDA/3_1/TOOLKIT/DOCS/NVIDIA_CUDA_C_PROGRAMMINGGUIDE_3.1.PDF](http://developer.download.nvidia.com/compute/cuda/3.1/toolkit/docs/nvidia_cuda_c_programmingguide_3.1.pdf)
- [HTTP://MESEEC.CE.RIT.EDU/551-PROJECTS/SPRING2015/3-2.PDF](http://measec.ce.rit.edu/551-PROJECTS/SPRING2015/3-2.PDF)
- [HTTPS://EN.WIKIPEDIA.ORG/WIKI/GRAPHICS_PROCESSING_UNIT](https://en.wikipedia.org/wiki/Graphics_processing_unit)
- [HTTPS://EN.WIKIPEDIA.ORG/WIKI/VIDEO_DISPLAY_CONTROLLER](https://en.wikipedia.org/wiki/Video_display_controller)
- [HTTPS://MEDIUM.COM/@SMALLFISHBIGSEA/BASIC-CONCEPTS-IN-GPU-COMPUTING-3388710E9239](https://medium.com/@smallfishbigsea/basic-concepts-in-gpu-computing-3388710e9239)
- [HTTPS://EN.WIKIPEDIA.ORG/WIKI/CUDA](https://en.wikipedia.org/wiki/CUDA)

MORE REFERENCES:

- [HTTPS://WWW.QUORA.COM/WHAT-IS-A-WARP-AND-HOW-IS-IT-DIFFERENT-FROM-A-THREAD-BLOCK-OR-WAVE-IN-CUDA](https://www.quora.com/What-is-a-Warp-and-how-is-it-different-from-a-thread-block-or-wave-in-CUDA)
- [HTTPS://ELECTRONICS.STACKEXCHANGE.COM/QUESTIONS/102695/HOW-DOES-A-GPU-CPU-COMMUNICATE-WITH-A-STANDARD-DISPLAY-OUTPUT-HDMI-DVI-ETC](https://electronics.stackexchange.com/questions/102695/how-does-a-gpu-cpu-communicate-with-a-standard-display-output-hdmi-dvi-etc)
- [HTTPS://EN.WIKIPEDIA.ORG/WIKI/VIDEO_CODEC](https://en.wikipedia.org/wiki/Video_codec)
- [HTTPS://ITSTILLWORKS.COM/VIDEO-CONTROLLER-1731.HTML](https://itstillworks.com/video-controller-1731.html)
- [HTTPS://WWW.NVIDIA.COM/EN-US/GEFORCE/20-SERIES/](https://www.nvidia.com/en-us/geforce/20-series/)