



The image shows two graphics cards. The top one is an NVIDIA GeForce RTX 2080 Ti with two fans and green accents. The bottom one is an AMD Radeon RX 5700 XT with a single large gold fan. The background is a blue gradient with circuit patterns.

THE GPU

- Very parallelized
- Very fast
- Let's go gamers 🕶️



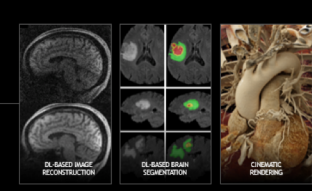
GPGPU APPLICATIONS

- Large matrix/vector operations (BLAS)
- Protein Folding/Molecular Dynamics
- Fast Fourier Transform
- Ray Tracing
- Physics Simulation
- Speech Recognition
- Databases
- Sort/Search
- Sequence matching
- And much more.

• Medical imaging →

CLARA – MEDICAL IMAGING SUPERCOMPUTER

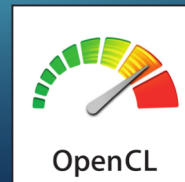
- GPU ACCELERATED KUBERNETES
- IMAGING & VISUALIZATION APPS
- CUDA | CUDNN | TENSORRT | OGL | RTX
- GPU CONTAINERS | VGPU
- NVIDIA GPU SERVER



LEFT TO RIGHT: <https://www.nvidia.com/en-us/deep-learning-ai/>, <https://www.nvidia.com/en-us/deep-learning-ai/>, <https://www.nvidia.com/en-us/deep-learning-ai/>

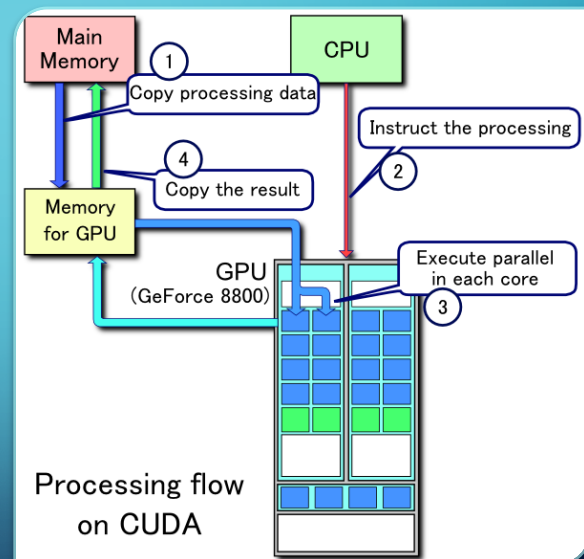
IS IT POSSIBLE TO LEARN THIS POWER?

- Not from AMD.
- The old days – OpenGL and DirectX
 - Had to convert data into graphical terms.
- The new hope – CUDA, OpenCL, MS DirectCompute



CUDA (NOT A SELLOUT BTW)

- Compute Unified Device Architecture – 2007
- Compatible with C, C++, and Fortran.
- Give direct access to the GPU's (if CUDA-enabled) virtual instruction set and parallel computational elements.
- Gone are the days of data conversion to graphics.



MY GPGPU EXPERIENCE - TL;DR

- OBJ: Move a nonlinear MRI data fitting process from the CPU to the GPU.
- Where? Loma Linda University Medical Center
- How? CUDA.
- Why? Because there are many data sets (30-100+), each taking up to 30 hours by CPU.
- Result: It got very fast.

	CPU Time (sec)	GPU Time (sec)	Approximate Acceleration
DCE – Patlak Model	124,645	20.7	6014x
DCE – Two Compartment Exchange Model (2CXM)	112,817	138.3	816x
Dixon – Three Parameter Fitting	417.7	43.5	9.6x

GPU VS CPU

CPU

- Fast caches (optimal data reuse)
- Fine branching granularity
- Lots of different processes/thread
- High performance on a single thread of execution

Conclusion: CPUs are optimal for **task** parallelism
GPUs are optimal for **data** parallelism

GPU

- Lots of math units; ALU-heavy
- Fast access to onboard memory
- Run a program on each fragment/vertex
- High throughput on parallel tasks

SOURCES

- <https://blogs.nvidia.com/blog/2018/03/28/ai-healthcare-gtc/>
- https://en.wikipedia.org/wiki/General-purpose_computing_on_graphics_processing_units
- <https://en.wikipedia.org/wiki/CUDA>
- <https://courses.cs.washington.edu/courses/cse471/13sp/lectures/GPUsStudents.pdf>
- <https://superuser.com/questions/308771/why-are-we-still-using-cpus-instead-of-gpus>
- https://graphics.stanford.edu/~mhouston/public_talks/R520-mhouston.pdf
- <https://www.tacc.utexas.edu/documents/13601/88790/8Things.pdf>
- <https://submissions2.mirasmart.com/ISMRM2020/ViewSubmissionPublic.aspx?sei=ZluGiDQgQ>

ANY QUESH'?