

The Graphics Processing Unit (GPU) revolution

Ramu Anandakrishnan

Outline

2

- The need for parallel processing
- Basic parallel processing concepts
- The GPU – a massively parallel processor

Lecture 1

-
- Overview of GPU hardware architecture
 - Introduction to GPU programming
 - Performance considerations
 - Homework assignment

Lecture 2

A typical scientific computing problem

3

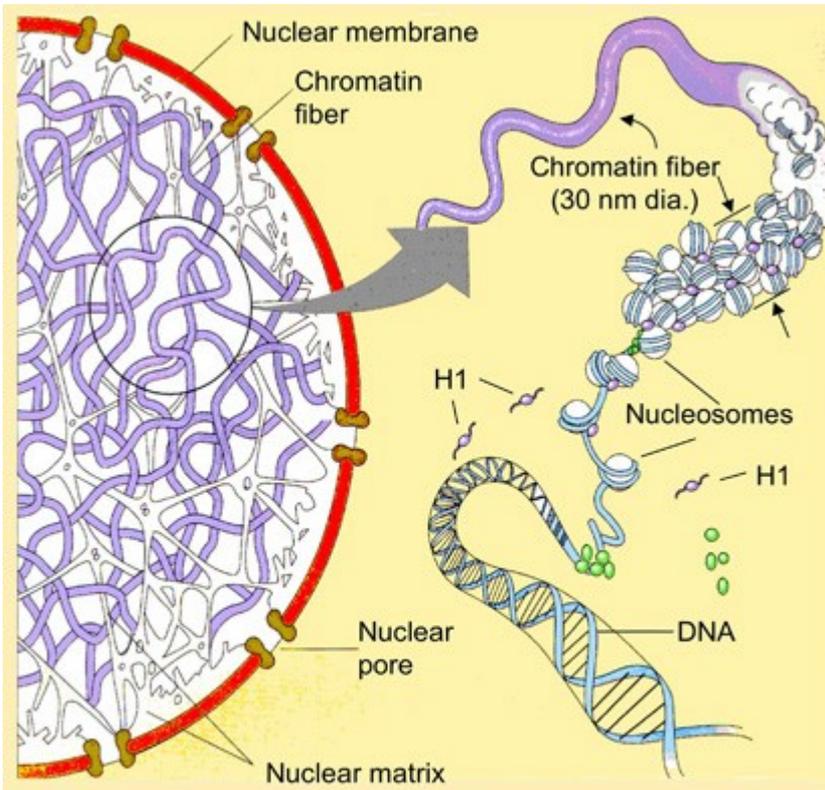
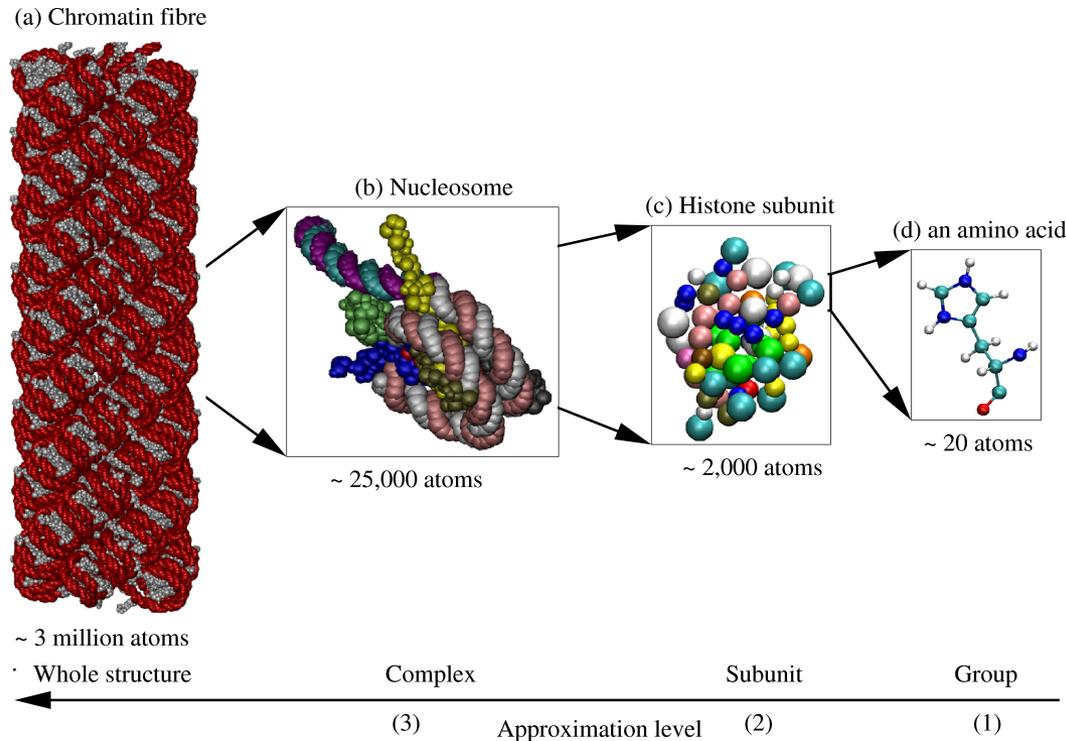


Fig. 1. Modifications of the histone components of nucleosomes help regulate DNA accessibility by promoting folding or unfolding of chromatin fibers, and by recruiting chromatin remodeling complexes and other factors to specific genomic loci.

- The DNA in the human genome is made up of ~300 billion atoms.
- How is this 3 m long DNA strand packed into the 3 μm nucleus of a cell?
- In this tightly packed environment how is a specific gene transcribed, as and when needed?
- Experimental tools are not yet sophisticated enough to answer such questions.

Molecular dynamics simulations can be used to study these systems

4



$$F_i = F_i^{elec} + F_i^{bond} + F_i^{vdw} = m_i a_i$$

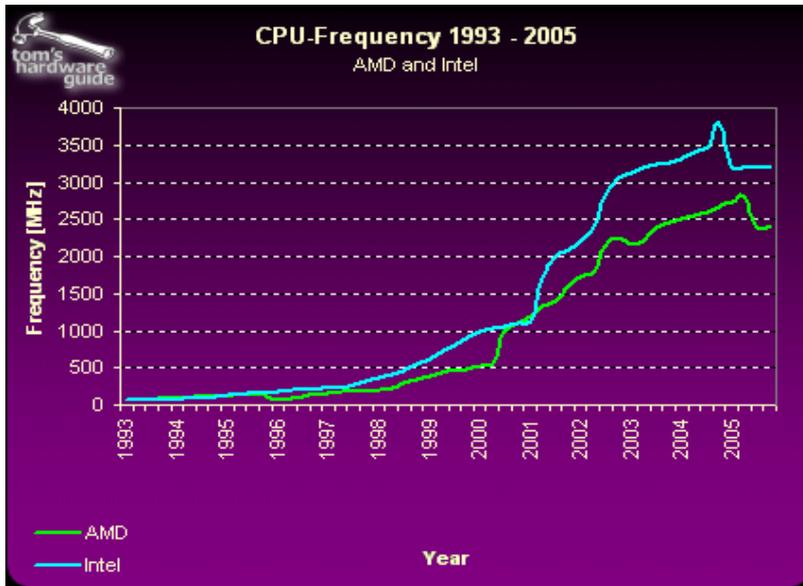
$$r_i(t + \Delta t) = r_i(t) + v_i(t + \Delta t) + \frac{1}{2} a_i t^2$$

$$F_i^{elec} = \frac{q_i}{4\pi\epsilon_0} \sum_{j=1}^N \frac{q_j}{r_{ij}^2}$$

- $N = 3 \times 10^6$ atoms
- $T = 10^{-15}$ sec per step (simulation time)
- $F = 10^{12}$ FLOPS (teraflop)
- Time simulated / year = $(T / N^2) \times F \times 3600 \times 8 \times 365$
- ~ 1 nanosecond
- Not long enough to “observe” meaningful activity

The “power wall” limits how much faster individual processors can run

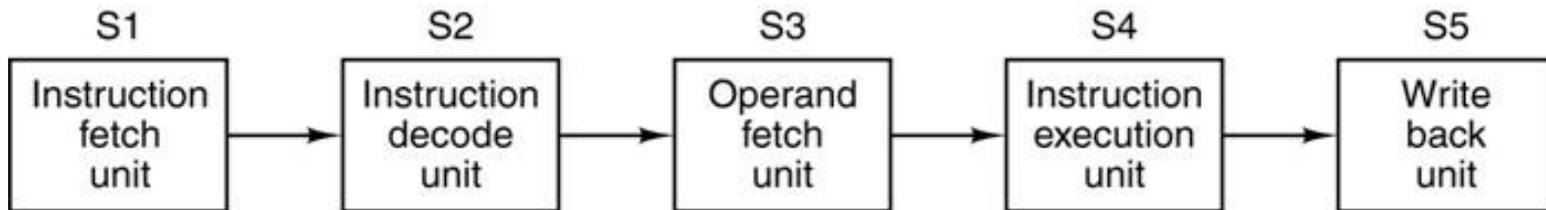
5



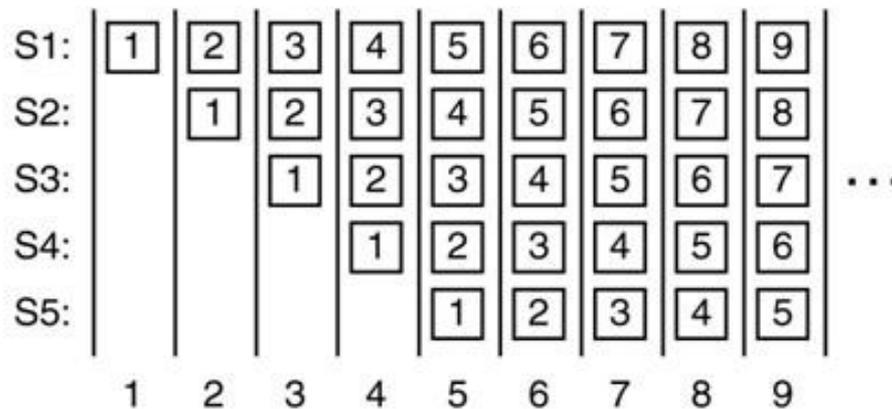
- In 1965 Intel founder Gordon Moore predicted that the number of components in a processor would double every year.
- To get a corresponding increase in computing power (FLOPS), you need to increase clock speed
- But the resulting increase in heat dissipation is limiting how fast you can run the processors - the power wall.
- The answer: parallel processing

Instruction level parallelism (ILP)

6



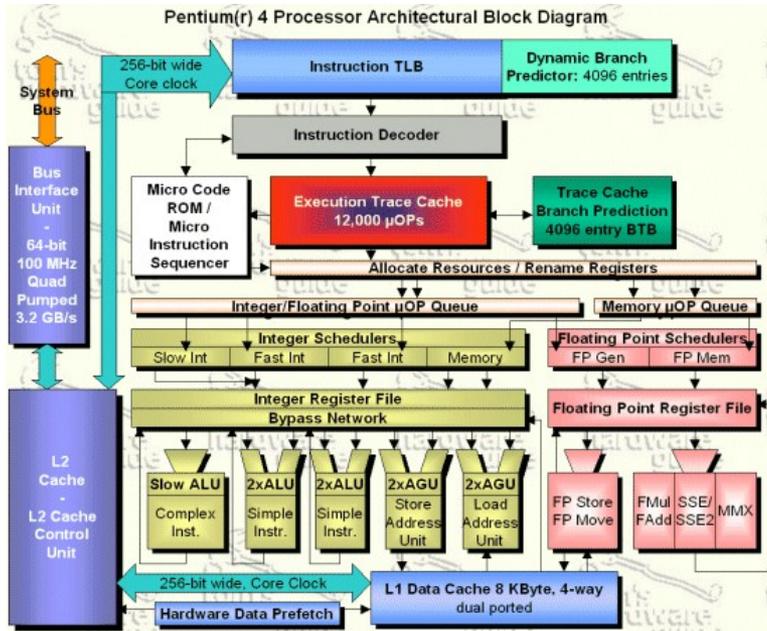
(a)



(b)

Limitations of instruction level parallelism

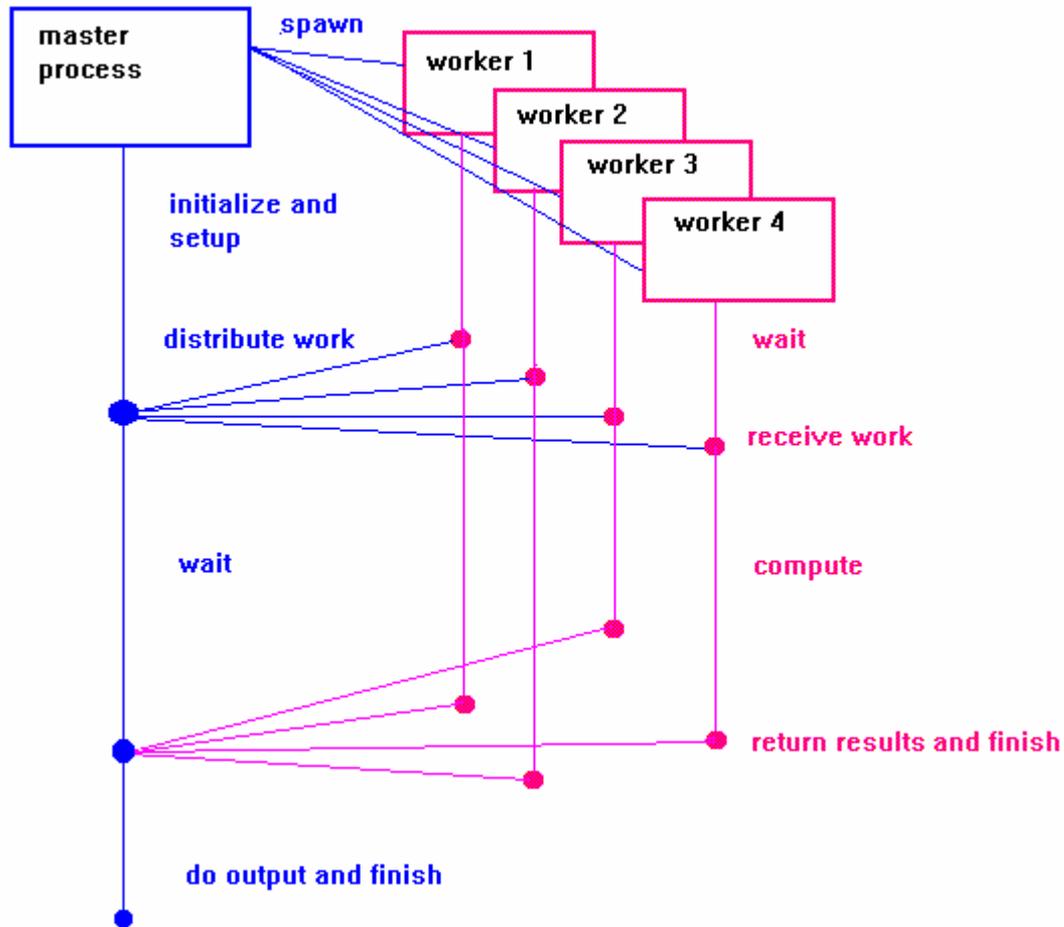
7



- The Pentium 4 has a 20 stage pipeline
- Typically every 5th or 6th instruction is a branch instruction – longer pipelines have a larger branch mis-prediction penalty
- Different instructions require different number of clock cycle and occur at different frequencies

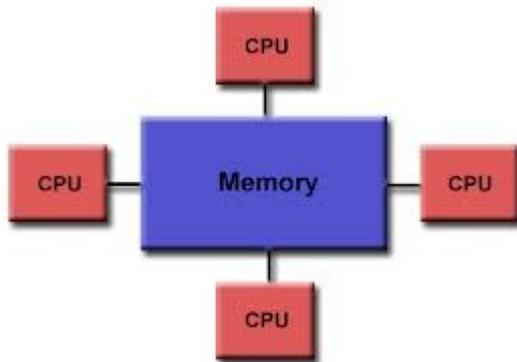
Application level parallelism

8



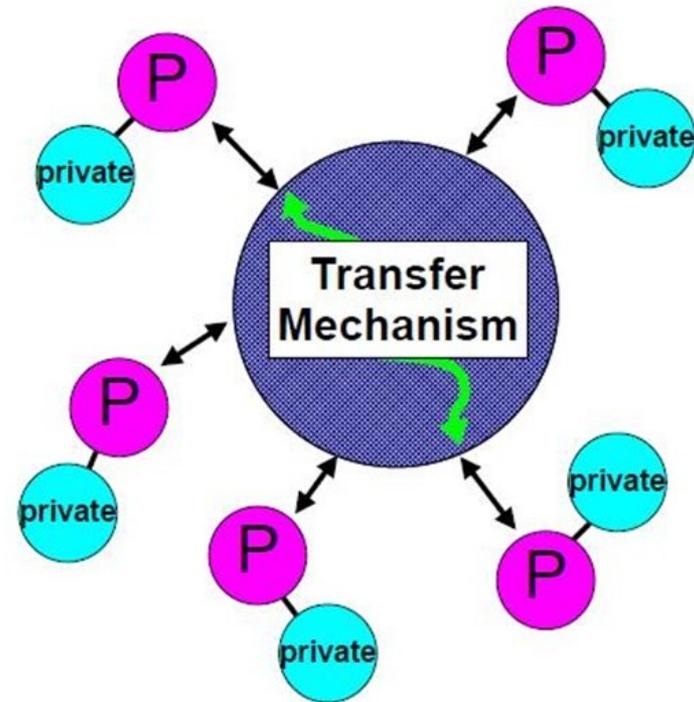
Application level parallelism - implementation models

9



Shared address space
(OpenMP)

HokieOne: 492 cores

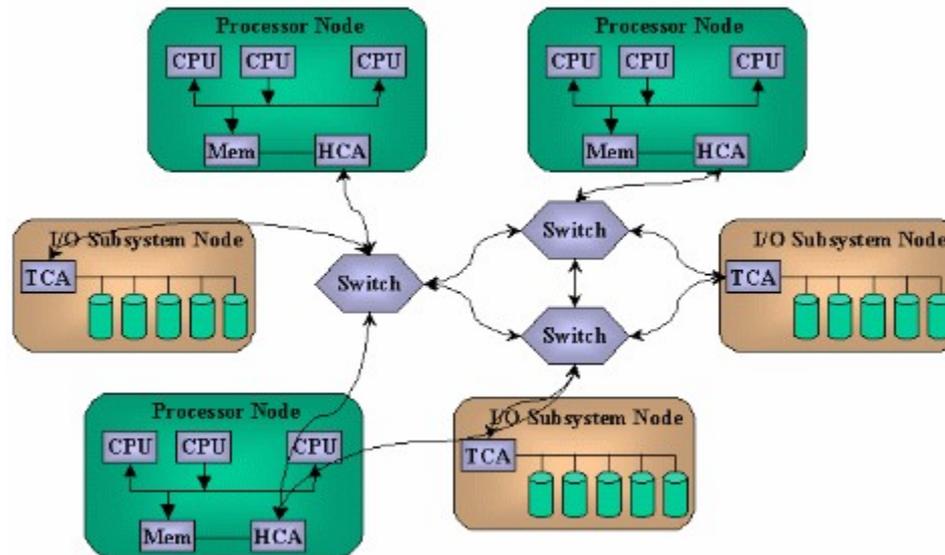


Message Passing
(MPI)

HokieSpeed: 2448 cores (204 nodes)

Limitations of supercomputer architecture

10



Gaming machines were way ahead of supercomputers when it came to parallel processing

11



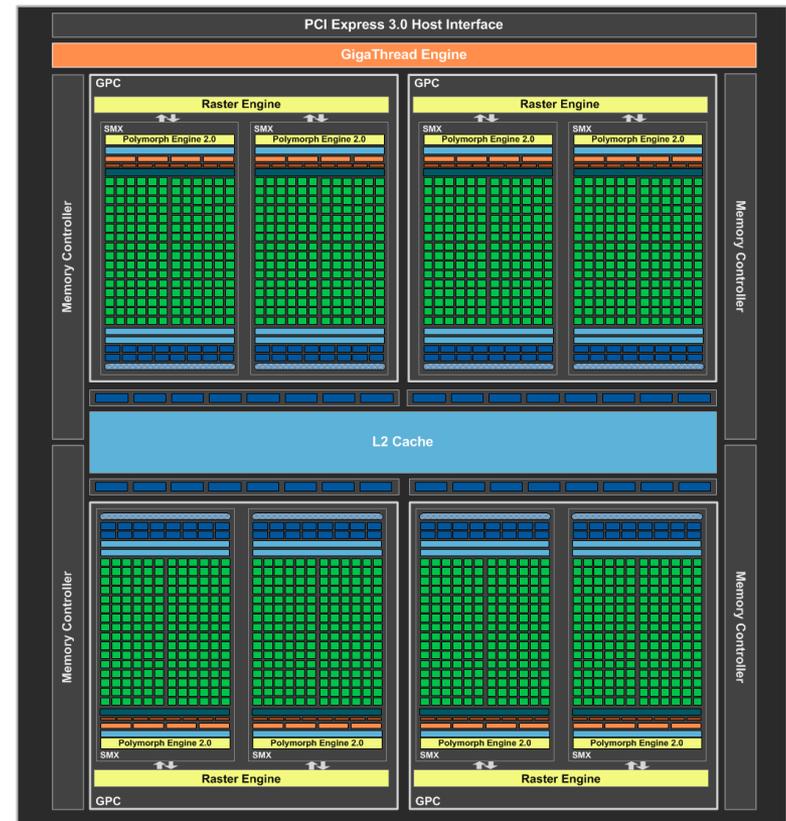
- General purpose central processing units (Intel/AMD) were just not up to the task of rapidly rendering realistic images
- Rendering graphics involves a large number of computations
- But each pixel can be computed more or less independently
- So the computer graphics were rendered using dedicated graphical processing units (GPUs)

GPUs – Massively parallel processors

12

GeForce GTX 690 Specifications

CUDA Cores	3072
Base Clock	915 MHz
Boost Clock	1019 MHz
Memory Config	4GB / 512-bit GDDR5
Memory Speed	6.0 Gbps
Power Connectors	8-pin + 8-pin
TDP	300W
Outputs	3x DL-DVI Mini-Displayport 1.2
Bus Interface	PCI Express 3.0



3 of the 5 fastest supercomputers in the world use GPUs (June 2012)

13

Rank	Site	Computer
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect Fujitsu
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 NUDT
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz Cray Inc.
4	National Supercomputing Centre in Shenzhen (NSCS) China	Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 Dawning
5	GSIC Center, Tokyo Institute of Technology Japan	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU , Linux/Windows NEC/HP

5 of the top 10 “green” supercomputers use GPUs (June 2012)

14

Green500 Rank	MFLOPS/W	Site	Computer
1	2026.48	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom
2	2026.48	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom
3	1996.09	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom
4	1988.56	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom
5	1689.86	IBM Thomas J. Watson Research Center	NNSA/SC Blue Gene/Q Prototype 1
6	1378.32	Nagasaki University	DEGIMA Cluster, Intel i5, ATI Radeon GPU , Infiniband QDR
7	1266.26	Barcelona Supercomputing Center	Bullx B505, Xeon E5649 6C 2.53GHz, Infiniband QDR, NVIDIA 2090
8	1010.11	TGCC / GENCI	Curie Hybrid Nodes - Bullx B505, Nvidia M2090 , Xeon E5640 2.67 GHz, Infiniband QDR
9	963.70	Institute of Process Engineering, Chinese Academy of Sciences	Mole-8.5 Cluster, Xeon X5520 4C 2.27 GHz, Infiniband QDR, NVIDIA 2050
10	958.35	GSIC Center, Tokyo Institute of Technology	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU , Linux/Windows
11	928.96	Virginia Tech	HokieSpeed, SuperServer 2026GT-TRF, Xeon E5645 6C 2.40GHz, Infiniband QDR, NVIDIA 2050