

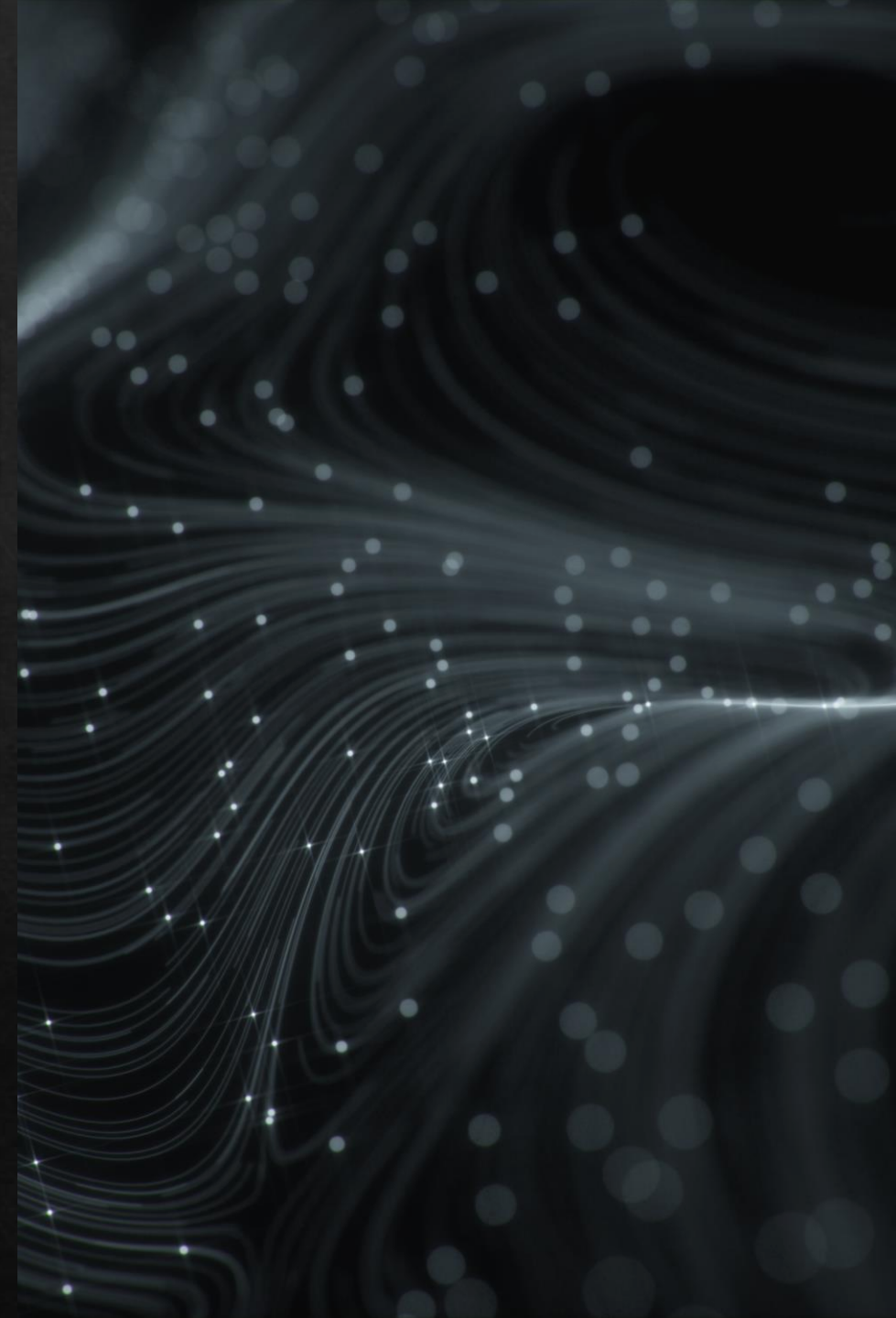
Discrete Space Hartree-Fock GPU Acceleration

ELEC 542 Project Presentation
Michel Kakulphimp

Introduction to GPUs

Introduction to GPUs

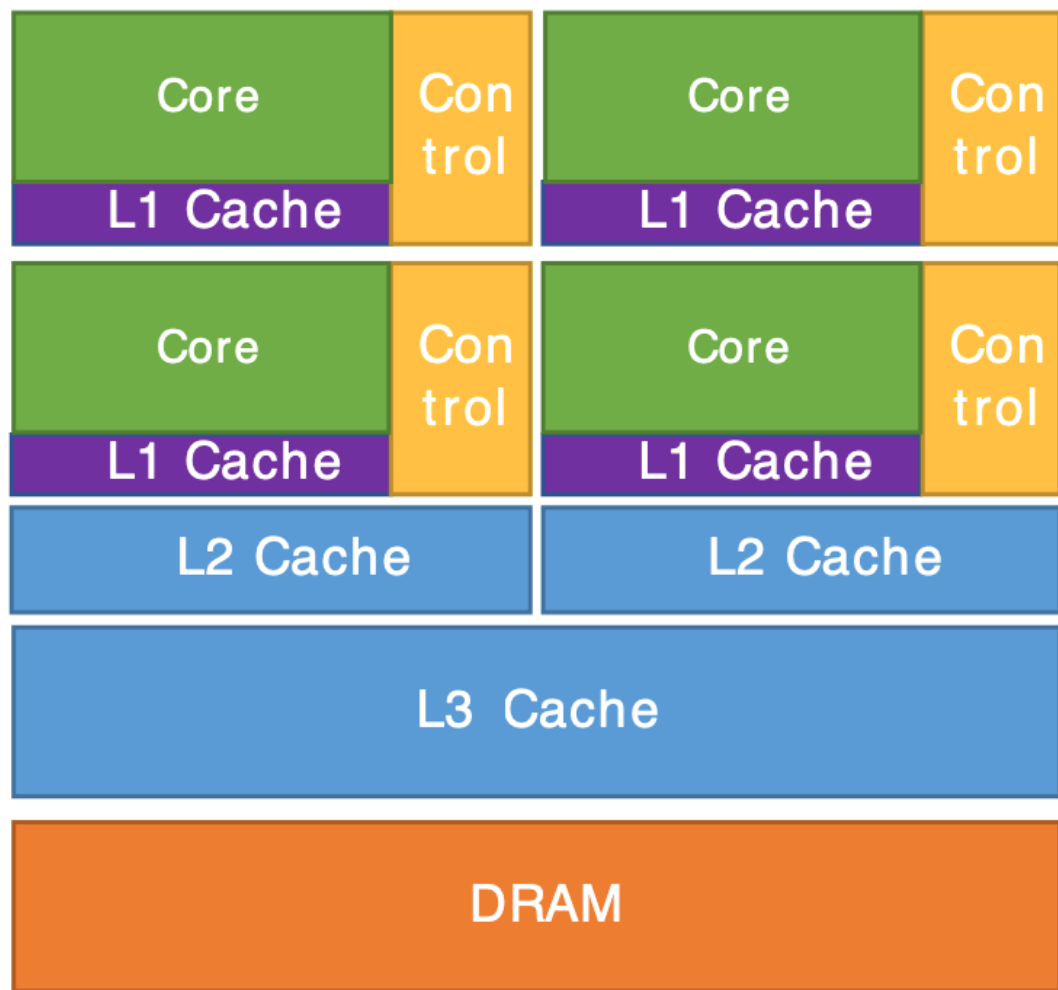
- ◆ Primarily designed for real-time computer graphics
- ◆ Well suited for highly parallelized computational workloads
- ◆ High memory bandwidth and many cores (RTX 3090 - 10496 cores!)
- ◆ Dependent on a host CPU, connected via high-speed high-bandwidth interconnect
- ◆ Different machine code than CPU, requires its own build tools
- ◆ Many GPU HW blocks are generic – just need to formulate problem differently
- ◆ Non-graphics workloads are a possibility and quite common
- ◆ Manufacturers provide APIs for development (NVIDIA CUDA used in this project)



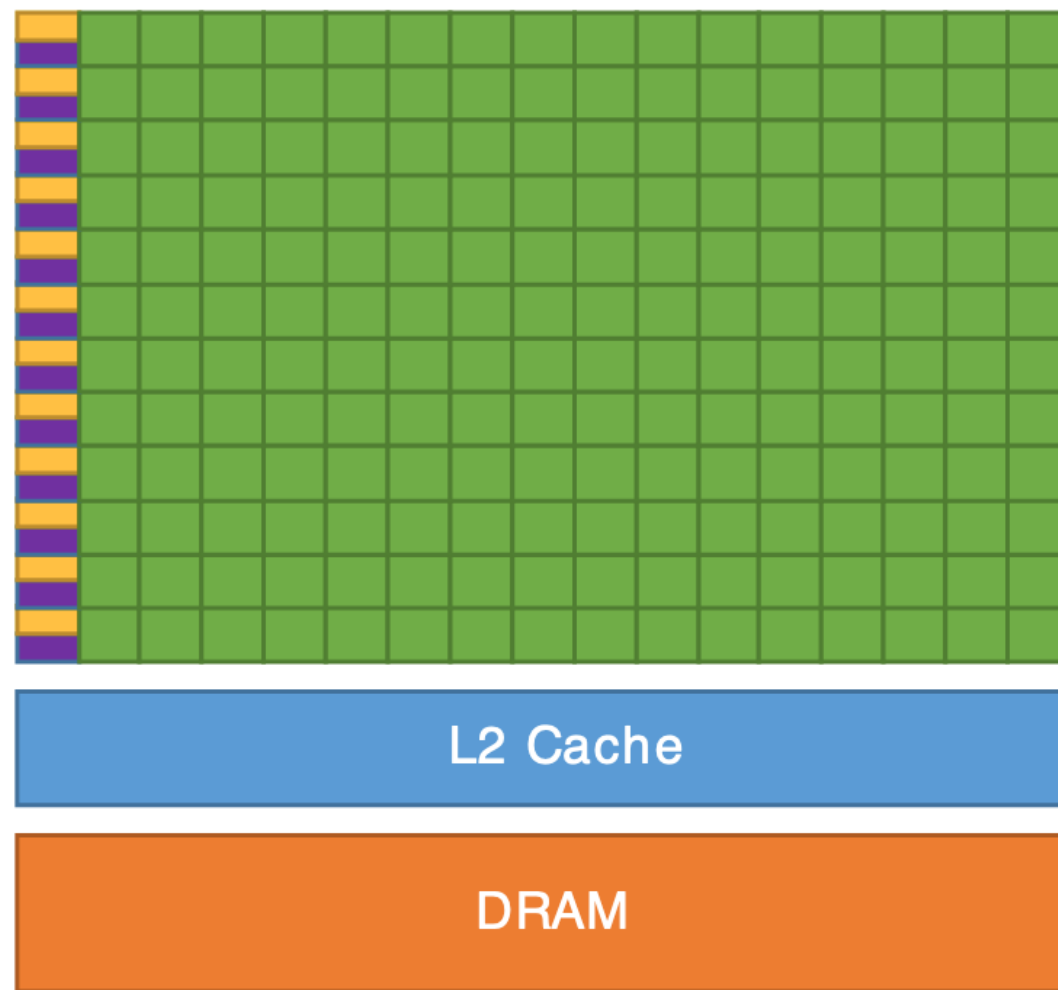
The background of the slide features a dark, abstract design. On the left side, there are numerous small, bright white dots arranged in a pattern that suggests a constellation or a network. These dots are connected by thin, flowing, wavy lines that create a sense of movement and depth. The lines and dots are set against a dark, textured background that transitions from black to a deep blue-grey. The overall effect is one of a futuristic or technological landscape.

GPU/CPU Comparison

- ◇ GPU is extremely parallel
 - ◇ Many more processing units (albeit weaker)
 - ◇ Optimized for SIMD instructions (i.e. one instruction – many data)
 - ◇ More transistors dedicated to ALU and FPU_s
- ◇ CPU more versatile
 - ◇ Larger instruction set
 - ◇ Better serial execution
 - ◇ Lower latency
 - ◇ Larger/more caches



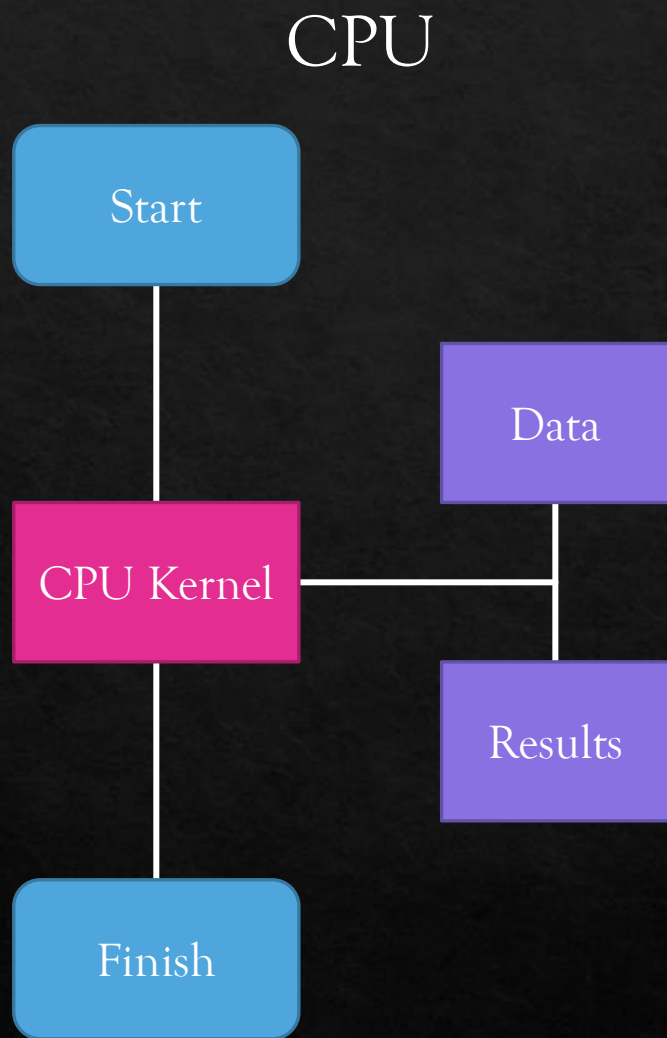
CPU



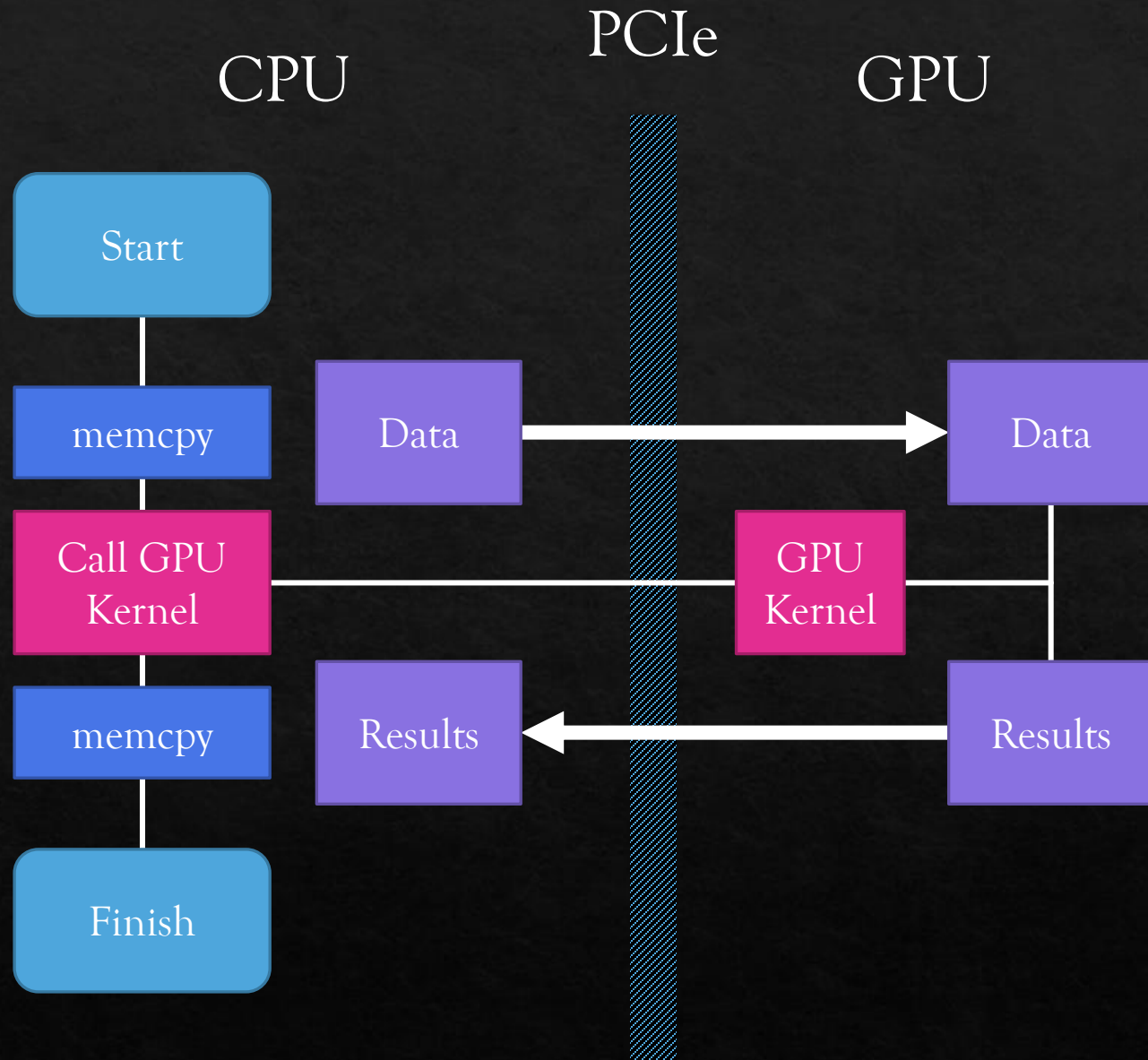
GPU

GPU Development Considerations

- ◆ Parallel Efficiency
 - ◆ Proper identification of parallelizable algorithmic components
 - ◆ Occupying as many GPU cores as possible
- ◆ Memory Throughput
 - ◆ Data must be well managed to ensure latency can be hidden
 - ◆ Moving data between CPU and GPU
- ◆ Instruction Throughput
 - ◆ Ensuring GPU cores are kept busy
 - ◆ Severe performance penalties from data dependencies (affects parallel efficiency as well)



CPU-only Execution



CPU-GPU Execution

GPU Accelerating Hartree-Fock

$$\hat{F}(\vec{r})\psi_n(\vec{r}) = \epsilon_n\psi_n(\vec{r})$$

$$\hat{F}(\vec{r}) = \hat{H}_{core}(\vec{r}) + \sum_{n=1}^{N/2} [2J_n(\vec{r}) - K_n(\vec{r})]$$

$$\hat{F}(\vec{r})\psi_n(\vec{r}) = \epsilon_n\psi_n(\vec{r})$$

$$\hat{F}(\vec{r}) = \hat{H}_{core}(\vec{r}) + \sum_{n=1}^{N/2} [2J_n(\vec{r}) - K_n(\vec{r})]$$

$$\hat{F}(\vec{r})\psi_n(\vec{r}) = \epsilon_n\psi_n(\vec{r})$$

$$\hat{F}(\vec{r}) = \hat{H}_{core}(\vec{r}) + \underbrace{\sum_{n=1}^{N/2} [2J_n(\vec{r}) - K_n(\vec{r})]}$$

Numerical Integration

$$\hat{F}(\vec{r})\psi_n(\vec{r}) = \epsilon_n\psi_n(\vec{r})$$

Eigensolver

$$\hat{F}(\vec{r}) = \cancel{\hat{H}_{core}(\vec{r})} + \underbrace{\sum_{n=1}^{N/2} [2J_n(\vec{r}) - K_n(\vec{r})]}$$

Numerical Integration

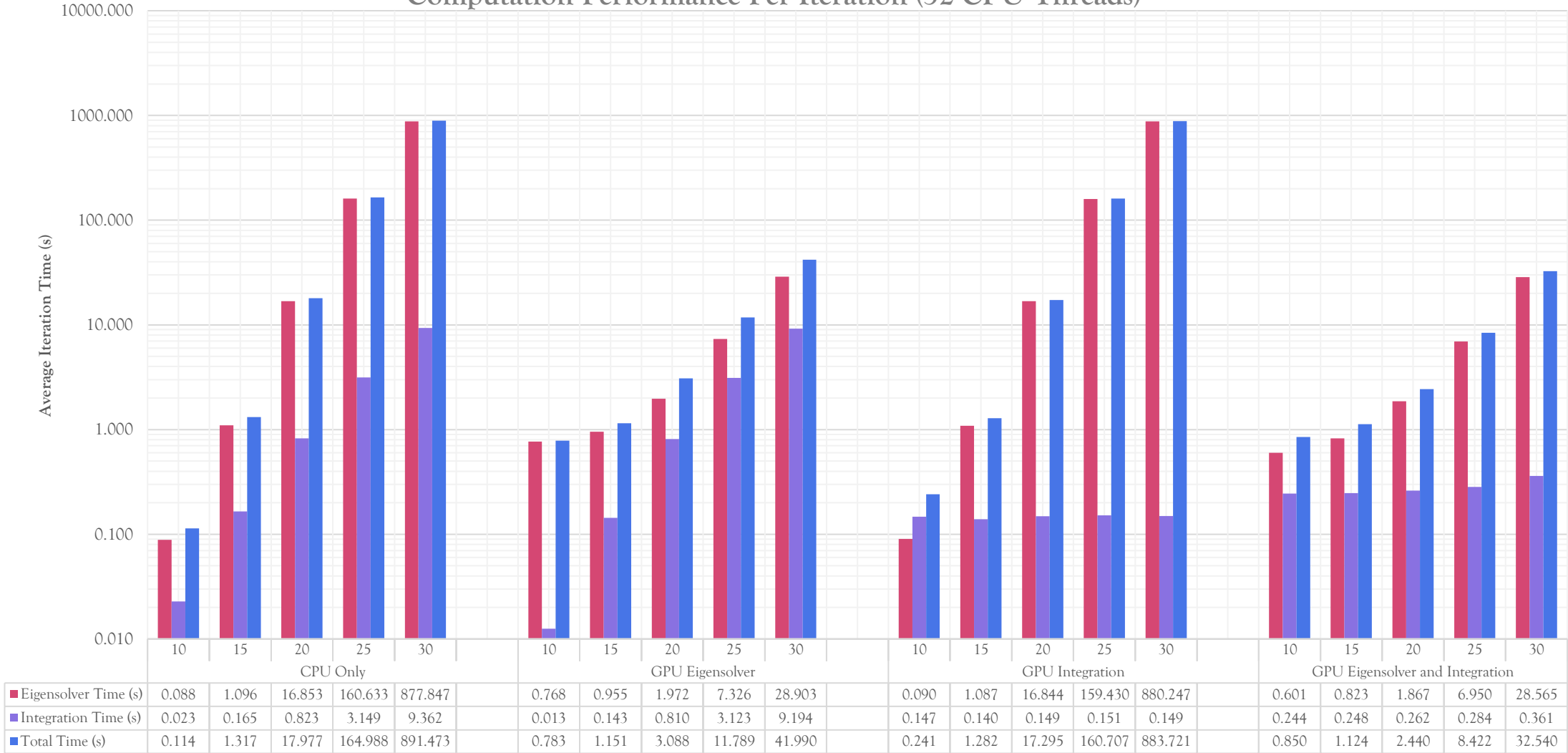
GPU Accelerating Hartree-Fock

- ◆ Fixed original Python implementation
- ◆ Ported Python implementation to C++ (BLAS/LAPACK used for eigensolver)
- ◆ Profiled C++ application
- ◆ Identified CPU-intensive sections (compounded by iterations)
- ◆ Re-implemented CPU-intensive sections as offloaded GPU workload
- ◆ Re-profiled with accelerated sections

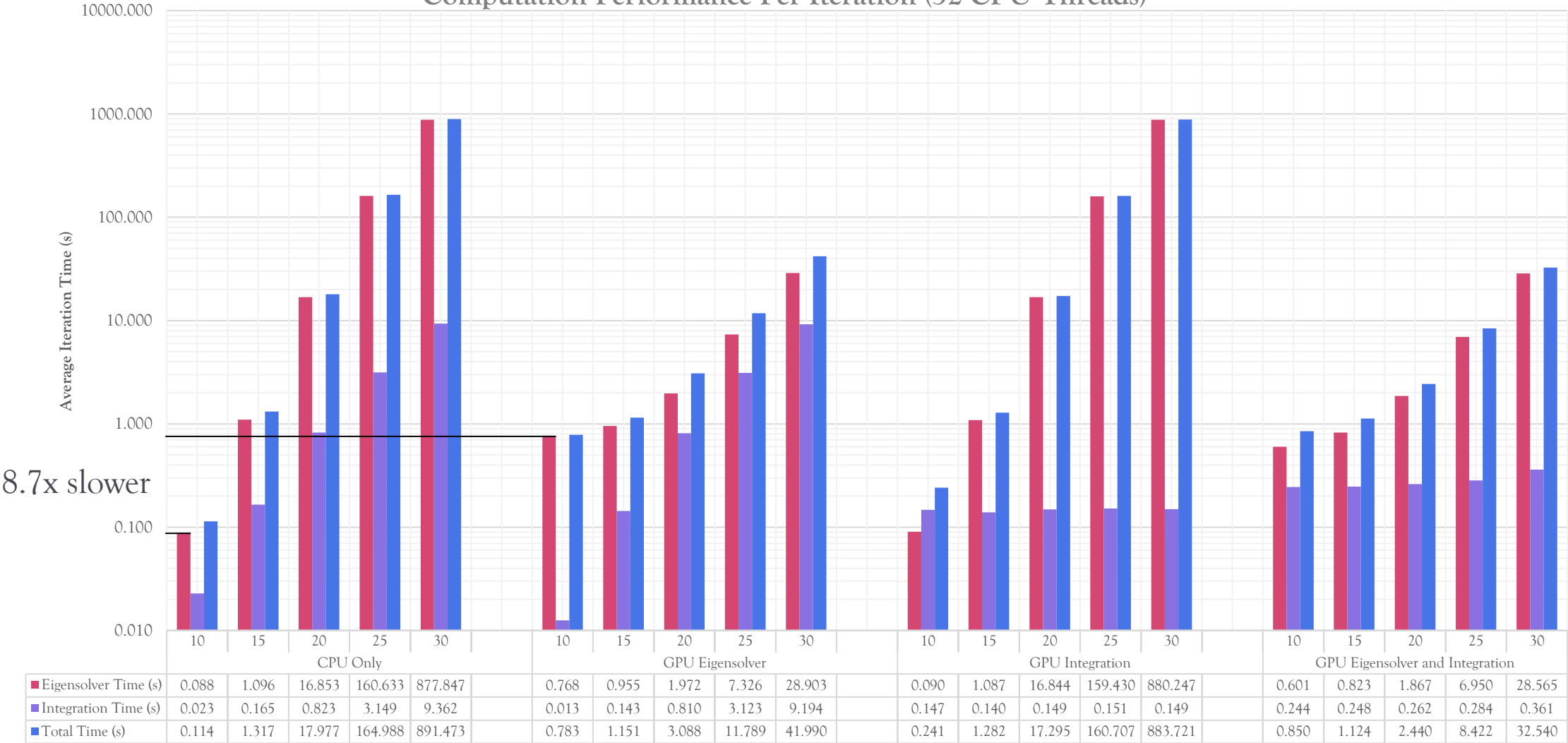
GPU Accelerating Hartree-Fock

- ◇ Numerical Integration
 - ◇ Large solution spaces require numerous integrations performed along the diagonal
 - ◇ Leverage GPU's massive parallelization capabilities
 - ◇ Ensure problem can be divided evenly into independent chunks (no memory dependencies)
- ◇ Diagonalization
 - ◇ Algorithm harder to perform parallel, but some speedups possible
 - ◇ Can leverage diagonalization routines provided by GPU vendor
 - ◇ CUDA cuSOLVER

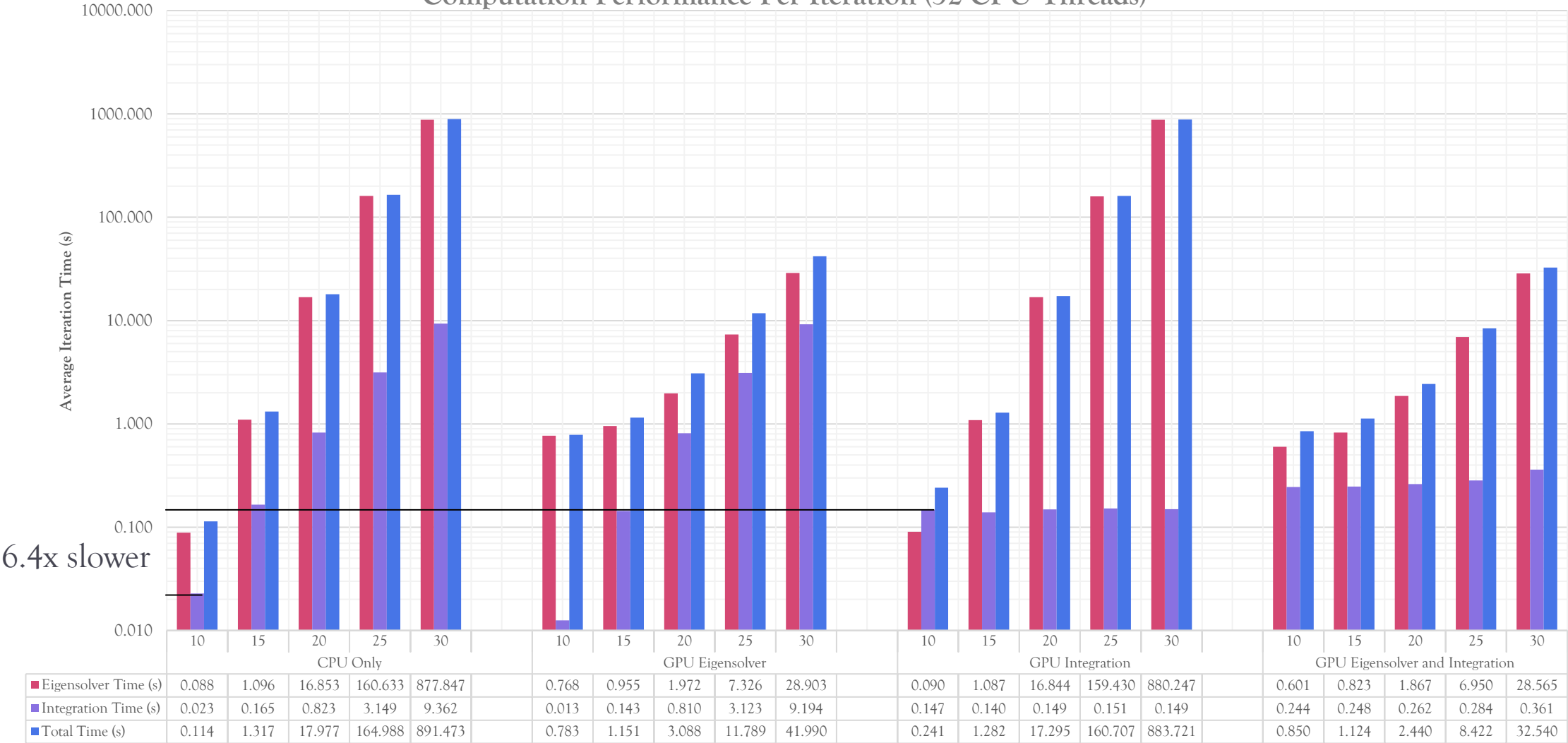
Discrete Space Hartree-Fock
Computation Performance Per Iteration (32 CPU Threads)



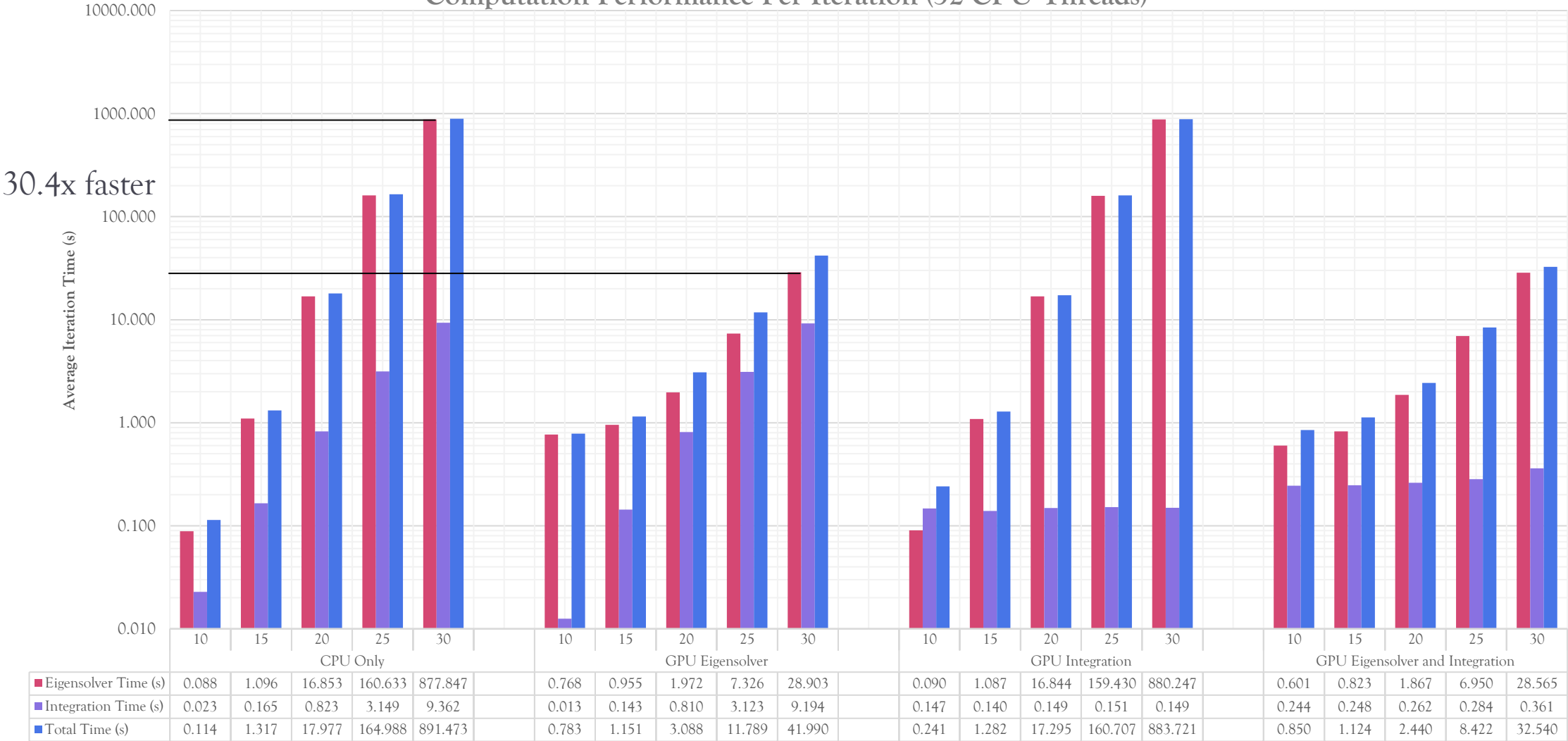
Discrete Space Hartree-Fock
Computation Performance Per Iteration (32 CPU Threads)



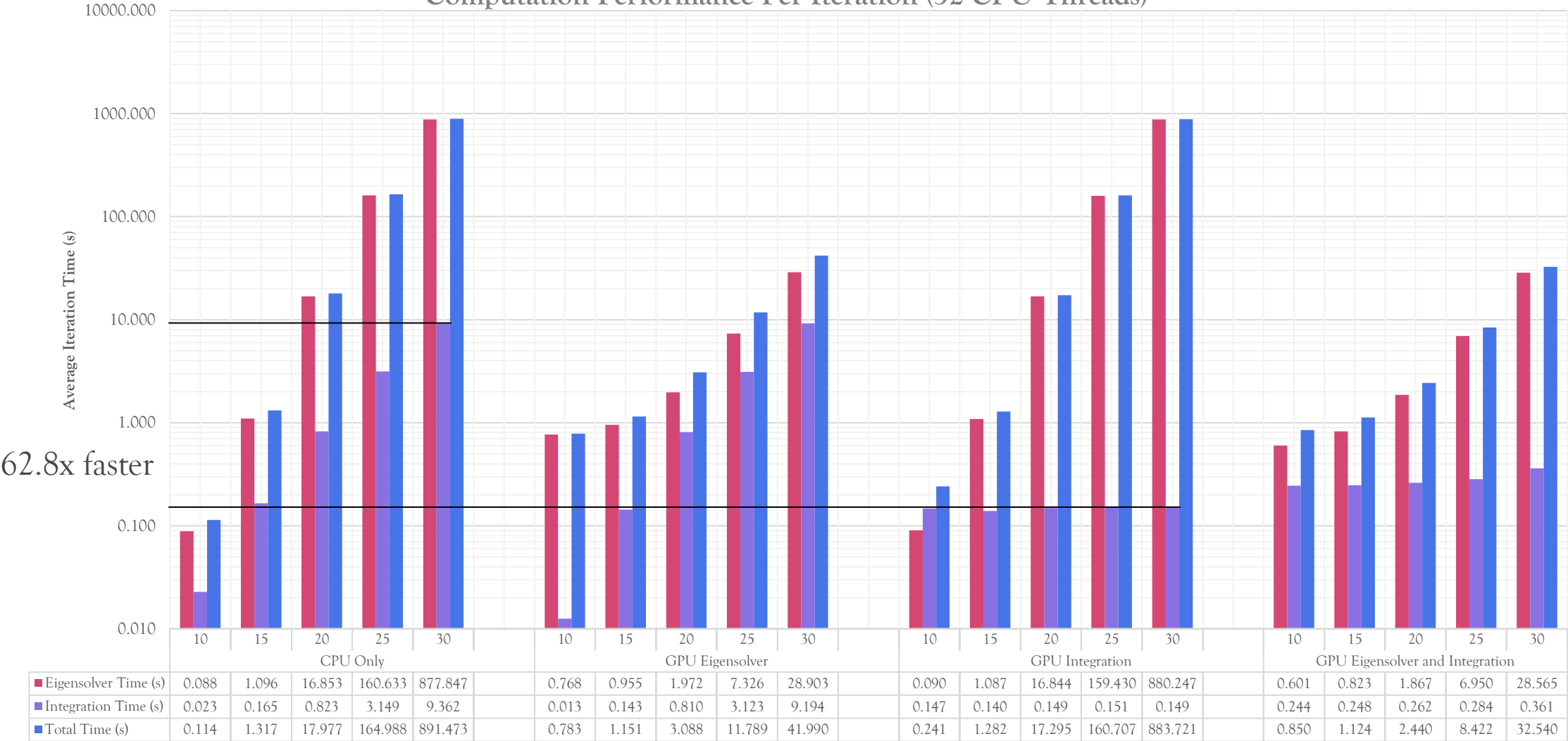
Discrete Space Hartree-Fock
Computation Performance Per Iteration (32 CPU Threads)



Discrete Space Hartree-Fock
Computation Performance Per Iteration (32 CPU Threads)



Discrete Space Hartree-Fock
Computation Performance Per Iteration (32 CPU Threads)



Future Considerations

- ◆ Current implementation naïve
 - ◆ Off-the-shelf eigensolver used, could break the algorithm down further and hand optimize (ROI?)
 - ◆ Numerical integration only parallelized per diagonal entry, integration itself could be parallelized as well
- ◆ Apply CPU acceleration techniques to basis-function HF and DFT
 - ◆ Other opportunities and challenges exist with both
- ◆ GPU Acceleration not trivial – requires decent understanding of architecture and algorithm for best results.

Questions?

Thank You!

Have a great summer ☺