



HYPERION RESEARCH

# Updates on the Intersection of AI and HPC

Tom Sorensen  
November 2023

# A New Category of HPC Workload

- HPC users adopting/integrating AI at high rates
- Many methods and models, LLMs draw attention
- ~90% of HPC users surveyed currently or plan to use AI methods for workloads
- AI methods introduce new demands on sites:
  - Hardware (processors, interconnect, data access)
  - Software (data management, queueing, dev tools)
  - Expertise (procurement strategy, maintenance, troubleshooting)
  - Regulatory (data provenance, privacy, legal)

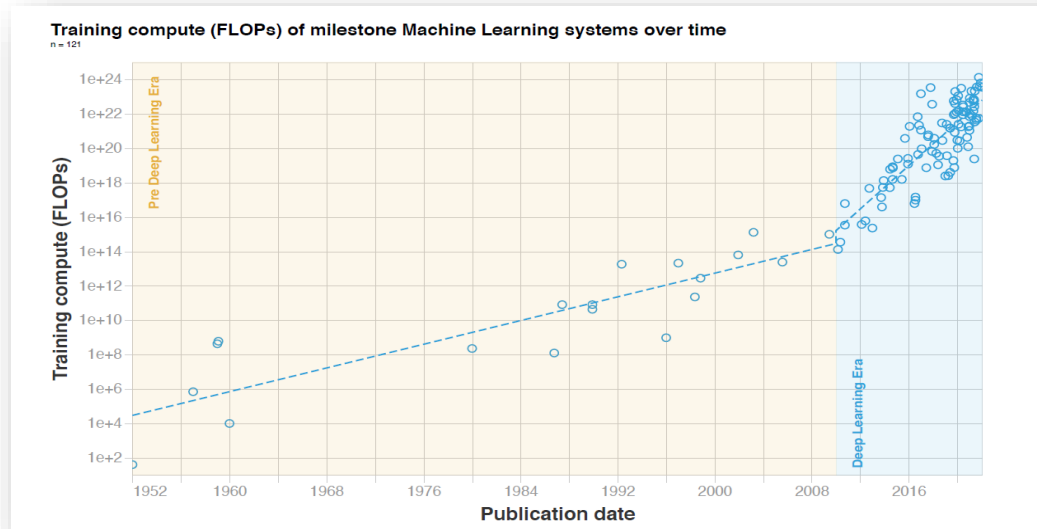
# Framing LLM/HPC Requirements

*Three elements dominate scaling of LLMs on HPCs*

- **Compute**: the absolute number of floating-point operations needed to train an LLM to a desired degree of accuracy
- **Dataset size**: input dataset used for training the LLM
- **Model size**: number of tokens or parameters
  - The larger the number of parameters, the more nuance in the model's understanding of each word's meaning and context
- **This scaling heuristic been called the ideal gas law of machine learning**
  - $PV = nRT$  encompasses a range of complex action
  - Scaling moves here as a  $f(C, D, M)$
- **LLM requirements ultimately define necessary HPC specifications**

# LLMs Consume Significant Flops

*LLM flops growth eclipses Top 500 growth*



- Pre 2010:
  - On the order of  $2 \times 10^{12}$  (200 Tflops)
  - Flops requirements doubling every 21.3 month
  - But not a lot of data points
- Post 2010 to Current:
  - Currently on the order of  $6 \times 10^{22}$  flops (60 Zettaflops)
  - Flops requirements doubling every 5.6 months
  - Roughly 11X faster than HPC Top 1 Linpack performance growth rate

See Compute Trends Across Three Eras of Machine Learning, arXiv:2202.05924

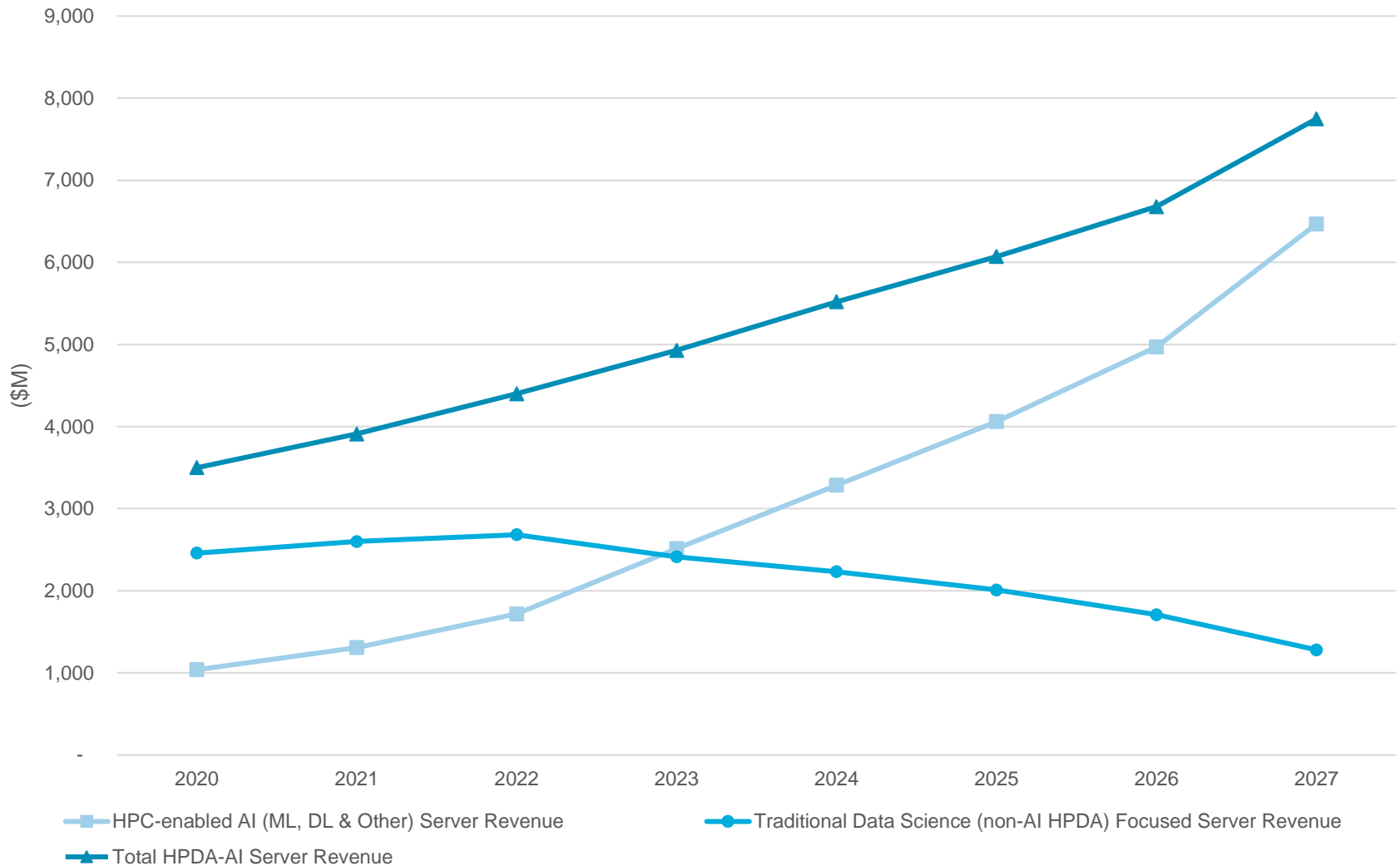
# Select Key Findings from AI Use Study

*From a survey of HPC users leveraging LLMs*

- **1:** LLMs are considered to be an important emerging asset for both current and planned HPC-related activity.
- **3:** Respondent organizations are looking at a broad set of LLM-related end uses.
- **5:** Many different HPC-related science and engineering algorithms were seen as viable for LLM enhancements.
- **9:** Open source is currently the most preferred option for accessing LLM software.
- **10:** Survey respondents looked to a wide range of LLM expertise to support the various stages of LLM development spanning foundation model construction, fine-tuning procedures, LLM integration into existing workloads, and supporting inference operations.

# Data-Intensive HPC Forecast (\$M)

Data-Intensive HPC Service Forecast Split (\$M)



# Putting This All Together

*Is this (another) new HPC architectural paradigm in the works?*

- **Based on a recent LLM analysis by Riken**
- **GTP variant flops requirements**
  - GPT-3.5 (ChatGPT):  $3 \times 10^{24}$  flops (estimated)
  - GPT-4.0:  $3 \times 10^{25}$  flops (estimated)
- **OpenAI System: Microsoft/Open AI collaboration**
  - Top 5 system when stood up
  - GPU-based BF16 312 Tflop/s x 25,000 = 7.8 Eflop/s TPP
  - GPT3.5 (ChatGPT): 4.5 days X 2
  - GPT-4.0 45 days X 2
- **Fugaku:**
  - FP32 6.76 Tflop/s X 158,976 = 1.07 Eflop/s (TPP)
  - GPT3.5 (ChatGPT): 32 days X 10
  - GPT-4.0 45 days X 2: 328 days X 10  $\sim$  8.9 years

Distributed Training of Large Language Models on Fugaku, <https://t.co/idofa7Tjyu>

# QUESTIONS?

[tsorensen@hyperionres.com](mailto:tsorensen@hyperionres.com)  
[ejoseph@hyperionres.com](mailto:ejoseph@hyperionres.com)

