

InPEX 2025 Japan: GenAI Breakout



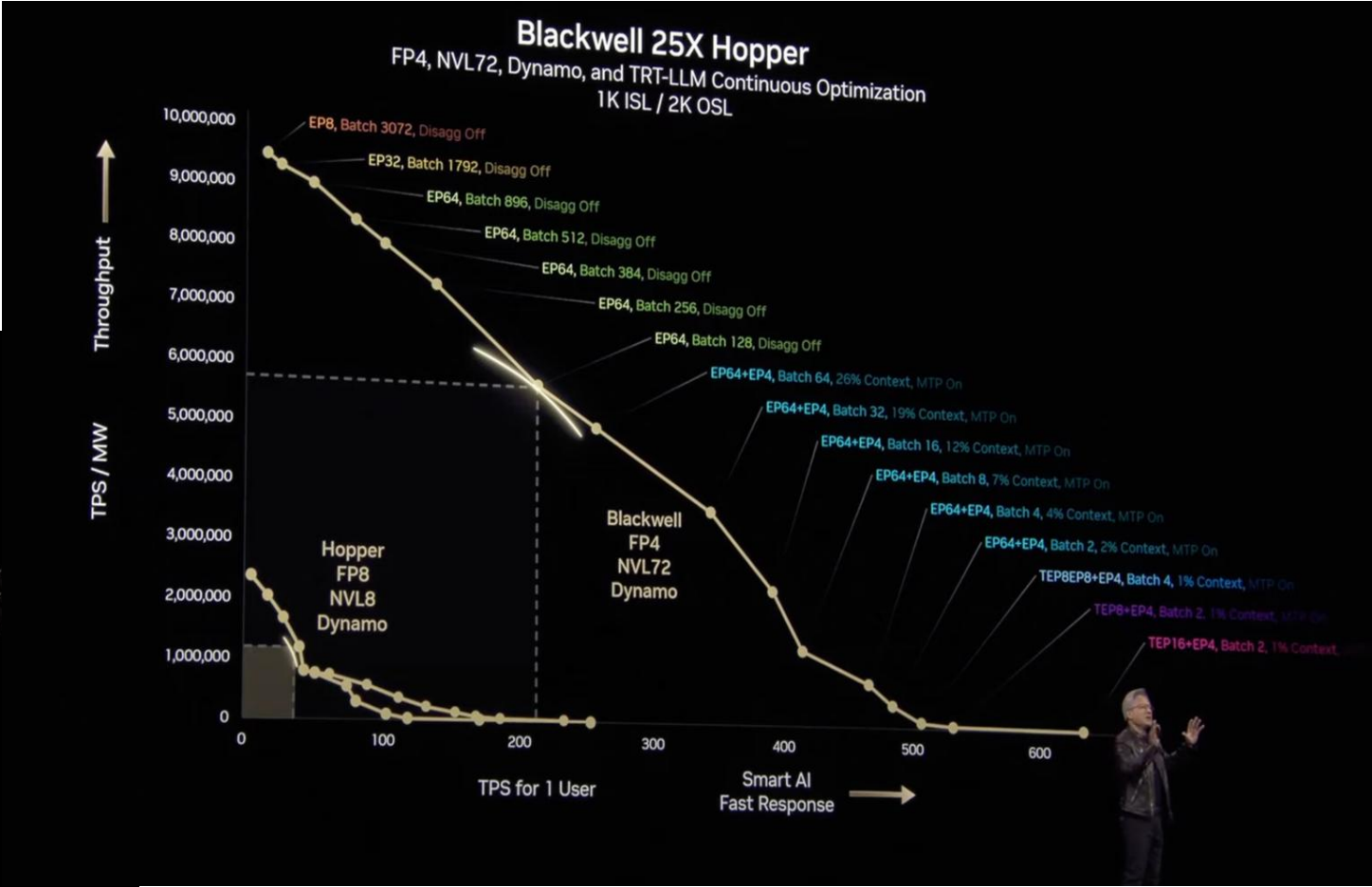
GenAI: Platforms

Blackwell 40X Hopper Inference Performance
NVLink Flops Per Watt ~ AI Factory Output



100 MW AI Factory	H100 NVL8	GB200 NVL72	
GPU Dies	45K	85K	
Racks	1,400	600	
Data Center Productivity	300M	12,000M	40X More Token Revenue

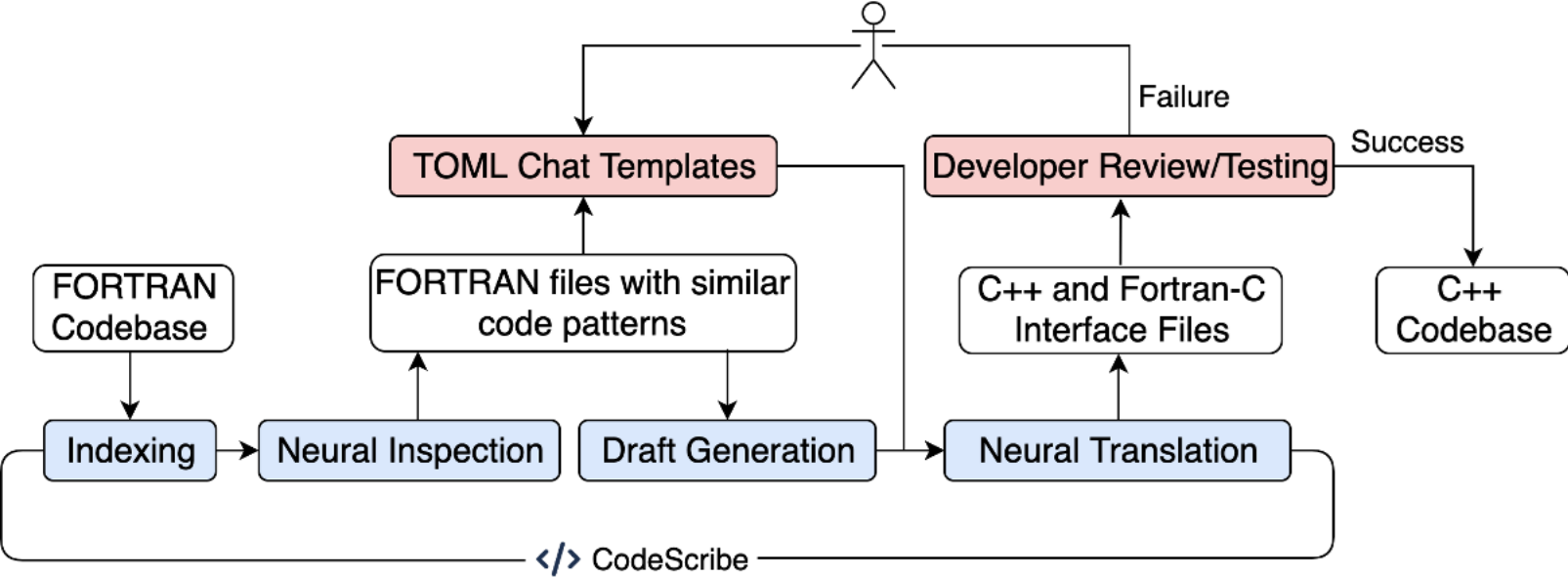
DeepSeek R1 Context and Generation, ISL=32K, OSL=8K @ Pareto optimal TPS/User



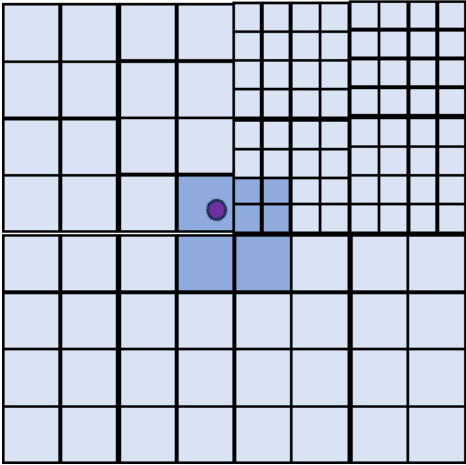
Key Insights: Future AI/HPC Architectures

- **Divergence vs. Integration?:** HPC remains unique due to precision and storage. AI brings low-precision inference and training demands—opportunities arise at the intersection.
- **Opportunities for Innovation:** Most innovation expected on the software stack, while hardware is vendor-driven. Emulated FP64, edge computing, and chiplet ecosystems were highlighted.
- **Portability & Sharing:** Emphasis on portability across platforms — especially for hybrid HPC+AI environments.
- **Operational Shifts:** Agents and inference are shaping new operational paradigms— we should address this with a working group
- **Global Coordination:** Strong sentiment: balance industrial AI trends and scientific HPC needs.

GenAI: Code



Monte Carlo code for LHC interactions



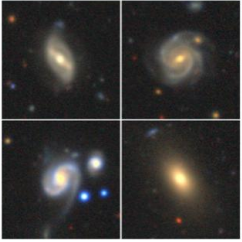
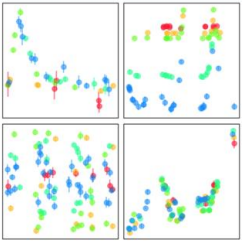
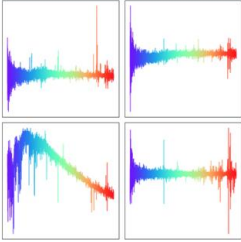
An algorithm for depositing from particle to mesh - decomposition, prompts and test-driven development, debugging – everything happening in (what should be) very precise English

Key Insights:

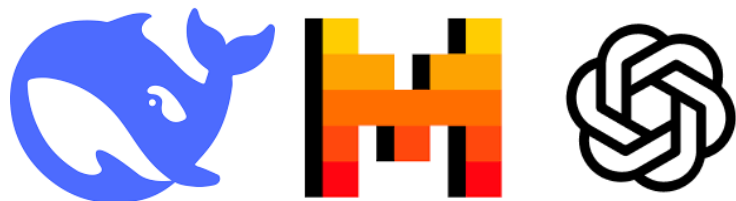
Code Generation, Precision, Debugging

- **AI Helps:** LLMs assist well with simple tasks and boilerplate code; fail more often on complex, hybrid, or algorithmic logic.
- **Prompting ≠ Programming:** Prompt engineering seen as inadequate pedagogy. Better for experienced developers, not a replacement for formal programming education.
- **Trust & Validation:** Consensus on the need for (more) robust test harnesses (unit tests, integration) and human-in-the-loop practices for validation.
- **Scientific Intuition:** (current) AI lacks intuition; safe precision reduction must remain scientist-guided, though AI can support validation and testing.
- **Path Forward: (FAST!)** International repositories (e.g., prompt libraries, agent systems) and shared fine-tuned models for HPC contexts.

GenAI: Data

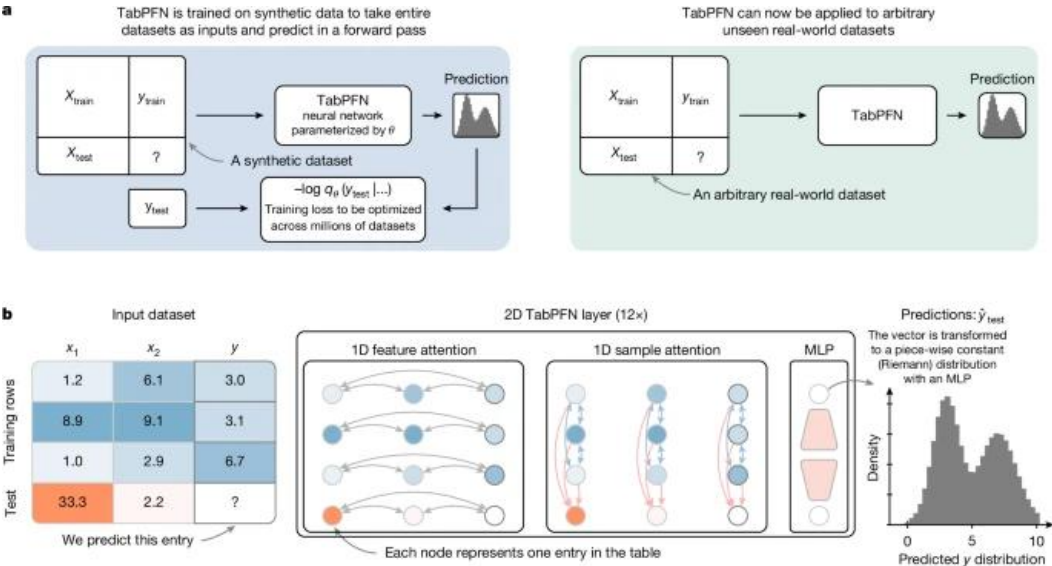
	Images	Time-Series	Spectra
# examples	140M	4.5M	225M
Description	images in a variety of wavelength ranges, including optical and infrared	multivariate time-series of flux + uncertainty in different wavelength ranges	flux as a function of wavelength
Tasks	galaxy classification, physical property estimation	time-series classification, redshift estimation	physical property estimation
Examples			

The Multimodal Universe Collaboration, 2024



>= 1T tokens

100M synthetic datasets



Hollman et al, 2024



The Well: 15TB of Physics Simulations

Ohana et al, 2024

Key Insights: Data Curation & Generation

- **Data Complexity Problem:** Each community has unique formats. There is a great need for AI to help make scientific data available for training.
- **Local Storage:** (a key difference in AI/HPC platforms): Supercomputing centers could support in-situ AI processing before discarding intermediate data.
- **Workflow Innovation:** Need for shared pipelines to load, preprocess, and stream data—especially for large GenAI models.
- **Access:** Social and political barriers limit data access.. This is difficult. We must continue to improve availability. Some fields (e.g., structural biology) are ahead; others lag behind.
- **Future:** Clear plan for improving, not just talking about, scientific data for AI

GenAI: Truth



AI FOR SCIENCE: 5 LESSONS FROM MY PHD

- #1: AI usage doesn't imply AI usefulness
- #2: AI will benefit some areas of science but not others
- #3: Evaluating whether AI is accelerating science is extremely difficult
- #4: Conflicts of interest and researcher degrees of freedom make AI-for-science overoptimistic
- #5: AI-for-science is often a solution looking for a problem

Nick McGreivy

Department of Astrophysics
Program in Plasma Physics
Princeton University

April 11th, 2025

Algorithmic Innovation & Entrepreneurship

Analysis | Published: 25 September 2024

Weak baselines and reporting biases lead to overoptimism in machine learning for fluid-related partial differential equations

[Nick McGreivy](#) ✉ & [Ammar Hakim](#)

[Nature Machine Intelligence](#) **6**, 1256–1269 (2024) | [Cite this article](#)

Key Insights: GenAI for Science

(Is AI Improving Scientific Progress?)

- **Impact Metrics:** Track *time to discovery, scientific output, human effort saved, community uptake, and cross-disciplinary impact*. Use both *quantitative* (e.g., citations, speedups) and *qualitative* (e.g., novelty, utility) measures.
- **Practical Benefits:** GenAI is best at supporting *tedious, complex, or repetitive tasks* like literature review, hypothesis generation, debugging, and code maintenance – freeing up human creativity.
- **ModSim vs AI:** GenAI can enable science in *data-rich domains* (e.g., astrophysics, protein folding) and assist where *first-principle models fall short*.
- **Best Practices:** Avoid biases (e.g., cherry-picking, weak baselines) via *provenance tracking, peer review, and international frameworks*.
- **Global Collaboration:** 1000 Scientist Jam – repeated in many places. Build **REAL** experience. Foster trust through transparency and reproducibility.

Thank You