# ADAPTING GENERAL-PURPOSE FOUNDATION MODELS FOR X-RAY PTYCHOGRAPHY IN LOW-DATA REGIMES

AI for Accelerated Materials Design: 39th Conference on Neural Information Processing Systems (NeurIPS 2025)

Robinson Umeike[1,*],   Neil Getty[2],   Xiangyu Yin[2]   and   Yi Jiang[2]

[1]The University of Alabama,  [2]Argonne National Laboratory

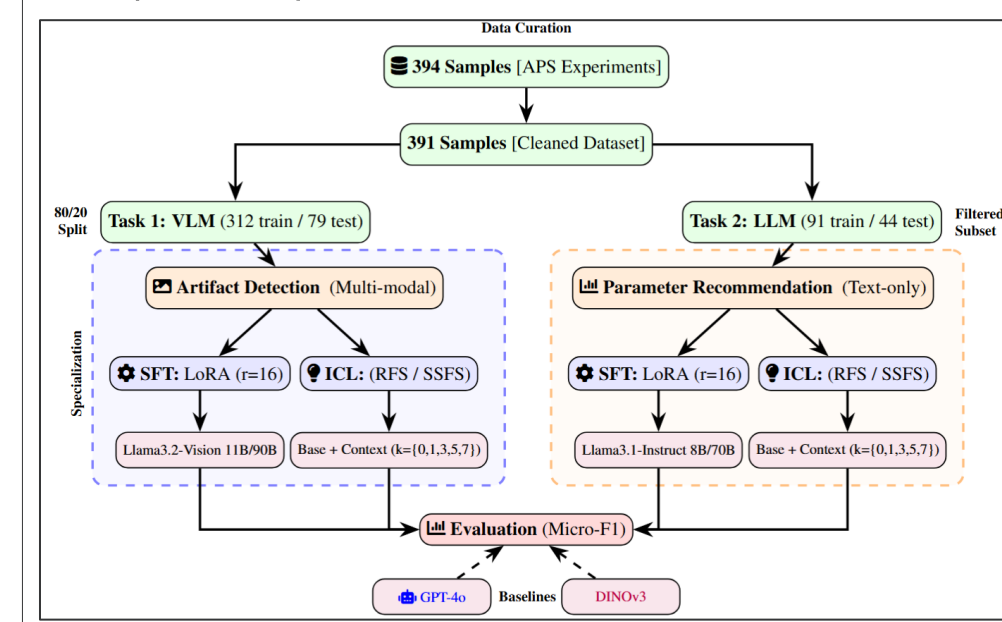* Work performed while at Argonne National Laboratory

## ABSTRACT SUMMARY

Adapting foundation models for specialized scientific tasks is critical, yet the optimal strategy remains unclear. We introduce **PtychoBench**, a multi-modal benchmark for ptychography, to systematically compare Supervised Fine-Tuning (SFT) versus In-Context Learning (ICL) in a data-scarce regime. Our findings reveal the optimal pathway is task-dependent. For visual tasks (VLMs), SFT and ICL are **complementary**, achieving the highest **mean performance (Micro-F1 0.728)**. Conversely, for textual tasks (LLMs), ICL on a base model is **superior (mean Micro-F1 0.847)**, outperforming a "super-expert" SFT model **(mean Micro-F1 0.839)**. Benchmarked against **GPT-4o** and **DINOv3**, our results on this dataset highlight that visual tasks benefit from hybrid specialization, while textual reasoning favors flexible base models in science-based agentic systems.

## CONTRIBUTION

❑ **Novel Benchmark:** 391 expert-annotated samples from the Advanced Photon Source (APS)

❑ **Dual Tasks:**
  ✓ **Task 1 (Visual):** Artifact Detection (VQA).
  ✓ **Task 2 (Textual):** Parameter Recommendation (QA).

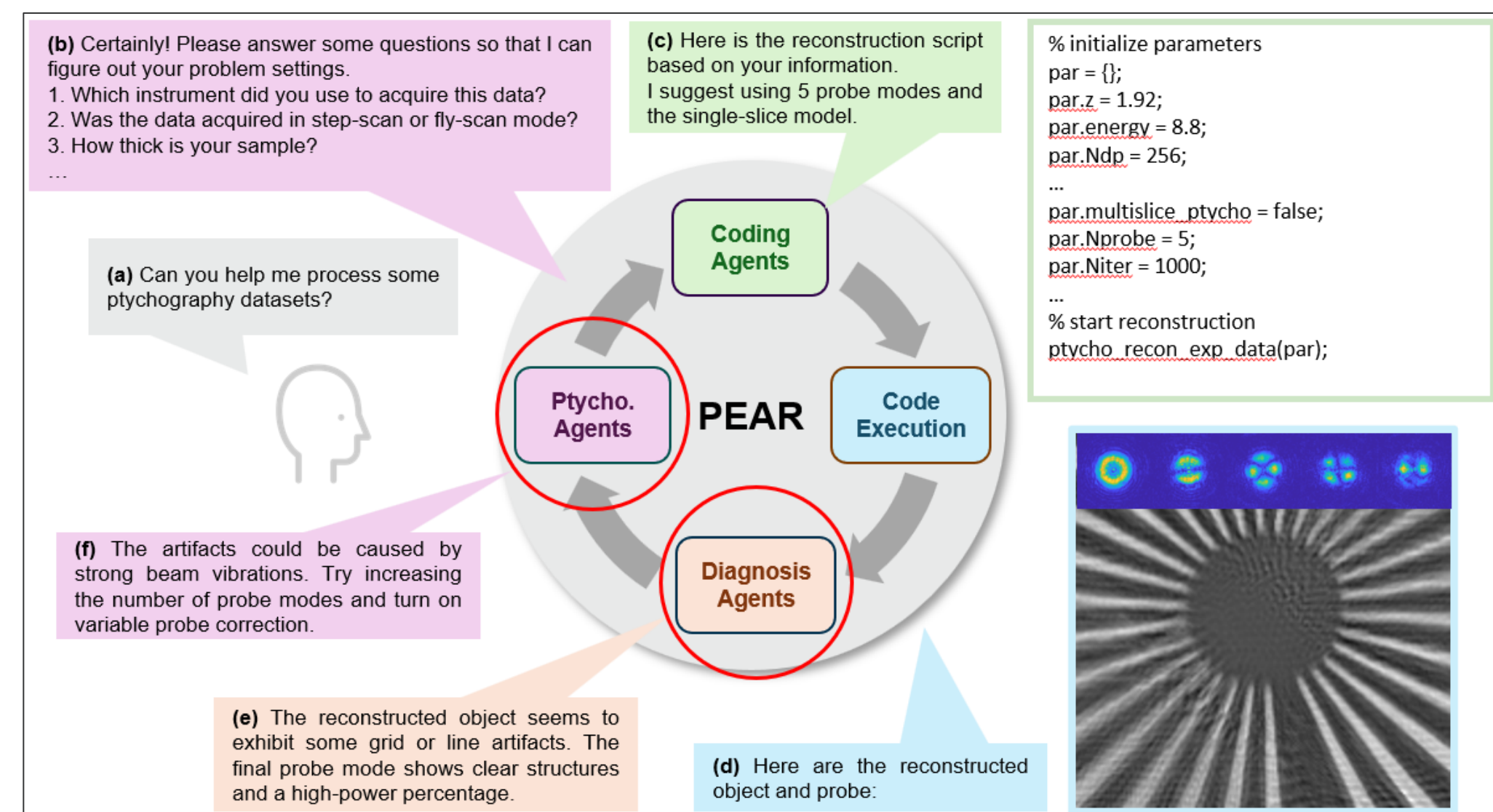❑ **Data-Scarce Regime:** Limited training data available, mimicking real experimental constraints.

## METHODS

❑ **Models:** Llama 3.2 (11B/90B *Vision*), Llama 3.1 (8B/70B *Text*).

❑ **ICL Strategy:** Compared Random Few-Shot (RFS) vs. Sample-Specific Few-Shot (SSFS).



## BACKGROUND & MOTIVATION

❑ **Context:** X-ray ptychography is a computational imaging technique that achieves high resolution but requires manual parameter tuning.

❑ **Problem:** General-purpose models struggle with this specialized scientific data, and _no standard benchmark exists_ to evaluate them.

❑ **Goal:** Automate manual analysis using "Agentic Workflows" (e.g., PEAR), with VLM (Diagnosis) and LLM (Recommendation) agents.

❑ **Approach:** We systematically compare two adaptation paths: **Supervised Fine-Tuning (SFT)** vs. **In-Context Learning (ICL)**.
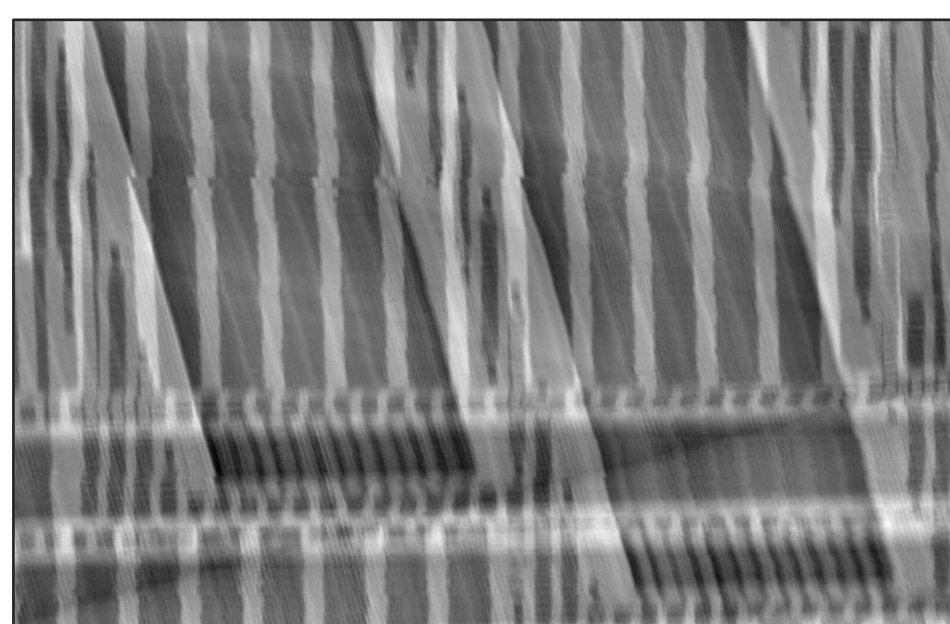


**Sample Info:**
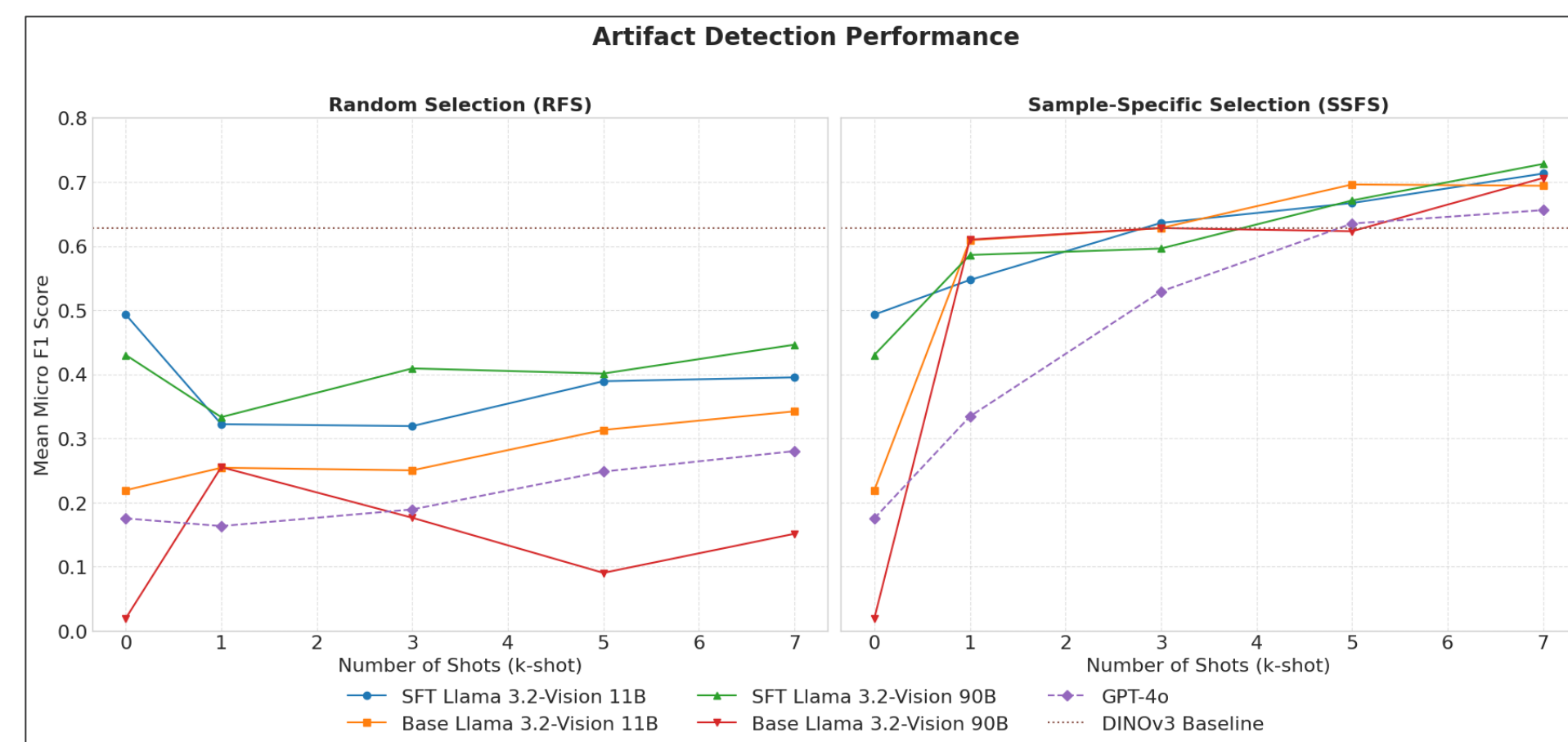  ✓ **Type:** Integrated circuit
  ✓ **Artifact:** Local Distortion

**Expert recommendation:**
  ✓ Increase the number of iteration.
  ✓ Reduce the size of diffraction patterns by a factor of 2
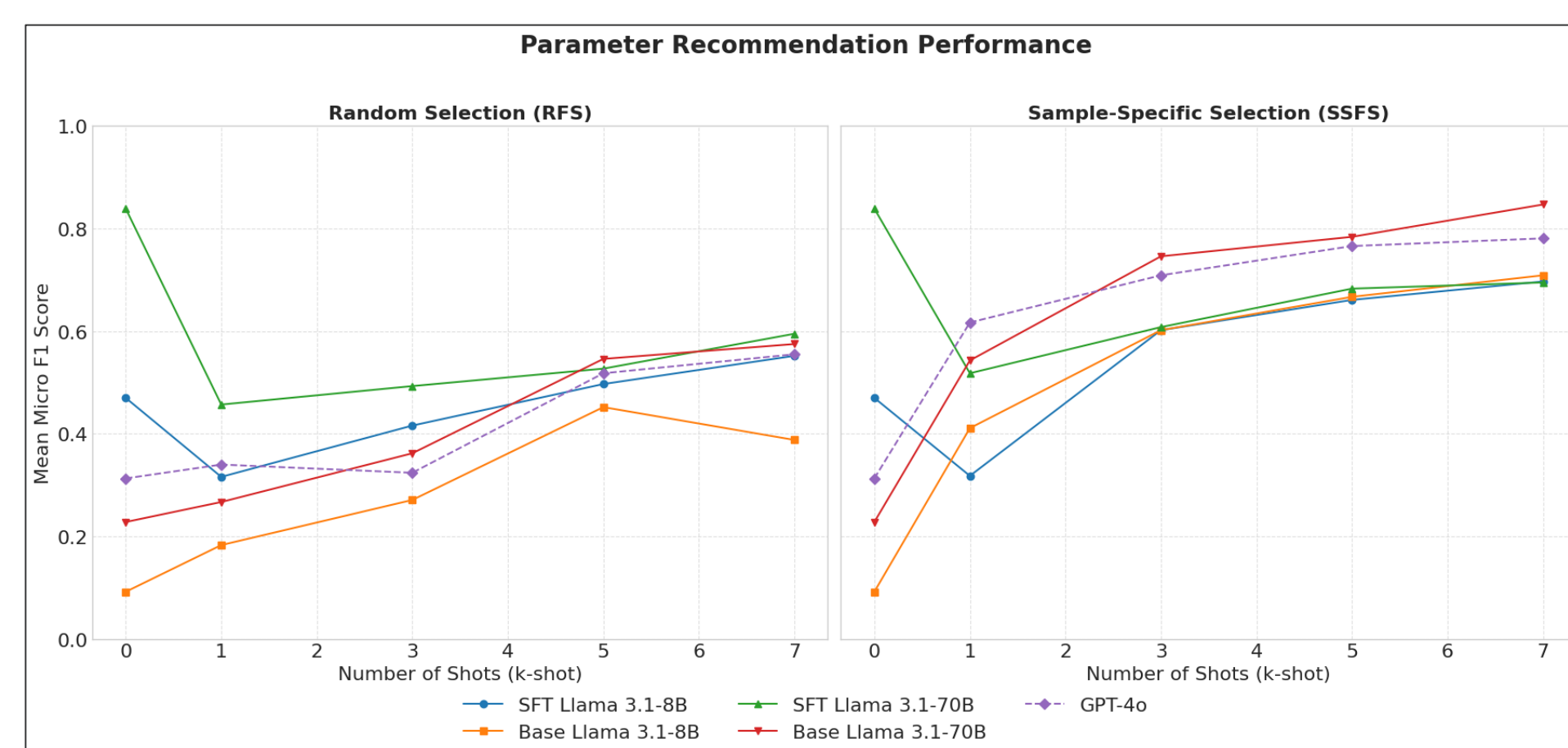  ✓ Turn off variable probe correction.

## RESULTS

❑ **Visual diagnosis:** The winning strategy **combines SFT and SSFS (0.728 F1)**. Fine-tuning builds a "_visual vocabulary_" that context refines, allowing the larger **90B model** to overcome its lower standalone baseline and outperform the smaller expert when unlocked by relevant examples.



❑ **Textual recommendation:** In contrast, textual recommendation favors **a Base 70B Model with SSFS (0.847 F1)**. While SFT created a "_brittle super-expert_" that degraded with context, the flexible base model effectively utilized examples to surpass the specialized expert.



## DISCUSSION & CONCLUSION

❑ **Contextual Interference:** We observed a consistent _contextual interference_ phenomenon across fine-tuned models. While relevant context boosts scores, **random examples actively confuse the expert models**, causing performance to plummet below the zero-shot baseline (e.g., a 35% drop for SFT-11B), highlighting that retrieval quality is a critical safety constraint.

❑ **Data Efficiency:** Our results also highlight remarkable data efficiency. The Base-90B model, using **only 7 inference-time examples**, achieved a **mean F1 score (0.706)** statistically comparable to a **DINOv3** vision classifier trained on the entire dataset **(0.628).** This proves that guided foundation models can match traditional supervised learning with a fraction of the data.

❑ **Conclusion:** Our dataset reveals no one-size-fits-all strategy. Visual tasks (**artifact detection**) benefit from SFT priming, whereas textual tasks (**parameter recommendation**) require the flexibility of base models.

## REFERENCES

❑ X. Yin, C. Shi, Y. Han, and Y. Jiang, (2024) *PEAR: A Robust and Flexible Automation Framework for Ptychography Enabled by Multiple Large Language Model Agents.*

❑ R. Umeike, N. Getty, F. Xia, and R. Stevens, (2025) *Scaling Large Vision-Language Models for Enhanced Multimodal Comprehension in Biomedical Image Analysis.* In 2025 IEEE 22nd International Symposium on Biomedical Imaging (ISBI), pp. 1--4.

**Have Questions? Please contact:**
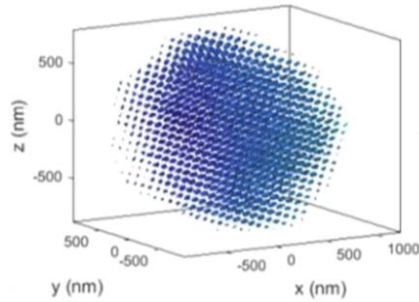  ✓ ngetty@anl.gov (Code)
  ✓ yjiang@anl.gov (Dataset)
  ✓ *Scan for Code & Data Access Guidelines* ➔
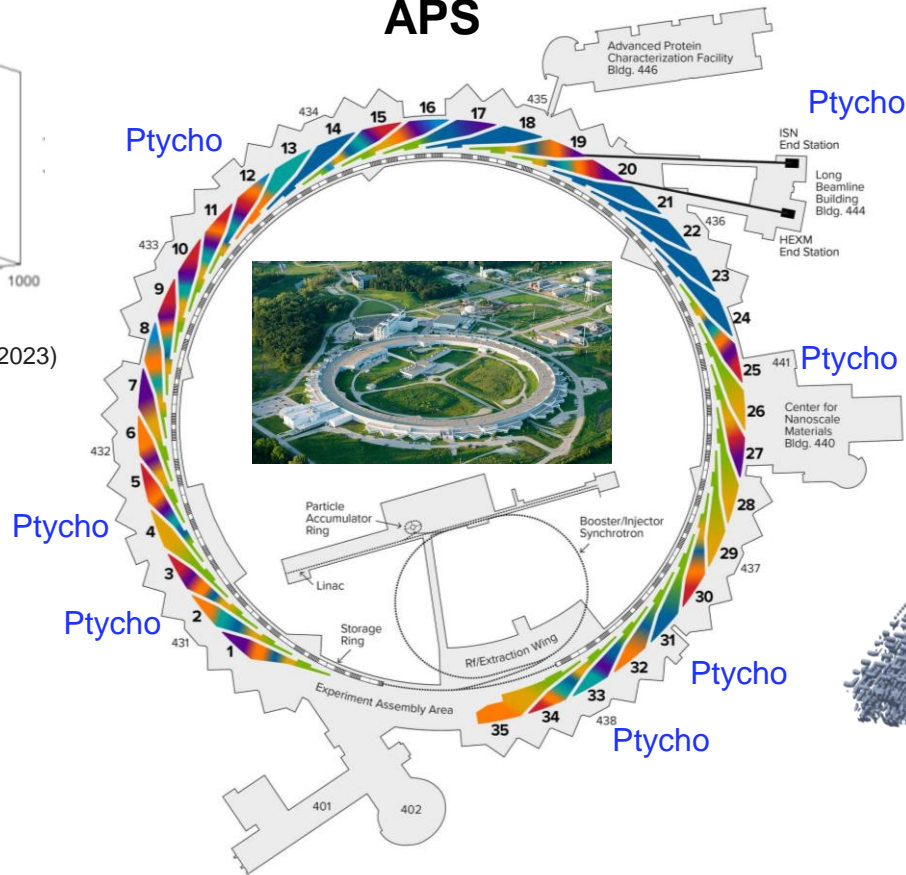
# Transmission X-ray Ptychography @ APS

supercrystals



H. Calcaterra et al. *ACS Nano.* (2023)

liquid-metal particles
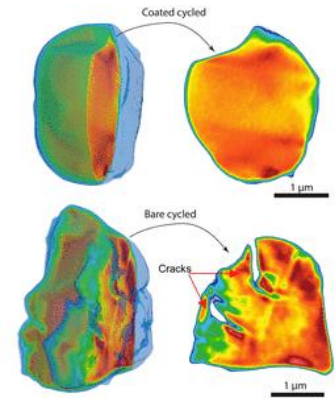


Y. Lin et al. *Nat. Chem.* (2022)

**APS**



Ptycho

Ptycho

Ptycho

Ptycho
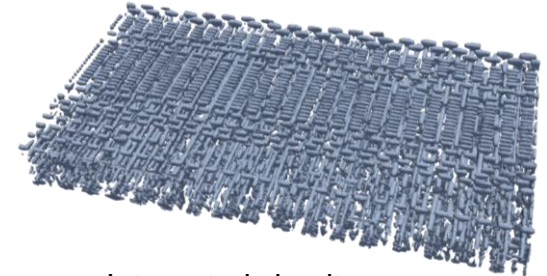
Ptycho

Ptycho

Ptycho

~10 instruments

Ni-Rich Cathodes



Q. Liu et al. *ACS Nano.* (2022)



Integrated circuit

Argonne
NATIONAL LABORATORY