

Virtual Summer School:
Machine Learning in Electron Microscopy

Sergei V. Kalinin and Gerd Duscher,
University of Tennessee, Knoxville

Maxim Ziatdinov and Rama Vasudevan,
Oak Ridge National Laboratory



Maxim
Ziatdinov



Rama
Vasudevan



Gerd
Duscher



Ayana
Ghosh



Tommy
Wong

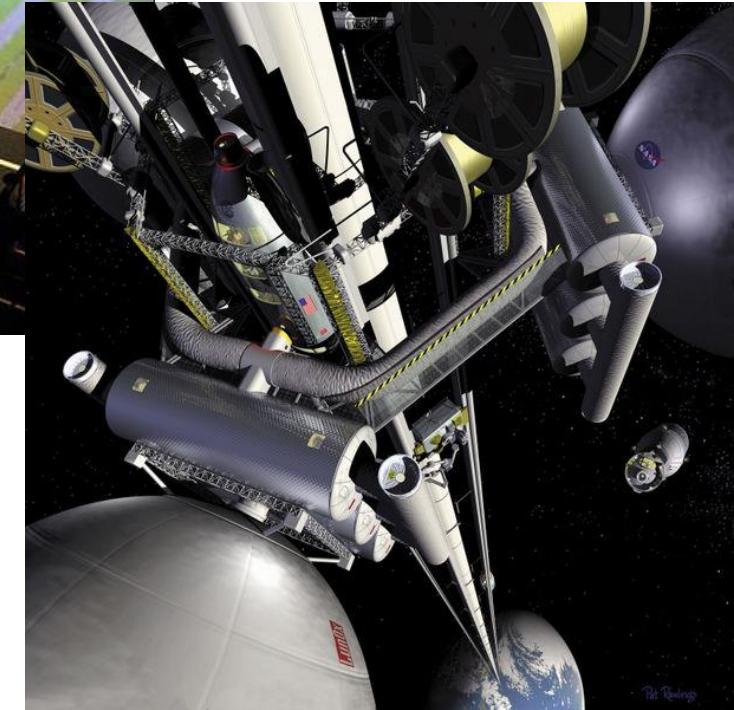
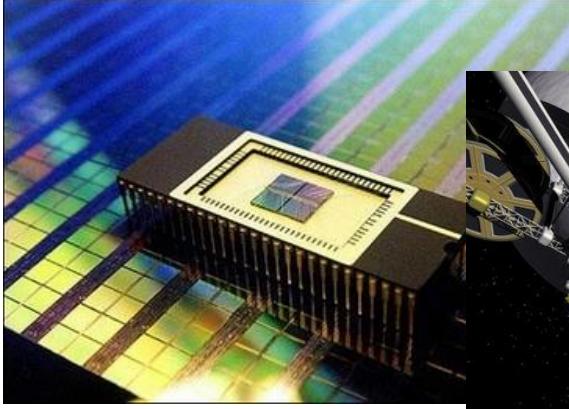


Kevin
Roccapriore



Thorfinn
Roccapriore

The World is Material Opportunity



Predicting crystal structure by merging
data mining with quantum mechanics

CHRISTOPHER C. FISCHER¹, KEVIN J. TIBBETTS¹, DANE MORGAN² AND GERBRAND CEDER^{1*}

¹Department of Materials Science and Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

²Department of Materials Science and Engineering, University of Wisconsin, Madison, Wisconsin 53706, USA

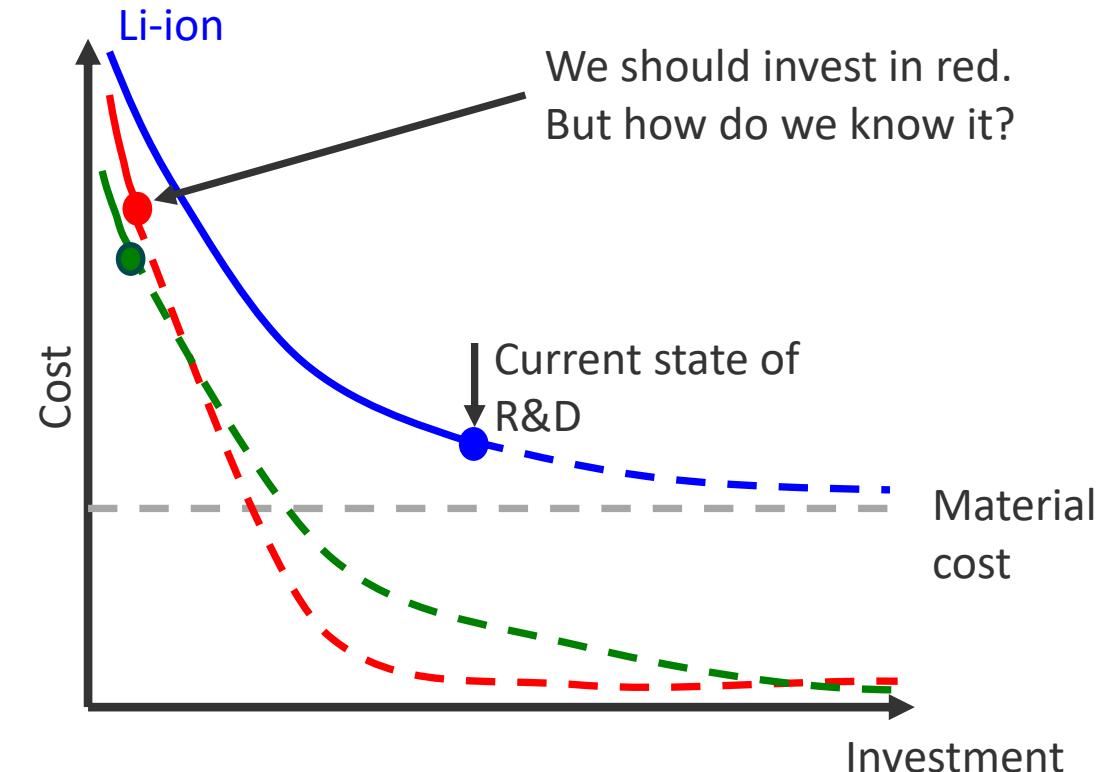
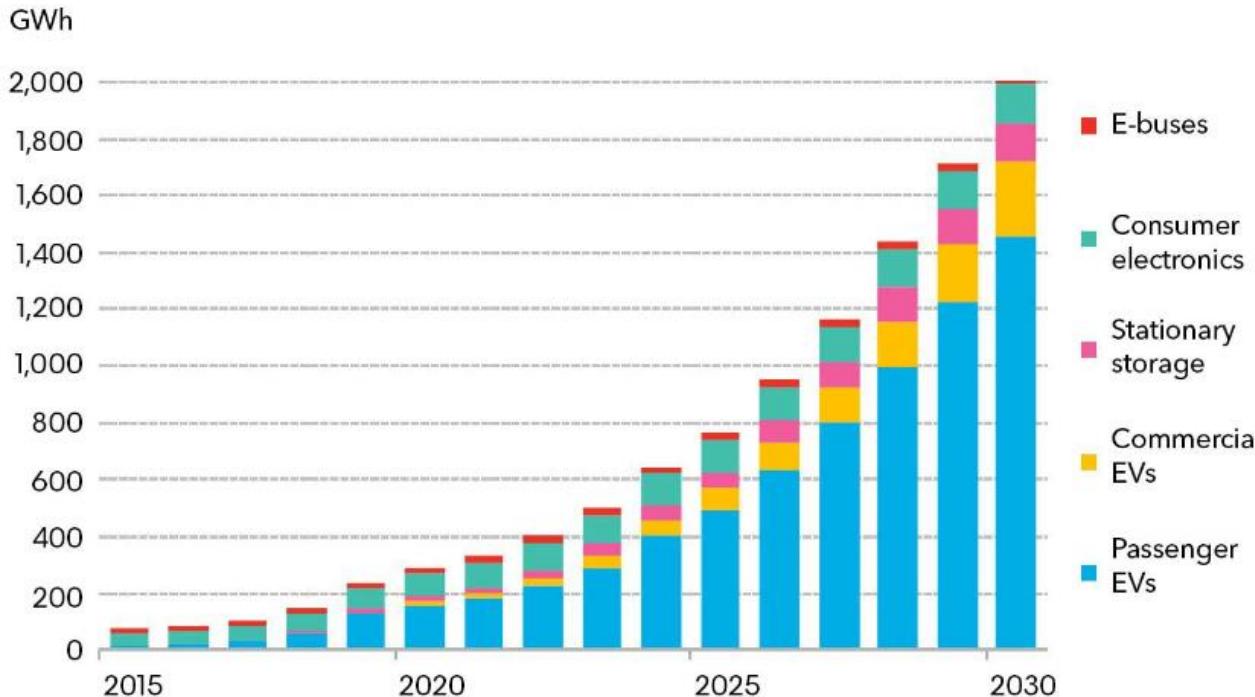
*e-mail: gceder@mit.edu

- “**Improve**”: Renewable energy, self-driving cars, transparent displays, new memory technologies
- “**Discover**”: Room temperature superconductivity, high mechanical stress materials
- “**Engineer**”: Quantum computing, single-atom catalysts, biomolecules

Functionality, manufacturability, cost

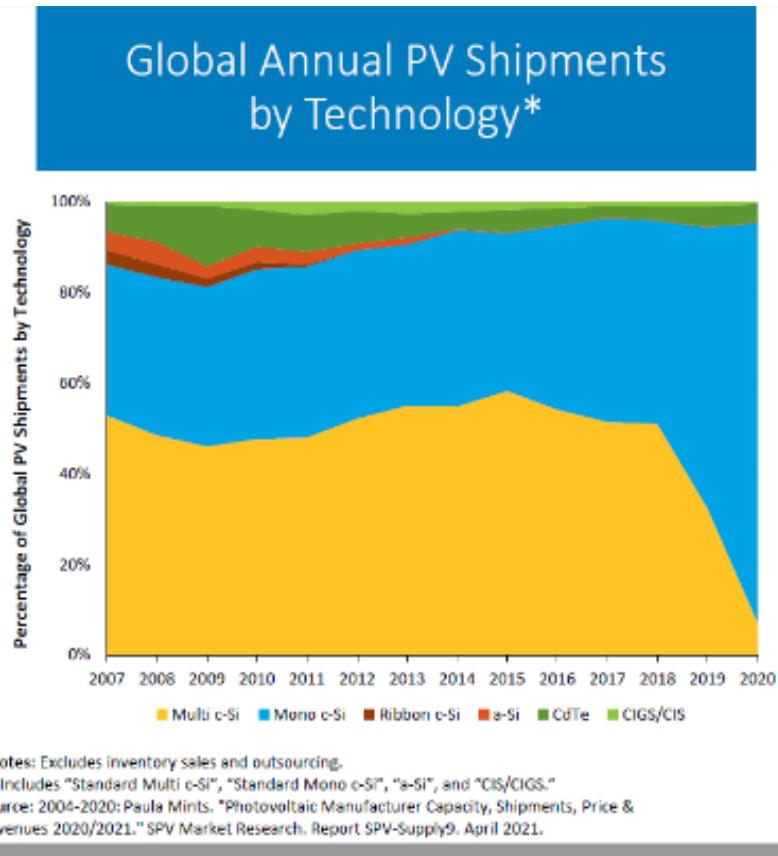
Batteries: Li-ion and Beyond

Annual lithium-ion battery demand



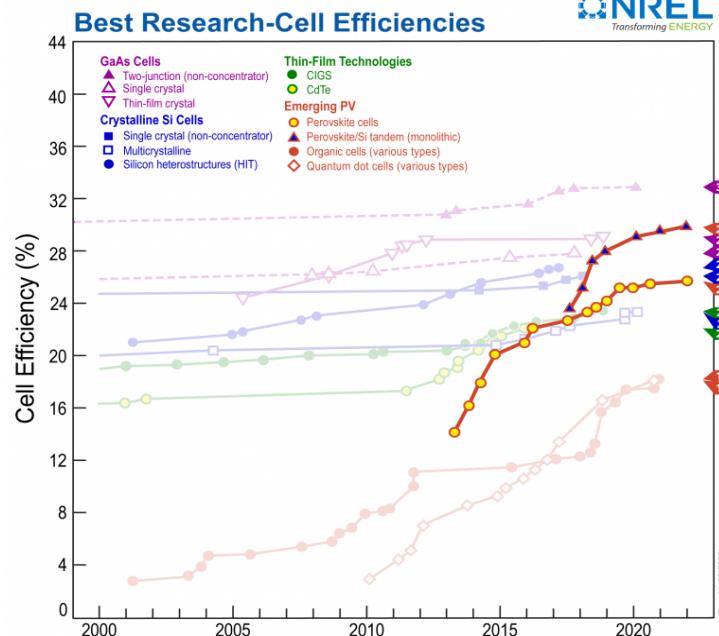
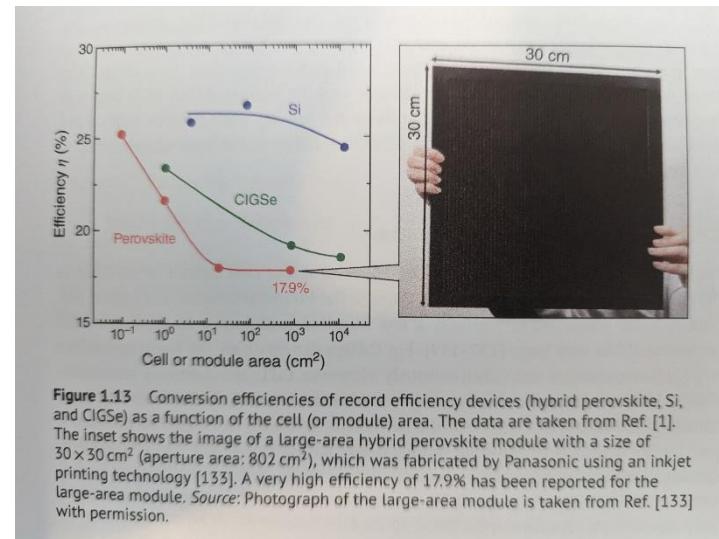
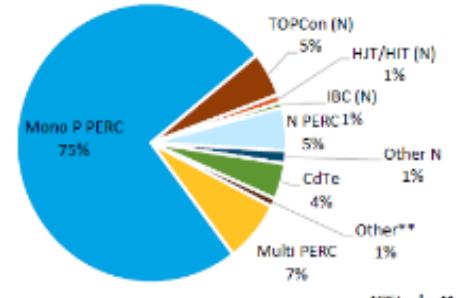
- Batteries are required element of energy transition (EVs, ESS, mobile devices)
- Currently Li-ion is the primary technology
- Optimization of Li-ion batteries takes years (even with same process on new Gigafactory)
- However, it is far from Goldilock zone for ESS or energy transport
- How can we optimize usage and safety for Li-ion batteries in EVs?
- How do we select beyond Li technologies for ESS?

Solar Energy: Will Silicon Ever Reign?



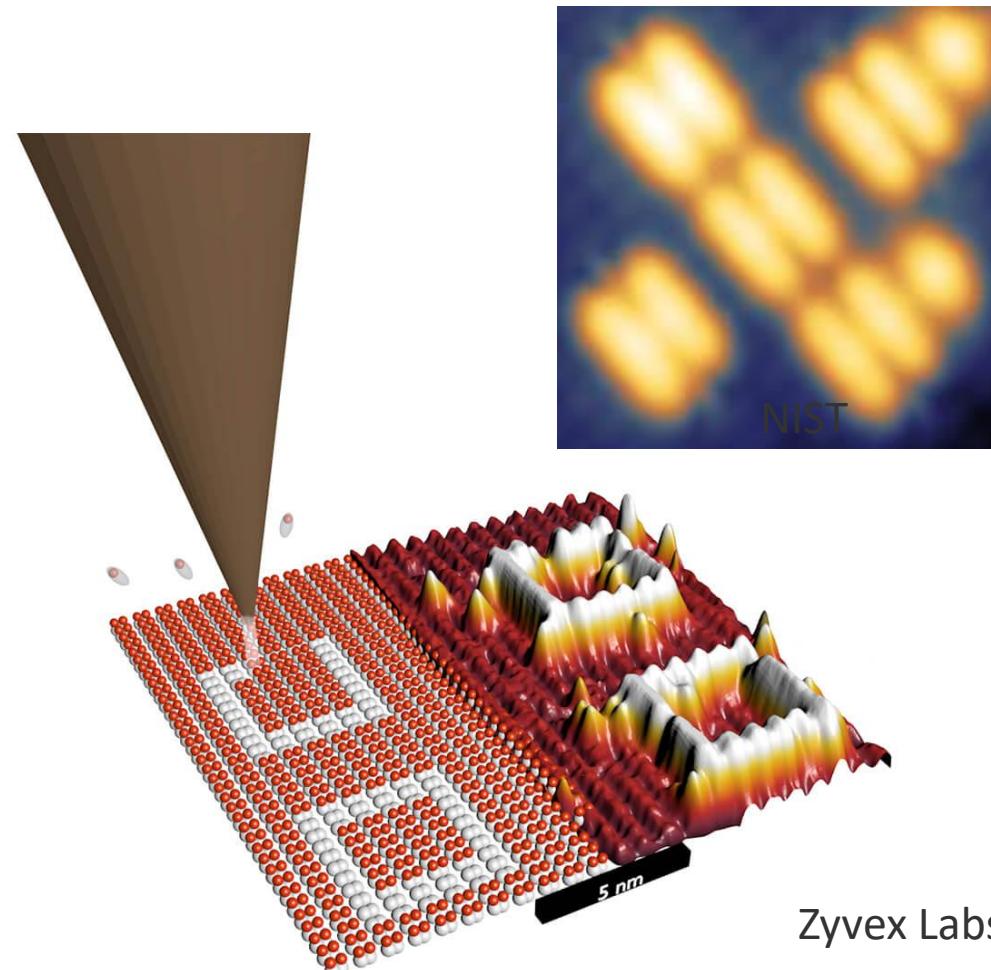
- In 2020, 88% of PV shipments were mono c-Si technology, compared to 35% in 2015 (when multi peaked at 58%).
- Mono P PERC was the dominant cell type in 2020, though n-type shipments grew 181%, y/y, to 13% of the market.

2020 Market Share by Cell Type

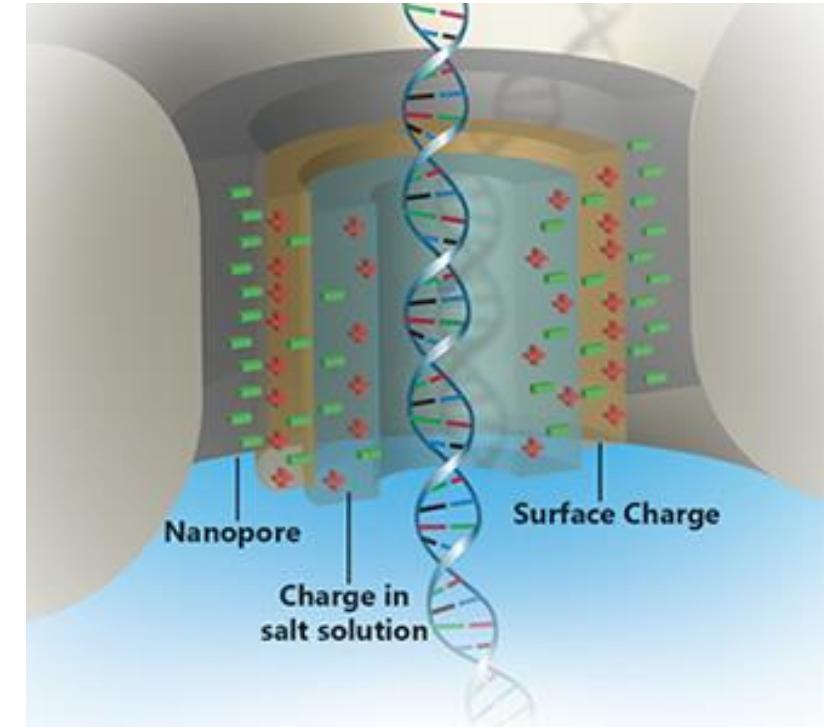


- Solar energy is the fastest growing energy sector
- Si is now reigning material – however, it is really not the optimal material for PV (heavy, expensive)!
- Hybrid perovskites can be used as ideal PV materials – if we can make them stable and scale manufacturing!

Quantum Computing and Single Molecule Bio



Zyvex Labs

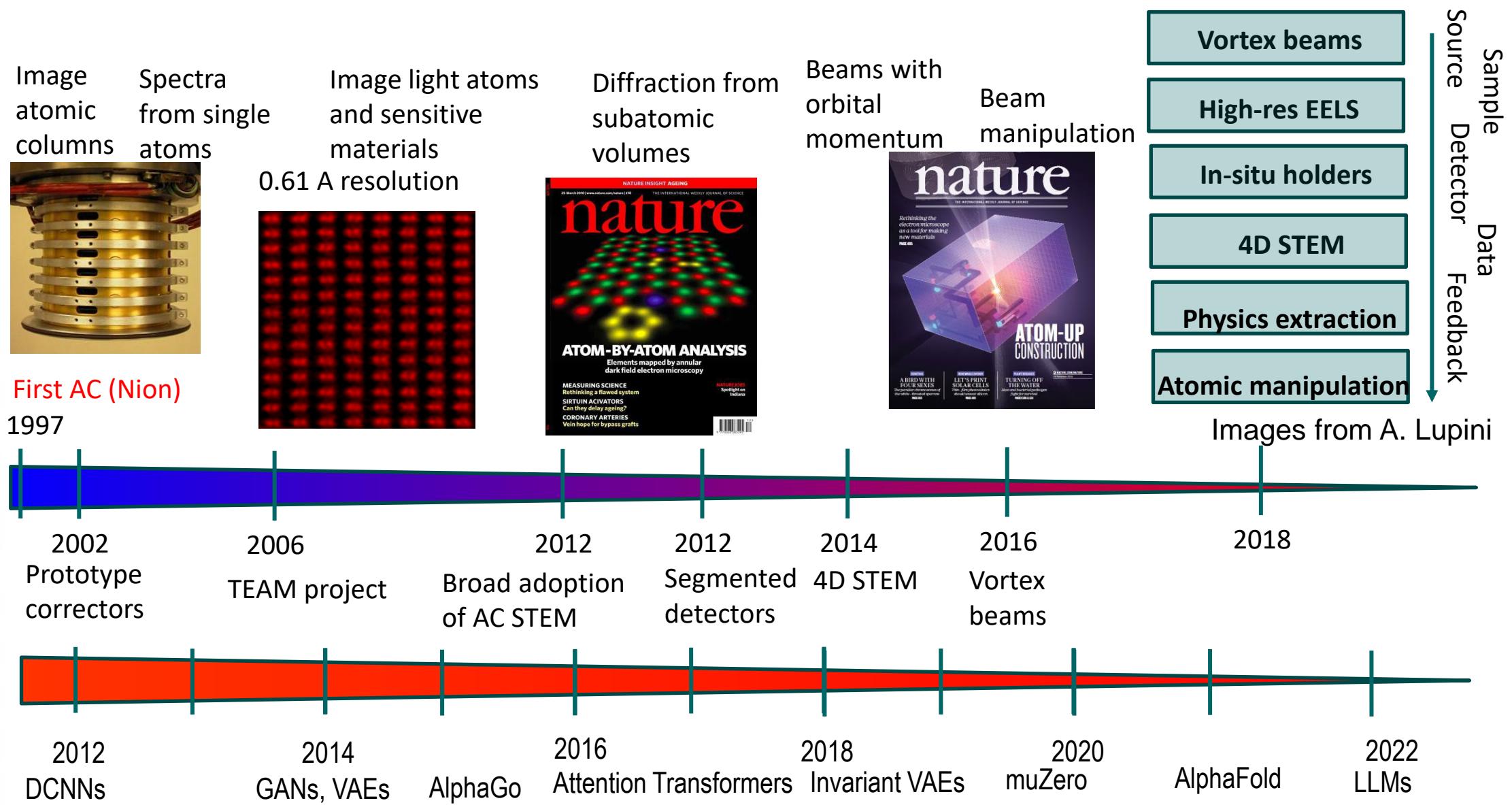


Oxford Nanopore

- Direct atomic fabrication: quantum communications and quantum computing, environmental sensing
- Single-molecule biological devices
- Success story 1: cryo-electron microscopy
- Success story 2: nanoelectron diffraction

We need to go small to go big!

The Lab on the Beam



However

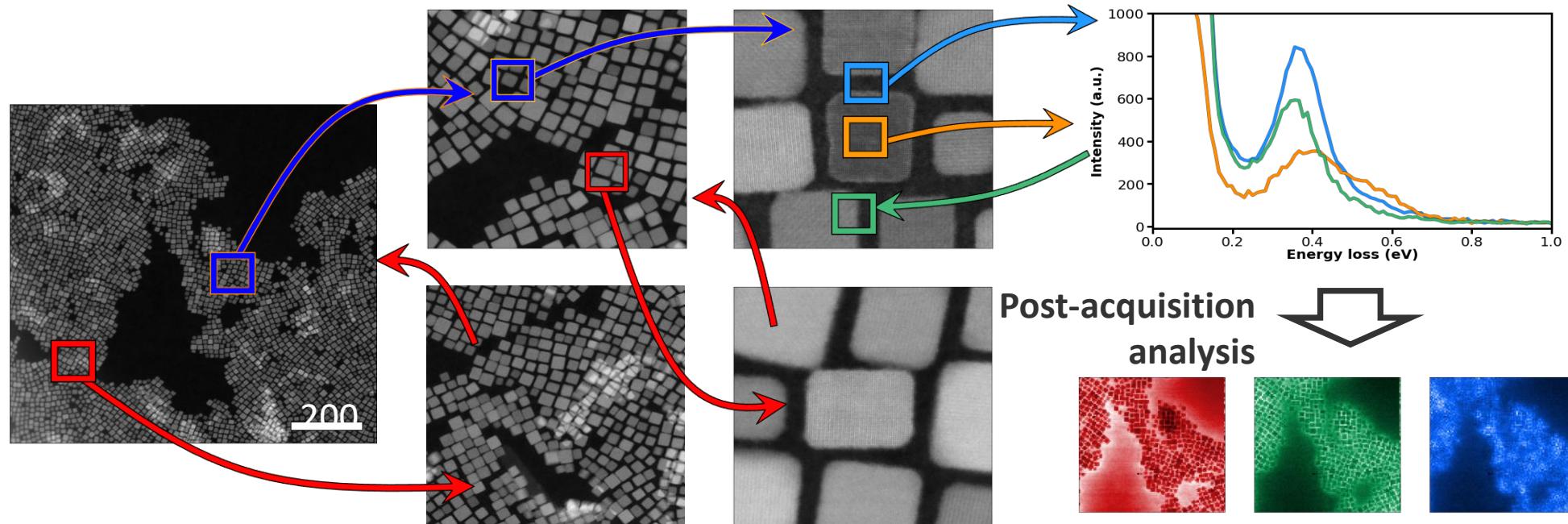
We run the instruments the same way as decades ago:

- Rectangular scans
- Point spectroscopy
- Hyperspectral imaging (spectroscopy on the grid)

Opportunities:

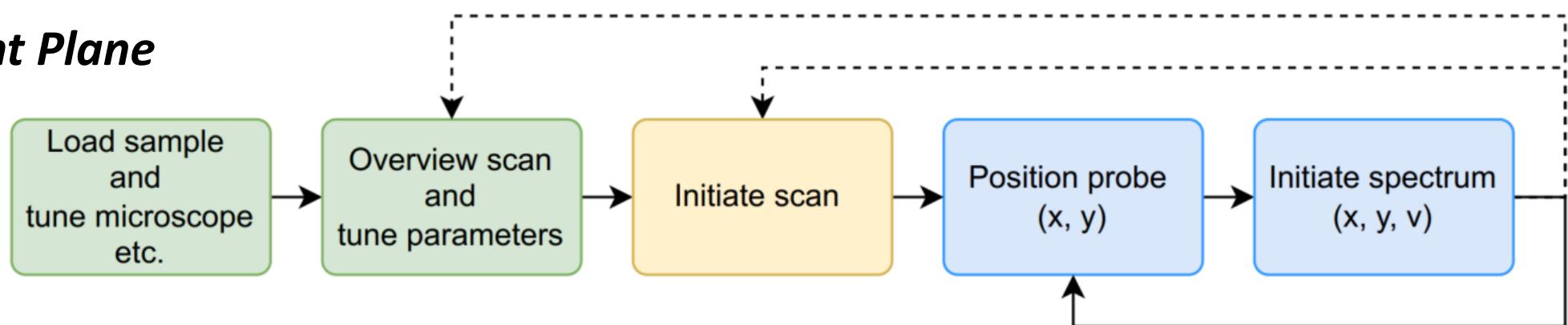
- Very high sampling images ($>2k$)
- Multidimensional spectroscopies
- Objects of interest are non-uniform in image plane
- Why microscopy?

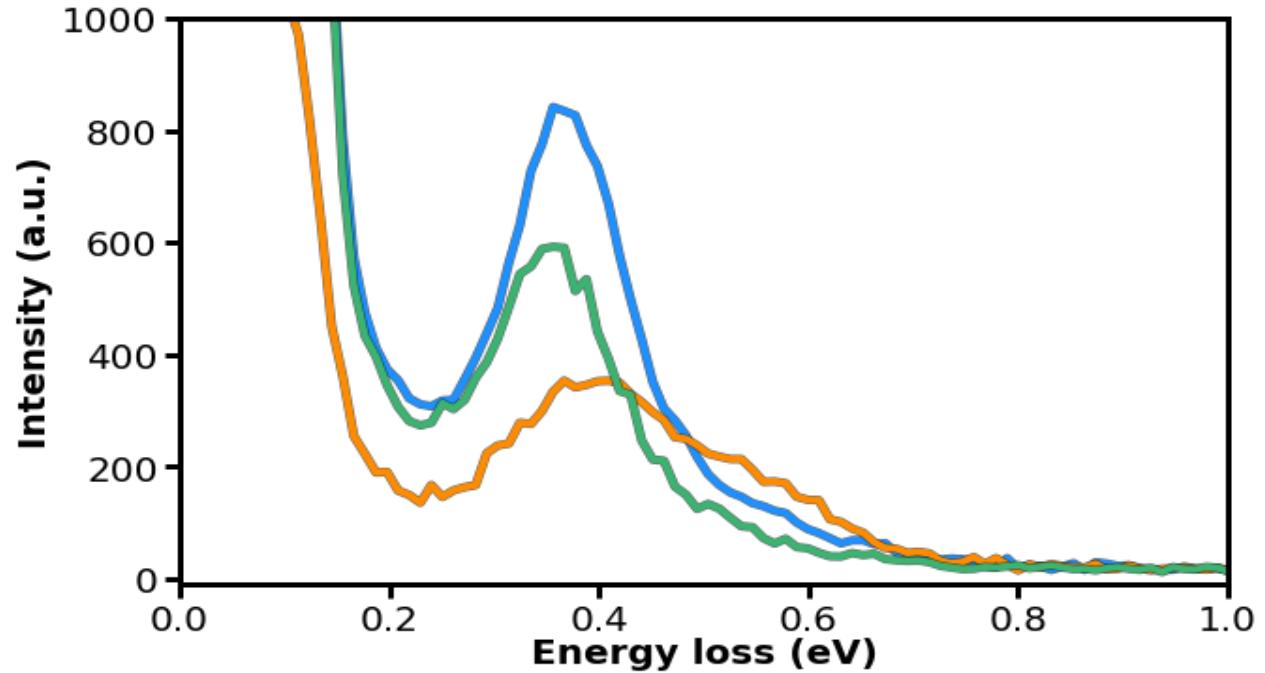
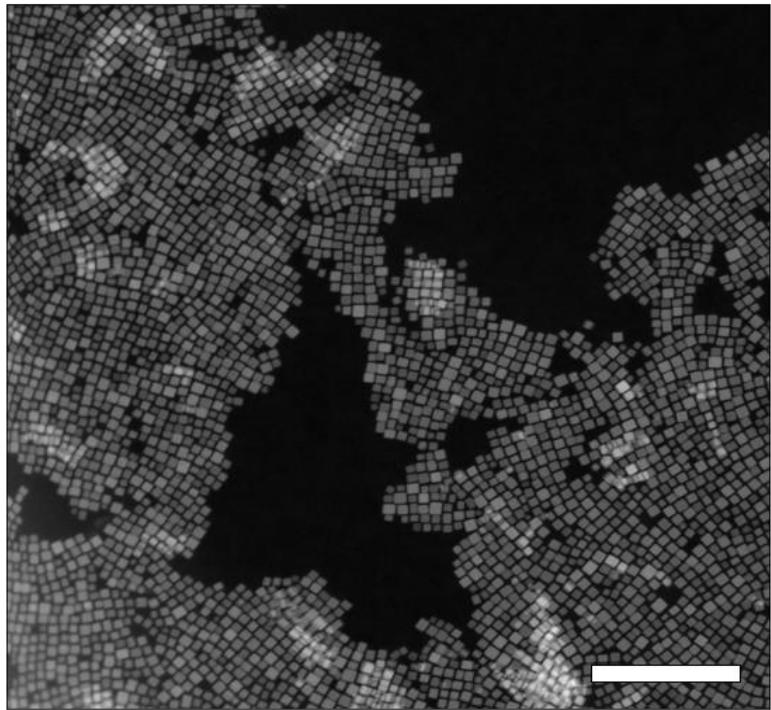
Workflows in STEM



Minimal instruction set control language

Instrument Plane





- What can we say about structure?
 - Interesting functionalities are expected at the certain structural elements
 - We can guess some; we have to discover others
 - Multiple goals while running experiment
-
- **Policy:** **what do we do depending on observation**
 - **Reward:** **what do we hope to achieve**
 - **Value:** **anticipated reward**

This workshop:

Learn how to transition from image
quantification to ML-enabled workflows

Tentative Program

1. Outline and structure of the course - June 6
2. Imaging in Scanning Transmission Electron Microscopy - June 9
3. Spectroscopy in STEM - June 13
4. Linear methods and dimensionality reduction for spectral data - June 16
5. High-resolution and Z-Contrast Imaging -June 20
6. Image registration methods - June 23
7. Linear methods and dimensionality reduction for imaging data - June 27
8. Diffraction and 4D STEM - June 30
9. Bringing Cloud and Edge to STEM: from tool to ecosystem - July 4
10. Image simulations - July 7
11. Deep convolutional networks - July 11
12. DCNN for image data - July 14
13. DCNN case studies - July 18

14. Gaussian processes and Bayesian Optimization - July 21

15. Bayesian Inference, Structured GP, and Hypothesis Learning - July 25

16. Variational Autoencoders – 1 - July 28

17. Variational Autoencoders – 2 - August 1

18. Encoders-decoders and structure-property relationships - August 4

19. Special topic: VAE for any tasks - August 8

20. Deep kernel learning: EELS and 4D STEM – August 11

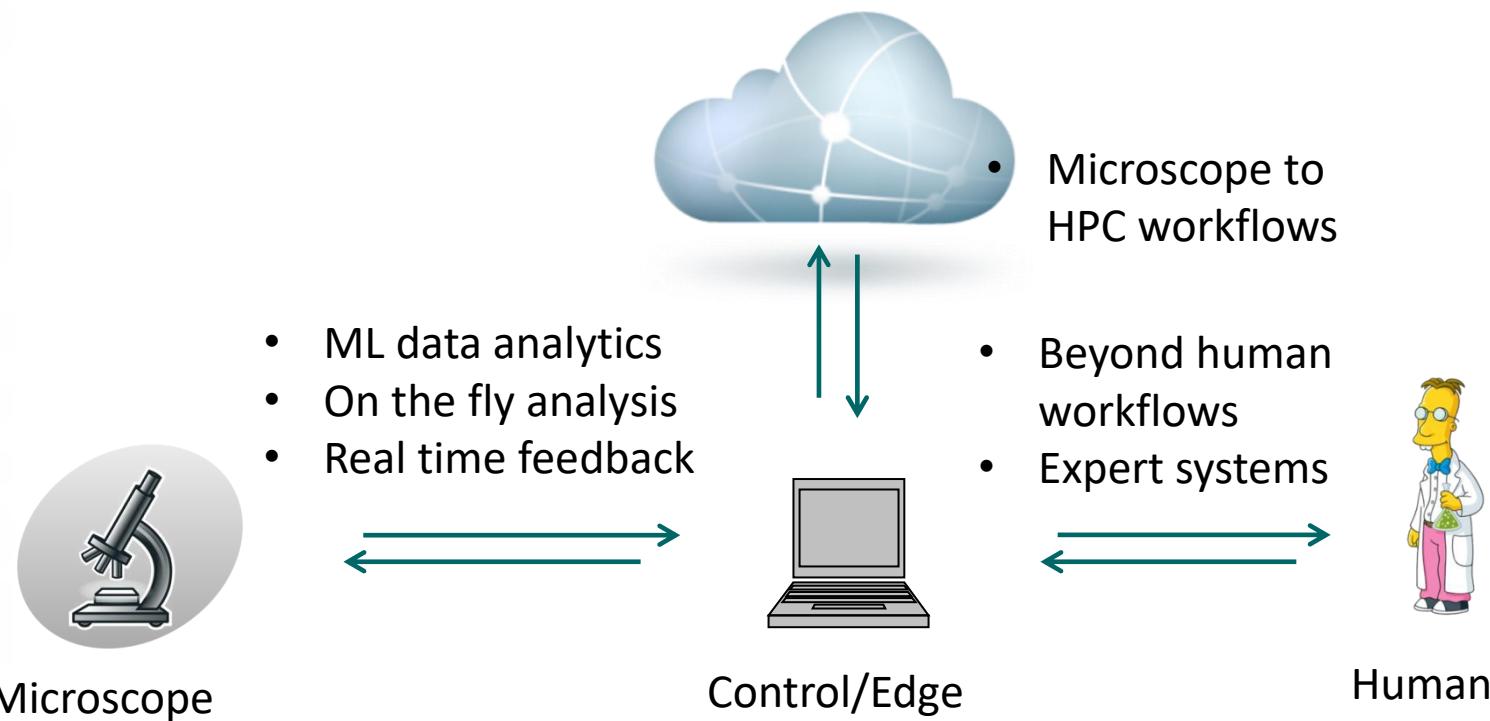
21. DKL forensics and human in the loop - August 15

22. Special topics: Reinforcement learning - August 18

23. Special topics: Learning physics from images - August 22

24. Special topic: Causality - August 25

1. Over 15-year experience in applying ML to experimental physical problems
2. Over 5-year experience in automated and autonomous microscopy via edge computing
3. Experience in ML for materials synthesis and synthesis/characterization workflow design
4. Leading teams developing data infrastructure infrastructure for scanning probe and electron microscopy

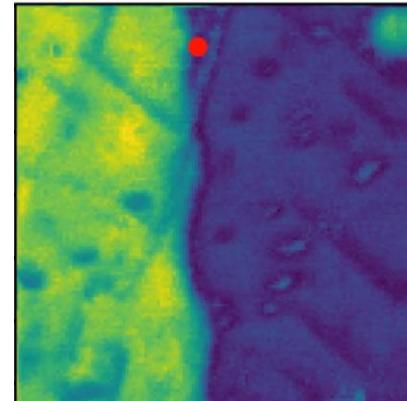


Sergei V. Kalinin

- Ph.D. 2002: U. Penn
- 2002-2022: Oak Ridge
- 2022-2023: Amazon, GC
- 2023: U. Tennessee, Knoxville

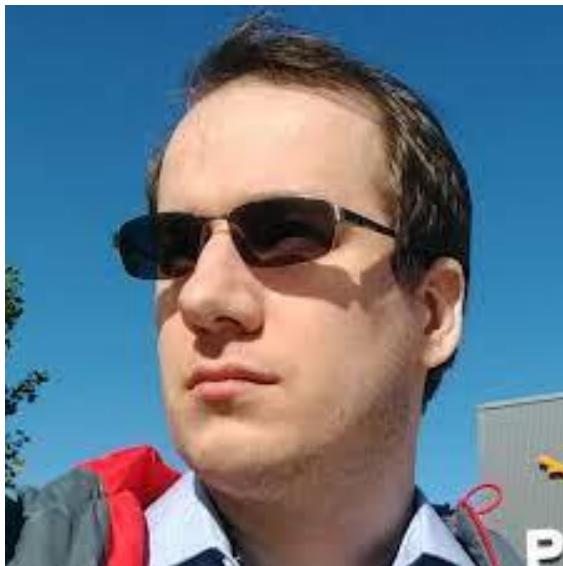
Rama Vasudevan

THE UNIVERSITY OF TENNESSEE  KNOXVILLE



- Science Interests: Autonomous labs, microscopy, open-source scientific software, machine learning and reinforcement learning
- Currently group leader of Data NanoAnalytics at CNMS/ORNL
- Running/Cricket/Gardening. Former opening bowler/now relegated to part time pie thrower
- Avid Gardener. Wannabe Farmer.

Maxim Ziatdinov



Dr. Maxim Ziatdinov, who has devoted his life's work to the masterful synthesis of machine learning, experimental practices, and theoretical principles, all aimed to expedite discoveries in the grand tapestry of physical sciences.

He has meticulously woven together science-informed machine learning algorithms, as detailed and intricate as the astrolabes of the Citadel, each informed and refined by prior realms of knowledge. His scholarly endeavors have also led to the creation of crucial pathways, akin to ancient Valyrian roads, connecting the latest instrumental platforms with the sheer might of high-performance computing facilities.

A master craftsman of the modern age, he has forged numerous software artifacts, each imbued with a unique purpose, much like the famed blacksmiths of Qohor. Among these are **AtomAI**, a tool as keen as Valyrian steel in the realm of deep learning applications for microscopy, **pyroVED**, a mystical sigil for invariant autoencoders in image and spectral analysis, and **GPax**, akin to the magical weirwood trees, it empowers physics-based active learning and Bayesian optimization in automated experiments. Each tool is shared freely in the community, like a lord sharing his feast with his people, contributing to a broader and deeper understanding of our world.

In the dominions of knowledge and discovery, Dr. Ziatdinov is not a solitary wanderer, but a cherished member of two esteemed brotherhoods. He holds vital roles in the noble houses of the CSED and CNMS at the towering fortress of ORNL, both renowned for their dedication to unraveling the mysteries of our universe.

Gerd Duscher



Research

- Electrical activity at interfaces and crystal defects
- Atomic and electronic structure at interfaces
- Atomic resolution Z-contrast imaging
- Atomic resolution electron energy-loss spectroscopy
- Atomic structure calculations with density functional theory

Dr. Duscher earned his master's degree (Diplom) in Physics from the University of Regensburg, Germany in 1990, and his Dr. Sci. in Materials Science and Engineering from the University of Stuttgart in 1996.

His dissertation research was performed at the Max-Planck Institute for Metals Research, Stuttgart, where he held a Fellowship. He used analytic TEM methods at both his masters and his doctoral research to investigate the relationships between the atomic structure and macroscopic properties. In his dissertation research he studied the influence of oxide precipitates on the electrical activity of grain boundaries in silicon.

In 1997, Dr. Duscher joined the group of Dr. Steve Pennycook in the Solid State Division at the Oak Ridge National Laboratory as a Postdoctoral Research Associate where he studied grain boundaries in oxides. In 1998, Dr. Duscher was a staff member of the Max-Planck Institute and in 1999 he accepted a position as a Research Assistant Professor for Theoretical Physics with Dr. Sokrates Pantelides at Vanderbilt University.

In 2001 Dr. Duscher joined the faculty at the North Carolina State University, Raleigh, and in 2006 he became the director of the Silicon Wafer Engineering and Defect Science a NSF IUCRC center. In 2008 Dr. Duscher joined the faculty at The University of Tennessee . He has authored or co-authored more than 100 technical papers and has contributed to over 50 (25 invited) technical presentations at national and international conferences.

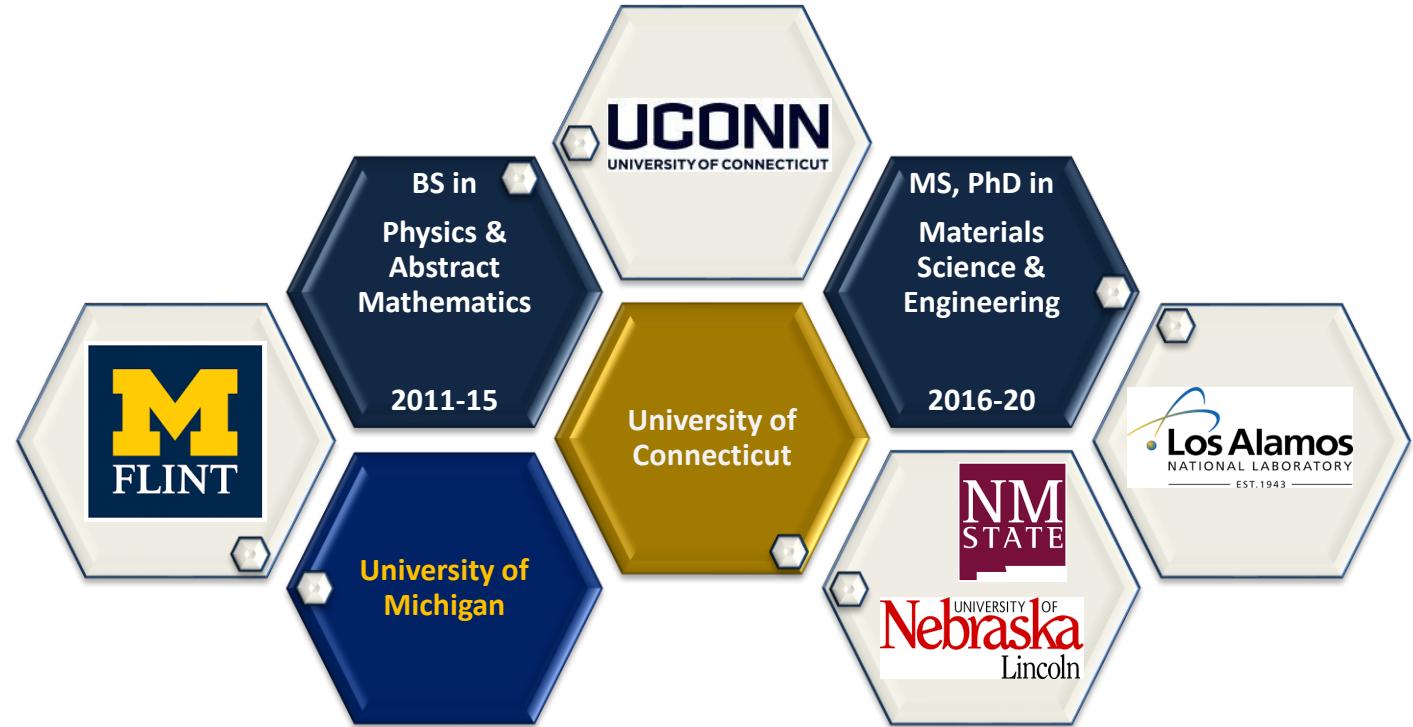
Ayana Ghosh



Postdoc at ORNL
(Sept. 2020 – Feb. 2023)

Research Scientist at ORNL
(March 2023 -- Present)

- Materials
- Physical Models
- Machine Learning & Causal analysis
- Combine Theory & Experiments



**Kevin Roccapriore:**

- Staff scientist at Oak Ridge National Laboratory's Center for Nanophase Materials Sciences (CNMS).
- PhD in Physics from the University of North Texas in 2018, which was nanofabrication for nanophotonics applications.
- Post-doc at ORNL (2018-2021), focusing on machine learning in the scanning transmission electron microscope (STEM) toward automated experiments, including STEM-EELS and 4D-STEM.
- Currently, focuses on ML on-the-fly in the STEM for precise atomic fabrication by site-specific atom targeting with the electron beam, automated experiments learning and utilizing structure-property relationships, and automated multidimensional data acquisition and analysis.

**Tommy (Chun Yin) Wong**

- Research interests: deep learning & computer vision, electron microscopy, radiation damage
- 3rd year Ph.D. student in Energy Science & Engineering, Bredesen Center, University of Tennessee
- M.S. in Materials Science & Engineering '22, University of Tennessee
- B.S.E. in Nuclear Engineering & Radiological Sciences '20, University of Michigan

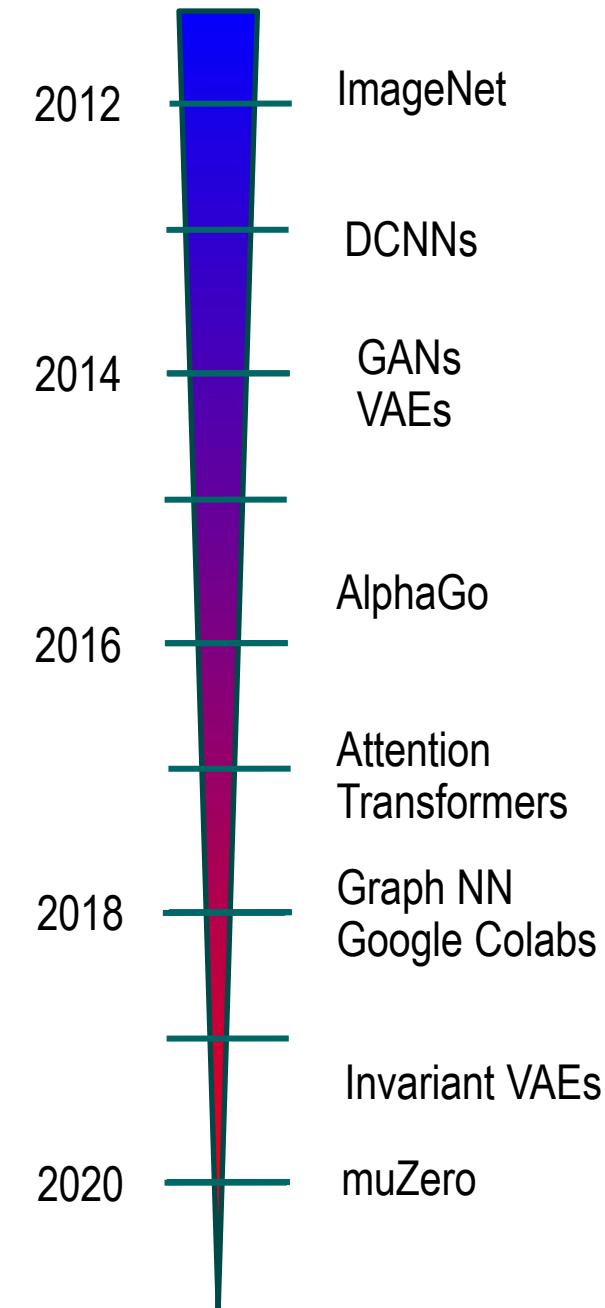
Why Machine Learning?

- Last decade has experienced an explosive growth of machine learning and artificial intelligence applications
- These developments have spanned areas from computer vision to medicine to autonomous systems and games
- However, the progress and impact as applied to experimental physical sciences has been minimal....

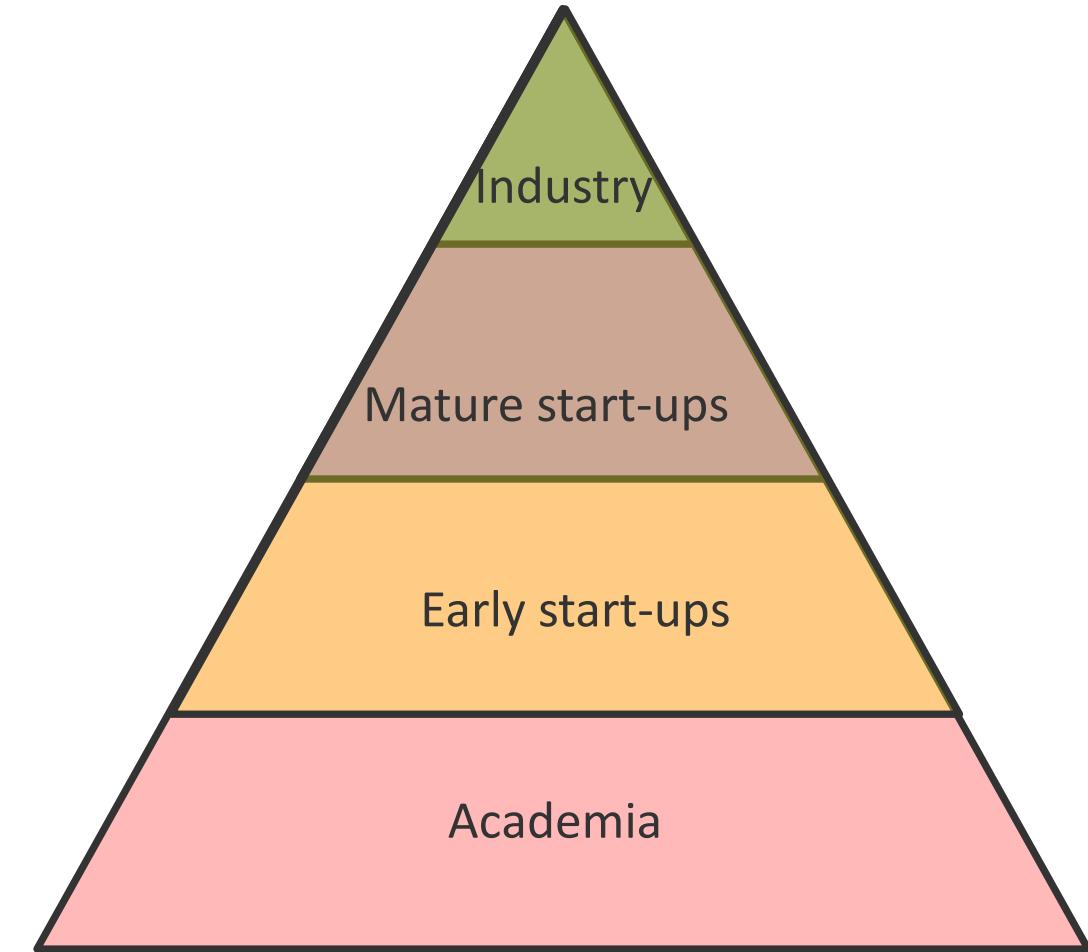
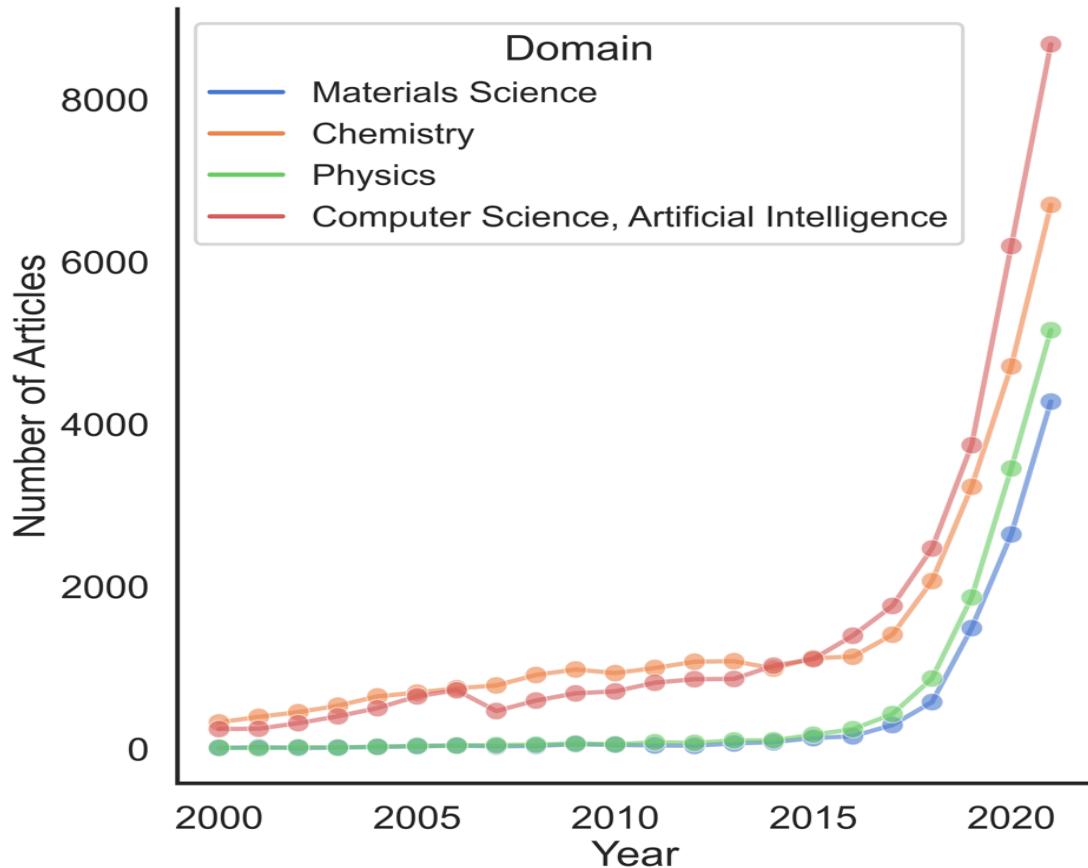
Why is it difficult?

- Requires domain expertise and domain-specific goals
- Deeply causal and hypothesis drive nature of domain sciences
- No single answer: culture, not a method
- Infrastructure, open code, open data
- **Most important:** active nature of scientific process

Microsoft: GitHub
Meta: Open Catalyst,
Meta: Papers with Code
Toyota: TRI
Google: AlphaFold
NVIDIA: protein folding



ML in Domain Sciences



Analysis by B. Blaiszik, Argonne

- The rapid adoption of ML in domain sciences and industrial R&D is a very recent trend
- Technologies and workforce emerge from academia into industry
- We can estimate potential growth rates comparing to cloud computing 15 - 20 years ago

“Eras” of ML in Industry

- **Before 2000:** It's all about IT (dotcoms, Amazon, etc)
- **2000 - 2010:** It's all about collecting and searching data (Facebook, Google, Uber, etc.)
- **2010 – 2020:** What do we learn from data (correlative era)
- **2020 – now:** Physics is the new data

[arXiv:2204.05095](https://arxiv.org/abs/2204.05095)

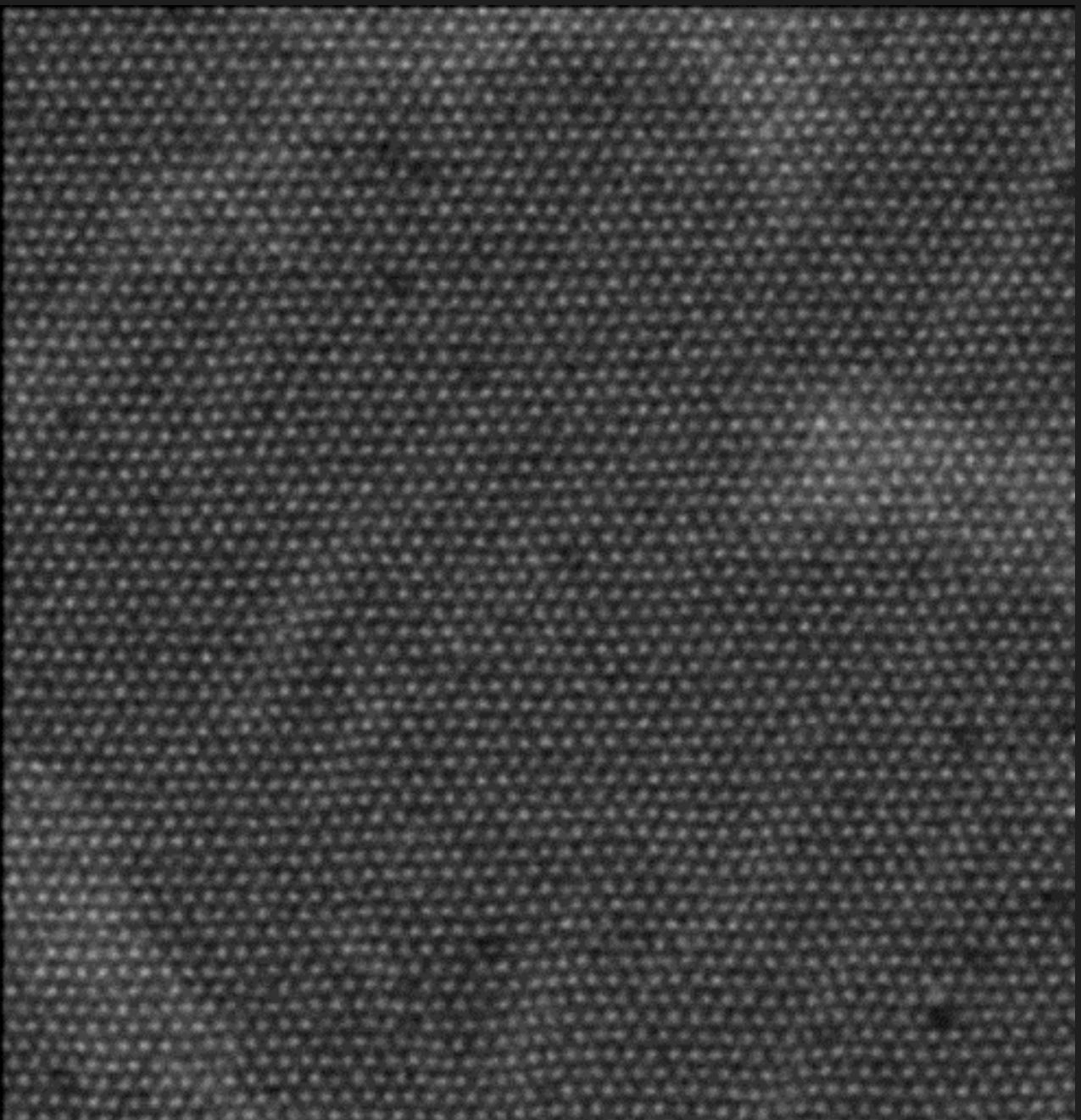
- Classical machine learning is underpinned by the existence of the large static data sets – from MNIST to emerging medical, bio, faces, etc.
- Real world problems are associated with the large distribution shifts, often small data sets, and presence of uncontrollable exogenous factors
- Also, real world problems are often active learning: we interrogate the data generation process and provide feedback, not deal with static data sets
- However, we often have extensive prior knowledge of past data, physical laws generalizing them, and strong set of inferential biases

[arXiv:2005.01557](https://arxiv.org/abs/2005.01557)

ML for real-world applications is different!

Can Machine Learning help us
analyze the acquired data?

Dynamic Atomic Changes



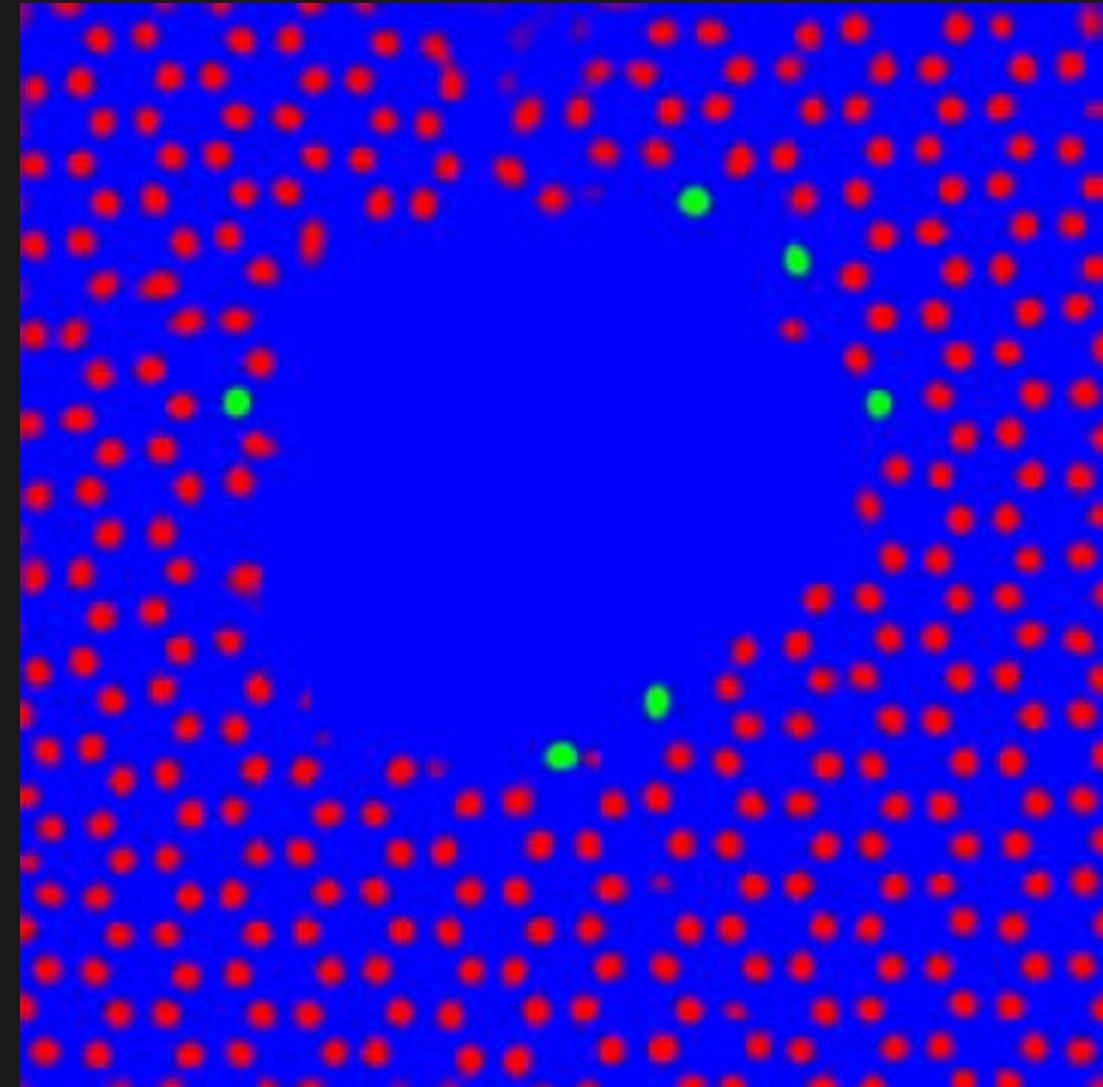
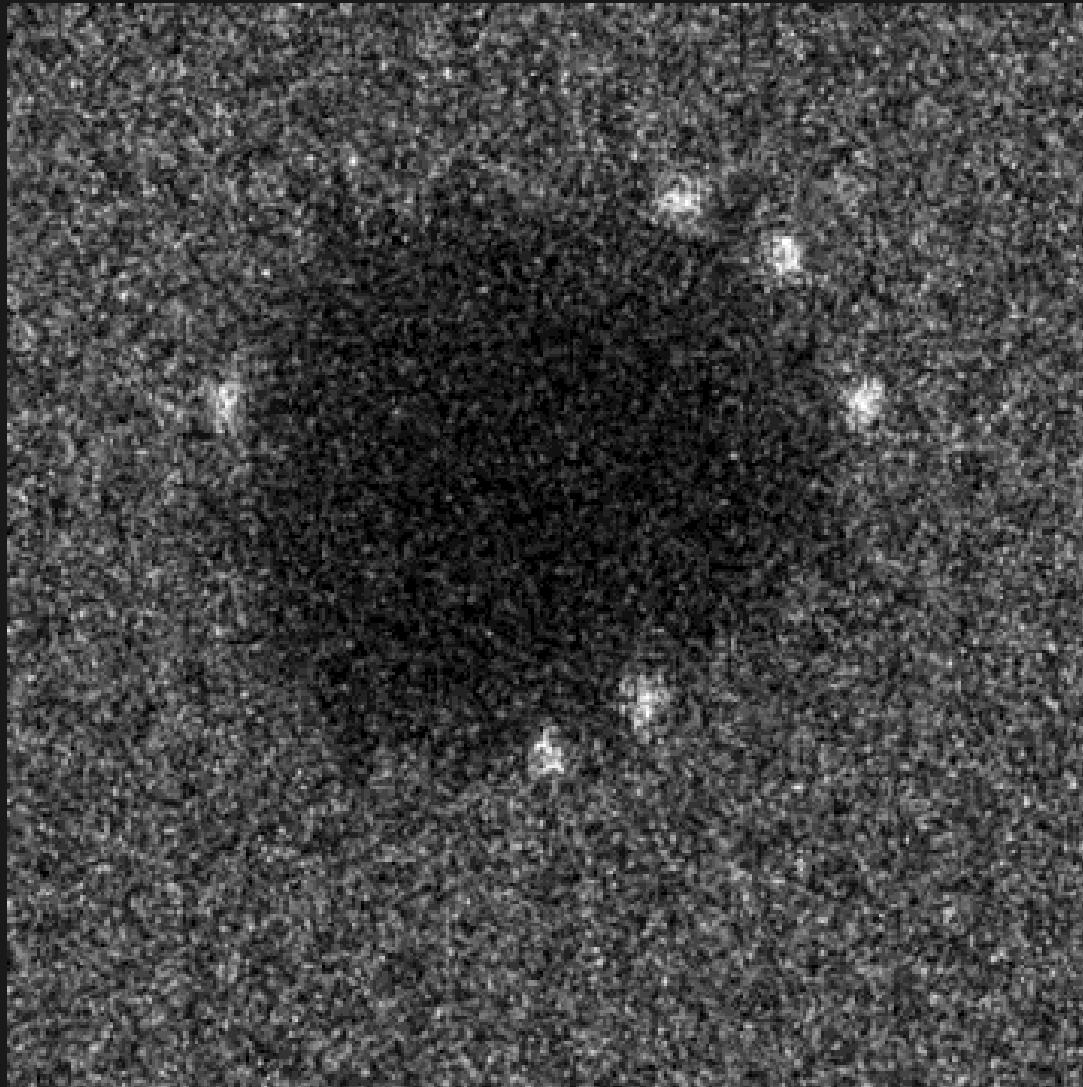
Observations of atomic dynamics induced by beam (or temperature, field, etc.) gives information on multiple atomic configurations as they form and evolve. Can we:

1. Minimize beam damage:
 - a. Low dose imaging
 - b. Low beam energy
 - c. Tailored measurement grid?
2. Can we learn the transformation mechanisms from atomic dynamics?
3. Can we learn force fields between the atoms?
4. Do we really have to take EELS and 4D STEM measurements everywhere?
5. Can we learn how to **direct** beam induced transformations atom-by-atom?

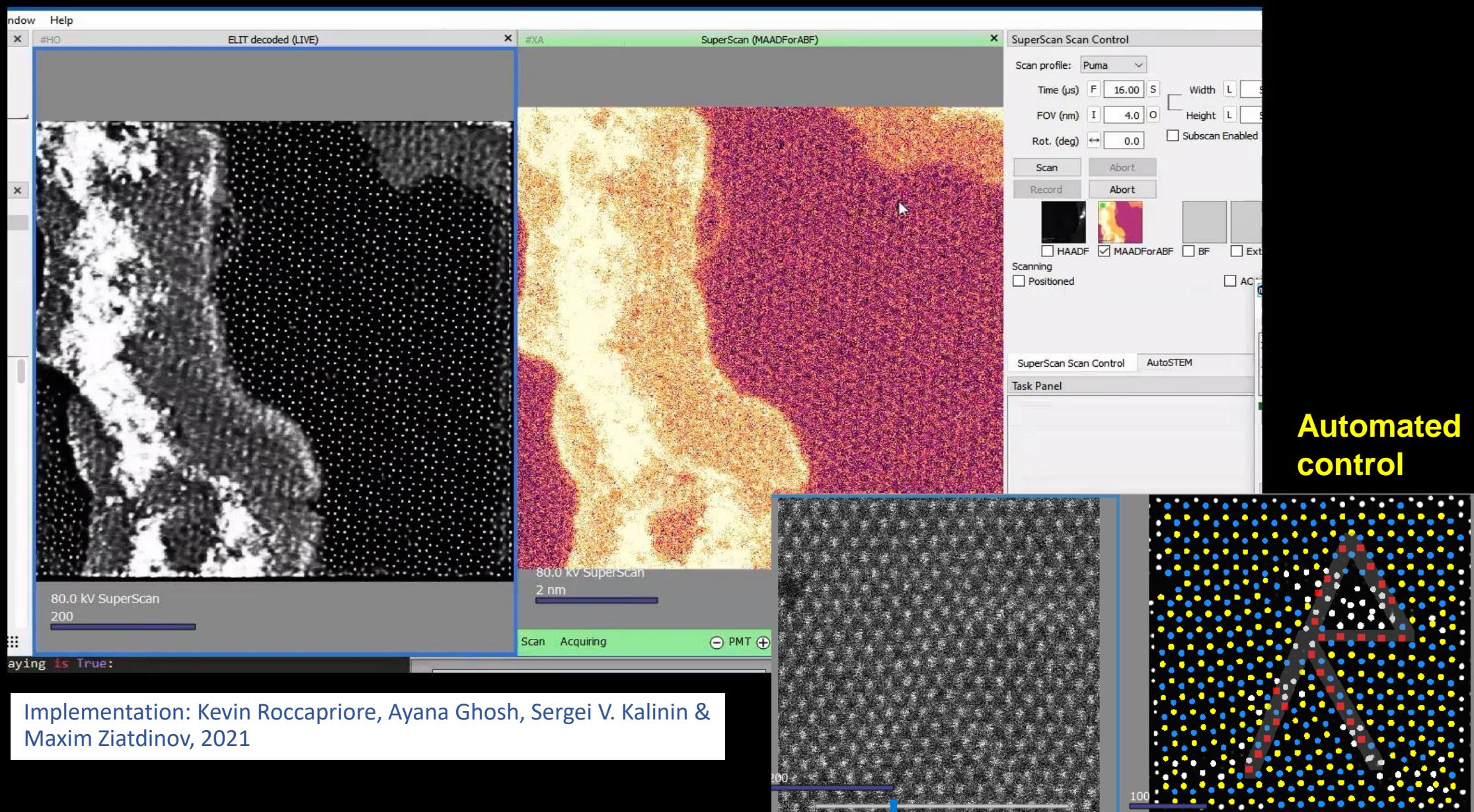
Deep learning works like a charm for:

- Drift correction
- Denoising
- Data processing/dimensionality reduction
- Feature finding (physics is in the training set)

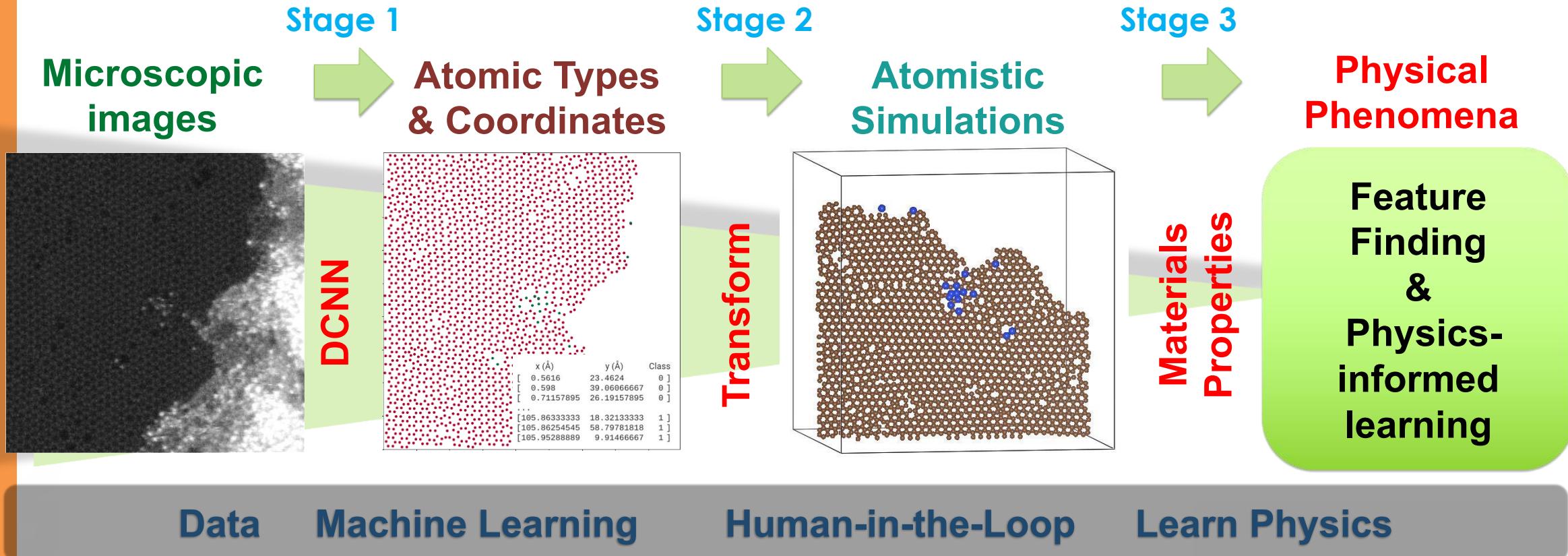
M. Ziatdinov et al, ACS Nano 2017



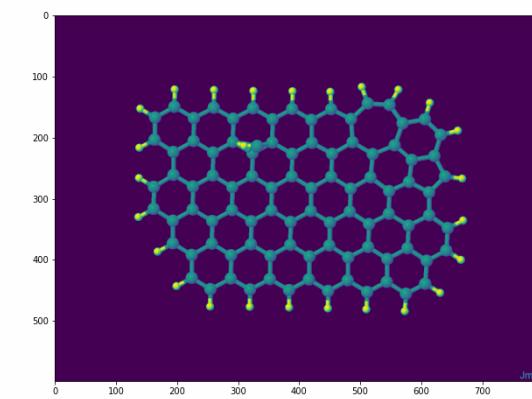
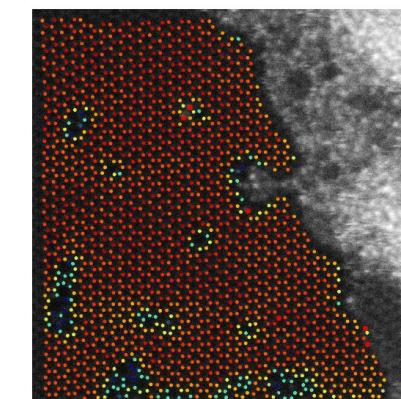
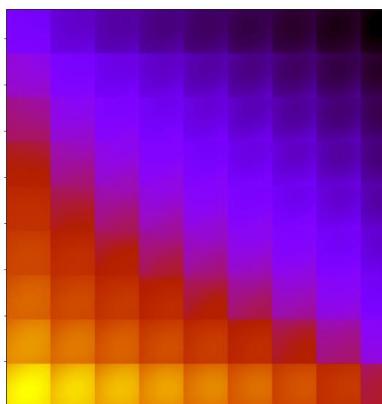
Automated Experiment:
almost easy.... If you know what you are looking for



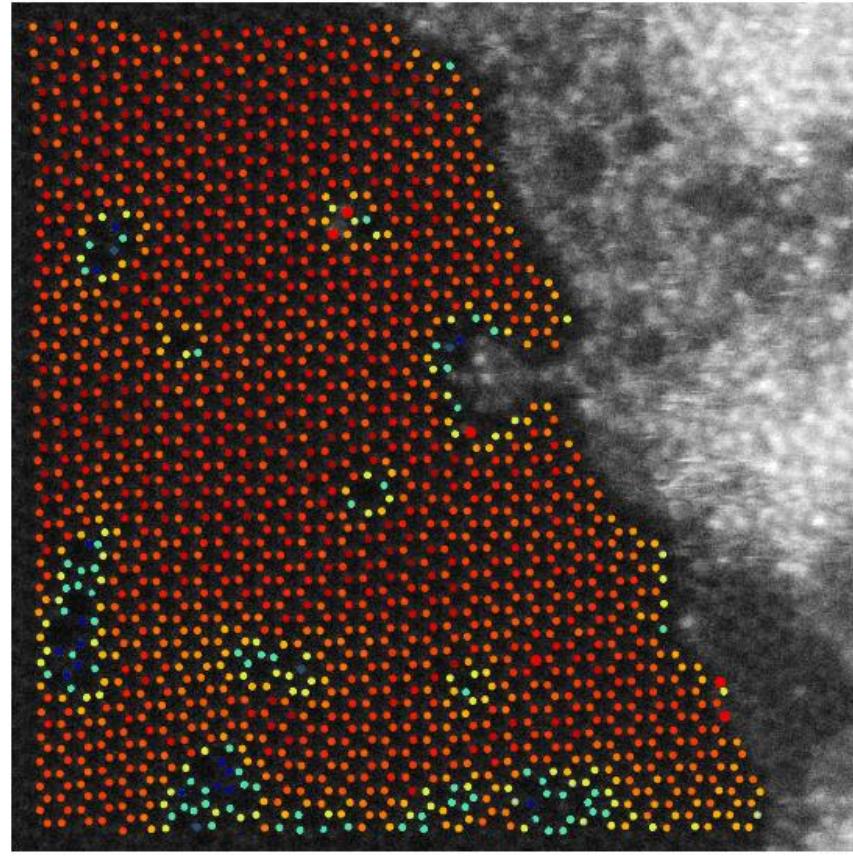
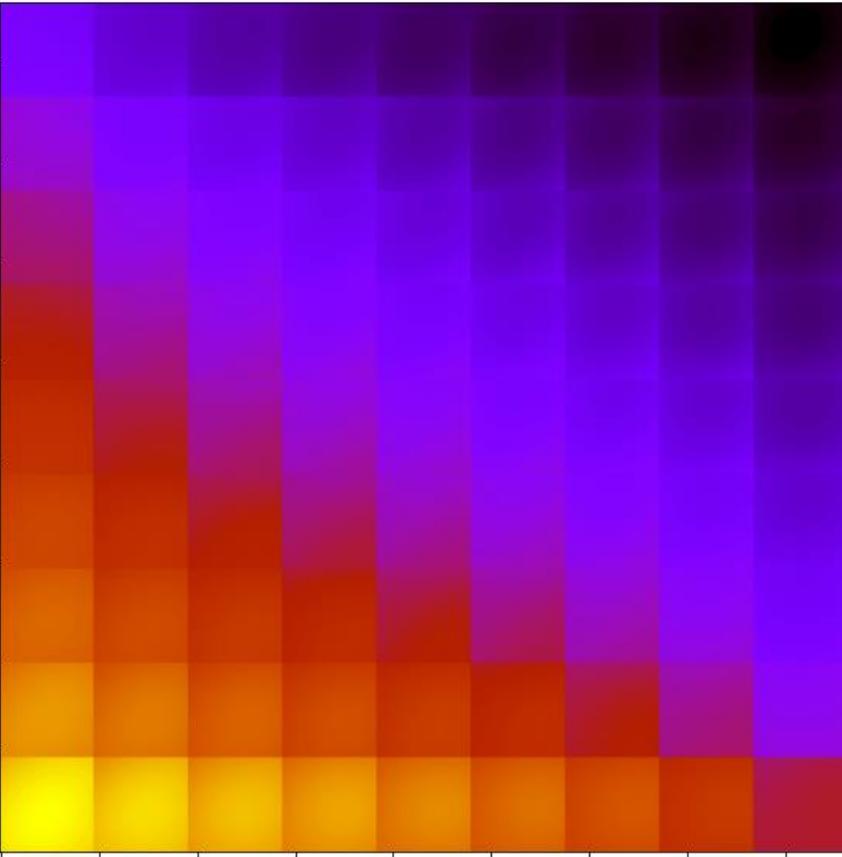
From Microscope to Simulation Environment



- DCNNs for atom finding
- Invariant VAEs for physics discovery
- Conditional VAEs
- GP reconstructions
- Piping data in DFT/MD environment

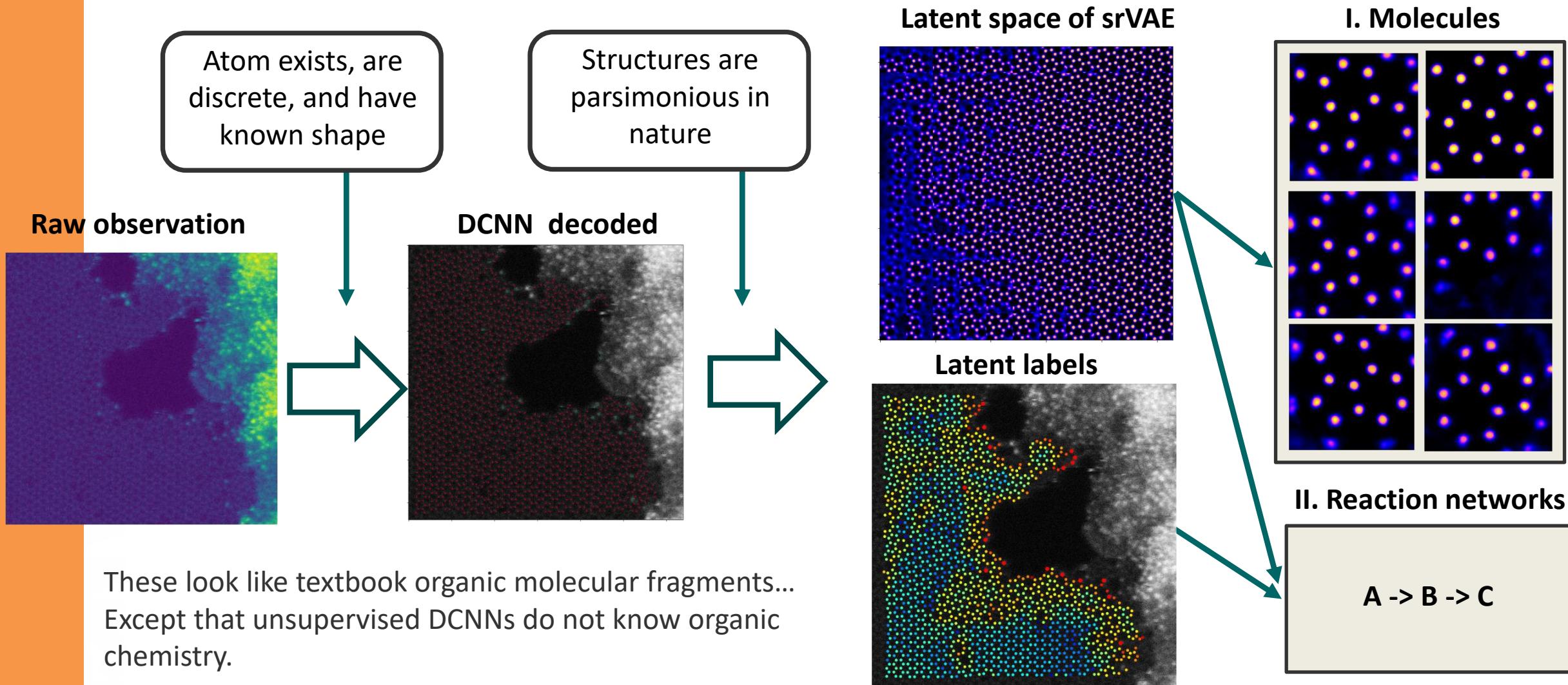


Discovering Building Blocks

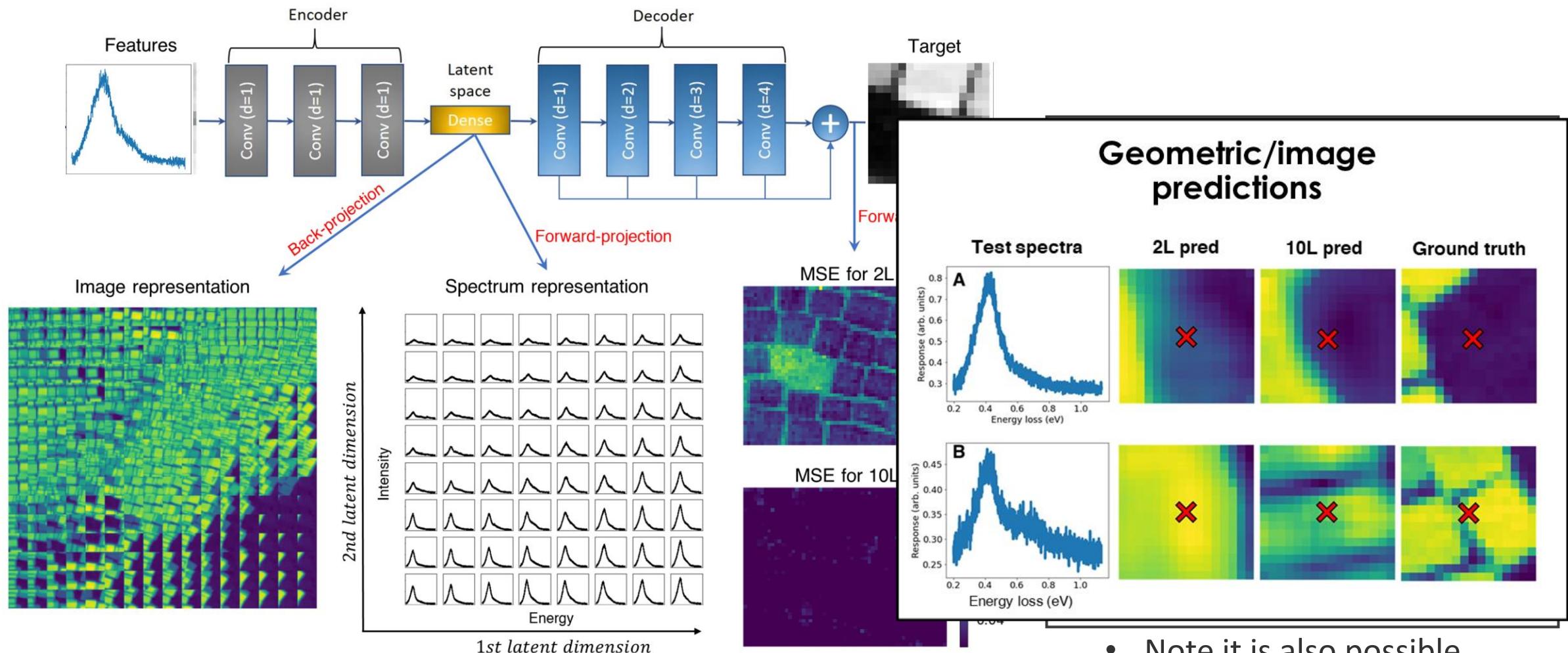


- Classical physical descriptions can be defined locally only in Bayesian sense
- We can argue that local descriptors are simple, if not necessarily known
- And the rules that guide their emergence are also simple, if not known
- rVAE discovers trends and (joint rVAE) building blocks irrespectively of orientation or shift in the image plane

Discovering Building Blocks



Predicting responses



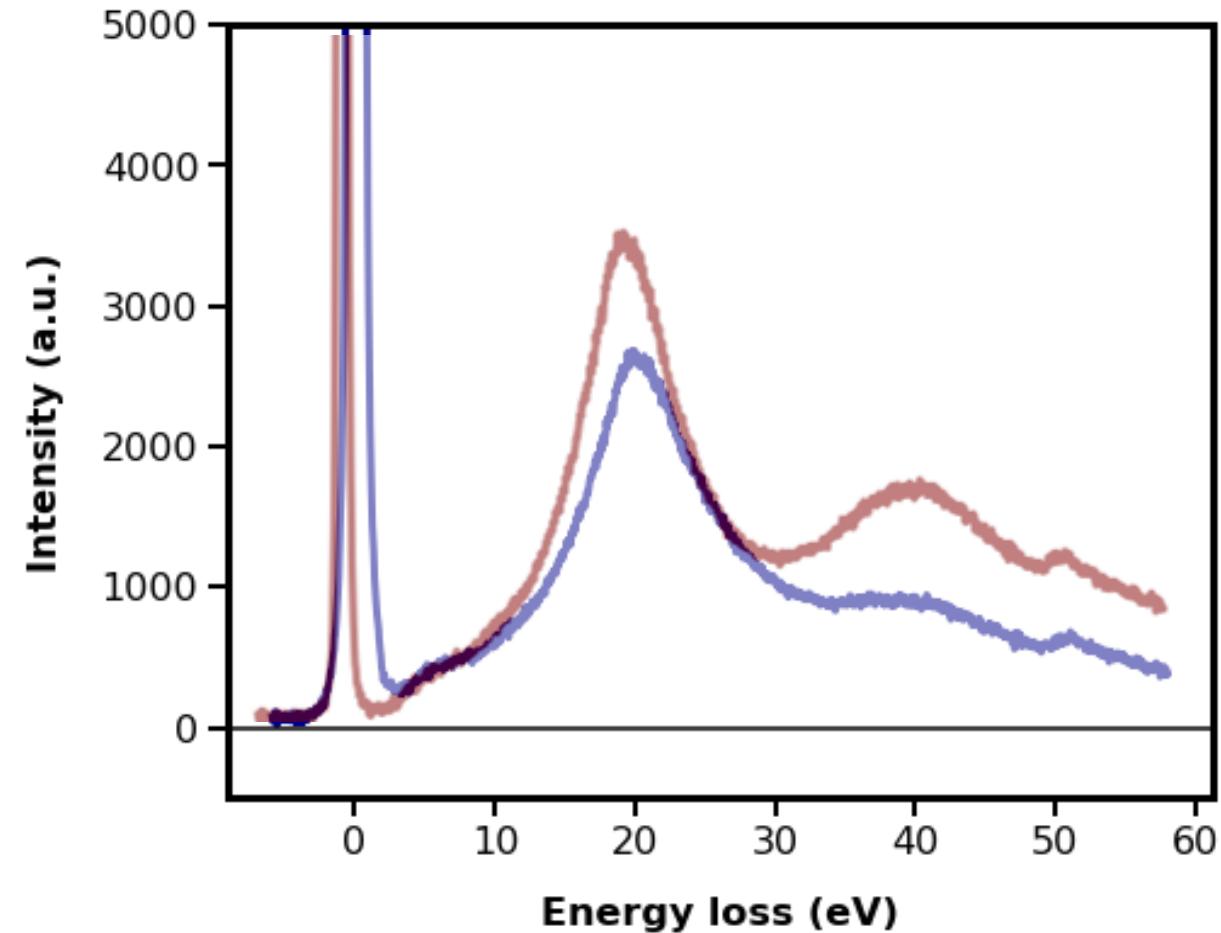
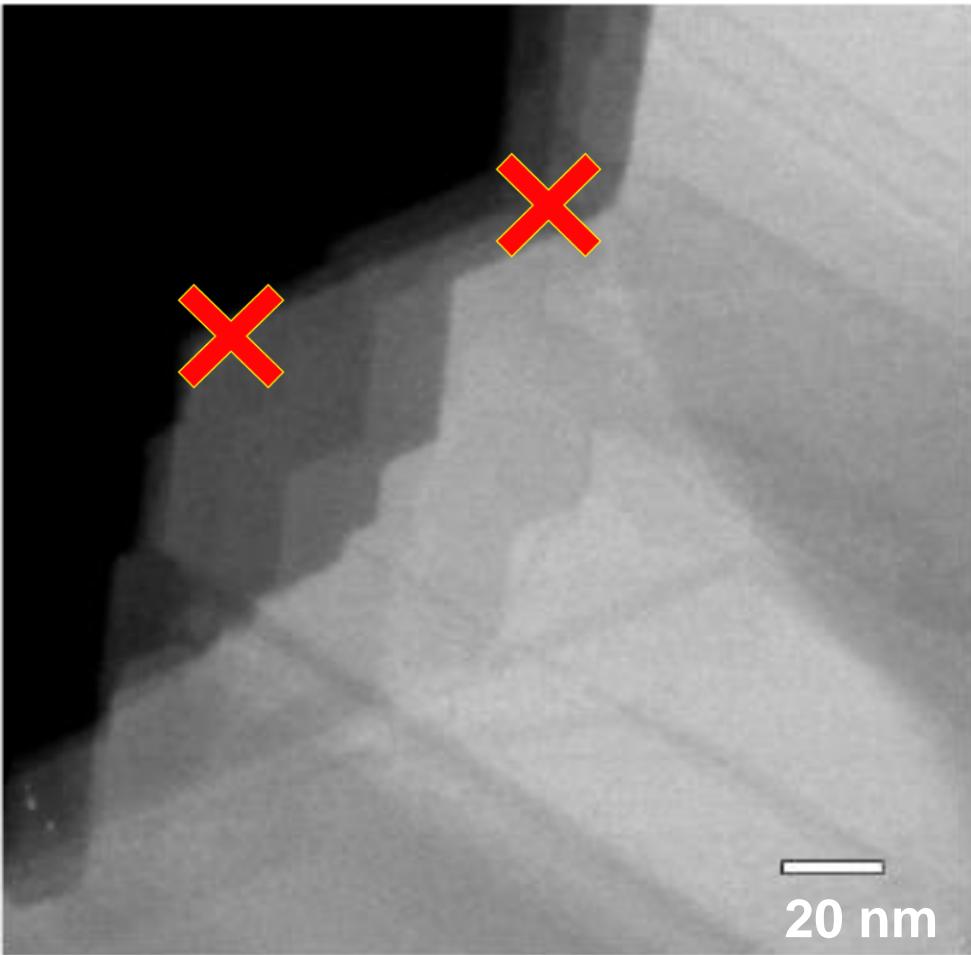
- Note it is also possible to do image prediction from spectrum (**reverse**)

Define **# of latent dimensions** to represent the data in a low dimensional space.

Automated Experiment:

... as a scientist...

Can ML run experiment as a scientist?

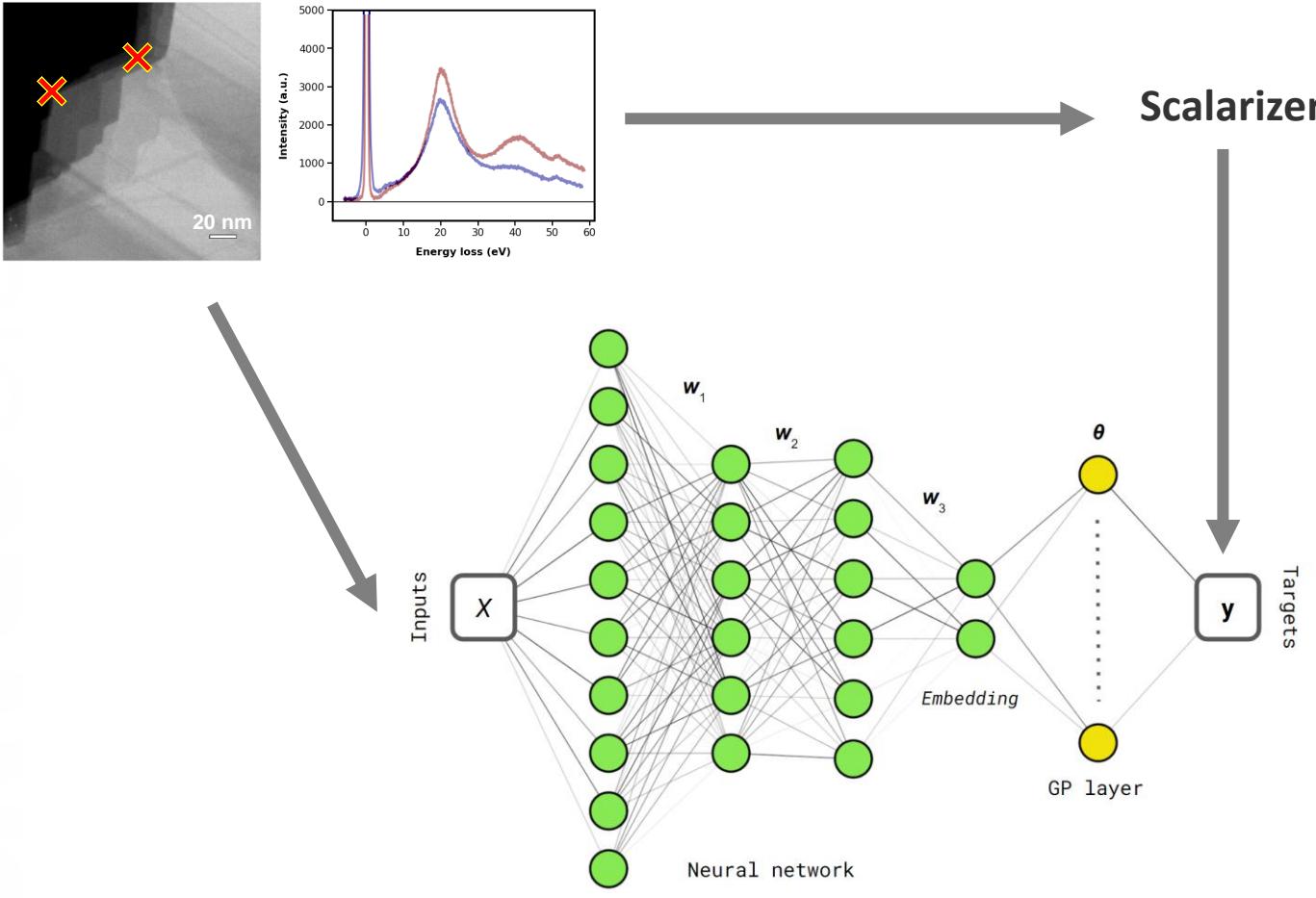


1. What if we have full access to structural information
2. And want to choose locations for (EELS, 4D STEM, CL, EDX) measurements
3. So as to **learn** relationship between structure and spectrum fastest
4. Or **discover** which microstructural elements give rise to specific **desired** spectral features?

Physics-based feature engineering: Deep kernel learning – Bayesian optimization

Specify physics criteria

Active learning



Acquire
structural data

Measure a
spectrum

Train DKL
model with new
data

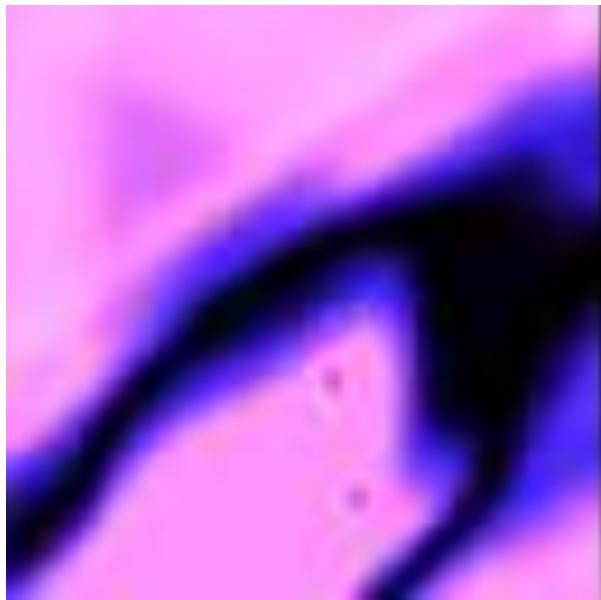
Decide next
position (optimize
physics criteria)

Allows navigation of the system to search for physics

Discovering Regions with Interesting Physics

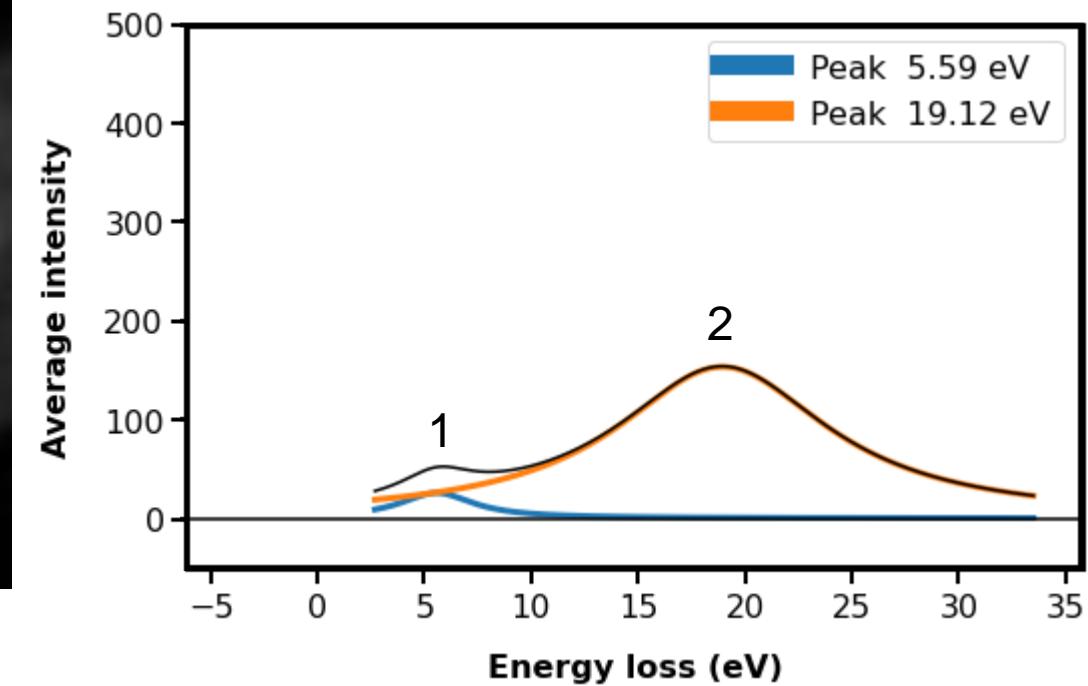
- Discovering physics in a “new” material MnPS_3
- Curve fitting to help enforce physical processes

“Acquisition function”



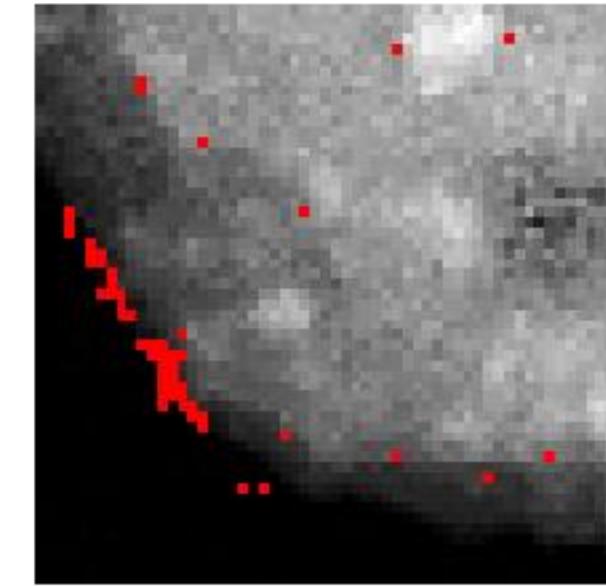
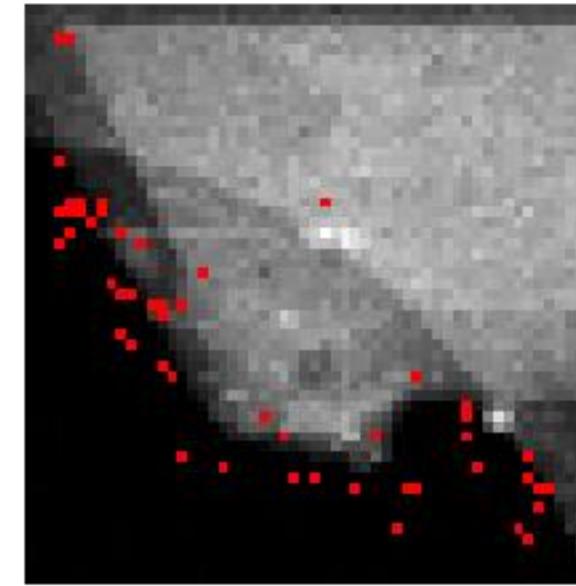
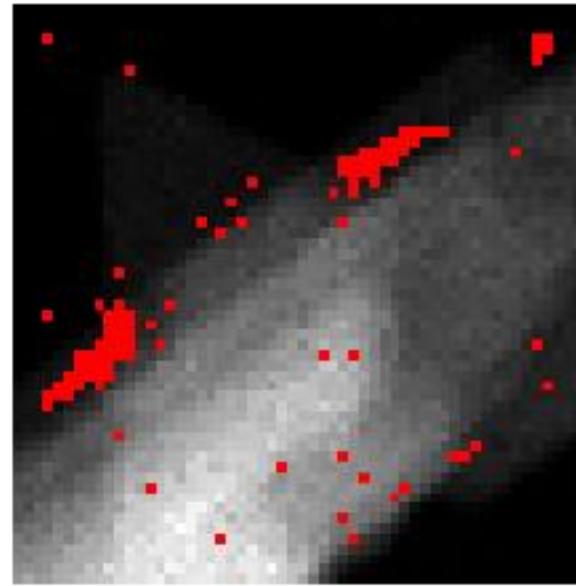
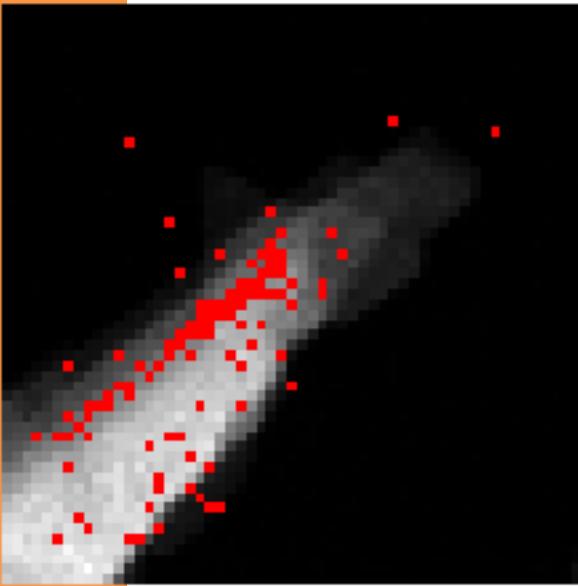
HAADF-STEM

Physics search criteria:
 $\text{Ratio} = \text{Peak 1} / \text{peak 2}$



More Examples of Physics Discovery

- Very similar behavior when searching for the same criteria!
- Success!

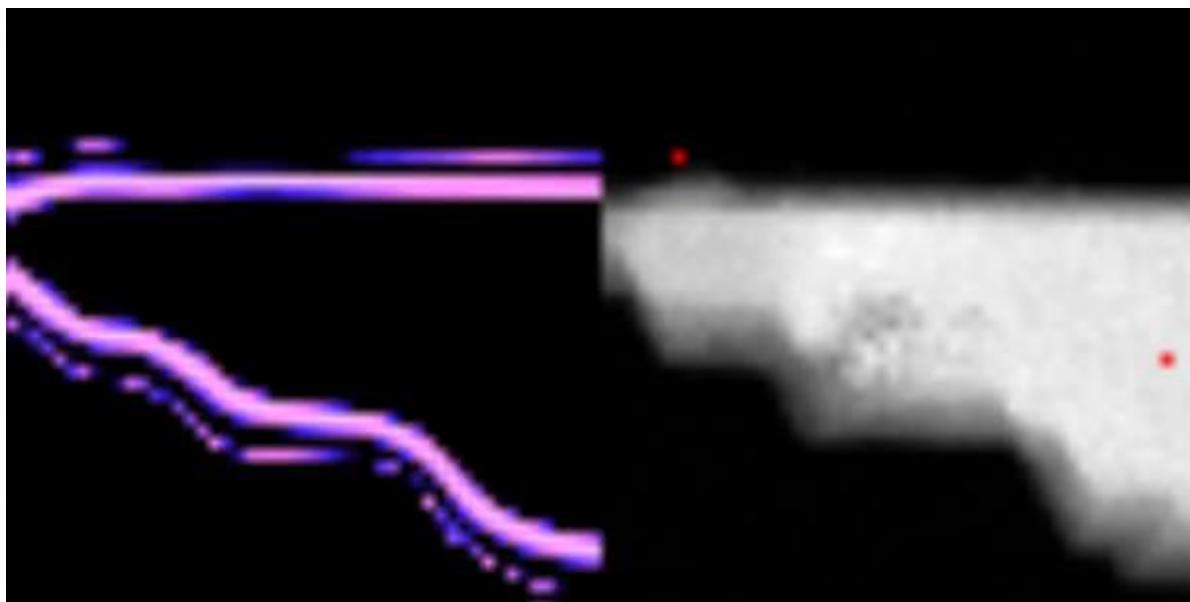


Discovery pathway depends on the reward structure (scalarizer that defines signature of physics we want to discover)!

Changing the Criterion

- (Same region) **Simple physics search:** peak max in selected region

“Acquisition function”

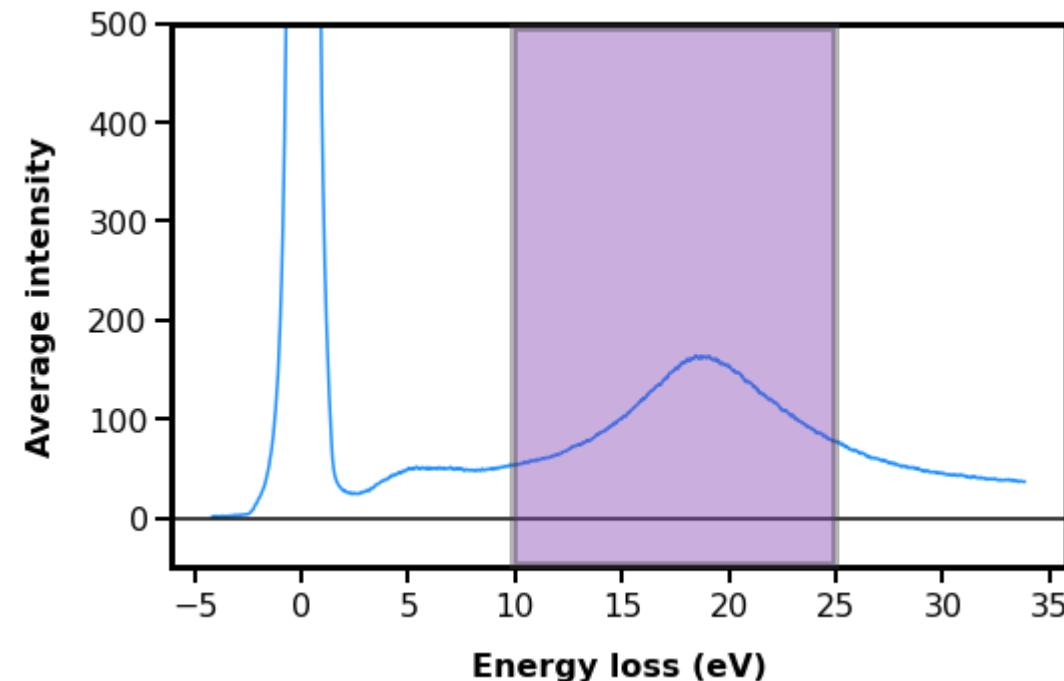


HAADF-STEM
+ points visited

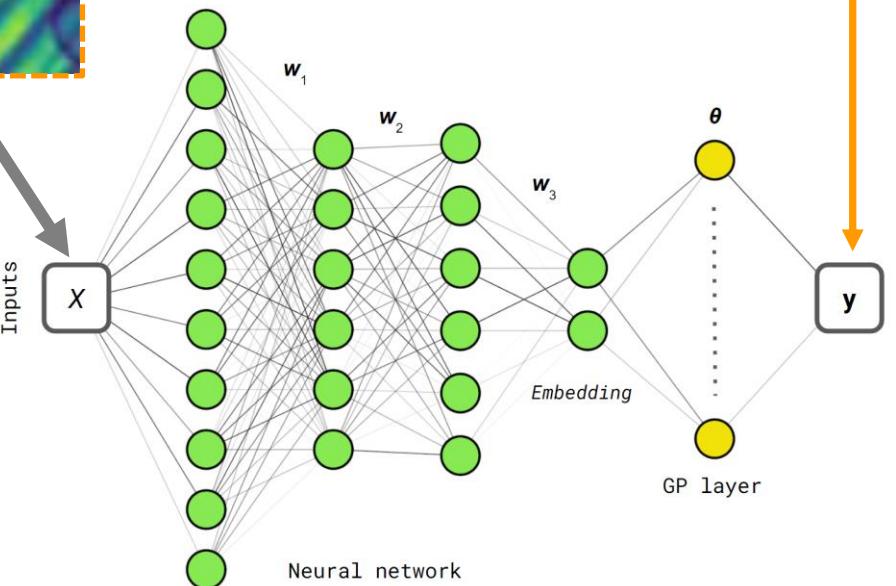
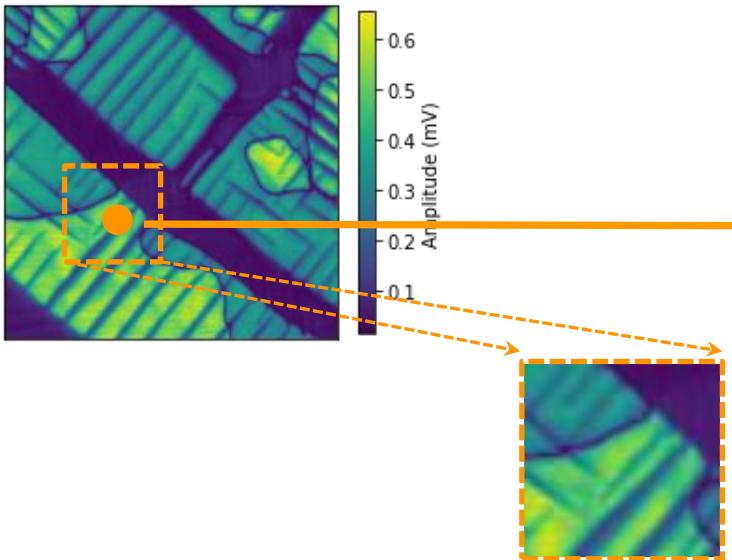
Physics search criteria:

$$\text{Maximize}(f)$$

(Specific peak intensity)



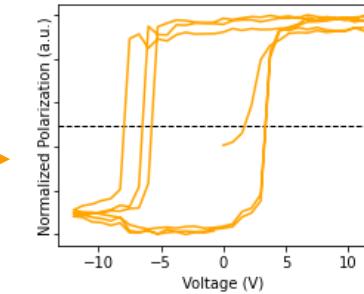
Deep Kernel Learning



- All patches are available in the beginning
- Spectra are made available sequentially
- We define what feature in spectrum are we interested in

Allows navigation of the system to search for physics

Specify physics criteria



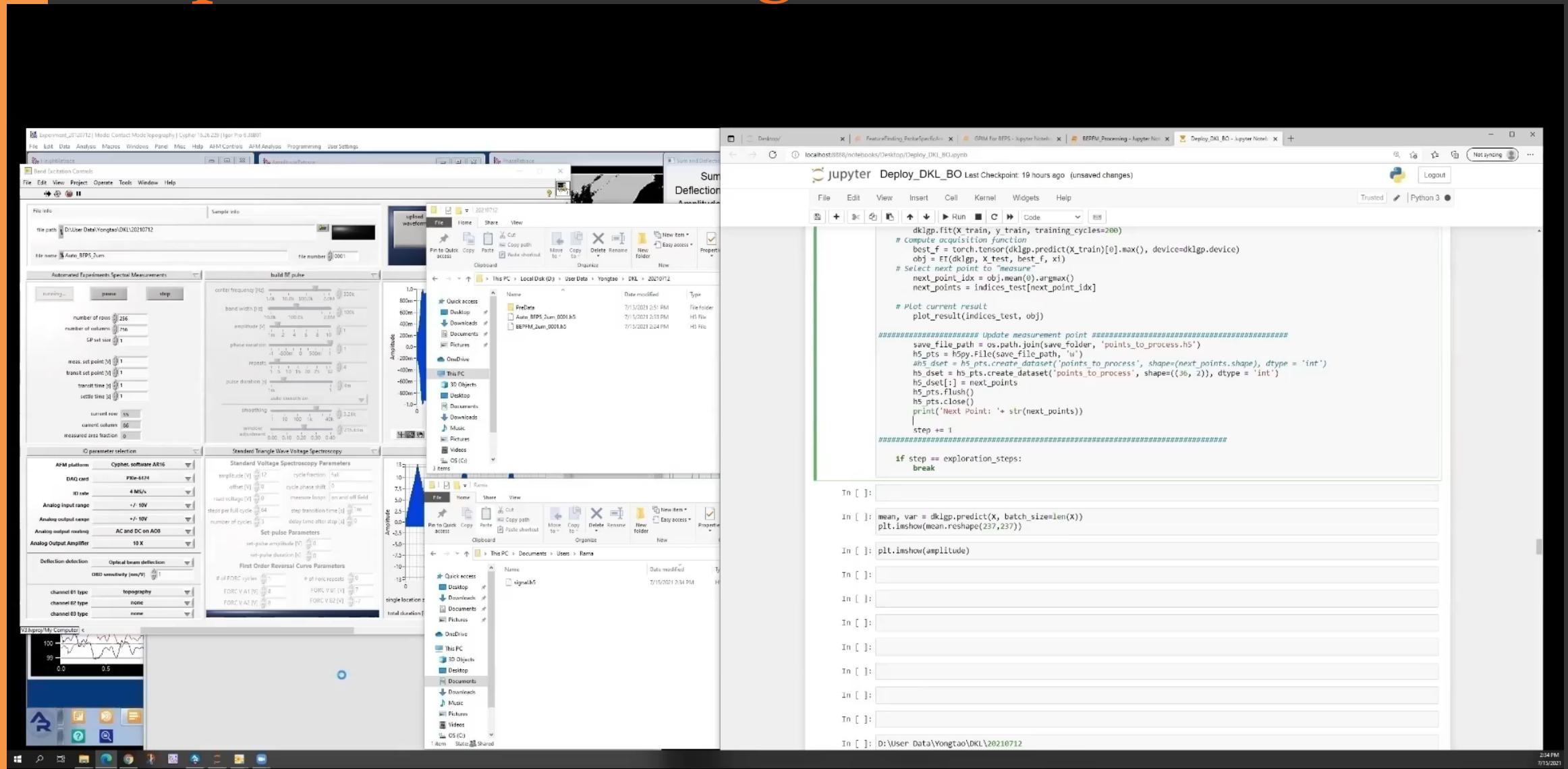
Acquire structural data

Measure a spectrum

Train DKL model with new data

Decide next position (optimize physics criteria)

Deep Kernel Learning AE

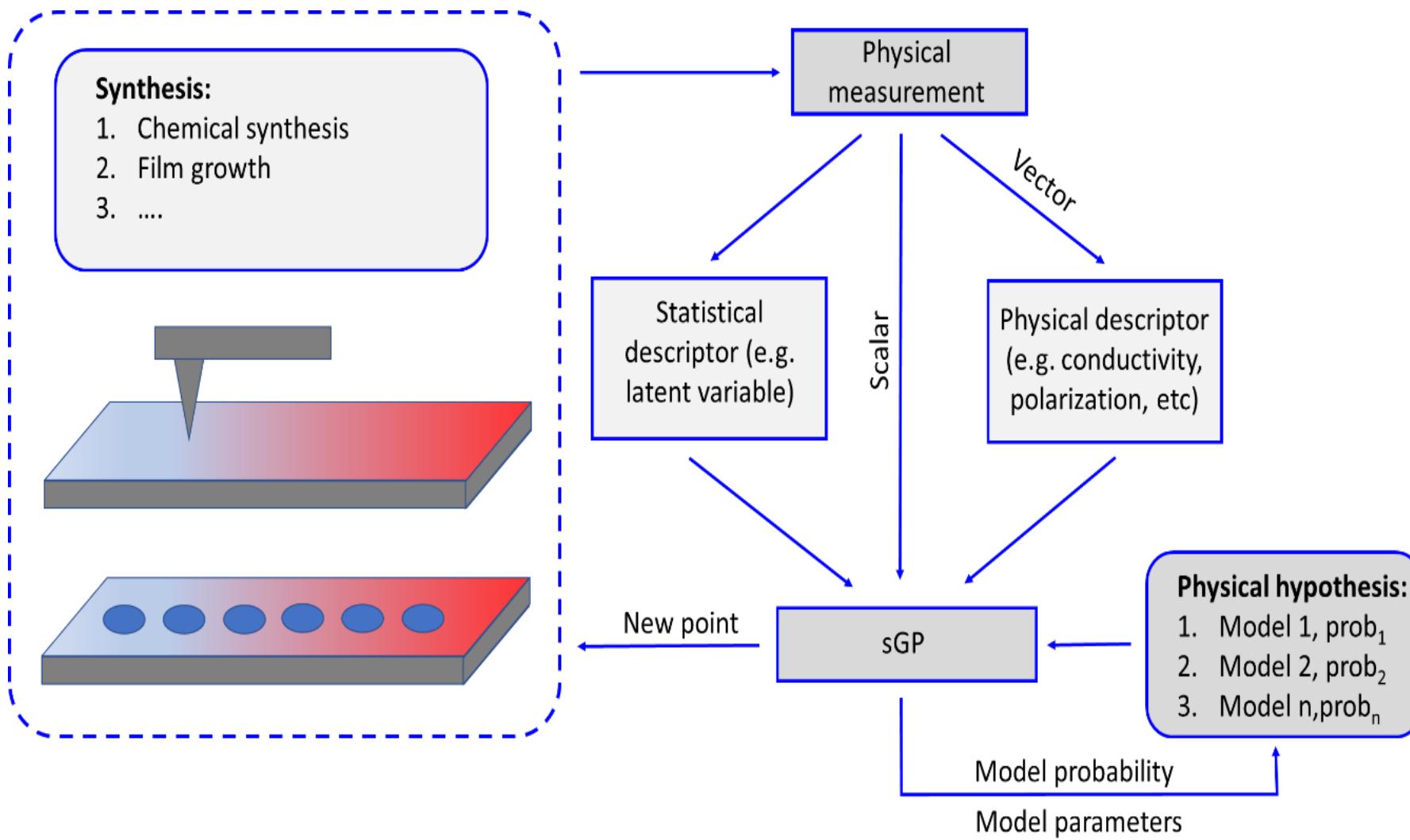


Automated Experiment: ... as a scientist...

Bayesian optimization:

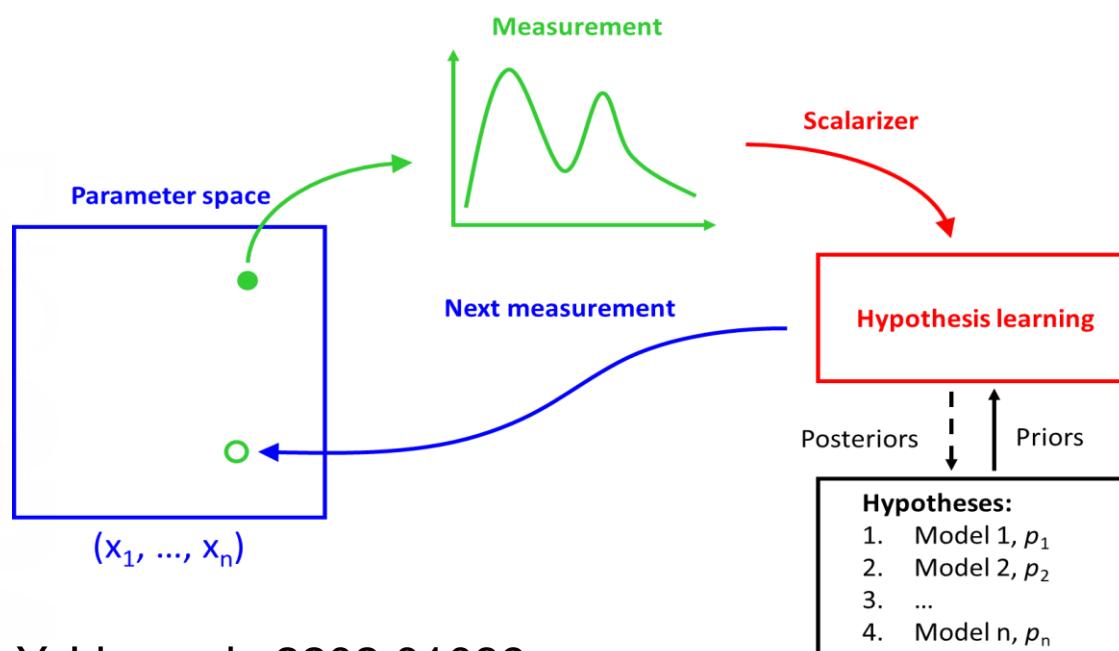
1. Works only in low-dimensional spaces
2. The correlations are defined by the kernel function (very limiting)
3. We do not use any knowledge about physics of the system
4. We do not use cheap information available during the experiment (proxies)

Hypothesis Active Learning



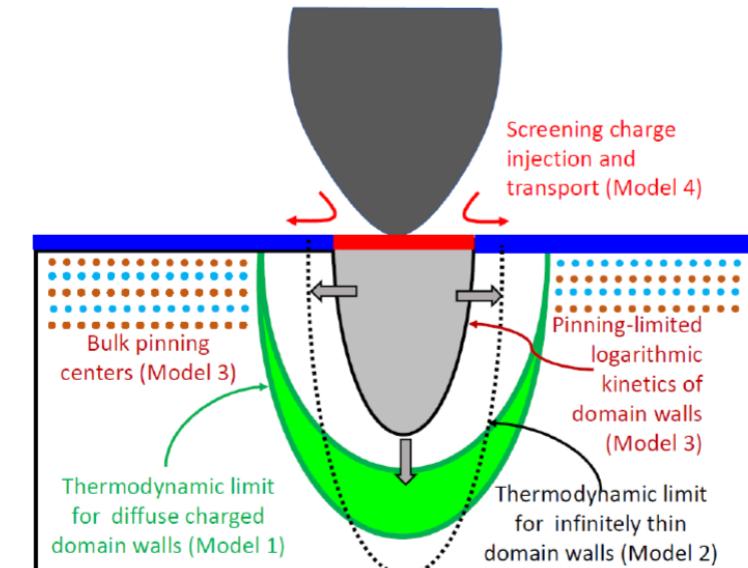
Hypothesis Learning

- Can ML algorithm think like a scientist?
- Yes – automated experiment can pursue hypothesis-driven science!



Y. Liu, arxiv 2202.01089

Y. Liu, arxiv 2112.06649



Model Equation

Thermodynamic 1

Model I

$$r(V) = r_{cr} + r_0 \sqrt{\left(\frac{V}{V_c}\right)^{2/3} - 1}$$

Thermodynamic 2

Model II

$$r(V) = r_{cr} + r_0 \sqrt[3]{\left(\frac{V}{V_c}\right)^2 - 1}$$

Wall pinning

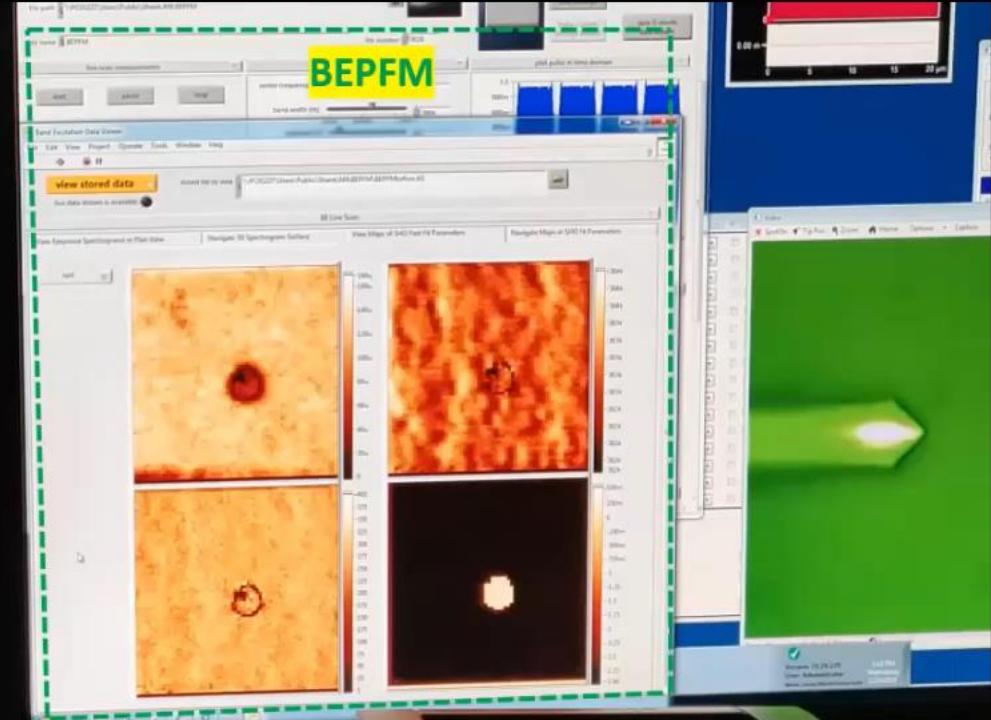
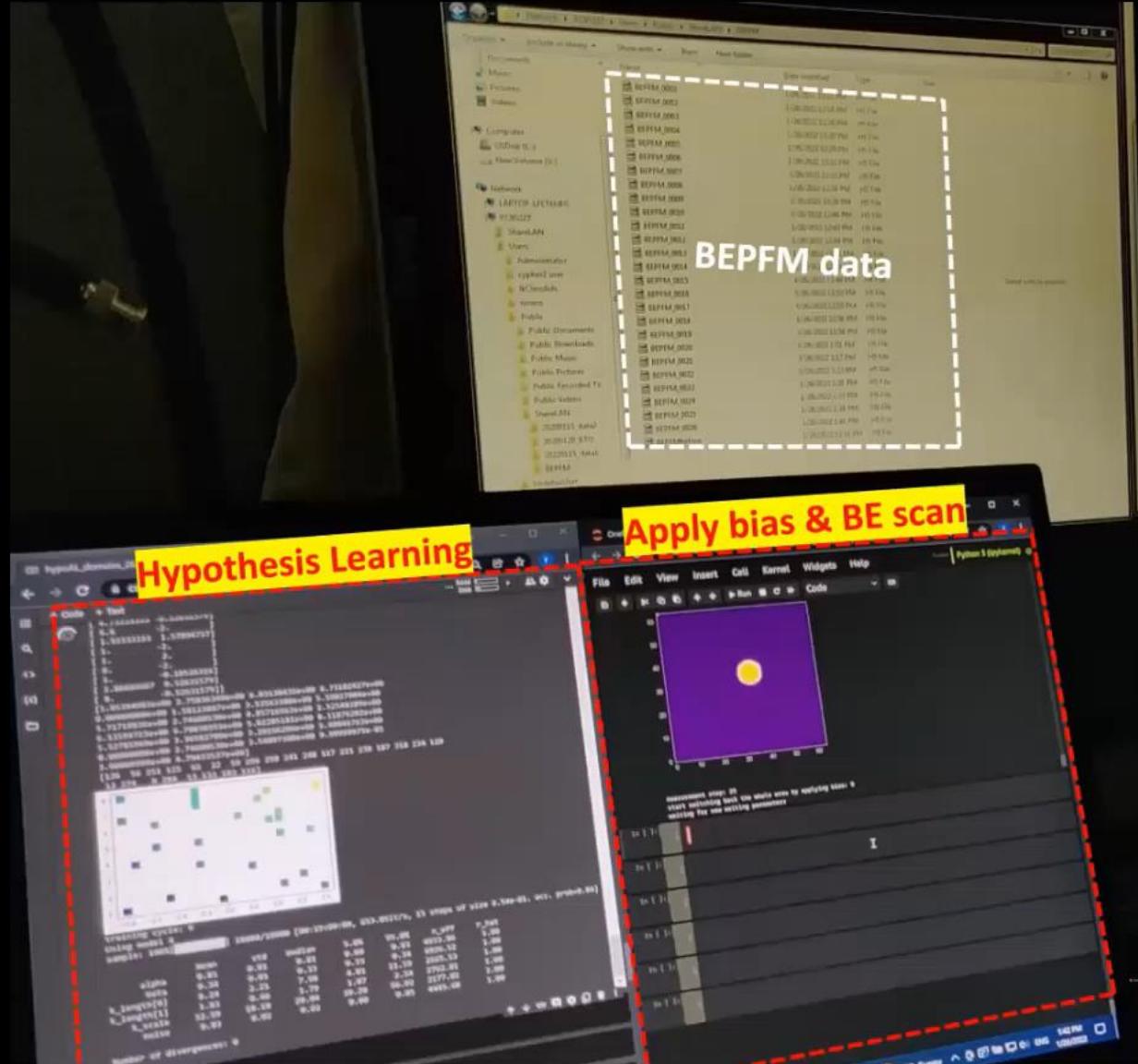
Model III

$$r(V, t) = V^\alpha \log \tau$$

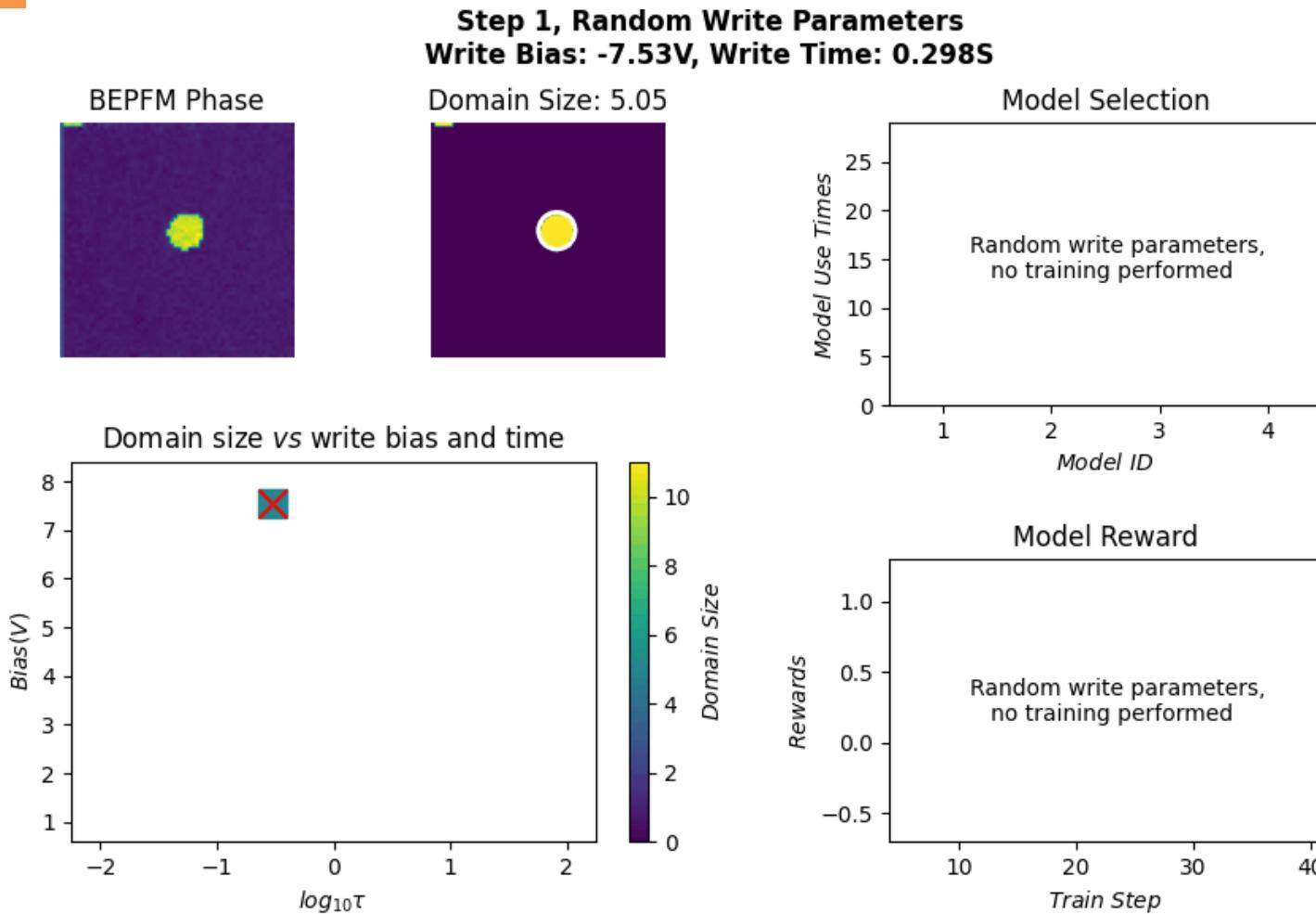
Charge injection

Model IV

$$r(V, t) = V^\alpha \tau^\beta$$



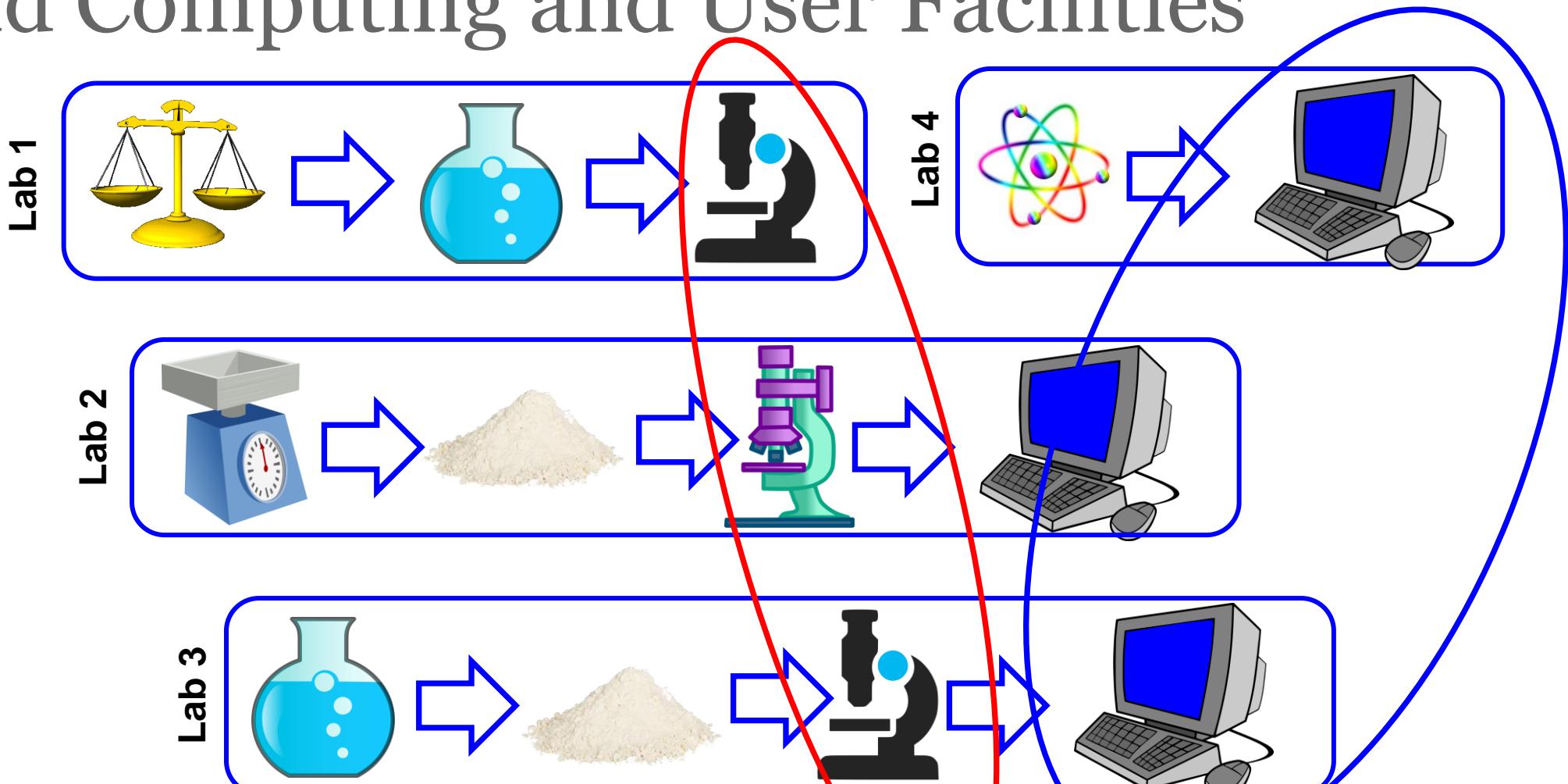
Hypothesis learning in action



- ML algorithm has 4 competing hypothesis on domain switching mechanisms
- These hypothesis represent full set of possibilities for this system
- The microscope chooses experimental parameters in such a way as to establish which hypothesis is correct fastest
- Important: the same approach can be implemented in synthesis and electrical characterization
- Machine learning meets hypothesis-driven scientific discovery!

Future of Instrumentation

Cloud Computing and User Facilities



- Big scientific tools (synchrotrons) are user facilities
- For ~20 years, medium-scale tools operated as user facilities
- Enterprise computing -> cloud computing
- Over last 5 years, cloud labs are emerging
- What about the workflows?

User facilities

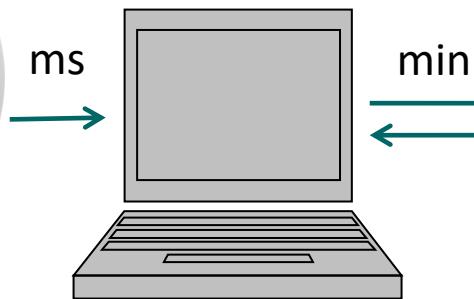
Cloud computing

Classical Instrumental Research

SPM: 30,000+ platforms worldwide:

Large weakly connected instrumental network

(S)TEM: ~100s top level machines,
much stronger integrated community



Instrument

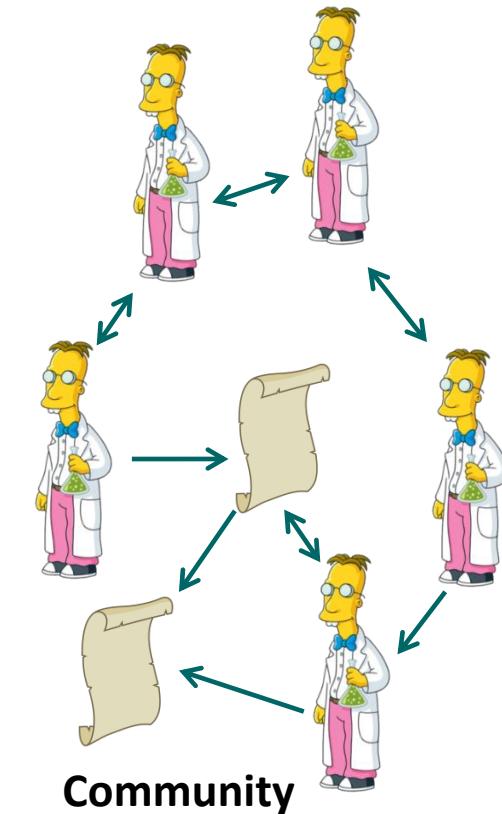
Control/data acquisition



Researcher

ms
min

Weeks -
months

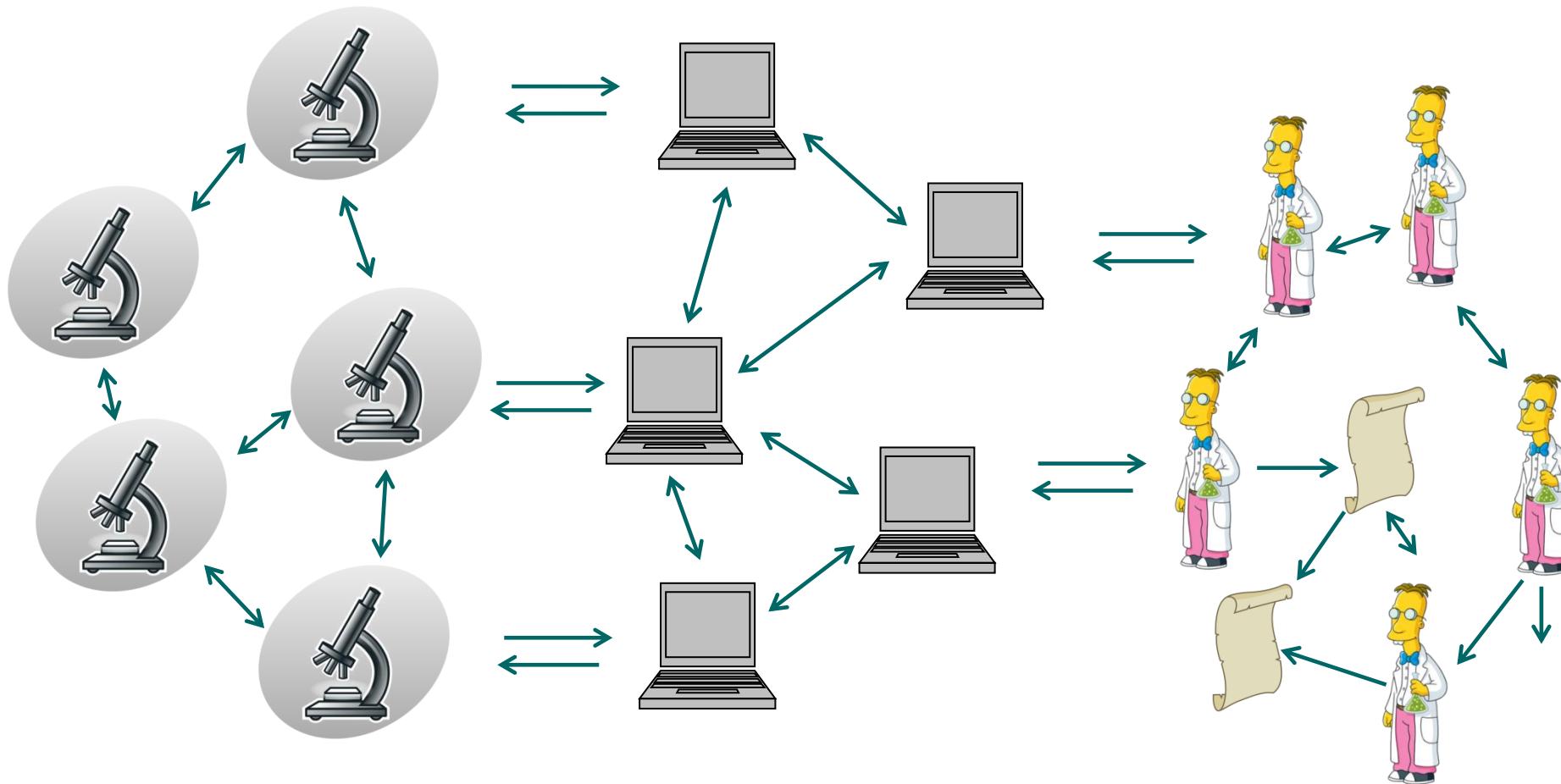


Community

- Social networking/education
- Publications/citations

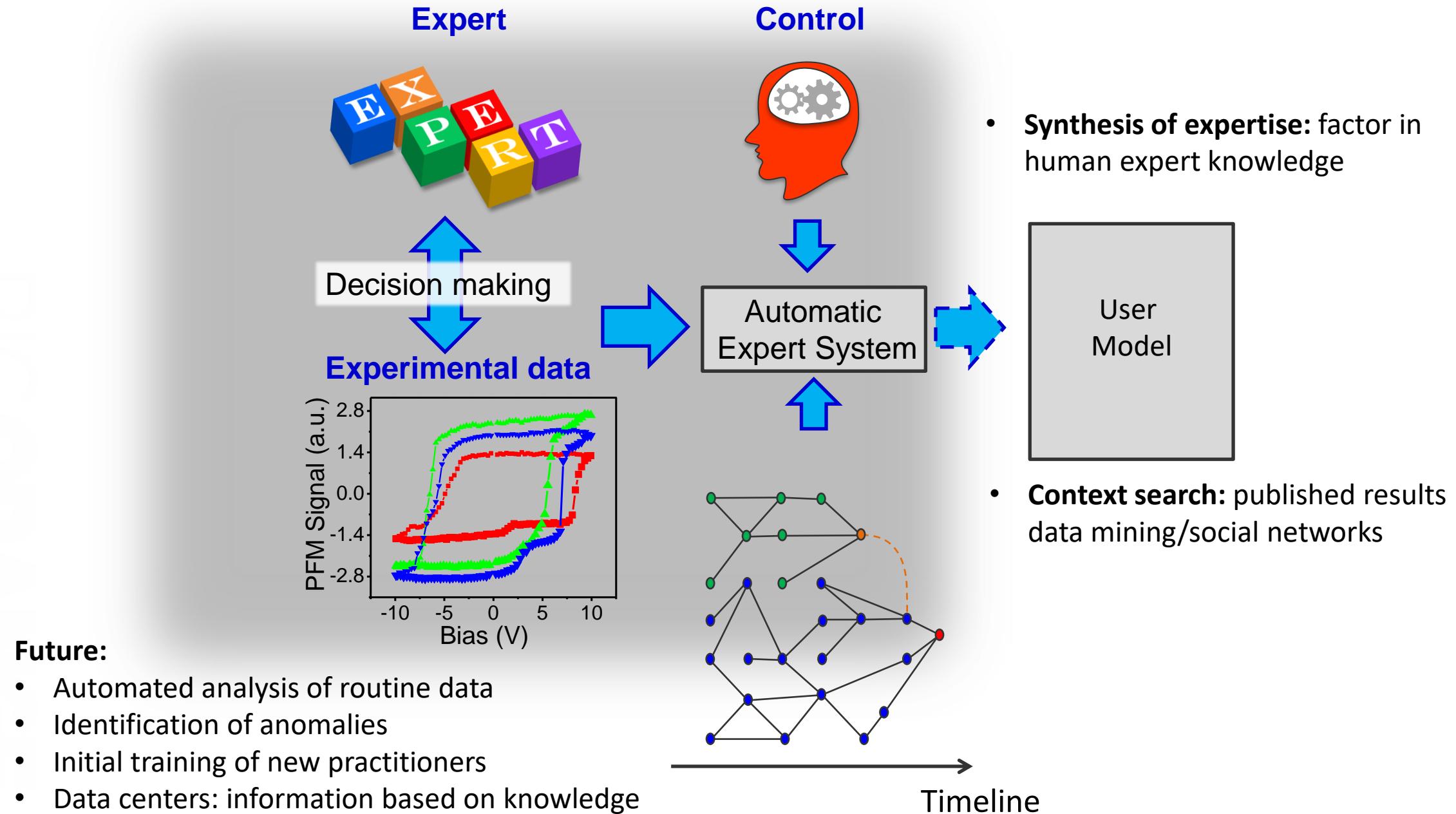
1. Only small fraction of data stream from the instrumentation is captured
2. Only small fraction of captured data is analyzed, interpreted, and put in the context
3. Human-machine interaction during acquisition is often slow and can be non-optimal
4. Human interpretation of data is limited: bias and ignoring serendipity
5. Information propagation and concept evolution in scientific community is slow and affected by non-scientific factors

Step 1: Cloud Integration

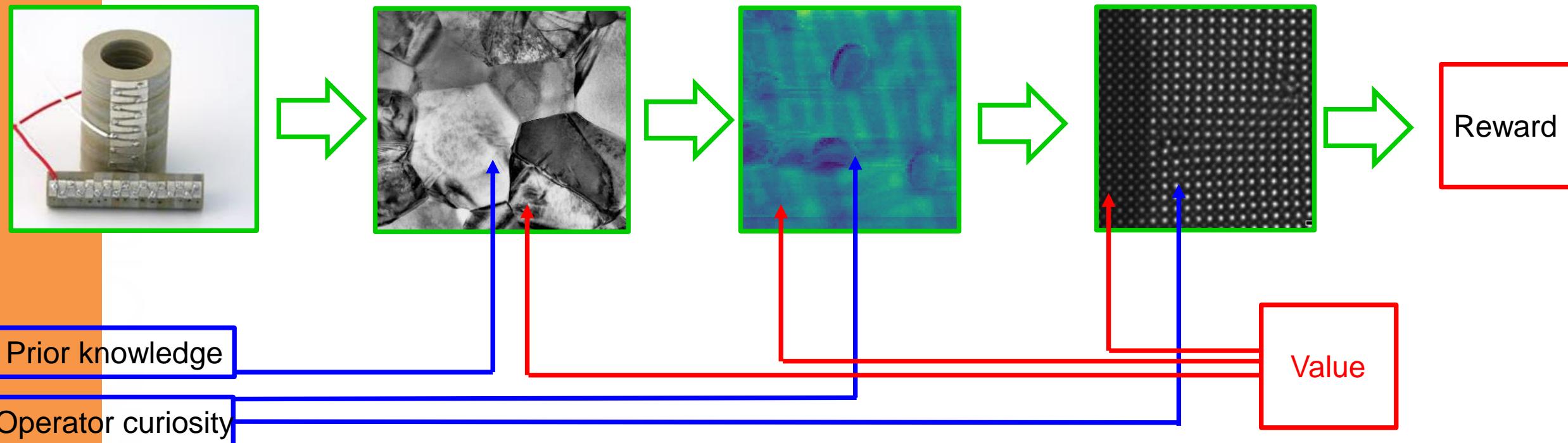


1. Multiple geographically-distributed data generation node
2. Full capture of instrumental data stream /compression/curation
3. Coordination of protocols and data/metadata across the cloud
4. Cloud-based processing and dimensionality reduction
5. Community-wide analytics

Step 2: Cloud Analytics



Step 3: Workflow Design



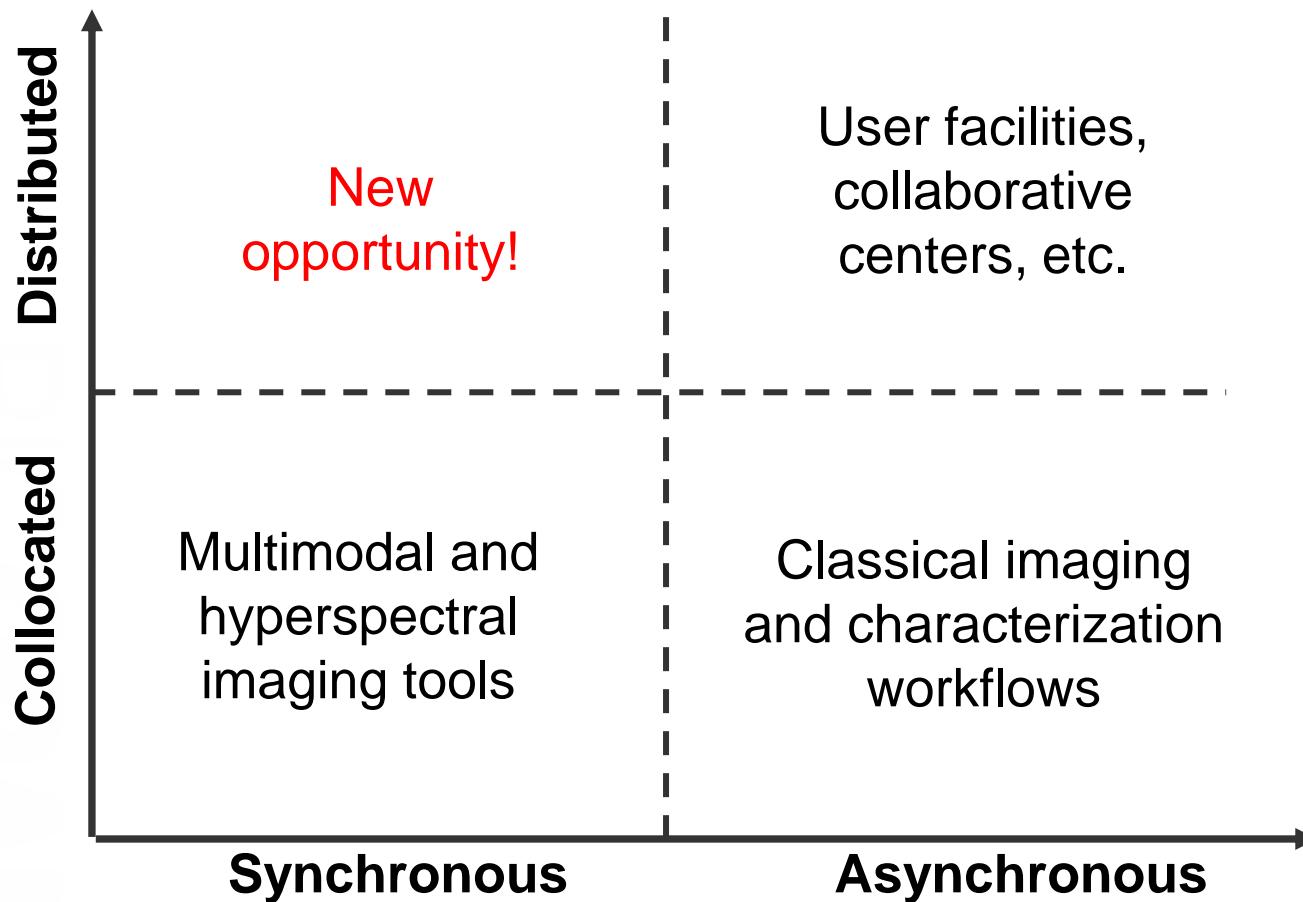
Traditional experiment:

1. Always based on workflows
2. Ideated, orchestrated, and implemented by humans
3. The “gain of value” during the workflow implementation is uncertain

Value of the step is key element:

- Either based on prior knowledge
- Or defined in a sense of the reinforcement learning Q-function

Cloud Labs: Facilities of the Future



Emerald Cloud Lab,
SF and CMU

1. Combined human-machine workflow implementation
2. Computer orchestrating agent
3. How would beyond human workflows be ideated?

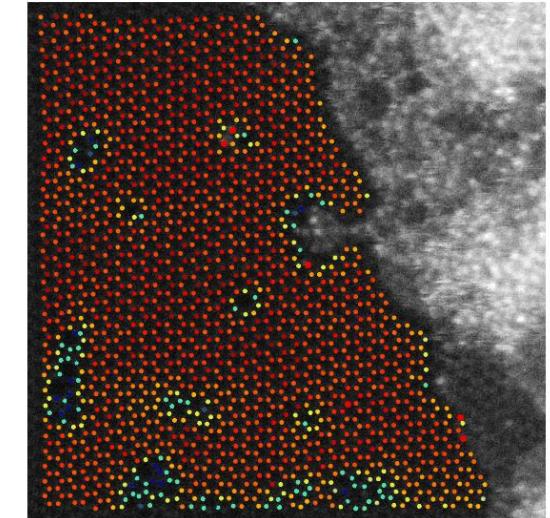
And technology

“New directions in science are launched by new tools much more often than by new concepts. The effect of a concept-driven revolution is to explain old things in new ways. The effect of a tool-driven revolution is to discover new things that have to be explained.”

Freeman Dyson

- Cloud connection and data infrastructure is necessary, but not enough
- Can we define workflow design, building, and optimization?
- Potential to change way R&D is performed: academia, start-ups, industry

- **AtomAI:** comprehensive toolbox for DCNN-based supervised exploration of STEM and SPM Data:
- **PyroVED:** building structure-property correlations and unsupervised and semi-supervised physical discovery
- **GPim:** Gaussian processing toolbox for image analytics and automated experiment
- **gpax:** hypothesis-driven structured Gaussian Processes
- **PyCroscopy:** General data formats, workflows, and image analytics



Welcome to Microscopy Project!