



Data-driven Scientific Discovery

July 12, 2021

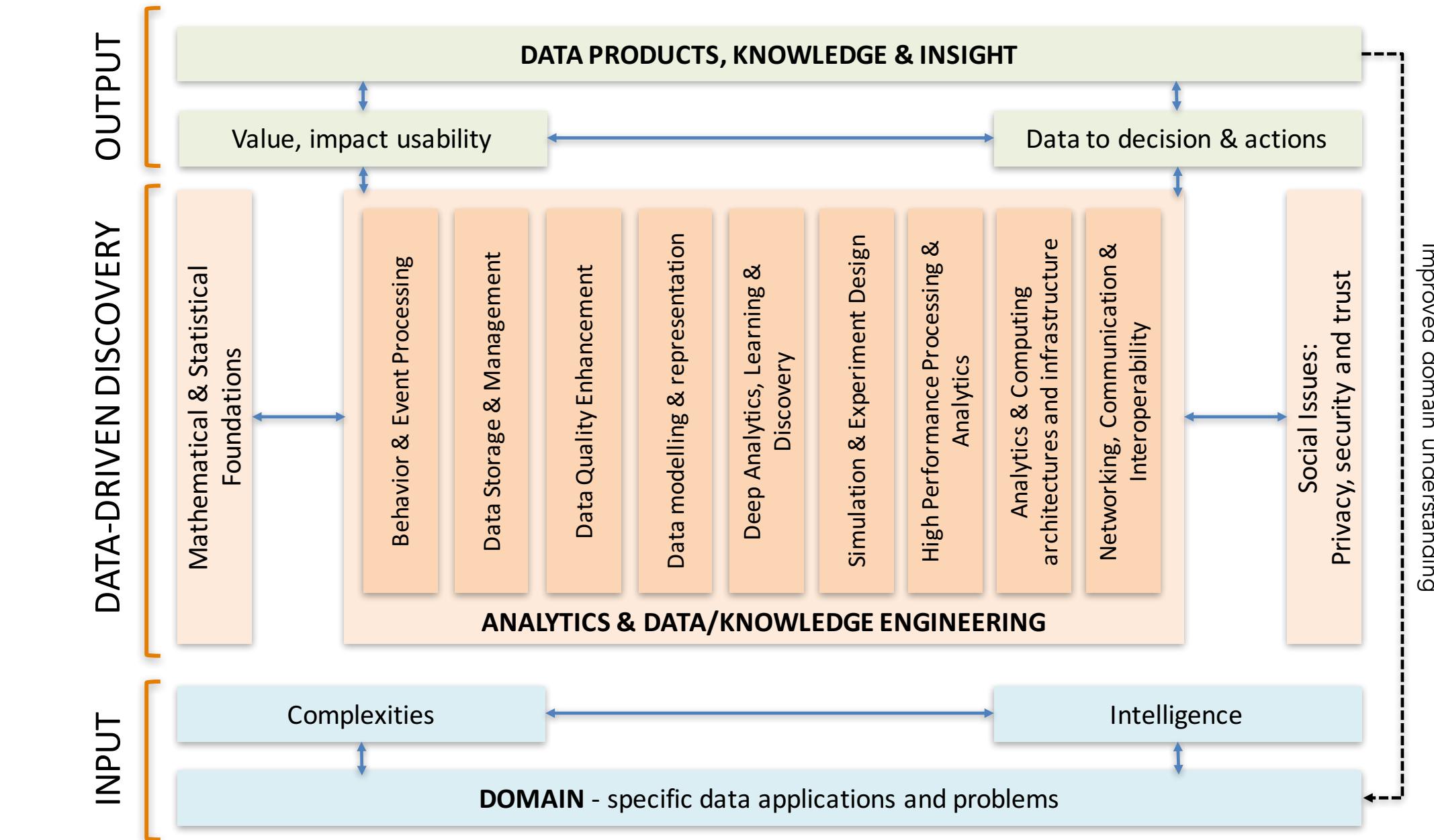
Robert Rallo
Director

Advanced Computing, Mathematics, and Data Division



PNNL is operated by Battelle for the U.S. Department of Energy

Data Science

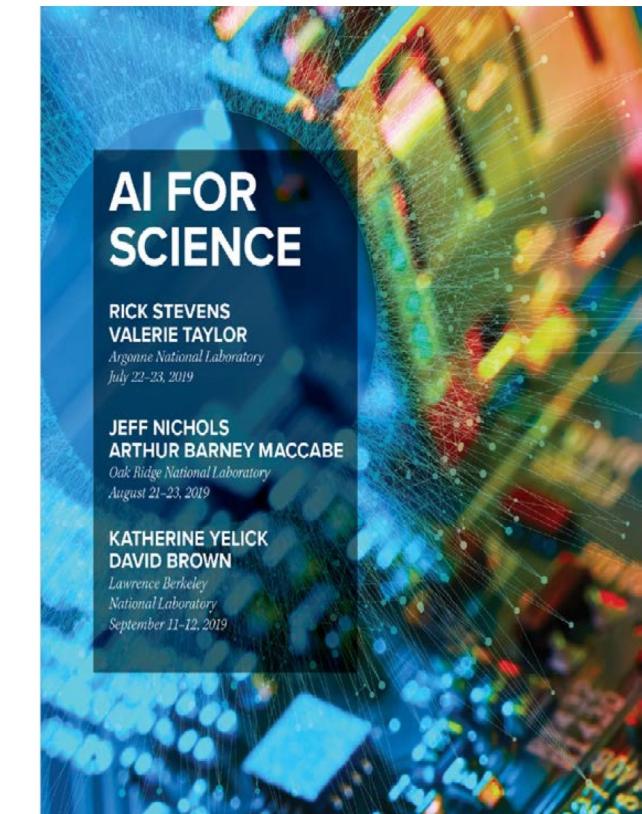
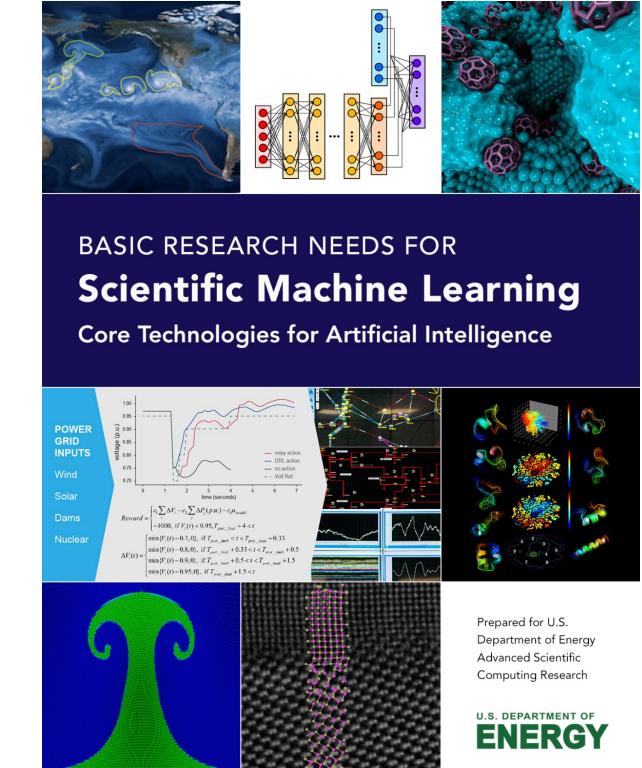




Pacific
Northwest
NATIONAL LABORATORY

Scientific Machine Learning

- Machine Learning algorithms should be:
 - Domain-aware
 - Interpretable
 - Robust
 - Data-intensive
 - Scalable, deployable
- Machine Learning must contribute to:
 - Enhance current modeling and simulation
 - Automation and decision support





Pacific
Northwest
NATIONAL LABORATORY

Data

- Data products
 - Beyond raw data
 - Models, workflows
- Data governance and data stewardship
 - Define process and procedures
 - Ensure compliance and security
 - Manages data, tools and processes
 - Comprehensive policies
- Data quality
 - Metrics
- Privacy preservation
 - Privacy-preserving analytics
 - Privacy-preserving AI/ML

DATAHUB
Pacific Northwest

Search by keywords, authors and much more... + About

Categories Datasets Data Sources Projects Publications People

DATASET
OmicsLHV-SHAE003

BIOLOGY HUMAN HEALTH

OMICS TRANSCRIPTOMICS MICROARRAY

Omics-LHV, SARS

Experiment SHAE003

The purpose of this SARS experiment was to obtain samples for transcriptome analysis in Human Airway Epithelial (HAE) cells infected with SARS-CoV, SARS deltaORF6, SARS BatSRBD mutants in a longitudinal study.

Overall Design: Primary Human airway epithelium (HAE) cells (resembling *in vivo* pseudo- stratified mucociliary epithelium) were infected with SARS-CoV, SARS deltaORF6, SARS BatSRBD mutant at an MOI of 2. RNA triplicates/quadruplicates are defined as 3/4 different wells, plated at the same time and using the same cell stock for all replicates. Time matched mocks done in triplicate from same cell stock as rest of samples. Culture medium (the same as what the virus stock is in) will be used for the mock infections. Time points used were: 0, 24, 48, 60, 72, 84, and 96 hrs post-infection.

Download

Projects (2)

Oomics-Lethal Human Viruses, SARS

Omics-LHV Profiling of Host Response to Severe Acute Respiratory Syndrome (SARS) Infection The Systems Virology project was one of four systems...

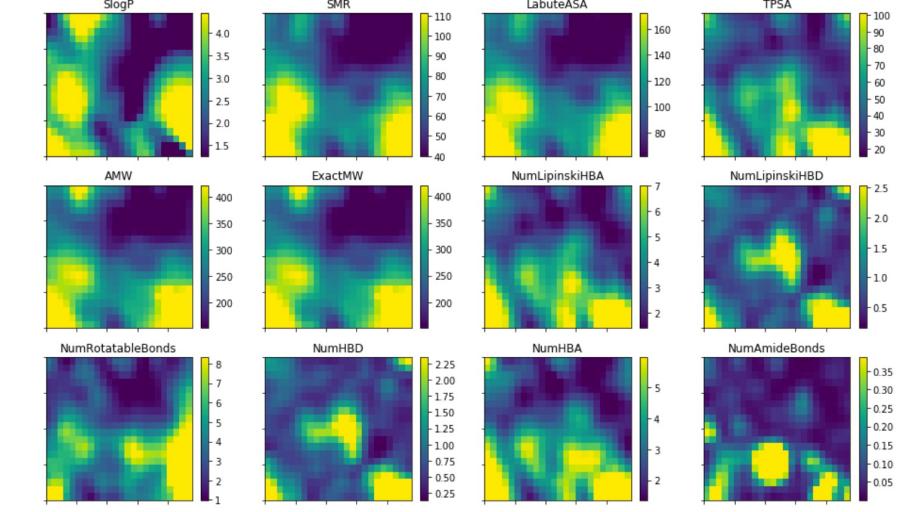
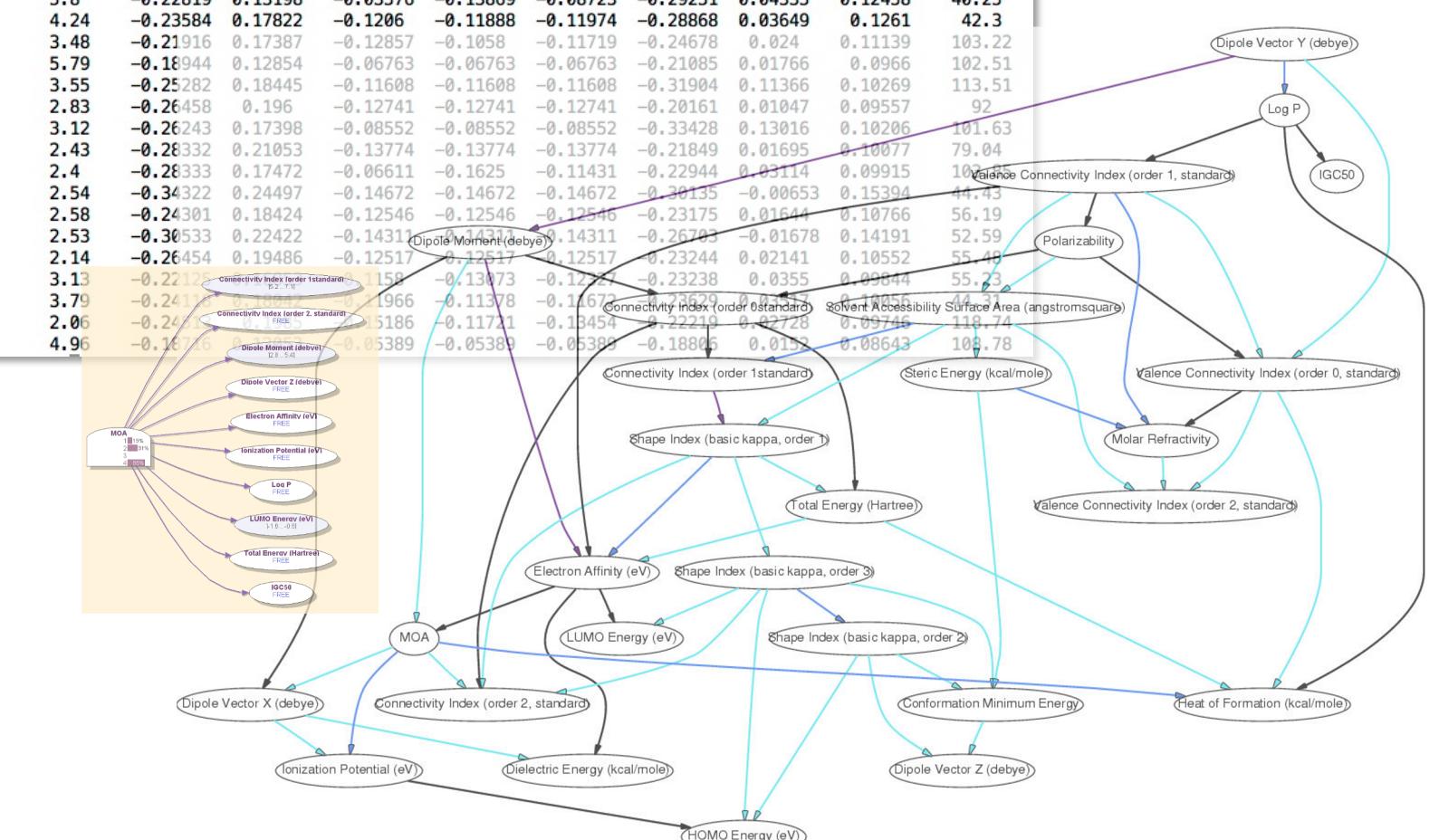
BIOLOGY HUMAN HEALTH



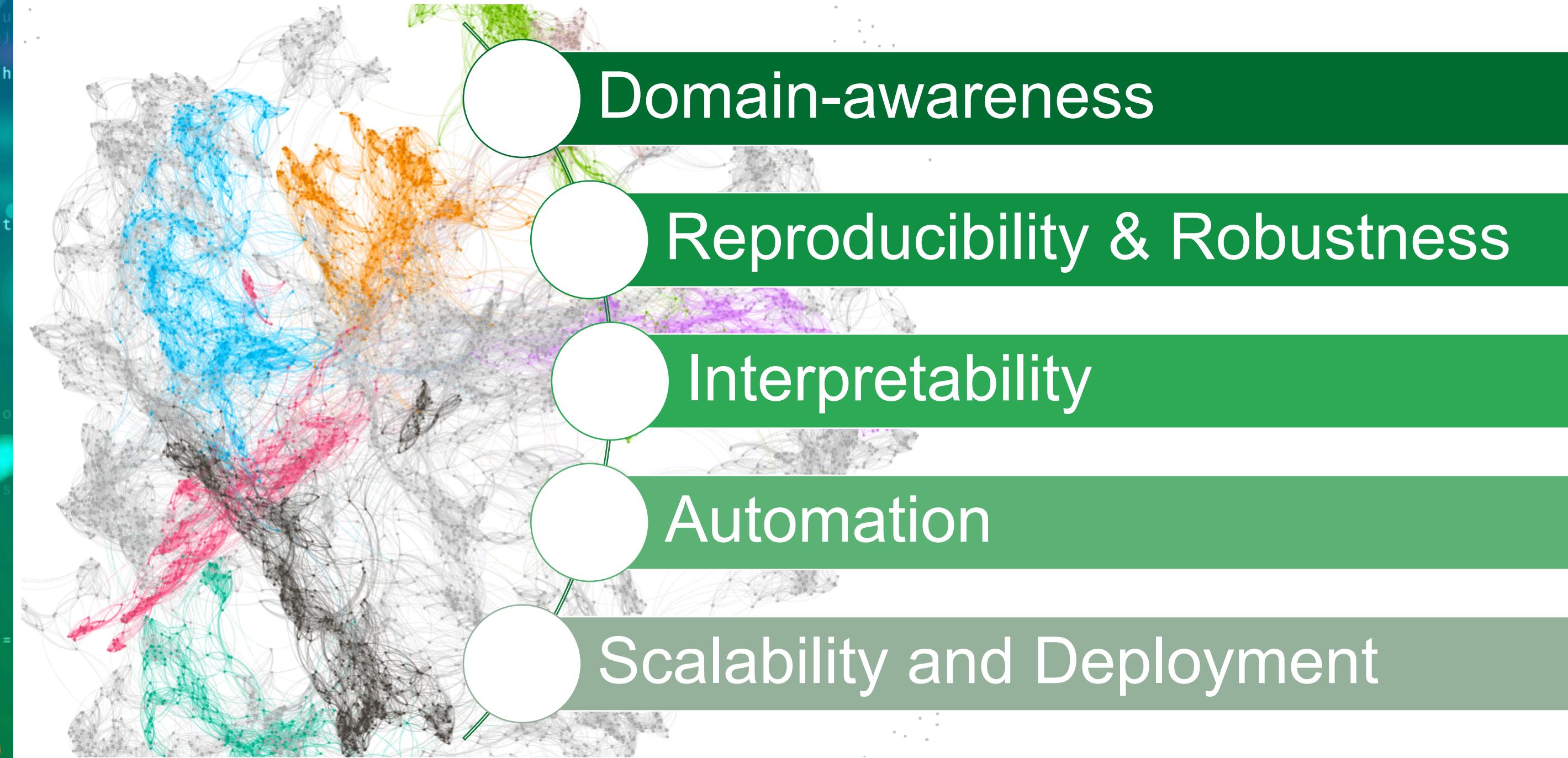
Data Representation

- Structured vs unstructured data
 - e.g., knowledge extraction from literature sources
- Efficient data structures
 - graphs
- Feature engineering
 - feature selection
 - feature synthesis
- Representation learning
 - interpretability

material	light_tox	dark_tox	HHOMO	LZELEHHO	LUMOA	LUMOB	ALZLUMO	MHOMOA	MLUMOA	QMELECT	Cp
ZnO	6.23	5.8	-0.22819	0.13198	-0.03576	-0.13869	-0.08723	-0.29251	0.04335	0.12458	40.25
CuO	5.71	4.24	-0.23584	0.17822	-0.1206	-0.11888	-0.11974	-0.28868	0.03649	0.1261	42.3
V2O3	3.78	3.48	-0.21916	0.17387	-0.12857	-0.1058	-0.11719	-0.24678	0.024	0.11139	103.22
Y2O3	5.84	5.79	-0.18944	0.12854	-0.06763	-0.06763	-0.06763	-0.21085	0.01766	0.0966	102.51
Bi2O3	4.02	3.55	-0.25282	0.18445	-0.11608	-0.11608	-0.11608	-0.31904	0.11366	0.10269	113.51
In2O3	3.48	2.83	-0.26458	0.196	-0.12741	-0.12741	-0.12741	-0.20161	0.01047	0.09557	92
Sb2O3	3.66	3.12	-0.26243	0.17398	-0.08552	-0.08552	-0.08552	-0.33428	0.13016	0.10206	101.63
Al2O3	2.75	2.43	-0.28332	0.21053	-0.13774	-0.13774	-0.13774	-0.21895	0.01695	0.10077	79.04
Fe2O3	2.54	2.4	-0.28333	0.17472	-0.06611	-0.1625	-0.11431	-0.22944	0.03114	0.09915	103.16
SiO2	2.92	2.54	-0.34322	0.24497	-0.14672	-0.14672	-0.14672	-0.30135	-0.00653	0.15394	44.43
ZrO2	3.04	2.58	-0.24301	0.18424	-0.12546	-0.12546	-0.12546	-0.23175	0.01644	0.10766	56.19
SnO2	3.24	2.53	-0.30533	0.22422	-0.14311	-0.14311	-0.14311	-0.26703	-0.01678	0.14191	52.59
TiO2	4.68	2.14	-0.26454	0.19486	-0.12517	-0.12517	-0.12517	-0.23244	0.02141	0.10552	55.49
CoO	3.33	3.13	-0.22128	0.158	-0.13073	-0.1227	-0.1227	-0.23238	0.0355	0.09844	55.23
NiO	3.87	3.79	-0.24165	0.1966	-0.11378	-0.11672	-0.11672	-0.22219	0.01717	0.10056	44.31
Cr2O3	2.06	2.06	-0.24124	0.15186	-0.11711	-0.13454	-0.13454	-0.22278	0.02728	0.09746	118.74
La2O3	5.56	4.96	-0.18726	0.05389	-0.05389	-0.05389	-0.05389	-0.18805	0.0158	0.08643	108.78



Opportunities for Accelerating Scientific Discovery with AI/ML



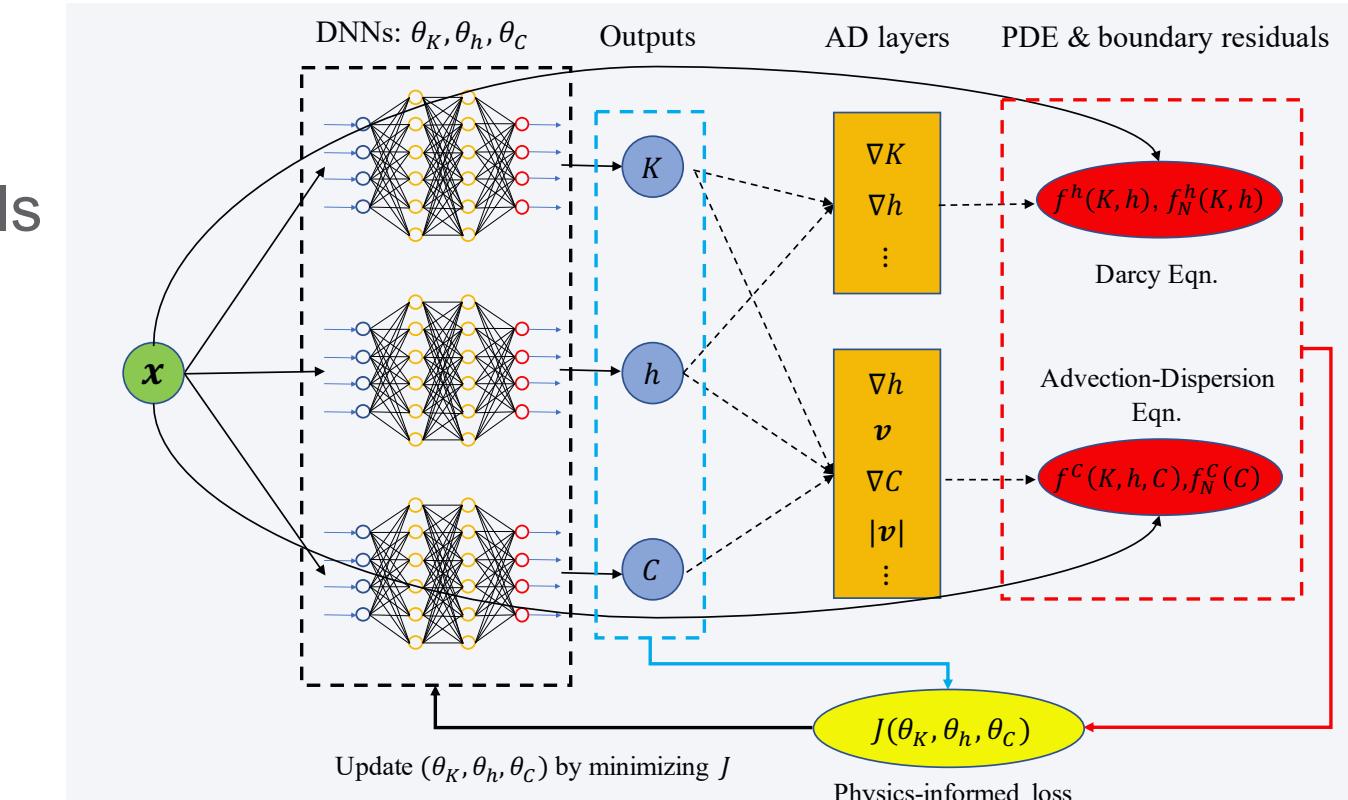
Domain-awareness

DATA

- Knowledge representations
 - Equations
 - Knowledge graphs
 - Causal models
 - Distributed knowledge models
- Domain-aware features
- Knowledge embedding
 - Constraints
 - Loss function
 - Coupling with simulations
 - Model structure

ALGORITHMS

Coupling limited experimental data with physics models improves accuracy and reduces the risk of model overfitting in subsurface transport modeling.

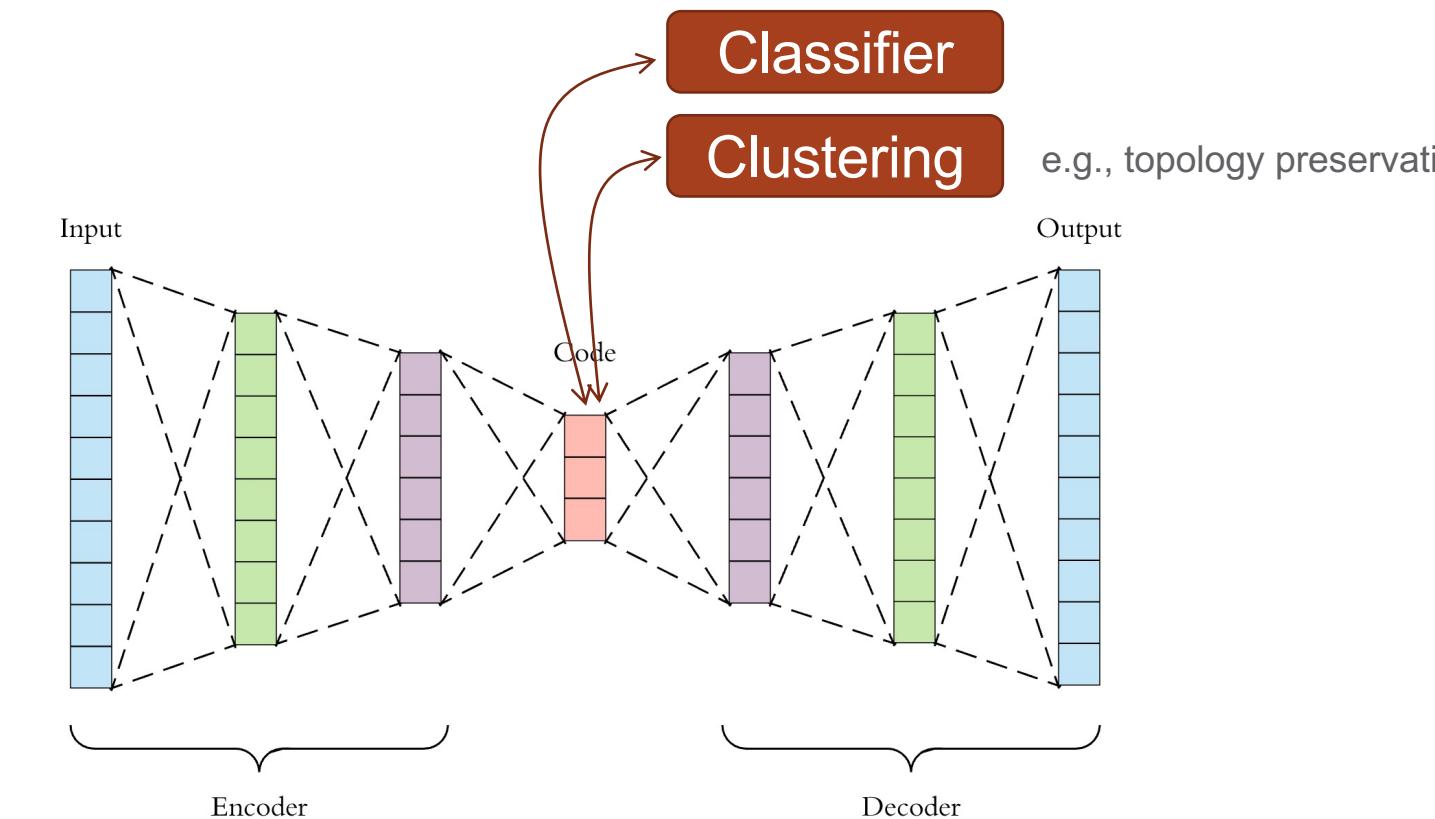


He Q et al.. 2020 Physics-informed neural networks for multi-physics data assimilation with application to subsurface transport. *Advances in Water Resources* (2020) 141.

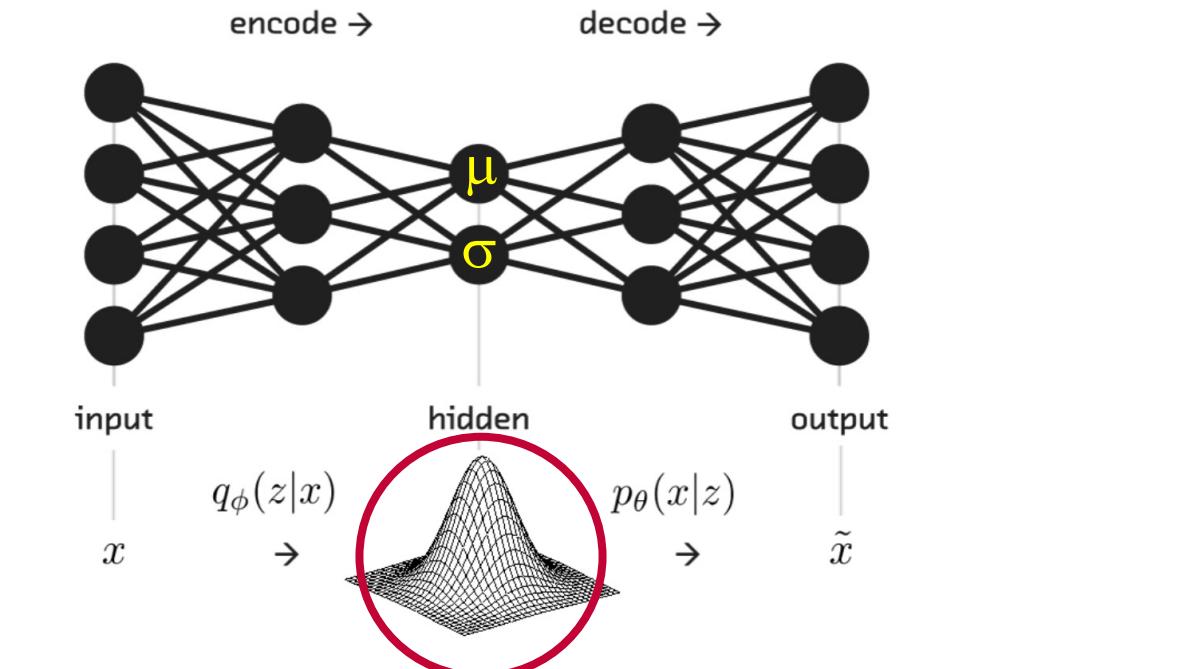
Domain-aware features

- Coupling autoencoders with:
 - Clustering algorithms
 - Classification algorithms

$$\mathcal{L}(x, \hat{x}) + \lambda_1 \sum_i |a_i| + \lambda_2 \mathcal{L}_{\text{clustering classifier}}$$



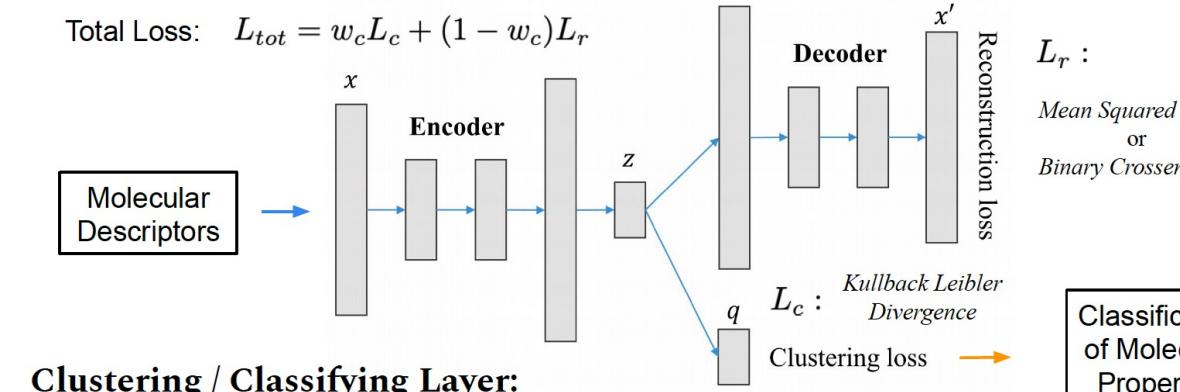
- Variational autoencoders with:
 - Predefined priors



- Domain-informed prior
- KL term in the loss function enforces the desired probability distribution

Biodegradability-aware embeddings

$$\text{Total Loss: } L_{tot} = w_c L_c + (1 - w_c) L_r$$



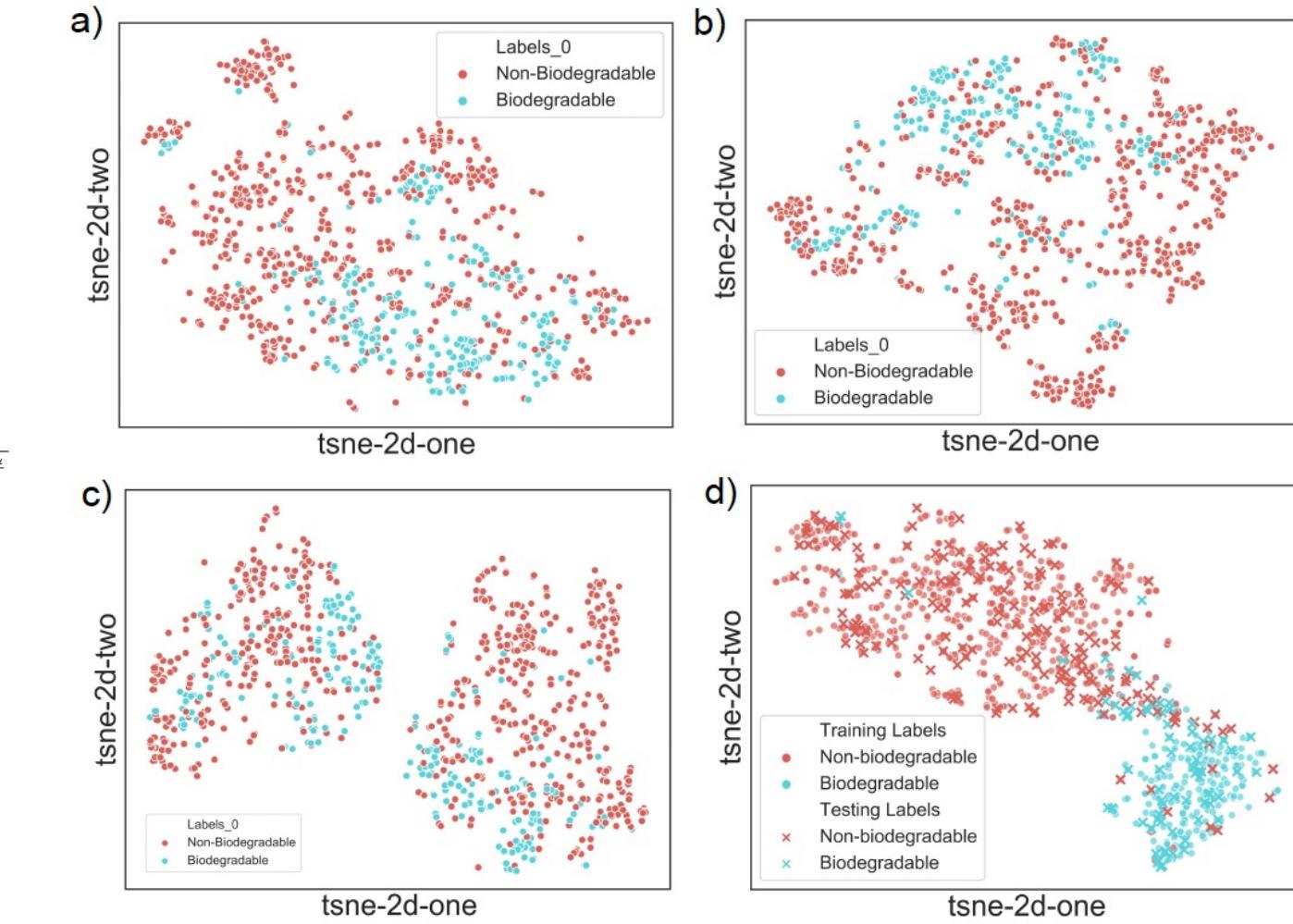
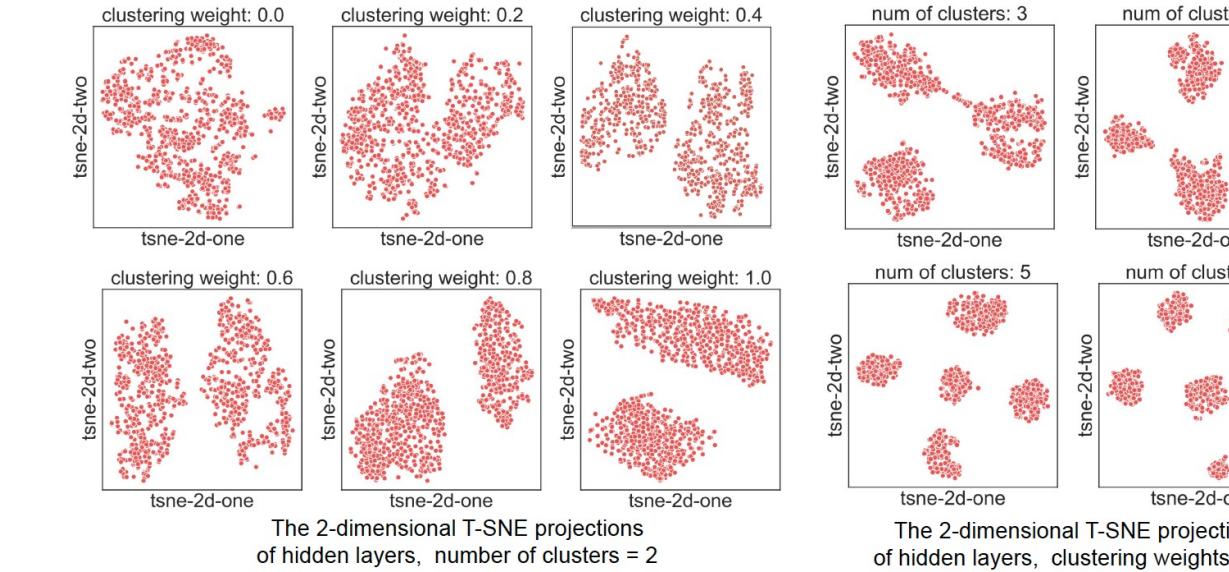
Clustering / Classifying Layer:

$$\text{Loss: } L_c = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$$\text{Label hardening for clustering layer: } p_{ij} = \frac{q_{ij}^2 / \sum_i q_{ij}}{\sum_j q_{ij}^2 / \sum_i q_{ij}}$$

$$\text{Soft label assigning point } z_i \text{ to cluster } \mu_j: q_{ij} = \frac{(1 + \|z_i - \mu_j\|^2)^{-\frac{1+\alpha}{2}}}{\sum_j (1 + \|z_i - \mu_j\|^2)^{-\frac{1+\alpha}{2}}}$$

$$\text{Using true label for classifying layer: } p_{ij} = \begin{cases} 1, & j = \text{true label} \\ 0, & \text{otherwise} \end{cases}$$

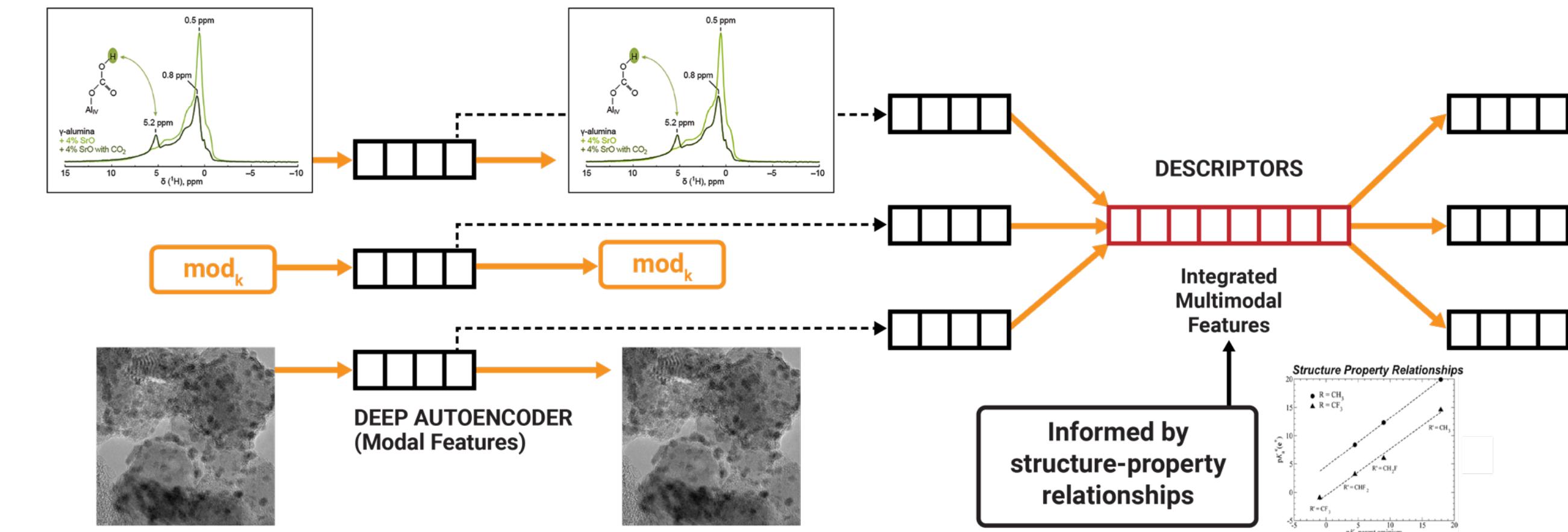


t-SNE projection: (a) raw data; (b) autoencoder; (c) clustering ; (d) classification



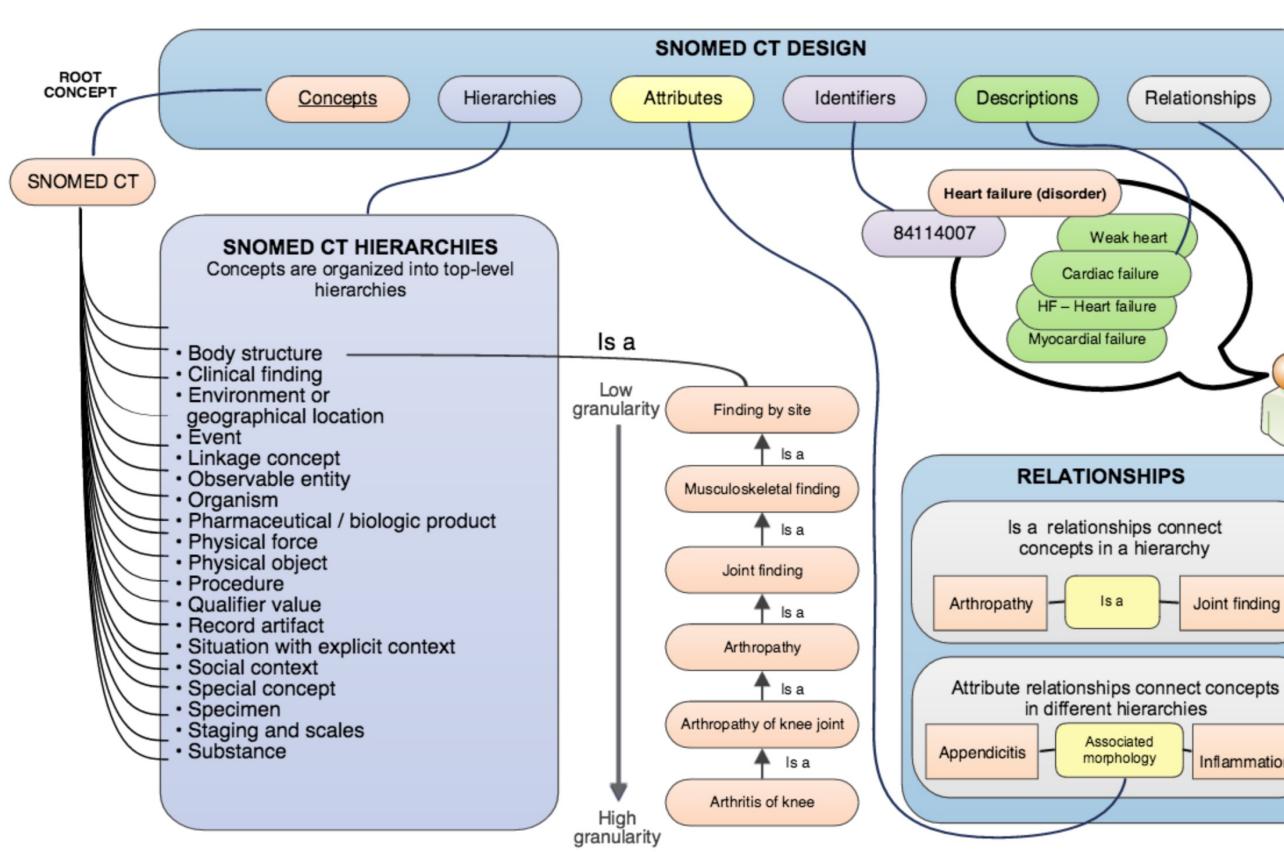
Integrated multi-modal features

- Heterogeneous data
- Domain-informed stacked autoencoders → integrated features

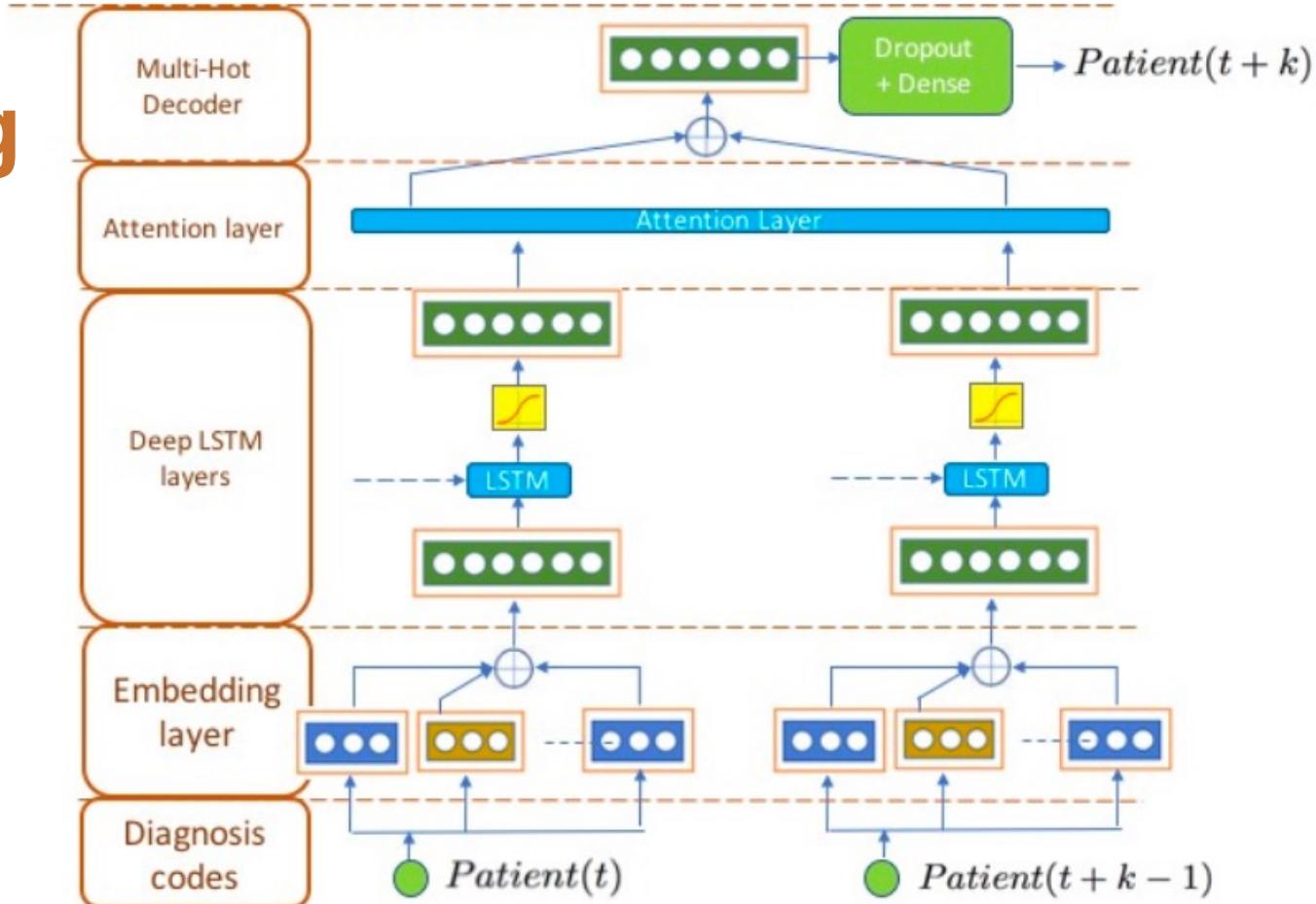


Knowledge embedding

- Medical knowledge / Knowledge graphs



Patient State Prediction (All Diagnosis)
Patient State Prediction (Frequent 20)
Patient State Prediction (Rare 20)

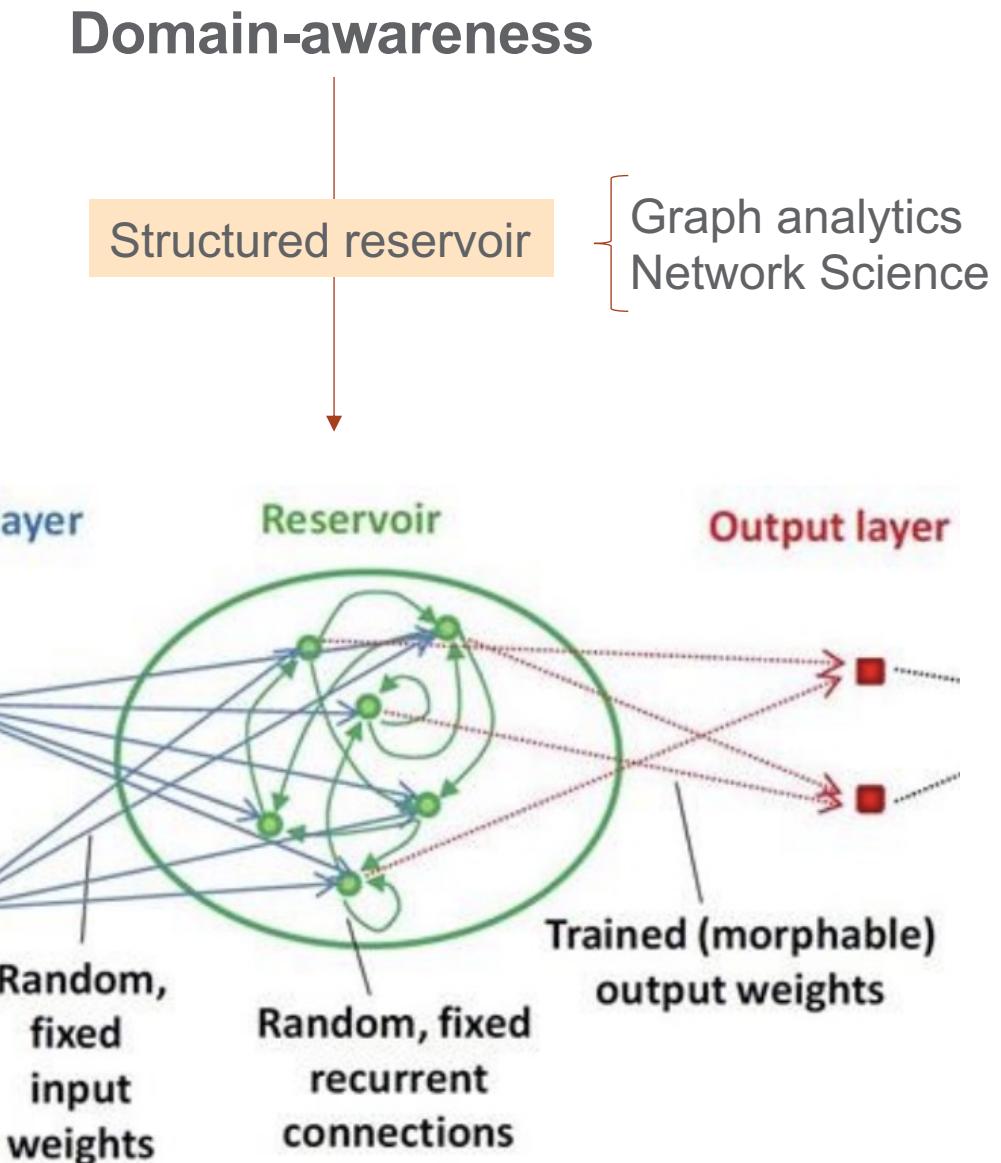


	Node2vec	Metapath2vec	Poincare	CUI2vec	Med2vec
Node Classification	0.817	0.3287	0.8579	0.5685	0.0409
Link Prediction	0.986	0.3988	0.7135	0.7222	0.8665
Concept Similarity (D1)	0.79	0.3	0.7	0.16	NA
Concept Similarity (D3)	0.90	0.46	0.31	0.15	NA
Concept Similarity (D5)	0.81	-0.32	-0.06	-0.01	NA
Patient State Prediction (All Diagnosis)	0.3938	0.3359	0.4197	0.3948	0.3881
Patient State Prediction (Frequent 20)	0.8465	0.9749	0.85	0.8035	0.7980
Patient State Prediction (Rare 20)	0.018	0.001	0.019	0.019	0.011



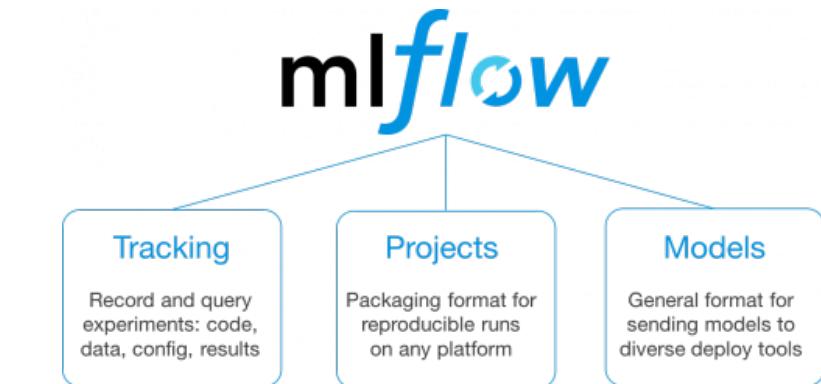
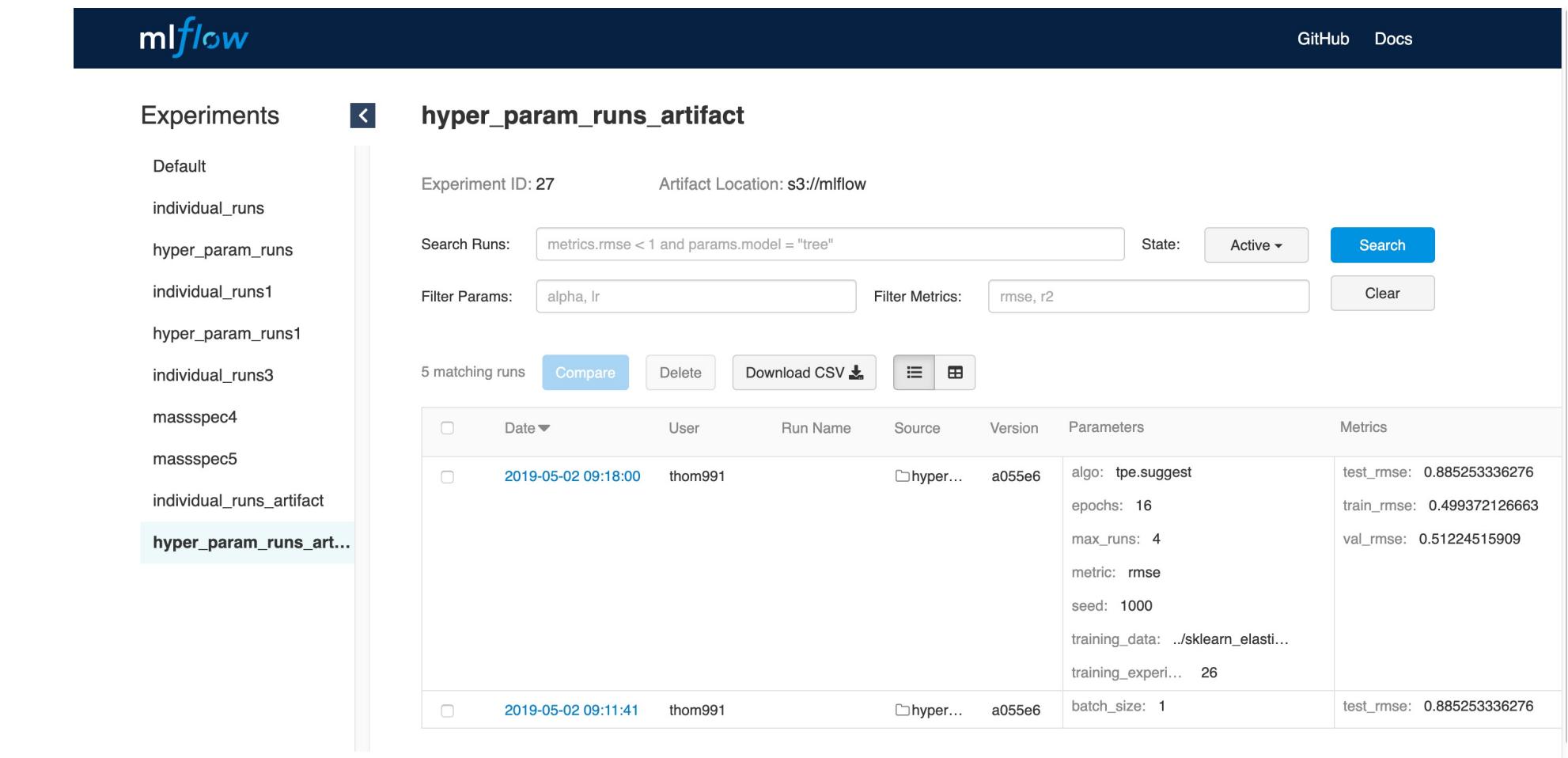
Model structure

- Reservoir computing
 - Generalization of recurrent neural networks
 - Dynamical systems
 - Maps inputs onto a high-dimensional space
 - Hardware implementation
- Elements
 - Input layer
 - ✓ random weights
 - Reservoir
 - ✓ Random sparse connectivity
 - ✓ Non-linear activation
 - Readout layer
 - ✓ Linear transformation of the reservoir state
 - ✓ Fast adaptation using ridge regression



Reproducibility

- Provenance tracking in AI/ML workflows
- Metadata on:
 - Task
 - Data
 - Algorithms
 - (Hyper)parameters
 - Performance
- Artifacts
- Packaging
- Deployment

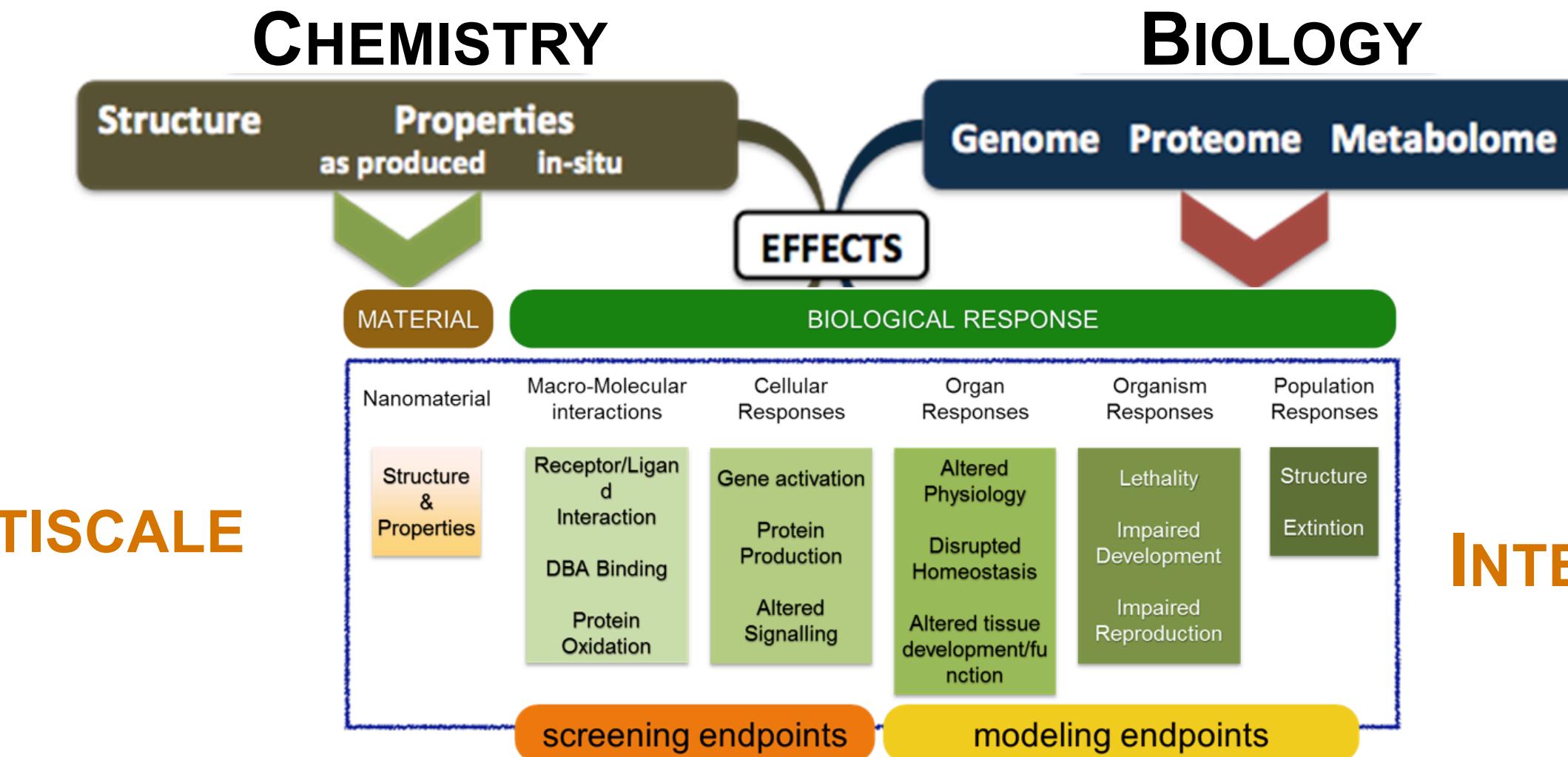
The screenshot shows the mlflow UI interface. At the top, there's a navigation bar with the mlflow logo, GitHub, and Docs links. Below the header, the page title is "hyper_param_runs_artifact". On the left, a sidebar lists various experiment runs: Default, individual_runs, hyper_param_runs, individual_runs1, hyper_param_runs1, individual_runs3, massspec4, massspec5, individual_runs_artifact, and hyper_param_runs_art... (which is highlighted). The main content area displays experiment details: Experiment ID: 27, Artifact Location: s3://mlflow. It includes search and filter fields: Search Runs: metrics.rmse < 1 and params.model = "tree", State: Active, Filter Params: alpha, lr, Filter Metrics: rmse, r2, and a "Search" button. Below these are buttons for Compare, Delete, Download CSV, and two icons. A table then lists the matching runs:

	Date	User	Run Name	Source	Version	Parameters	Metrics
<input type="checkbox"/>	2019-05-02 09:18:00	thom991	<input type="checkbox"/> hyper...	a055e6		algo: tpe.suggest epochs: 16 max_runs: 4 metric: rmse seed: 1000 training_data: ./sklearn_elasti... training_experi... 26	test_rmse: 0.885253336276 train_rmse: 0.499372126663 val_rmse: 0.51224515909
<input type="checkbox"/>	2019-05-02 09:11:41	thom991	<input type="checkbox"/> hyper...	a055e6		batch_size: 1	test_rmse: 0.885253336276

Interpretability: Integrated Domain-informed frameworks

MULTISCALE

DATA
INTEGRATION



Multiple levels of biological organization

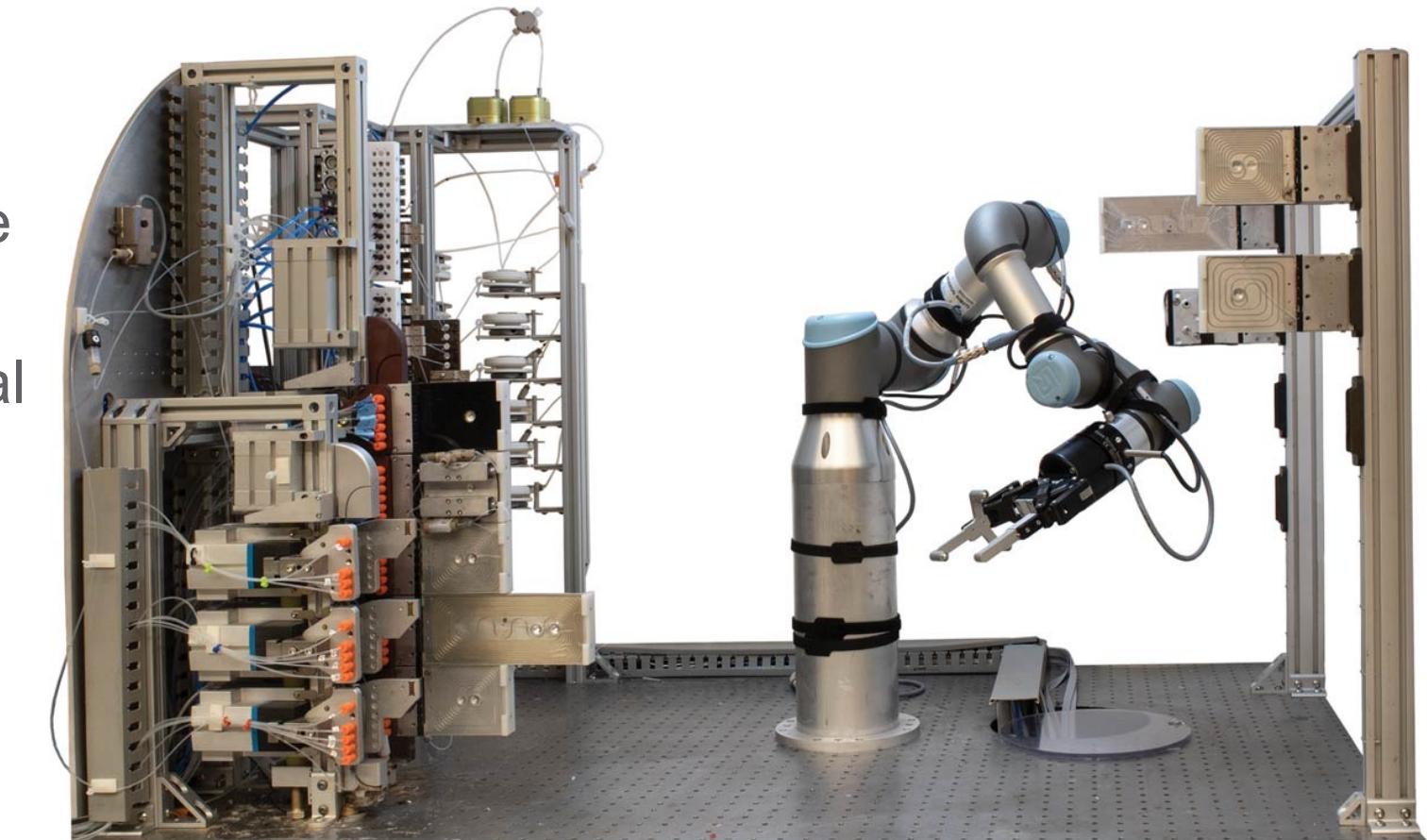
Automation

Interfacing with scientific instruments

- Edge computing for in-situ processing.
- Opportunity to leverage existing real-time control and optimization capabilities.
- Optimal exploration of large combinatorial spaces.

Beyond automation: autonomy

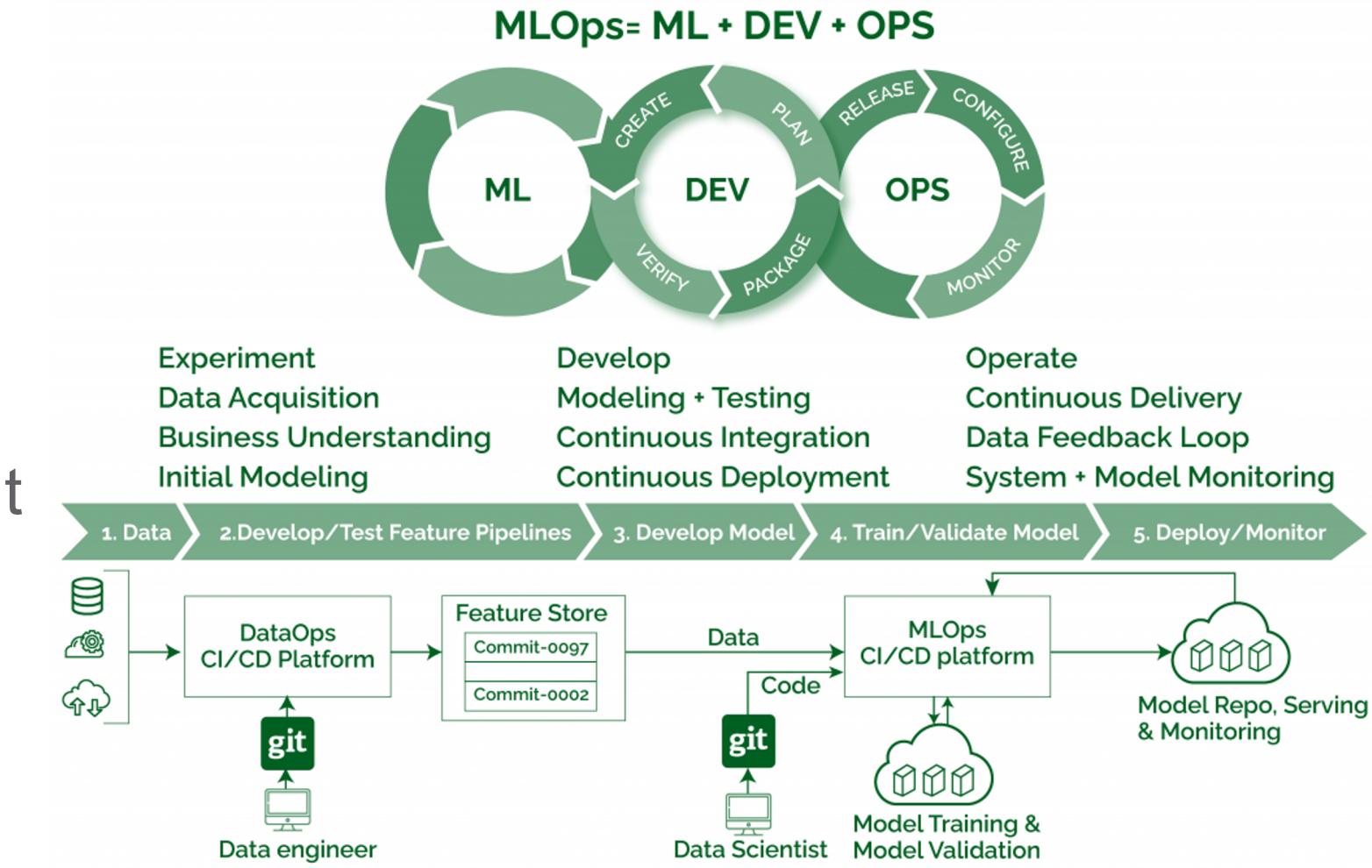
- Hypothesis generation, validation and refinement.
- From machine learning to machine reasoning





Scalability and Deployment

- Rethinking software engineering
 - Composable software systems
 - Heterogeneous computing
- Connecting “Software 1.0” with “Software 2.0”
- Managed and scalable deployment of AI/ML models
 - Lifecycle management
 - Reproducible results
 - Model governance
 - Monitoring and operational use



Adapted from:

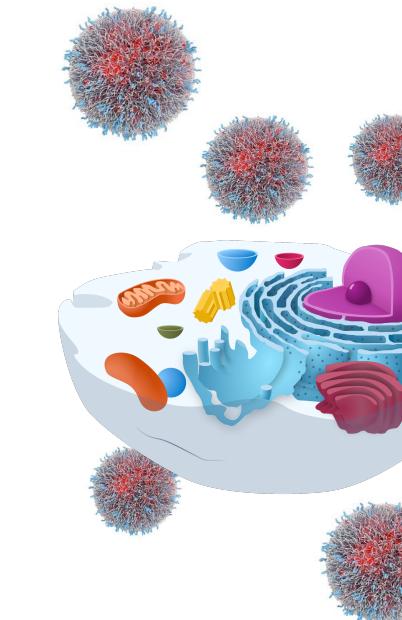
[The fundamentals of MLOps – The enabler of quality outcomes in production environments](#)



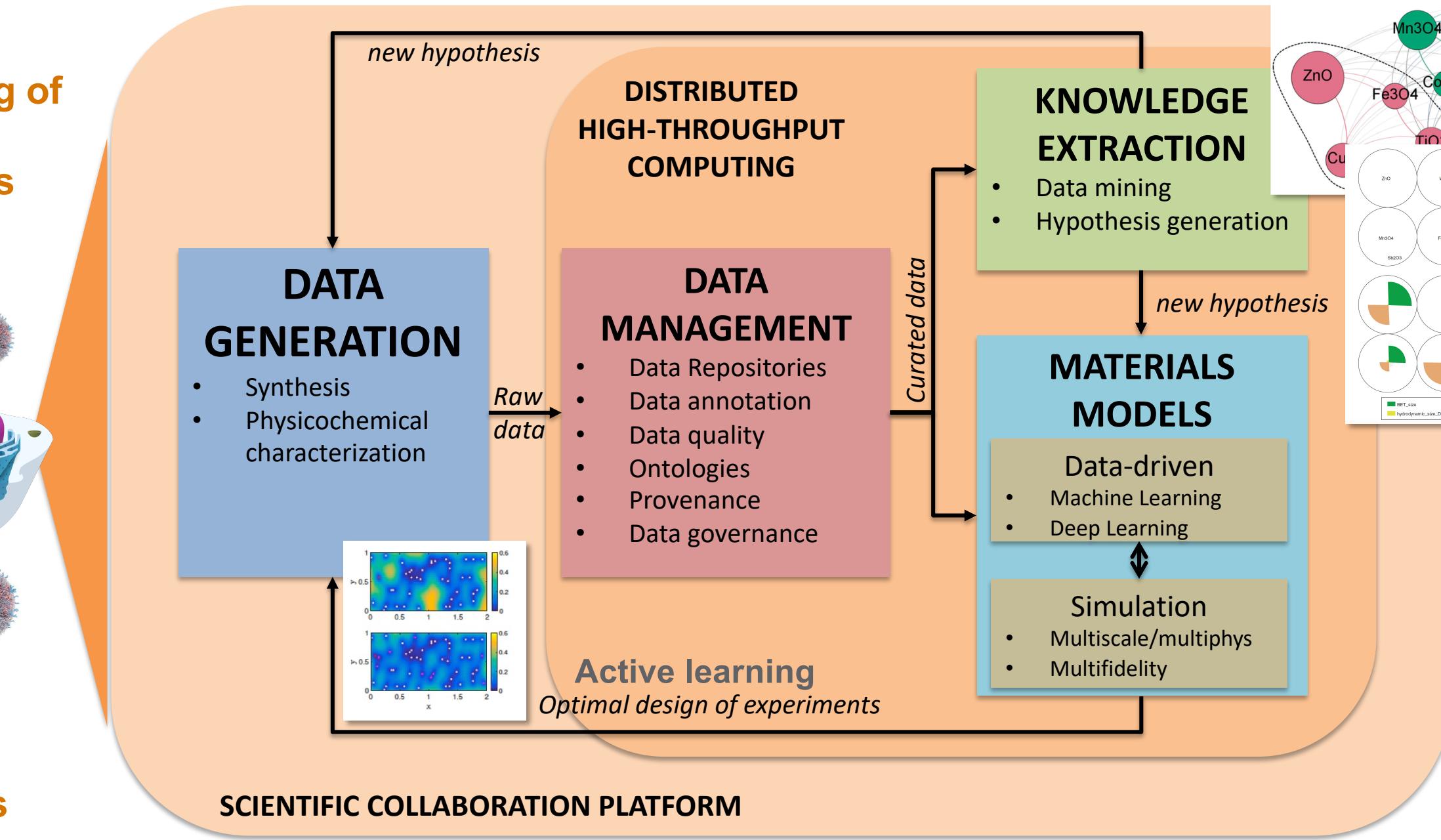
Pacific
Northwest
NATIONAL LABORATORY

Summary

Better
understanding of
bio-nano
interactions



Enhanced
predictive
capabilities





Pacific
Northwest
NATIONAL LABORATORY

Thank you

