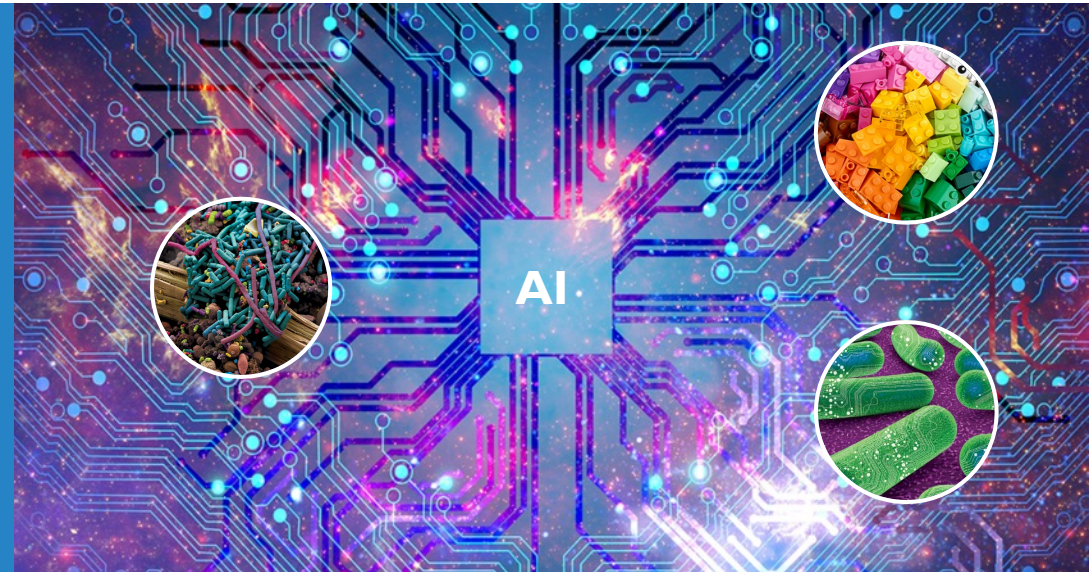# Foundation models for complex biological systems + design
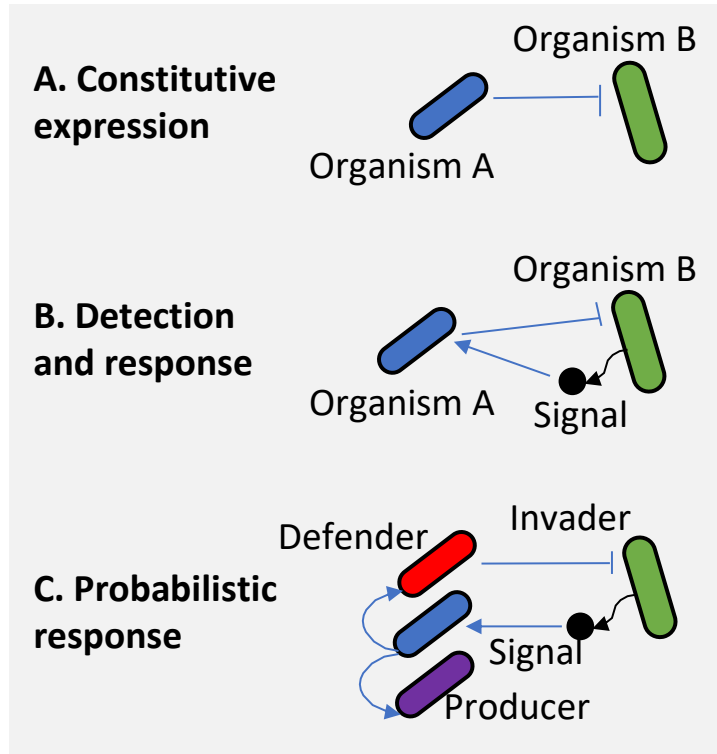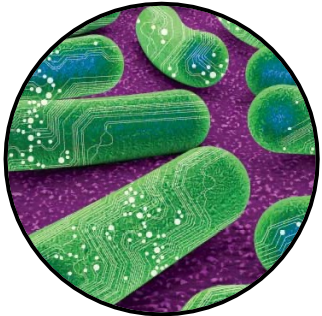
**Arvind Ramanathan/ ramanathana@anl.gov**

**Argonne National Laboratory/ University of Chicago Consortium for Advanced Science and Engineering (CASE)**

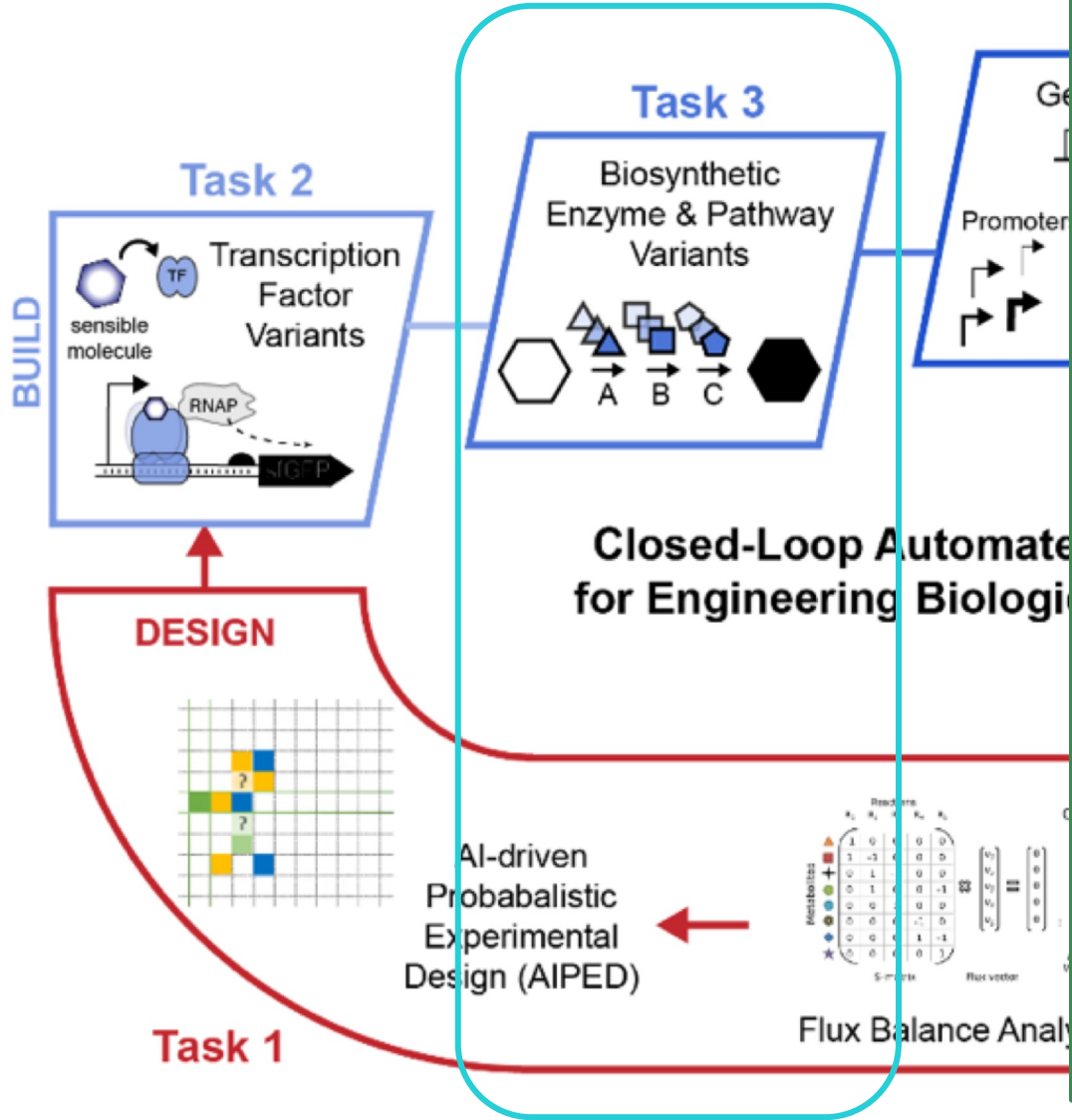**Northwestern-Argonne Institute for Science and Engineering (NAISE)**

# Project overview



**A. Constitutive expression**

Organism B
Organism A

**B. Detection and response**

Organism B
Organism A    Signal

**C. Probabilistic response**

Defender    Invader
Signal
Producer

- Building microorganisms
  - System for generating **biological modules** (*in vitro* closed loop system)
  - **Biocontainment** development and implementation (CRISPR)
  - System for **coupling modules with biocontainment** (CRAGE)
  - System for **evaluating integrated systems** (experimental; modeling; ML- and NLP-based modeling)
- Building communities of microorganisms
- Concept of biocontainment

# An integrative platform for rapid engineering of biological parts
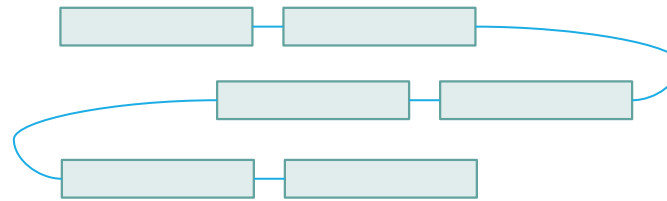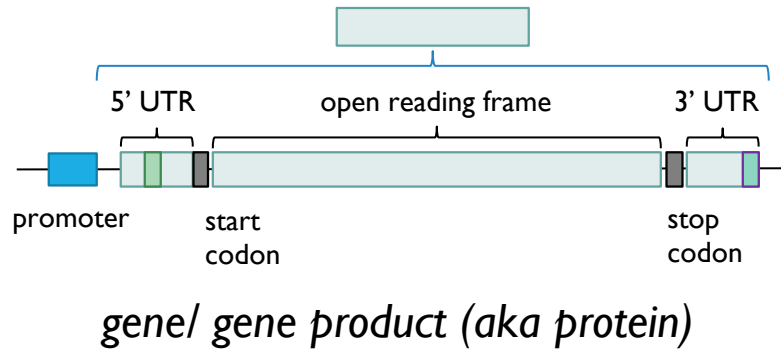


**Poster 2: Priyanka Setty**

- How to think about biological hierarchy and information?
  - implicit and explicit representation learning
  - genome-scale language models
- How to understand hierarchy of biological information?
  - individual gene/protein
  - pathways
  - genome-scale
- How do we generate new examples?
  - designing new genes (MDH as an example)
  - understanding how SARS-CoV-2 is evolving
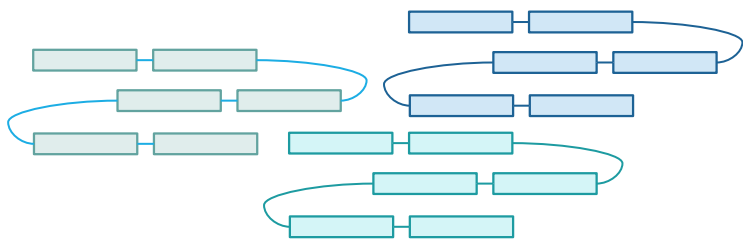  - integrating with biophysics (experiments + simulations)
- Future work/perspectives

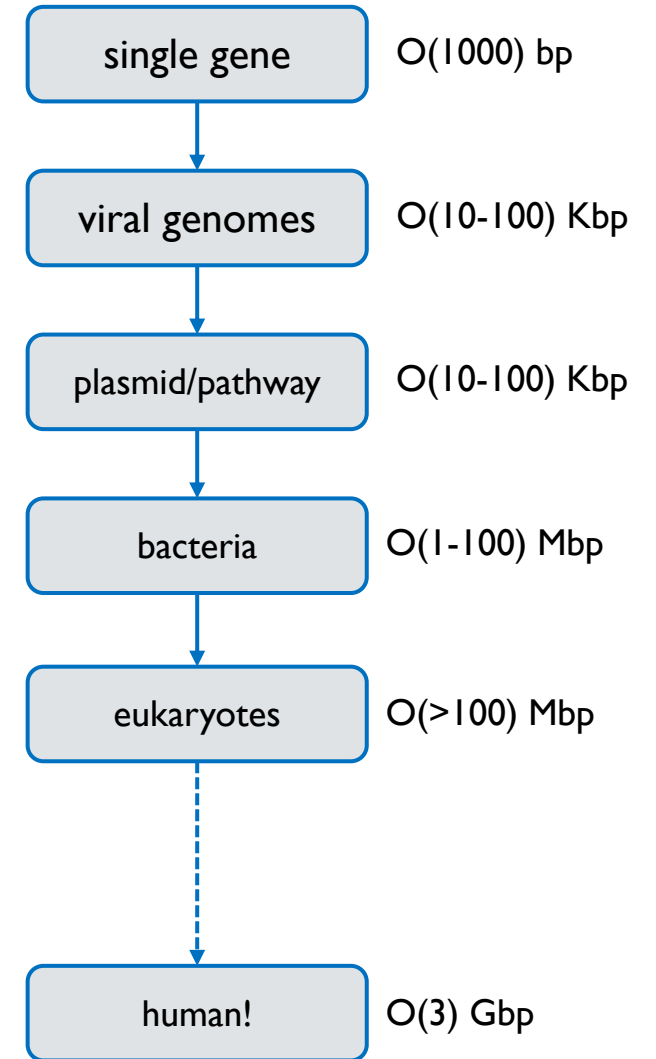# Biological information and hierarchy

# Hierarchical information within '-omics' data



*gene/ gene product (aka protein)*

5' UTR

open reading frame

3' UTR

promoter

start codon

stop codon

*collection of genes (either as "contigs" or ORFs)*

*entire genomes*

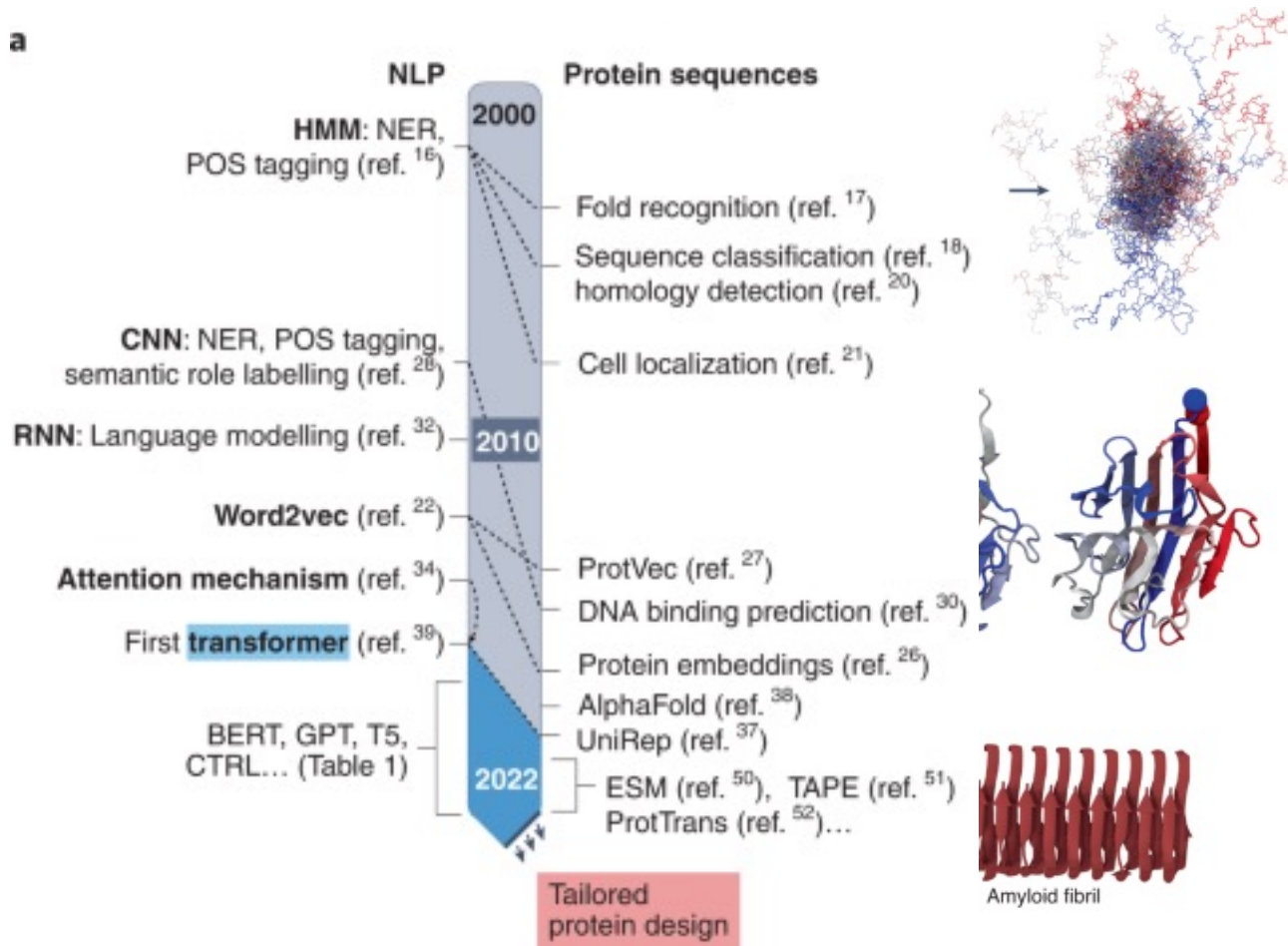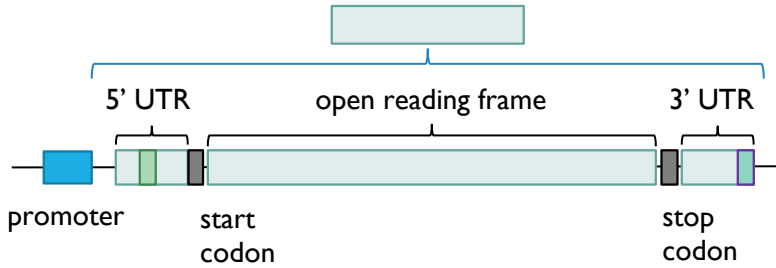| single gene | O(1000) bp |
| viral genomes | O(10-100) Kbp |
| plasmid/pathway | O(10-100) Kbp |
| bacteria | O(1-100) Mbp |
| eukaryotes | O(>100) Mbp |
| human! | O(3) Gbp |

# Traditional approach: use protein language models…



- Transformers have evolved to be good constructs for capturing protein "language"
  - consists of 20 natural amino acids chained together
  - attention mechanism capture interactions across amino-acids
- Considerable challenges in translating NLP based approaches to protein language
  - sequence length
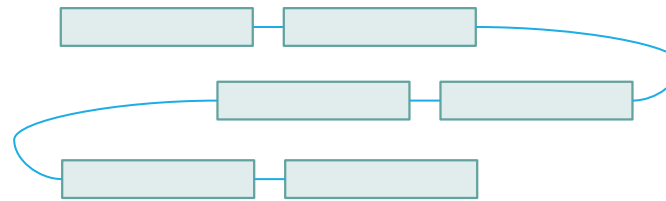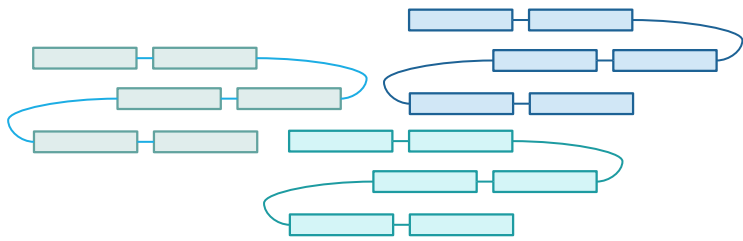  - attention mechanism

Controllable protein design with language models, N. Ferruz, B. Hocker, Nature Machine Intelligence (https://www.nature.com/articles/s42256-022-00499-z)

# Genome-scale language models (GenSLM)



Enzyme/transcription factor sequences - codon level tokenization, GPT models

*gene/ gene product (aka protein)*

Open reading frames – codon level tokenization, Reformer model

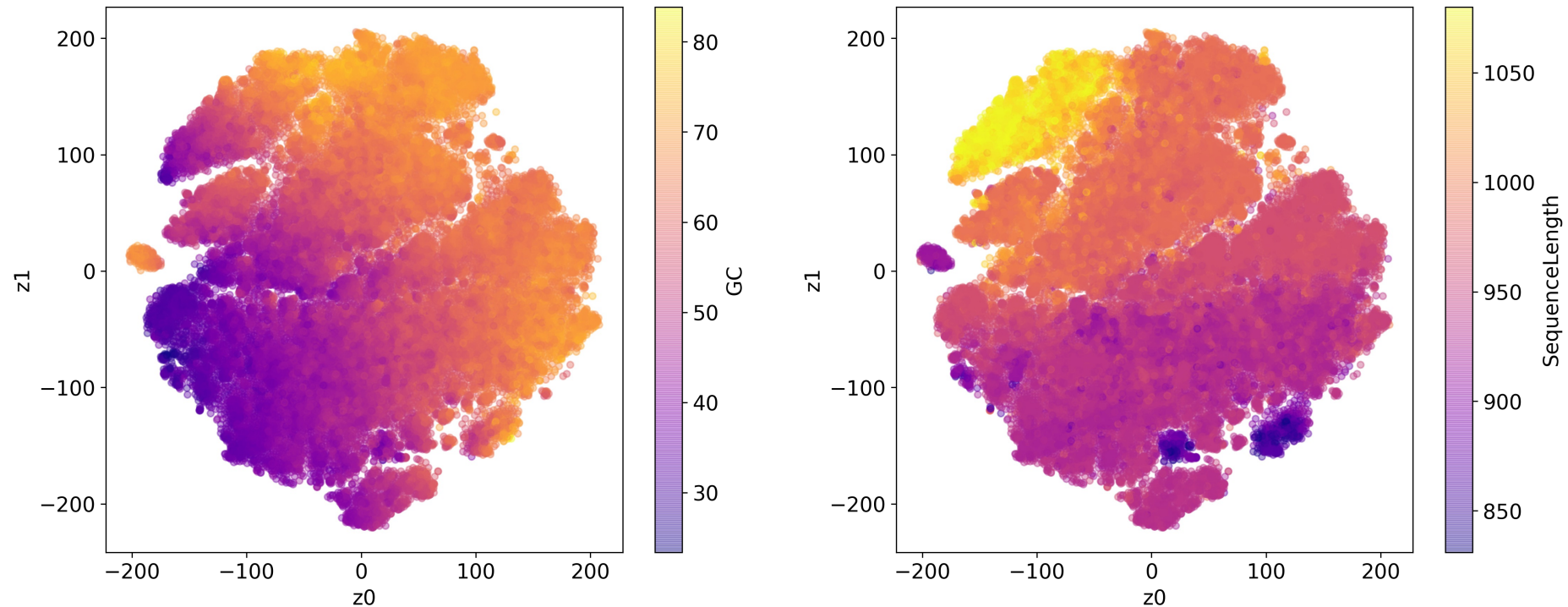*collection of genes (either as "contigs" or ORFs)*

Full genome sequences - BPE Encoding, cannot currently generate full genomes
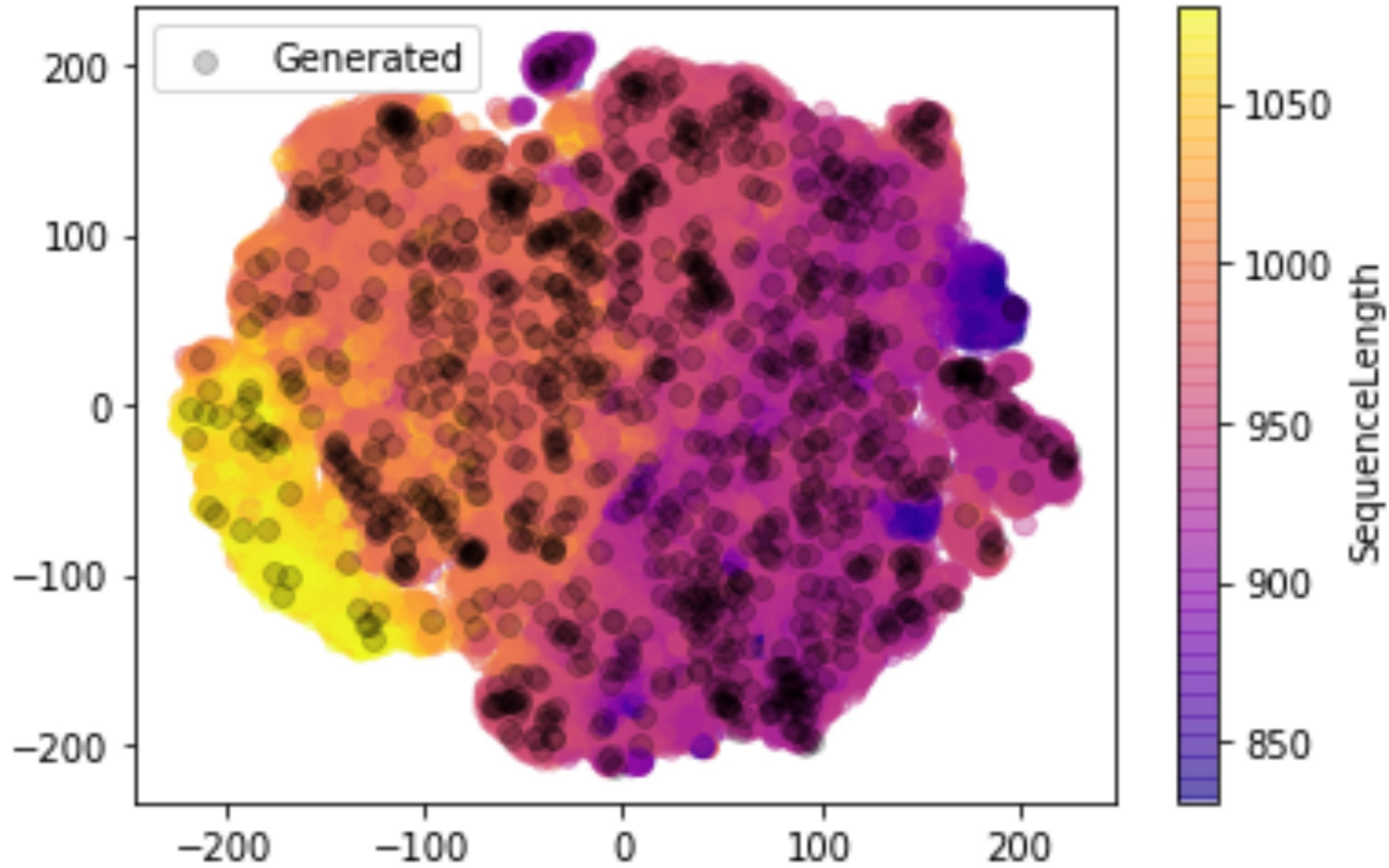
*entire genomes*

- Go beyond traditional k-mer models:
  - variable length issues
- At each level of hierarchy maintain information learned at the lower levels (gene → collection/cluster → full genomes)
- Scale at each level but "tie" it together with stable diffusion models

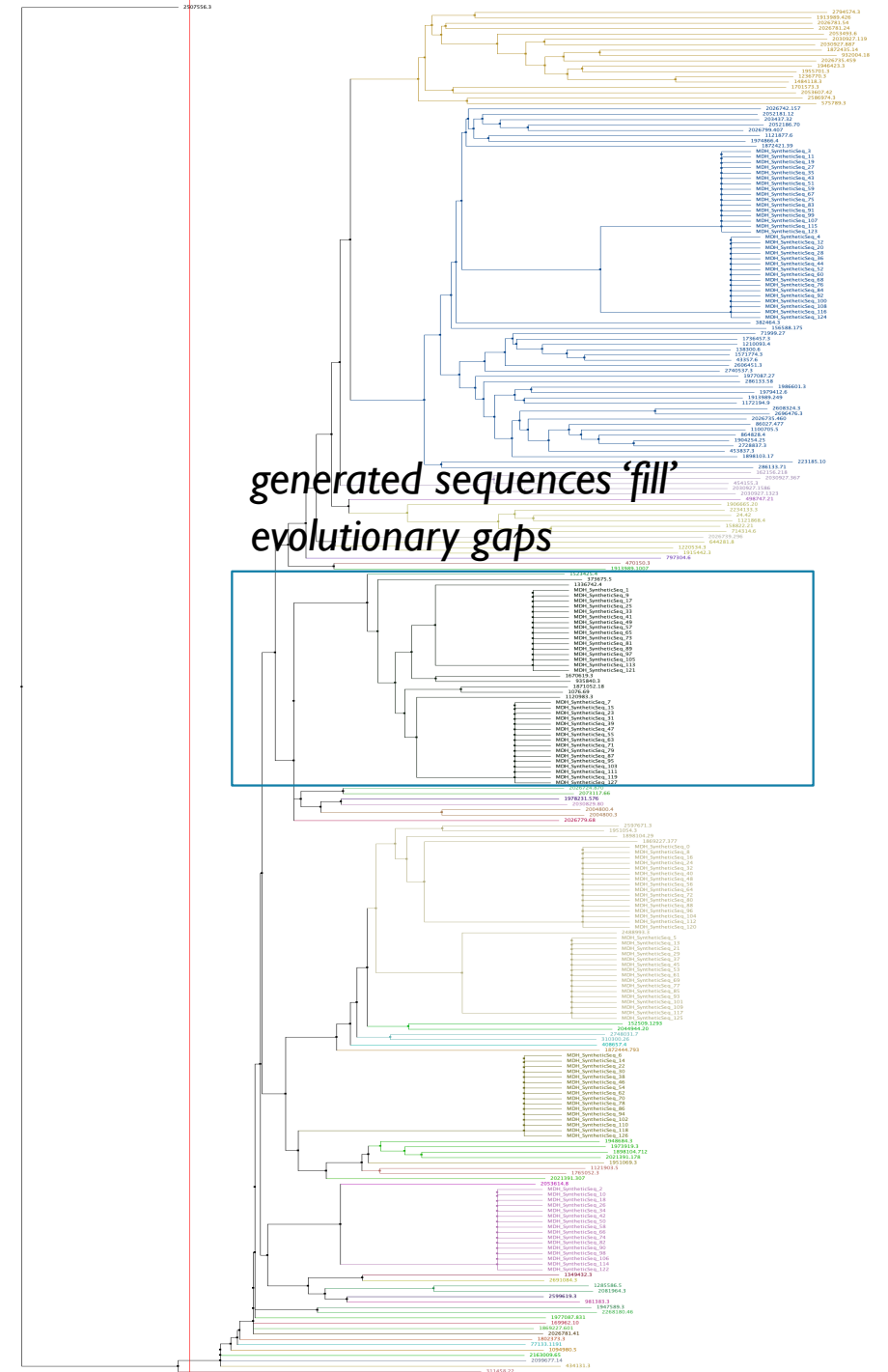# GenSLMs capture individual gene-level information



- Unsupervised learning on 36K malate dehydrogenase (MDH) sequences

- Embedding of the latent space using t-SNE (or even UMAP) reveals characteristic features:
  - GC content, sequence length variations, molecular weight…

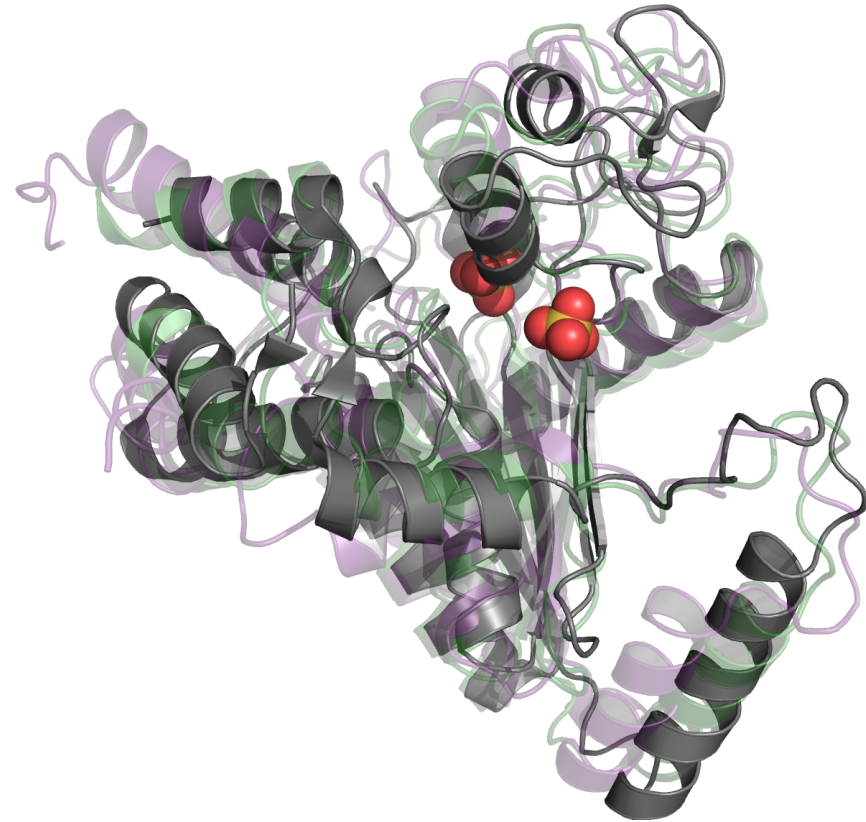- These individual gene-level models can be used to interpolate and sample

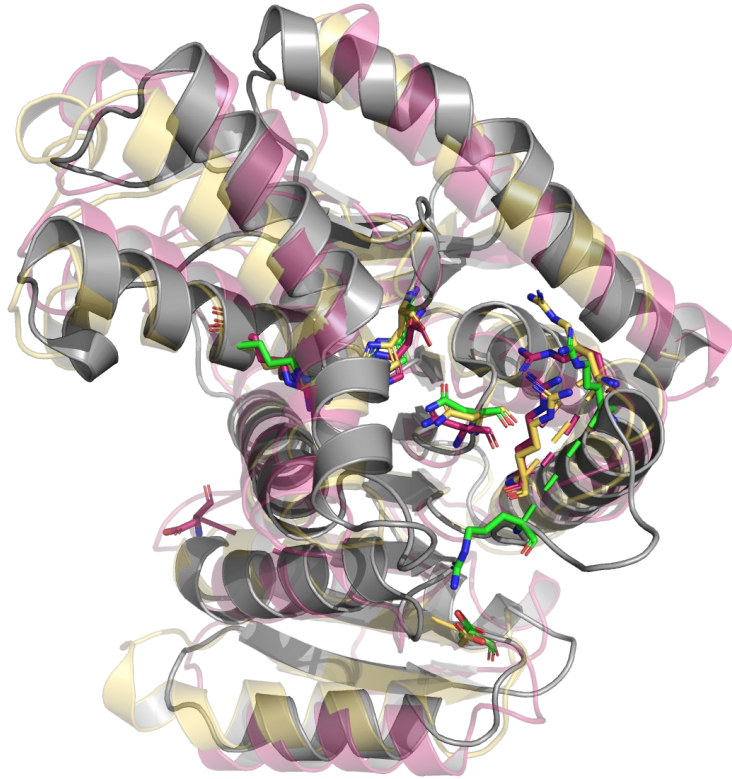# GenSLMs are accurate enough to generate gene sequences…



UMAP embeddings of generated sequences agree with learned embeddings using GPT-2



*generated sequences 'fill' evolutionary gaps*

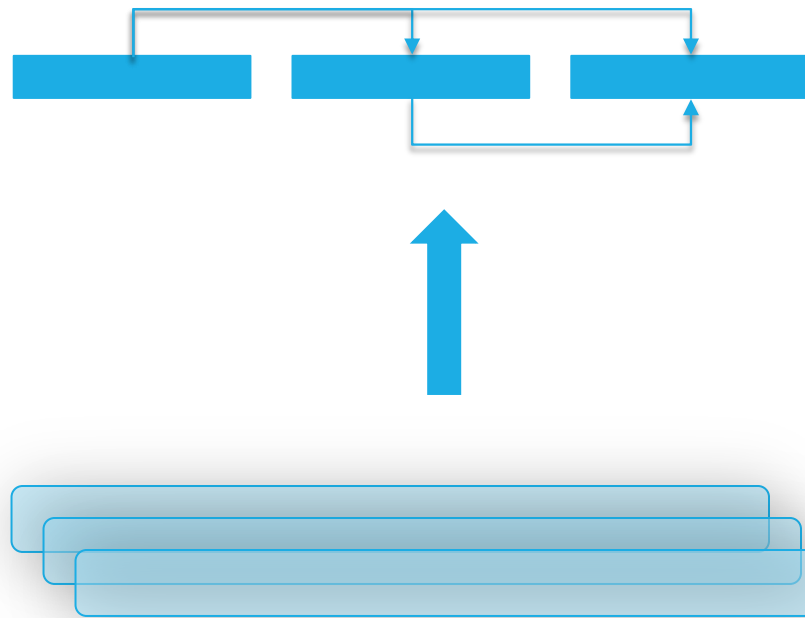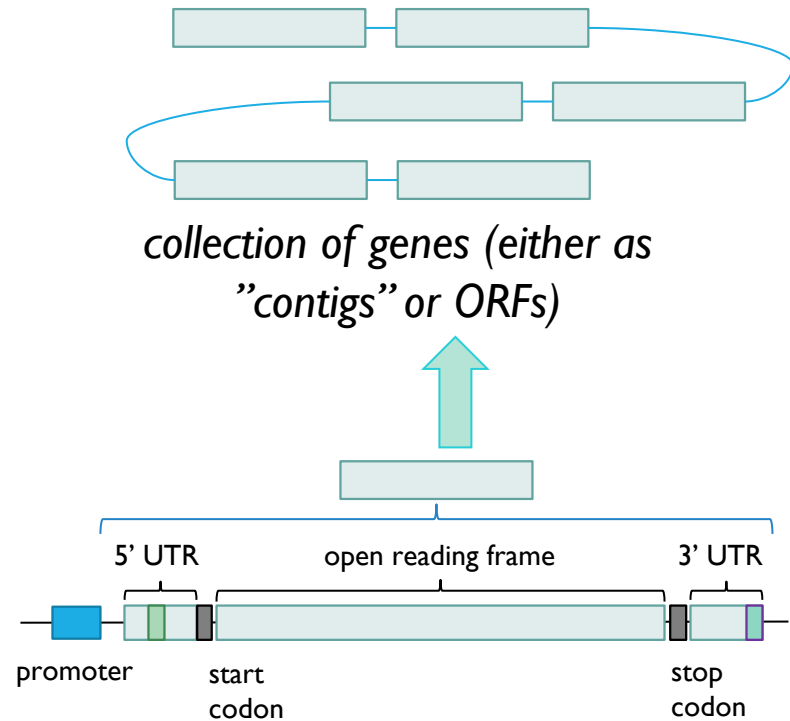# … genes → proteins share MDH similarity at key sites as predicted via OpenFold



GenSLMs learn the two distinct isoforms for MDH and within each isoform we find conservation of key residues and placement of binding sites

# Extending GenSLMs to model SARS-CoV-2 evolutionary dynamics

# Overcoming length limitations of GenSLMs

collection of genes (either as "contigs" or ORFs)

5' UTR    open reading frame    3' UTR

promoter    start codon    stop codon

Stable diffusion models (read and learn context amongst ORFs)

GPT-2 like (8 layers + 8 attention heads + 10240 x 10240 positional encoding (no ordering in ORFs)

- Implicit representation of hierarchy by integrating LLMs with stable diffusion models

# A Foundation Model approach for SARS-CoV-2 genomes ...

**TRAINING**

Train on all available sequences of SARS-CoV-2

- Periodically retrain on new variants sequenced across specific time window
- **Performance**: CS-2, Frontier, Polaris, Perlmutter

**PREDICTION WORKFLOW**

**DETECTION WORKFLOW**

# reveals intrinsic evolutionary patterns of SARS-CoV-2



- Variant clustering

- Semantic distance in the latent space can classify variants of interst

# We can generate new sequences that look like SARS-CoV-2

# We can prioritize sequences that are novel and perceived VOCs

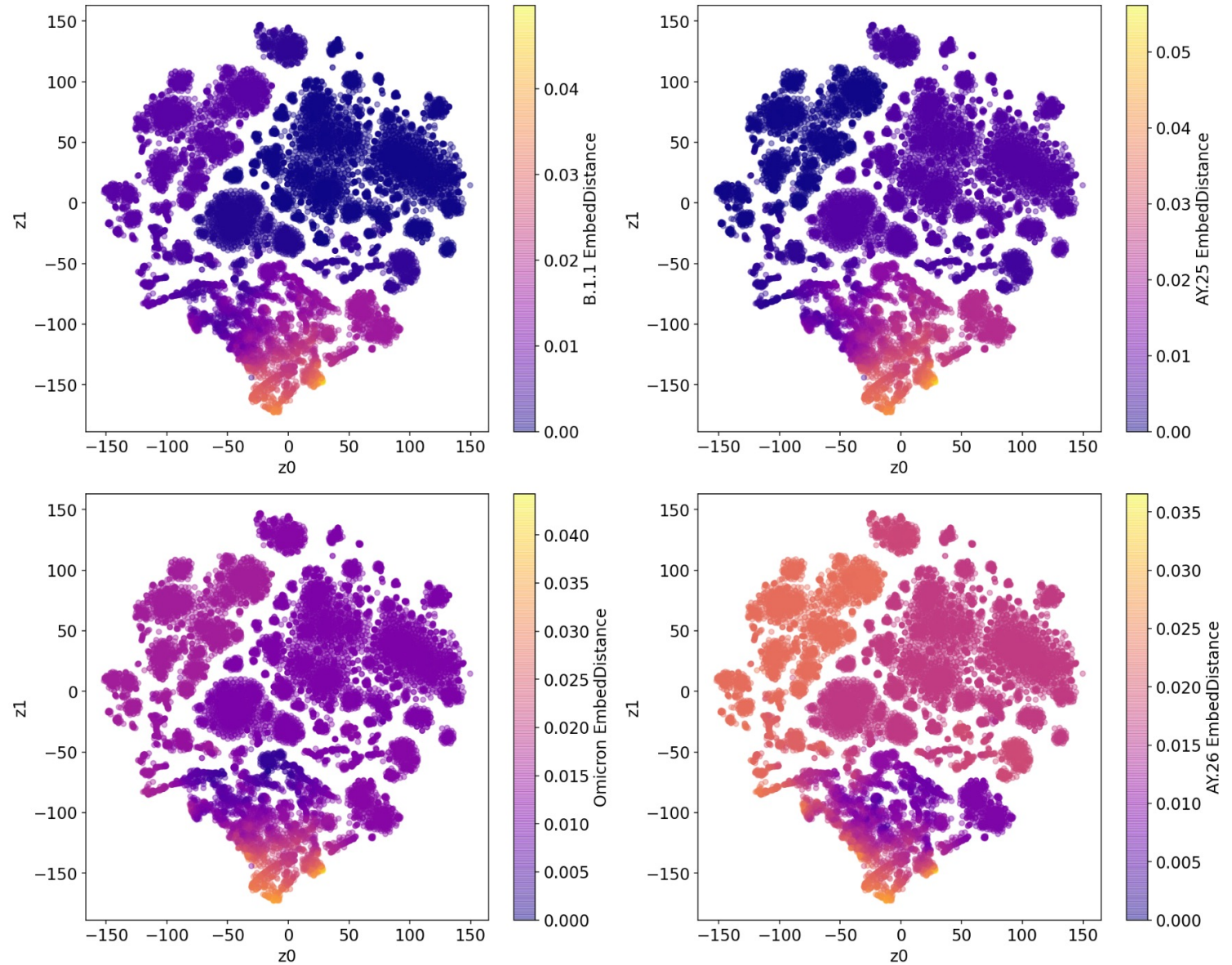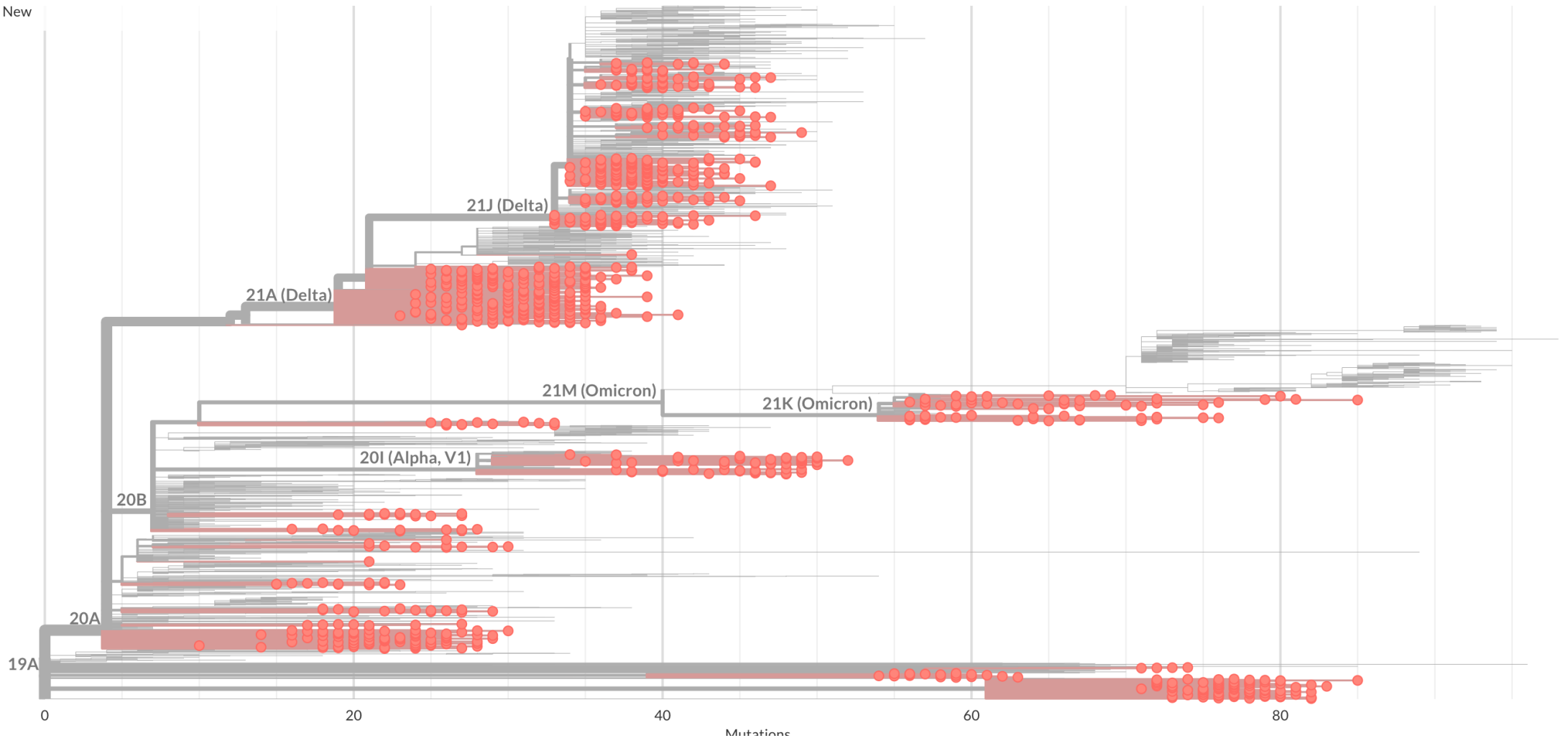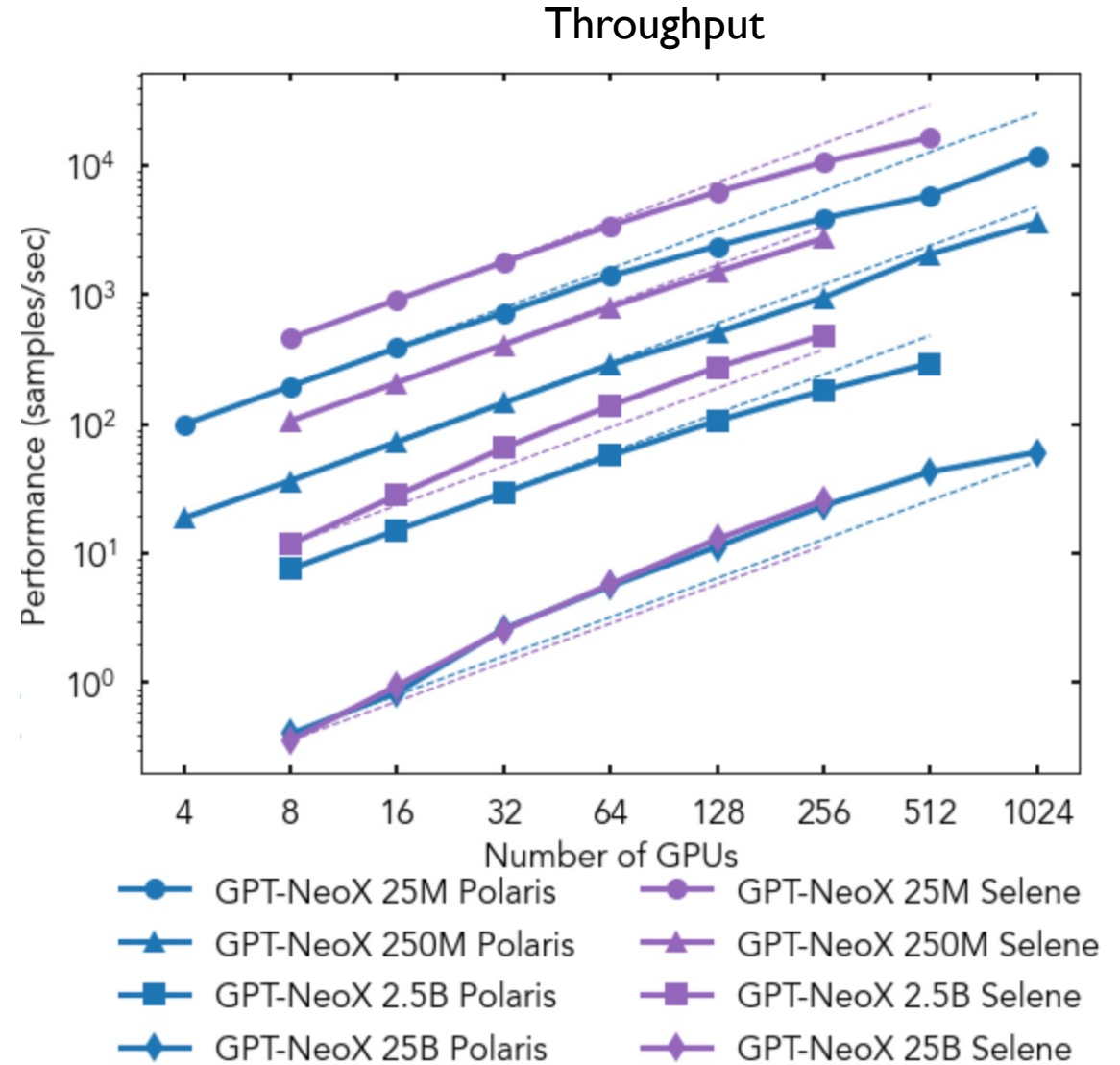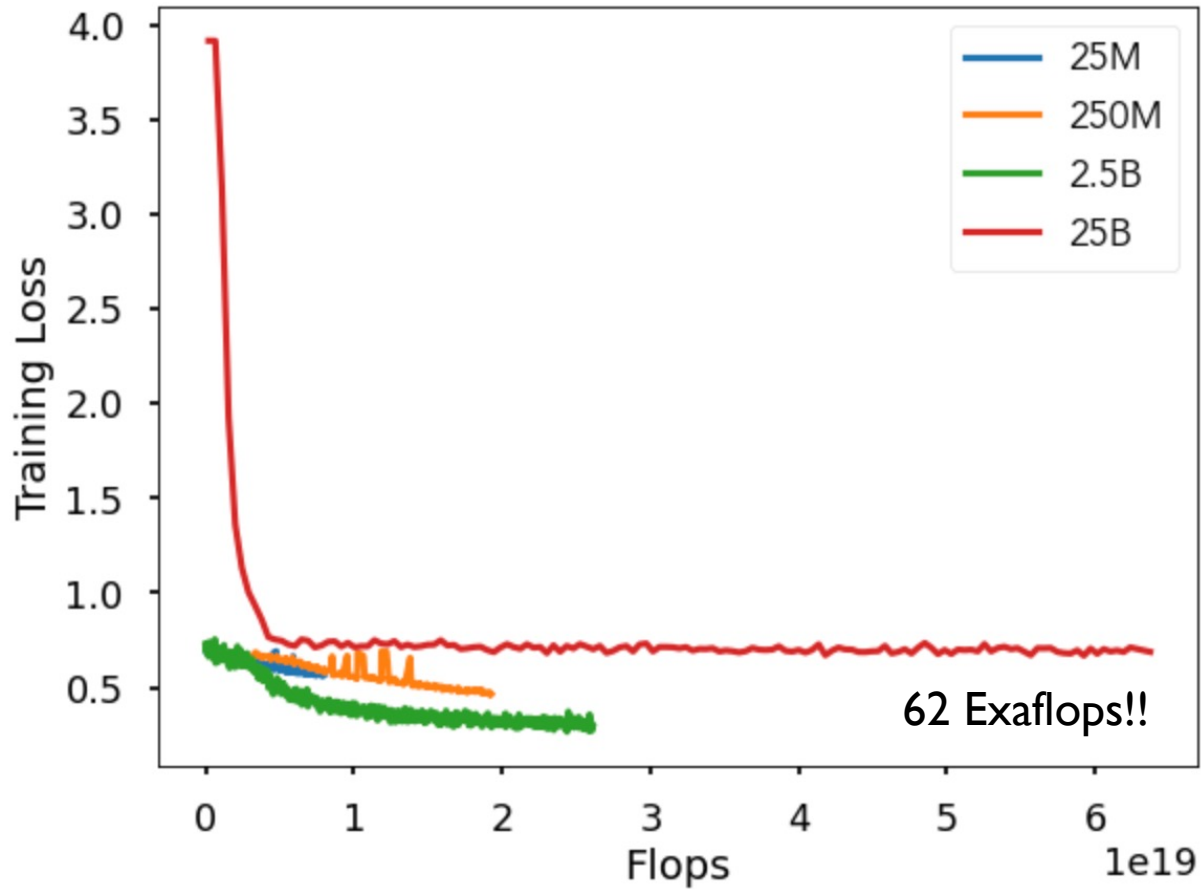| | Genome_ID | Of_Interest | Predicted_Variant | Distance_to_Reference | Neighbors_win_X | K_Neighbors_Variant_Dictionary |
|---|---|---|---|---|---|---|
| **212** | top-p0-9-0212 | True | B.1 | 19.0 | 538 | {'B.1': 16, 'B.1.206': 3, 'B.1.596': 1} |
| **295** | top-p0-9-0295 | True | B.1 | 28.0 | 274 | {'B.1': 16, 'B.1.206': 3, 'B.1.596': 1} |
| **313** | top-p0-9-0313 | True | B.1 | 19.0 | 538 | {'B.1': 16, 'B.1.206': 3, 'B.1.596': 1} |
| **349** | top-p0-9-0349 | True | omicron | 76.0 | 298 | {'omicron': 20} |
| **398** | top-p0-9-0398 | True | B.1 | 28.0 | 274 | {'B.1': 16, 'B.1.206': 3, 'B.1.596': 1} |
| **416** | top-p0-9-0416 | True | B.1.1.7 | 56.0 | 67 | {'B.1.1.7': 17, 'B.1.1': 2, 'None': 1} |
| **438** | top-p0-9-0438 | True | B.1.1.7 | 49.0 | 71 | {'B.1.1.7': 13, 'B.1.1': 5, 'None': 2} |
| **540** | top-p0-9-0540 | True | omicron | 76.0 | 298 | {'omicron': 20} |
| **544** | top-p0-9-0544 | True | B.1.1.7 | 56.0 | 67 | {'B.1.1.7': 17, 'B.1.1': 2, 'None': 1} |
| **715** | top-p0-9-0715 | True | B.1.1.7 | 49.0 | 71 | {'B.1.1.7': 13, 'B.1.1': 5, 'None': 2} |
| **807** | top-p0-9-0807 | True | B.1 | 10.0 | 650 | {'B.1': 16, 'B.1.206': 3, 'B.1.596': 1} |

# Scaling laws for GenSLMs…



Throughput

62 Exaflops!!

Designing a robust automated protein engineering platform…

# System architecture of a protein engineering platform

- Diverse robotic platforms that don't necessarily talk to each other:
  - Liquid handling, sealer, peeler
  - no standardized interface
  - sample movement
- Lack of a programming environment for even relatively simple protocols
- An engineering platform that can inter-operate across diverse robots
- Scalable, open Python API:

https://github.com/AD-SDL

- Coordinates many different systems (UR-3 arms, Hudson platform, etc.)

# Making it all work together

# Digital twins enabled beamline experiments can enable reproducibility

# Towards the automation of scientific experiments…

# Summary

- Foundation models (GenSLMs) can learn from genome-scale datasets:
  - applications for function annotation workflows, completing metagenome data, and many more
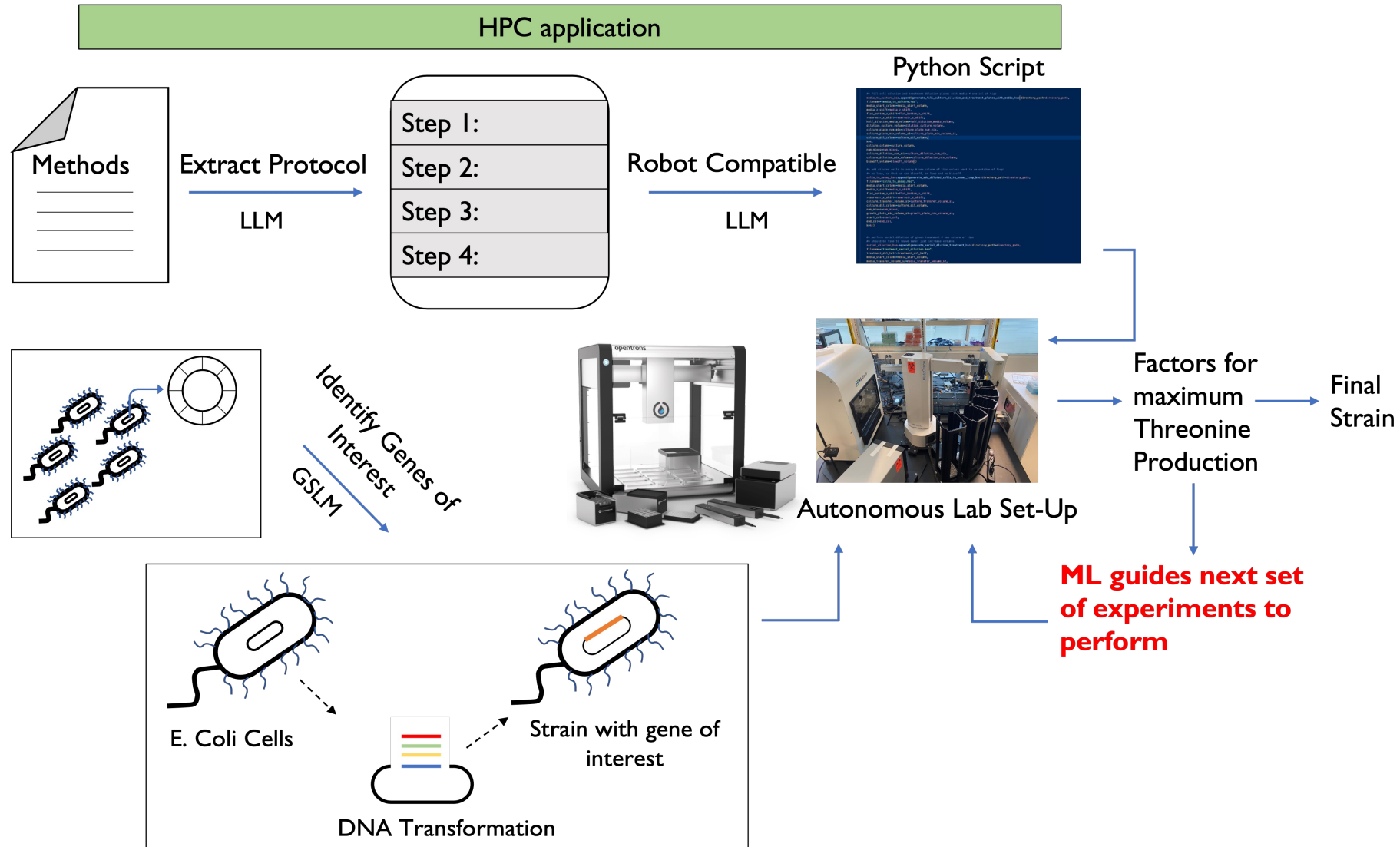  - generative models provide suggestions for experiments → can be integrated with fitness optimization tools for design workflows
  - scaling on whole genomes faces many challenges
- AI-driven simulations will be important for advances in integrating experiments and theoretic understanding of complex bio-systems:
  - Hierarchical AI models enable multi-scale simulations
  - Acceleration with emerging AI hardware
- Automated platforms for experimental design and workflows

# Acknowledgements



Tom Brettin
**Alexander Brace**
NVIDIA
Argonne/
University of
Chicago Mysore

**NVIDIA**
**Maxim Zvyagin**
Argonne

Anima
Anandkumar *
Caltech/
NVIDIA

Hyunseung Yoo
Argonne
**Zongyi Li**
**Caltech**

Ji Yin
**Michael Salim**
Argonne

Aristedis Tsaris
ORNL

Austin Clyde
Argonne/
University of
Chicago

Srinivas
University of
**Pittsburgh**

John McCalpin
TACC
University of
Chicago

Lei Huang
**TACC**

Joseph Insley
Jessica Liu
Argonne
**Cerebras Inc.**

Silvio Rizzi
Subbiah
Argonne
Cerebras Inc.

Sarah Harris*
University of Leeds
**Heng Ma**
Argonne

**Anda Trifan**
UIUC/ Argonne

Geoffrey
Venkatram
Vishwanath*
Argonne

**Defne Gorgun**
UIUC/ Argonne

Rick Stevens
Argonne/
University of
Chicago

**Visualization Aids**

Hardy
**Victor Mateevisti**
Argonne

Tom Burnley
Scientific
Facilities and
Technology
Council

**Anda Trifan**
UIUC/ Argonne

Noah Trebesch

Murali Emani
Argonne/ University
of Chicago

Emad Tajkhorshid*
UIUC

Jim Philipps
UIUC

Janet Knowles
Argonne

# Acknowledgements

## Funding

- DOE- National Virtual Biotechnology Laboratory (NVBL)
- Exascale Computing Project Cancer Deep Learning Environment (CANDLE)
- Exascale Workflows Project (ExaWorks)
- Codesign for Online Data Reduction and Analytics (CODAR)
- DOE Codesign for multimodal AI approaches
- NSF MRI: Multi-modal imaging

## Computing Time

- Argonne Leadership Computing (Theta/ Theta-GPU/ AI-testbed)
- Oak Ridge Leadership Computing (Summit)
- Texas Advanced Computing Center (Frontera/Longhorn)
- National Energy Research Computing Center (Perlmutter)

## Colleagues

- Dan Stanzione
- Paul Lim
- Jack Dillespie
- Dwight Nissley
- Shantenu Jha
- Mike Papka
- Katherine Riley
- Bronson Messer
- Harry Petty
- Rommie Amaro

## Questions/Comments

*ramanathana@anl.gov*