

Biological Data in the Era of Foundation Models, Embodied AI and Autonomous Laboratories



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)

The Traditional vs. Emerging Paradigms of Scientific Discovery

- The current paradigm for discovery is inefficient, time-consuming, labor intensive, and costly
- The level of specialization required to explore a large hypothesis space limits researchers
- Data sets are clustered, sparse, lessinterpretable, and incomplete
- Reproducibility crisis

ARTIFICIAL INTELLIGENCE GUIDED, ROBOTICALLY EXECUTED EXPERIMENTS



The Traditional vs. Emerging Paradigms of Scientific Discovery

- Accelerate the discovery process
- Elevate human creativity to higher level goals
- Democratize biological/materials synthesis and characterization
- Unbiased data collection and evaluation

ARTIFICIAL INTELLIGENCE GUIDED, ROBOTICALLY EXECUTED EXPERIMENTS



Autonomous Discovery @Argonne

The vision

- A system that starts with a high-level description of a hypothesis and autonomously carries out computational and experimental workflows to confirm or reject that hypothesis
- Use of AI in robotics and simulations to close the loop on planning, execution, and analysis of experiments
- Builds on
 - Al approaches to planning (multiple steps), and integration of results, causality, etc.
 - Machine learning/simulation to design and predict exp properties and outcomes
 - Automation of experimental protocols (robotic steps and workflows)
 - Active Learning or RL for selection of next experimental targets, etc.

ARTIFICIAL INTELLIGENCE GUIDED, ROBOTICALLY EXECUTED EXPERIMENTS



https://github.com/anl-sdl/ https://www.cs.uchicago.edu/~rorymb/

Outline

How biological data management is rapidly evolving?



Design of antimicrobial peptides

An antimicrobial peptide (AMP) is a short (typically 12 to 50 amino acid) molecule that can target and kill viruses, bacteria, fungi, and other pathogens

Challenge: Design an AMP that can kill specified bacterial strains without harming host cells

With 20 possible amino acids, there are $20^{20} = 10^{26}$ AMPs of length 20

A rational design approach might combine knowledge of bacterial cell membrane composition and structure, AMP molecular and structural properties, host cell membrane characteristics and intracellular pathways—knowledge that may be gained by database/literature search, simulation, experiment



L. T. Nguyen, E.F. Haney, H.J Vogel, The expanding scope of antimicrobial structures and their modes of action, Trends in Biotechnology, 20 (9): 464-472





Embodied Agent for Automated Lab Code Generation

Task: Prepare the master mix for the PCR reaction.

Candidate Code





Outline

How biological data management is rapidly evolving?



Genome-scale Language Models (GenSLMs)



Model	Seq. length	#Parameters	Dataset
GenSLM- Foundation	2048	25M, 250M, 2.5B, 25B	110M
GenSLM	10240	25M, 250M, 2.5B, 25B	I.5M
GenSLM- Diffusion	10240	2.5B	I.5M

- Scaling LLMs with 25B parameters:
 - O (L²) complexity in the attention computation
 - overcome communication overheads, parameters, checkpointing
- Variation within SARS-CoV-2 sequences can be small (< 1% overall variation)
 - Need foundation model to accommodate
 diversity
- One of the largest foundation model trained on raw nucleotide sequences

Infrastructure of GenSLM Foundation Models





license Apache-2.0 website online

release v2.0.0

Q



=

import torch import numpy as np from torch.utils.data import DataLoader from genslm import GenSLM, SequenceDataset

model = GenSLM("genslm_25M_patric", model_cache_dir="/content/gdrive/MyDrive")
model.eval()

Input data is a list of gene sequences
sequences = [
 "ATGAAAGTAACCGTTGTTGGAGCAGGTGCAGTTGGTGCAAGTTGCGCAGAATATATTGCA"

"ATTAAAGATTTCGCATCTGAAGTTGTTTGTTAGACATTGGCGCAGAGTTATGCCGAAGGT",

dataset = SequenceDataset(sequences, model.seq_length, model.tokenizer)
dataloader = DataLoader(dataset)

Compute averaged-embeddings for each input sequence embeddings = [] with torch.no_grad(): for batch in dataloader: outputs = model(batch["input_ids"], batch["attention_mask"], output_hidden_states=True) # outputs.hidden_states shape: (layers, batch_size, sequence_length, hidden_size) emb = outputs.hidden_states[0].detach().cpu().numpy() # Compute average over sequence length emb = np.mean(emb, axis=1) embeddings.append(emb)

Concatenate embeddings into an array of shape (num_sequences, hidden_size)
embeddings = np.concatenate(embeddings)
embeddings.shape
>>> (2, 512)



GenSLM Foundation models reveal new biological insights on gene-level organization

GenSLMs also reveal function level organization of genes



Energy

- Miscellaneous
- **RNA** Processing

Embeddings produced by GenSLM 1.3B, MSL 10,240



Fast time-to-solution: Train from scratch with MSL 10,240 in less than half a day with 4 CS-2s



	GenSLM 123M		GenSLM 1.3B	
	1 CS-2	4 CS-2	1 CS-2	4CS-2
Training steps	5,000	3,000	4,500	3,000
Training samples	165,000	396,000	49,500	132,000
Time to train (h)	4.1	2.4	15.6	10.4
Validation accuracy	0.9615	0.9625	0.9622	0.9947
Validation perplexity	1.031	1.029	1.031	1.025

15

Designing enzymes by incorporating experimental feedback (aka ChatGPT for protein design)



- Need general framework that enables generative design of proteins by incorporating experimental feedback
- Genome-scale language models (GenSLMs)¹ provide a means to incorporate generative modeling for gene sequences:
 - complementary to protein language models
- Rewards for the model:
 - intrinsic sequence specific (e.g., GC content for environmental adaptation)
 - extrinsic functional annotation/ enzyme activity measured via experimentation

M. Zvyagin, et al, Genome-scale language models map the evolutionary ¹⁶trajectories of SARS-CoV-2 (SC'22 Gordon Bell Prize)

Multi-objective RL for generative design allows greater sequence diversity across MDH sequences



Range 17: 9 to 244 GenPept Graphics

Vext Match A Previous Match & First Match

Score		Expect Meth	od		Identities	Positives	Gaps	
89.4 bit	ts(220)	5e-14 Com	positional	matrix adjust	. 72/247(29%)	115/247(46%)	12/247(4%)
Query	9218	VAVTGAAG	QIGYSLLF	RIASGSMFGPD	QPVVLHLIEIEP	ALPALQGVVMELE	DCAFPLLK	9277
Sbjct	9	LVIVGAGG	MIGSNM	-VQSALMLGLT	PNICLYDI	FEPGVHGVFDEIQ	QCAFPGVN	61
Query	9278	GIVPTASL				GKIFVGQGKAIAA	NAAKDVRI	9337
Sbjct	62	-VTYTVNP	EEAFTGAK	YIISSGGAPRK		CKIAAEFGDNIKK	YCPEVEHV	120
Query	9338	LVVGNPCN	TNCLIAMN	NAADVPRDRWF	AMTRLDENRAKA		TNMTIWGN	9397
Sbjct	121	VVIFNPAD	VTALTALI	HSGLKP-NQLT	SLAALDSTRLQQ	ALALEFGVQQDKV	TGAHTYGG	179
Query	9398	HSATQYPD	FYNAHING	RPANEV-IHDE	AWLKGDFITTVQ	RGAAIIKARGLS	SAASAANA	9456
Sbjct	180	HGEQMAVF	ASQVKVDG	KPLAEMGLSDE	W + I C RWEEIKHHTV(QGGSNIIKLRGRS	S S A SFQSPAYN	237
Query	9457	IIDTVKS	9463					
Sbjct	238	AVKMIEA	244					

- We can generate new sequences with varying degrees of sequence identity + positive matches
- We can also generate minimal sequences that have functional 17 domains and can function as a productive enzyme

Why existing supercomputers may not be prepared for self-driving laboratory workflows?



- Exploring even top 1% (1,000 variants x 20 simulation windows = 2,000 simulations) from the embedding space using simulations can overcome the limits on nodes (for a single iteration of RL-based finetuning)
- Labeling productive designs and ranking \rightarrow large compute requirements across multiple computing sites/ facilities

Outline

How biological data management is rapidly evolving?



We can simulate even beamlines @APS



Summary and some take aways...

- HPC workflows are evolving to include experimental data in new and novel ways:
 - direct feedback to existing data (augmentation)
 - collection of new data to drive additional experiments
- Infrastructure support for such workflows need new thinking:
 - Al-based representation for access and relationship discovery across heterogenous datasets
 - rethinking of data fabric: experimental data, simulations, and theory within the same infrastructure
- Foundation models:
 - for empirical data, foundation models may be ideal (e.g., biology)
 - for grounded data (physics, chemistry), multimodal foundation models with theoretical underpinnings are necessary
- Automation and robotics:
 - key progress drivers; need "operating systems" that support effective integration across standards, vendors, and operators



Acknowledgements

Funding

- DOE- National Virtual Biotechnology Laboratory (NVBL)
- Exascale Computing Project Cancer Deep Learning Environment (CANDLE)
- Exascale Workflows Project (ExaWorks)
- DOE Codesign for multimodal AI
- NSF MRI: Multi-modal imaging
- DOE-MEDAL (RENEW) project for workforce training

Computing Time

- Argonne Leadership Computing (Theta/ Theta-GPU/ Al-testbed)
- Cerebras/Nvidia
- o NERSC

Data/ Code/ Models

- o https://github.com/ramanathanlab/genslm
- Access to model weights will also be available via API

Colleagues

- **Richard Scheuermann**
- James Olds

0

- Wesley Scott
- Anda Trifan
- Ashka Shah
- Ozan Gokdemir
- Mike Tynes

Questions/Comments

ramanathana@anl.gov

ZIdsu