# Generative Foundation Models for Sciences

Tung Nguyen, CS261 Winter 2024



### Science Advances by Al

AlphaFold: a solution to a 50-year-old grand challenge in biology



GraphCast: Al model for faster and more accurate global weather forecasting



Accelerating fusion science through learned plasma control



Predicting gene expression with AI



Al has the capability to discover and simulate complex patterns!

### **General Problem**

Learn a surrogate model  $f_{ heta}$  to approximate y = f(x)



Chemical molecules



Reactivity, Stability



Formula, structure

### **General Problem**

Learn a surrogate model  $f_{ heta}$  to approximate y = f(x)







### **General Problem**

Learn a surrogate model  $f_{ heta}$  to approximate y = f(x)









### Surrogate Modeling Use Cases

#### Prediction

How to use  $f_{ heta}(x)$  to **forecast/simulate** f(x)?



Weather forecasting



PDE modeling



Protein structure prediction

### Surrogate Modeling Use Cases

#### Prediction

How to use  $f_{\theta}(x)$  to **forecast/simulate** f(x)?



Weather forecasting



PDE modeling



Protein structure prediction

#### Black-box Optimization

How to use  $f_{ heta}(x)$  to **optimize** f(x)?



Material design



Molecular optimization



Hyperparameter optimization

#### Prediction

How to use  $f_{\theta}(x)$  to **forecast/simulate** f(x)?



Weather forecasting



PDE modeling



Protein structure prediction



How to use  $f_{ heta}(x)$  to **optimize** f(x)?



Material design



Molecular optimization



Hyperparameter optimization

### **Black-box Optimization**

General Problem: **Optimize** a function **without** its functional form or **gradient** information



Limited labeled data but plenty unlabeled data!

### Task-specific Surrogate Models



#### **Foundation Models**



### **GPT-3** Pretraining



#### **GPT-3** In-context Learning



Emergent in-context learning!





#### **Foundation Models for Sciences**

#### Part 1: Learning to Learn In-context

How to train a model that generalizes to unseen functions via in-context conditioning?

#### Part 2: Cross-domain In-context Learning

How to leverage a pretrained LLM for in-context learning in an unseen domain?

**Question**: How to train a model to **generalize** to unseen **functions/tasks?** 

□ In-context learning is an *emergent capability* of GPT-3 at a certain scale



Can we explicitly train a model to generalize to unseen functions?

# Generalizing over Functions

**Question**: How to train a model to **generalize** to unseen **functions/tasks?** 

Data Points Generalization

Learning over data points drawn from the same function

 $\mathbb{E}_{x \sim p_{ ext{train}}(x), y = f(x)}[\log p(y \mid x)]$ 

The model generalizes to unseen data points

 $x \sim p_{ ext{test}}(x), y = f(x)$ 

# Generalizing over Functions

**Question**: How to train a model to **generalize** to unseen **functions/tasks?** 

Data Points Generalization

Learning over data points drawn from the same function

 $\mathbb{E}_{x \sim p_{ ext{train}}(x), y = f(x)}[\log p(y \mid x)]$ 

The model generalizes to unseen data points

$$x \sim p_{ ext{test}}(x), y = f(x)$$

**Functional Generalization** 

- $\Box \quad \text{Learning over functions, given finite context of labeled points} \\ \boxed{\mathbb{E}_{f \sim \mathcal{F}_{\text{train}}, (x, y), C \sim f}[\log p(y \mid x, C)]}$
- The model generalizes to unseen functions given the context

$$f \sim \mathcal{F}_{ ext{test}}, (x,y), C \sim f$$

#### Learning to Learn In-context

- lacksquare Sample  $f\sim\mathcal{F}_{ ext{train}}$  and  $x_{1:N},y_{1:N}\sim f$
- Observe context points  $\{x_i, y_i\}_{i=1}^m$  and make predictions for a set of target points  $\{x_j\}_{j=m+1}^N$ .



#### Model Architecture

U We instantiate the framework with a **transformers** architecture called **TNP** 



$$= \mathbb{E}_{f \sim \mathcal{F}_{ ext{train}}, x_{1:N}, y_{1:N} \sim f, m} \left[ \sum_{i=m+1}^N \log p(y_i \mid x_i, x_{1:m}, y_{1:m}) 
ight]$$

Nguyen, Tung, and Aditya Grover. "Transformer neural processes: Uncertainty-aware meta learning via sequence modeling." ICML 2022.

# **Properties of TNP**

#### Property 1. Context invariance.

The model's predictions do not depend on the permutation of the context points.



#### Property 2. Target equivariance.

Whenever we permute the target inputs, the predictions are permuted accordingly.



### Variants of TNPs

#### We introduce TNP-A and TNP-ND to improve expressivity

#### Autoregressive TNP (TNP-A)

Predict the target jointly with an autoregressive factorization:

$$egin{aligned} p_{ heta} \left( y_{m+1:N} \mid x_{1:N}, y_{1:m} 
ight) \ &= \prod_{i=m+1}^{N} p_{ heta} \left( y_i \mid x_{1:i}, y_{1:i-1} 
ight) \end{aligned}$$

#### Non-Diagonal TNP (TNP-ND)

Predict the target jointly using a multivariate Gaussian distribution with a non-diagonal covariance matrix:

$$p_{ heta}\left(y_{m+1:N} \mid x_{1:N}, y_{1:m}
ight) \ = \mathcal{N}\left(y_{m+1:N} \mid \mu_{ heta}\left(x_{1:N}, y_{1:m}
ight), \Sigma_{ heta}\left(x_{1:N}, y_{1:m}
ight)
ight)$$

Nguyen, Tung, and Aditya Grover. "Transformer neural processes: Uncertainty-aware meta learning via sequence modeling." ICML 2022.

#### **TNPs for 1D Regression**

□ Train on functions generated from **GPs** with an **RBF kernel** 

$$f\sim \mathcal{GP}(0,\mathcal{K}), \;\; \mathcal{K}(x,x')=\sigma^2 \expigg(-rac{(x-x')^2}{2\ell^2}igg)$$

**D** Test on functions generated from **GPs** with **different kernels** 



Nguyen, Tung, and Aditya Grover. "Transformer neural processes: Uncertainty-aware meta learning via sequence modeling." ICML 2022.

#### **TNPs for 1D Regression**

□ Train on functions generated from **GPs** with an **RBF kernel** 

$$f\sim \mathcal{GP}(0,\mathcal{K}), \;\; \mathcal{K}(x,x')=\sigma^2 \expigg(-rac{(x-x')^2}{2\ell^2}igg)$$

**D** Test on functions generated from **GPs** with **different kernels** 



#### TNP generalizes well to unseen functions!

# **TNPs for Image Completion**

- $\Box \quad \text{Train to map from pixel coordinates} \rightarrow \text{pixel values}$
- Each image is a 2-dimensional function



Image completion on unseen classes from 100 pixels



Samples from 20 context points.

### Discussion

- TNP learns to learn in-context **from scratch** 
  - Lt's often desired in deep learning to start from a pretrained model
- There exists powerful in-context models in other domains, e.g., large language models.

Can we transfer in-context learning from LLMs to new domains?

#### Part 1: Learning to Learn In-context

How to train a model that generalizes to unseen functions via in-context conditioning?

#### Part 2: Cross-domain In-context Learning

How to leverage a pretrained LLM for in-context learning in an unseen domain?

### Molecular Optimization

Problem: find the molecule that optimizes a certain property.

 $x^* = rgmin_{x \in \mathcal{X}} f(x)$ 



□ Important domain and distinguished from language.

#### **TNPs for Molecular Optimization**

First attempt: TNPs for in-context molecular property prediction





Solution: Use a more capable model, i.e., a pretrained LLM

#### LLMs for Molecular Optimization

□ Motivation: Repurpose a pretrained LLM for in-context learning in a new domain.



# LLMs for Optimization

#### Existing works directly prompt a pretrained LLM in the text space.

- Not applicable to non-textual domains ×
- Do not generalize to underrepresented domains X



*Liu, Tennison, et al. "Large Language Models to Enhance Bayesian Optimization." arXiv preprint arXiv:2402.03921 (2024).* 

```
Your task is to generate the instruction <INS>. Below are some previous instructions with their scores.

The score ranges from 0 to 100.

text:

Let's figure it out!

score:

61

text:

Let's solve the problem.

score:

63

Generate an instruction that is different from all the instructions <INS> above, and has a higher score

than all the instructions <INS> above. The instruction should begin with <INS> and end with </INS>.
```

The instruction should be concise, effective, and generally applicable to all problems above.

Yang, Chengrun, et al. "Large language models as optimizers." arXiv preprint arXiv:2309.03409 (2023).

### **BOLM Architecture**

Learn separate embedding and prediction layers for the new domain

- Domain-agnostic
- lacksquare Computationally efficient  $\checkmark$



# **BOLM Training**

Train BOLM to perform in-context learning prediction



### **Training Data Generation**

- $\label{eq:constraint} \Box \quad \mbox{Ideal family of functions } \tilde{\mathcal{F}}$ 
  - **Given and Close to downstream functions**
  - Diverse



### **Training Data Generation**

- $\label{eq:constraint} \Box \quad \mbox{Ideal family of functions } \tilde{\mathcal{F}}$ 
  - lacksquare Close to downstream functions  $\checkmark$
  - Diverse



### **Training Data Generation**

- $\label{eq:constraint} \Box \quad \mbox{Ideal family of functions } \tilde{\mathcal{F}}$ 
  - lacksquare Close to downstream functions  $\checkmark$
  - 🗅 🛛 Diverse 🗸

Semi-synthetic training!



#### **Prediction Results**

**Test BOLM on in-context molecular property prediction** 



### **BOLM for Black-box Optimization**

- We ultimately care about optimization
- We consider the online setting

Algorithm 1 Black-box optimization with BOLM **Require:** objective f, fine-tuned LLM surrogate  $f_{\theta}$ , budget B, candidate pool size C, acquisition function  $\alpha$ , batch size k Initialize  $\mathcal{D}_{obs} = \{\}$ while  $|\mathcal{D}_{obs}| < B$  do **Evolutionary** Generate a set of candidates  $\{x_i\}_{i=1}^C$ for each candidate  $x_i$  do BOI M Predict  $\mu_i, \sigma_i = f_{ heta}(x_i, \mathcal{D}_{ ext{obs}})$ Compute utility score  $u_i = \alpha(\mu_i, \sigma_i)$ UCB end for Select k candidates with the highest utility scores for each selected candidate  $x_i$  do Evaluate  $x_i$  using the actual objective  $y_i = f(x_i)$ Add  $(x_i, y_i)$  to the observation dataset  $\mathcal{D}_{obs}$ end for end while

#### **PMO Benchmark**

U We test (a single) BOLM on 21 benchmark optimization tasks from PMO



### Experiments

U We test (a single) BOLM on 21 benchmark optimization tasks



#### BOLM achieves the best overall scores and ranking!

#### **Ablation Studies**

□ Importance of semi-synthetic training



#### Semi-synthetic is better than intrinsic or synthetic alone

#### **Ablation Studies**

□ Importance of language instruction



#### Language instruction is crucial to BOLM

# Summary

- We can train a model to learn in-context
  - □ In-context learning can **generalize** from **synthetic to real** functions
  - □ In-context learning can be **transferred across domains**
- **L** End goal: have a GPT-like model that scientists can use to optimize **arbitrary objectives**



# **Future Work**

- □ Improve generality
  - **G** Foundation models that can perform cross-domain optimization
- □ Improve performance
  - **G** Scale up data and model size
  - Parameter-efficient finetuning
- **G** Foundation models for end-to-end optimization
  - **Candidate generation, exploration, surrogate modeling, etc.**

Thank You! Q&A