# SNIP: Bridging Mathematical Symbolic and Numeric Realms with Unified Pre-training



*Kazem Meidani     *Parshin Shojaee     Chandan K. Reddy     Amir Barati Farimani
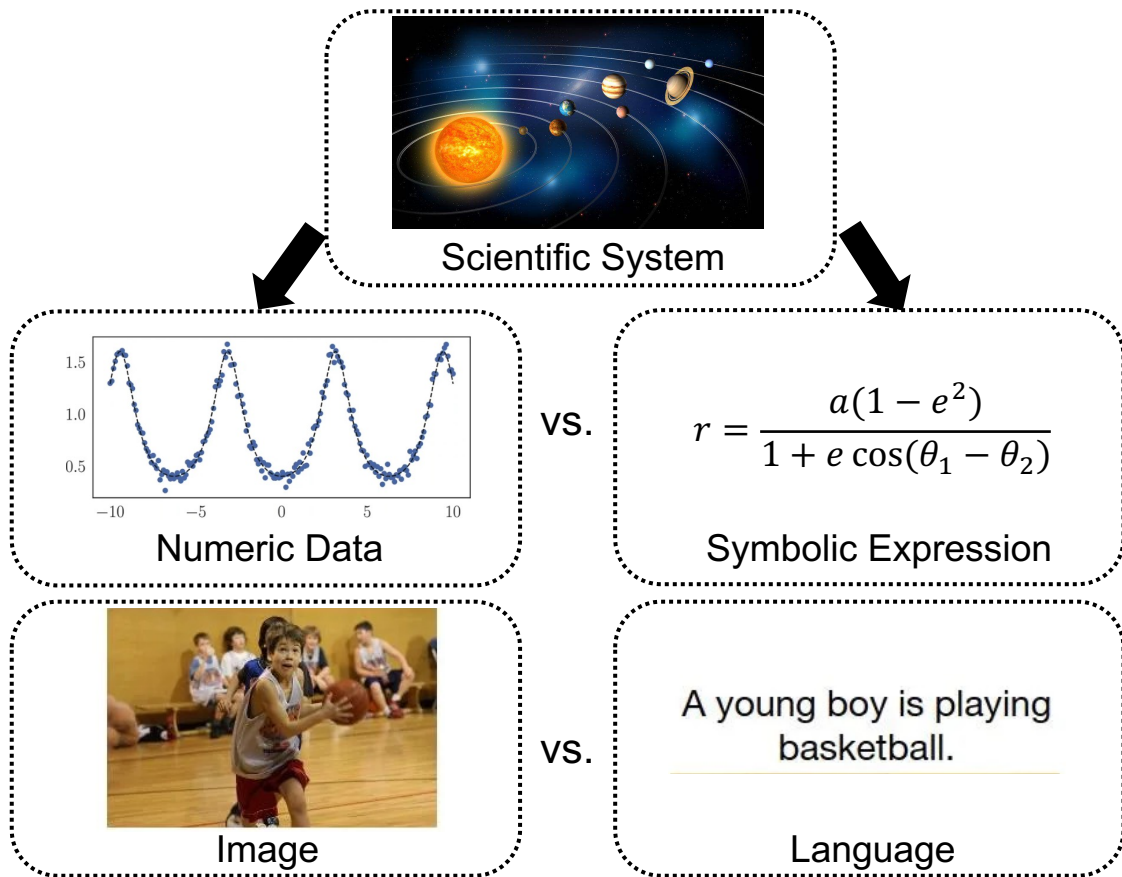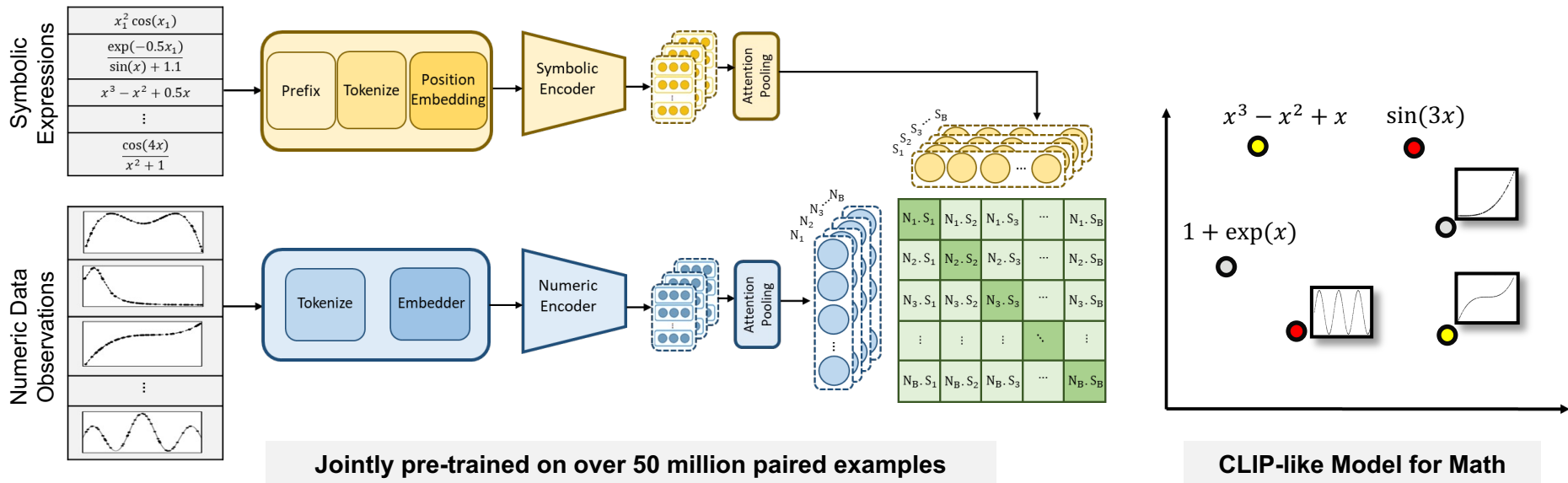
*Equal contribution

ICLR 2024 **Spotlight** Presentation

# Math: Language of Science

- Symbolic Mathematics has been **unreasonably effective** for understanding, predicting, and controlling various scientific systems.

- Obtaining mathematical equations from data is an essential part of **scientific discovery**.

- Each system can be represented by two modalities of **Numeric Data Observations** and **Symbolic Mathematical Expressions**

- **Multi-modal representation learning** has shown success in many domains including vision-language models.


Scientific System


Numeric Data

vs.

$$r = \frac{a(1 - e^2)}{1 + e \cos(\theta_1 - \theta_2)}$$

Symbolic Expression


Image

vs.

A young boy is playing basketball.

Language

# SNIP: Symbolic Numeric Integrated Pre-training



**Jointly pre-trained on over 50 million paired examples**

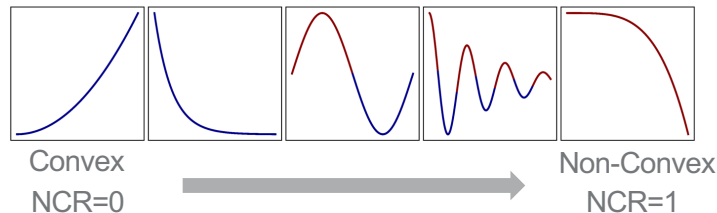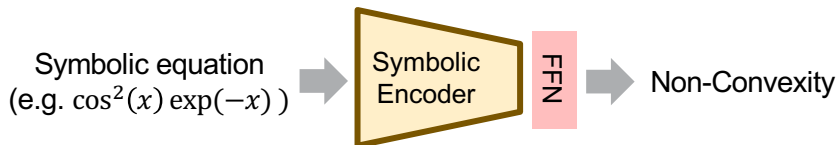**CLIP-like Model for Math**

- Data samples are synthetically generated.

- Transformer Encoders are used for both symbolic and numeric encoders.

- Contrastive loss is used for pre-training joint embeddings.

$$\mathcal{L} = - \sum_{(v,s) \in B} \left( \log \text{NCE}(\boldsymbol{Z}_S, \boldsymbol{Z}_V) + \log \text{NCE}(\boldsymbol{Z}_V, \boldsymbol{Z}_S) \right)$$

$$\text{NCE}(\boldsymbol{Z}_S, \boldsymbol{Z}_V) = \frac{\exp\left(\boldsymbol{Z}_S \cdot \boldsymbol{Z}_V^+\right)}{\sum_{\boldsymbol{z} \in \{\boldsymbol{z}_V^+, \boldsymbol{z}_V^-\}} \exp\left(\frac{\boldsymbol{Z}_S \cdot \boldsymbol{Z}}{\tau}\right)}$$
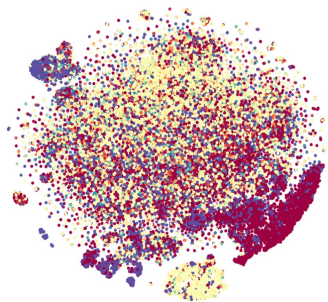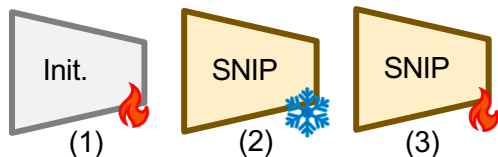
# Task 1: Cross-Modal Property Prediction

- Predicting numeric properties from symbolic input and vice vera.
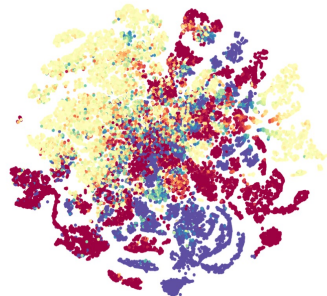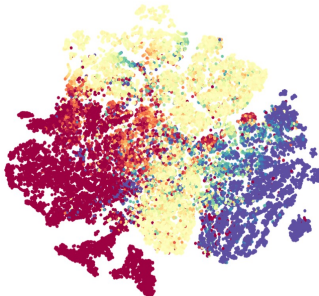- Example: Predicting convexity ratio (numeric property) of a function based on its symbolic expression

Symbolic equation
(e.g. $\cos^2(x)\exp(-x)$ ) → Symbolic Encoder → FFN → Non-Convexity

Convex NCR=0 → Non-Convex NCR=1

Using a predictor head, we compare:

(1) Supervised Model
(2) SNIP (frozen)
(3) SNIP (finetuned)

Init. (1)   SNIP (2)   SNIP (3)

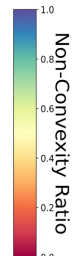| Model | Non-Convexity Ratio | | Upwardness | |
|---|---|---|---|---|
| | $\uparrow R^2$ | $\downarrow$ NMSE | $\uparrow R^2$ | $\downarrow$ NMSE |
| Supervised | 0.4701 | 0.5299 | 0.4644 | 0.5356 |
| SNIP (frozen) | 0.9269 | 0.0731 | 0.9460 | 0.0540 |
| SNIP (finetuned) | **0.9317** | **0.0683** | **0.9600** | **0.0400** |



(a) Without Pretraining

(b) Pretrained, Before Finetuning

(c) Pretrained, After Finetuning

Non-Convexity Ratio

R²  # of Training Samples

Supervised
SNIP (Frozen)
SNIP (Finetune)

# Task 2: SNIP for Symbolic Regression



(a) Training

Pretrained SNIP Encoder: Embedder → Numeric Encoder → Attention Pooling → Latent $Z_y$ → Mapping Layers → Pretrained E2E Decoder (Decoder) → $\hat{f}(x)$

Symbolic Expression Space
× Discrete / Combinatorial
× Very Large
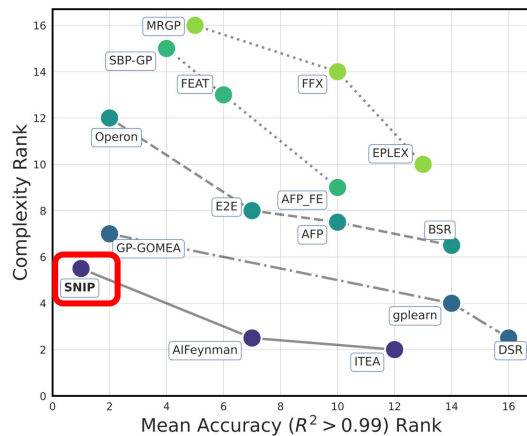× Numeric-Ignorant

Decode → $\hat{f}(\cdot)$ → $(X, \hat{y})$

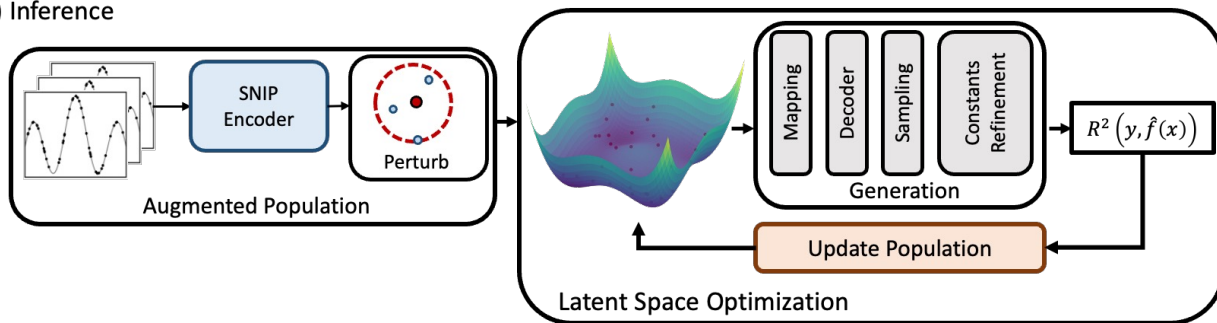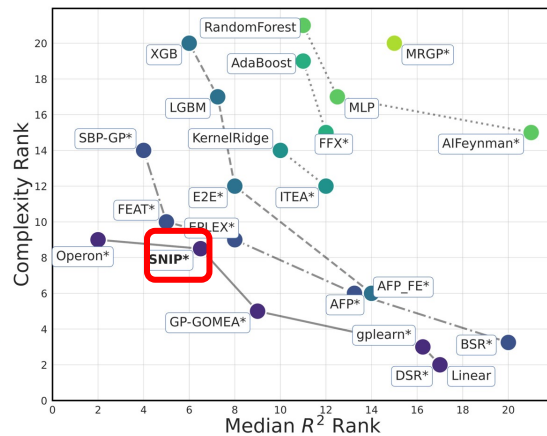**Latent Space**
✓ Continuous
✓ Low-dimensional
✓ Interpolatable
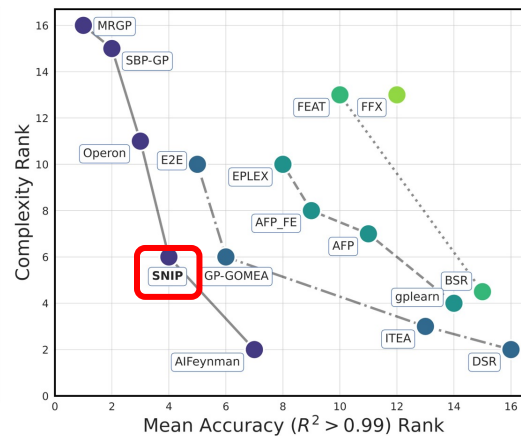✓ Numeric-informed

# Task 2: SNIP for Symbolic Regression



(a) Strogatz datasets

(b) Black-box datasets
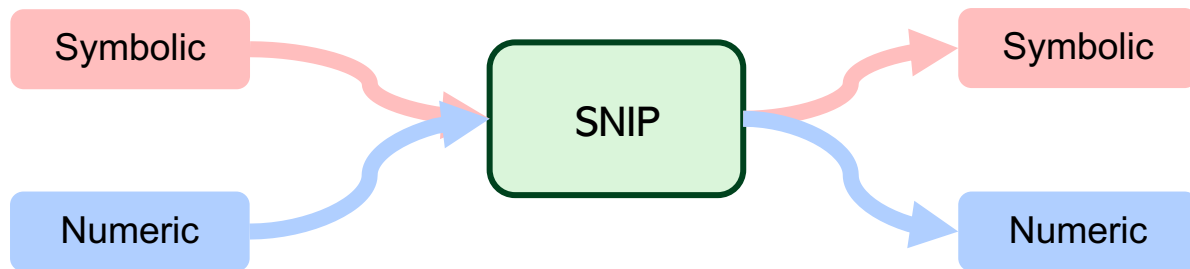
(c) Feynman datasets

# Future Work

Symbolic → SNIP → Symbolic

Numeric → SNIP → Numeric

Multi-Modal
Pre-training

Multi-Modal
Symbolic Regression

Carnegie
Mellon
University

Contact: mmeidani@andrew.cmu.edu

VIRGINIA TECH