# A Unified Framework for Deep Symbolic Regression

Thirty-sixth Conference on Neural Information Processing Systems
Tue Nov 29th through Dec 1st, 2022. New Orleans, Louisiana, US

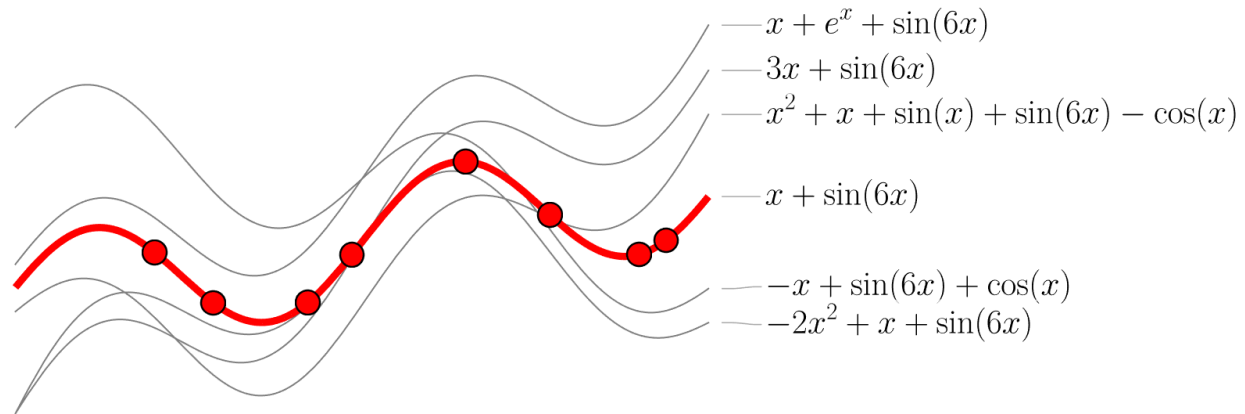

NEURAL INFORMATION
PROCESSING SYSTEMS

**Mikel Landajuela**, Chak Lee, Jiachen Yang, Ruben Glatt, Claudio P. Santiago, Ignacio Aravena, Terrell N. Mundhenk, Garrett Mulcahy, Brenden K. Petersen

Lawrence Livermore
National Laboratory

# Symbolic Regression: A classical Problem

Given a dataset $(X, y)$, where each point $X_i \in \mathbb{R}^n$ and $y_i \in \mathbb{R}$, find an analytic expression $f : \mathbb{R}^n \to \mathbb{R}$ such that $f(X_i) \approx y_i$
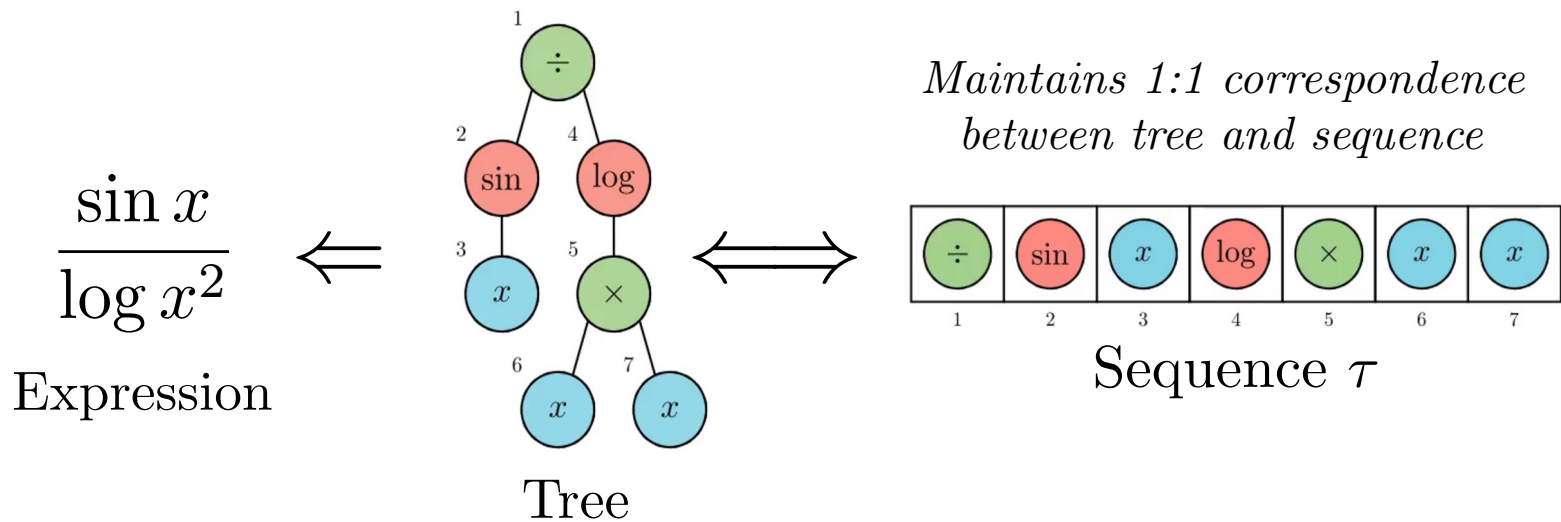


- $x + e^x + \sin(6x)$
- $3x + \sin(6x)$
- $x^2 + x + \sin(x) + \sin(6x) - \cos(x)$
- $x + \sin(6x)$
- $-x + \sin(6x) + \cos(x)$
- $-2x^2 + x + \sin(6x)$

- Symbolic Regression (SR) leads to interpretable models with high performance and generalizability, even in the small dataset regime Brøløs et al., 2021; Wilstrup et al., 2021

- SR has received lot of attention in recent years Cranmer et al., 2020; Udrescu et al., 2020; Petersen et al., 2021; Landajuela et al., 2021; Biggio et al., 2021 ; Kamienny et al., 2022 ; …

# Symbolic Regression as Discrete Optimization

- Using expression trees, the problem becomes a discrete optimization one:

$$\underset{n \leq N, \tau_1, \ldots, \tau_n}{\arg\max} \; [R(\mathtt{ET}(\tau_1, \ldots, \tau_n))] \; \text{ with } \tau_i \in \mathcal{L} = \{+, \ldots, \sin, \ldots, x_1 \ldots\}$$

$$\frac{\sin x}{\log x^2} \Longleftarrow$$

Expression



Tree

*Maintains 1:1 correspondence between tree and sequence*



Sequence $\tau$

- **Exponentially large** search space $|\mathcal{L}|^N$. SR is **NP-hard** (Virgolin et al., 2022), i.e., the search for the best solution can be **intractable**.

# Solution Strategies for Symbolic Regression

- Over the last few years, there are now several quite different approaches to SR:

  - **Problem Simplification**
    Udrescu et al., 2019 and 2020

  - **Neural-guided Search**
    Bello et al., 2017; Petersen et al., 2021

  - Genetic Programming
    Koza, 1994; Mundhenk et al., 2021; …

  - Large Scale Pre-training
    Biggio et al., 2021; Kamienny et al., 2022; …

  - Linear Regression
    Legendre,1805; Brunton et al., 2016; …

Exploits $(X, y)$ data to simplify a SR problem into lower-dimensional sub-problems.



Neural network learns to search over time, with the ability to incorporate in situ constraints.



Rapidly explores the search space via genetic operators.



Leverages big data, learning from many other problems by conditioning on the $(X, y)$ data.



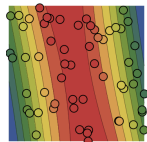Quickly learn sparse coefficients of a linear combination of basis functions.

# uDSR: A Unified framework for Deep Symbolic Regression

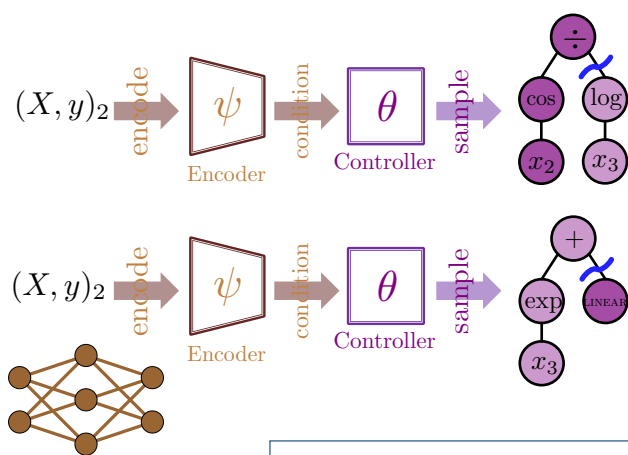# uDSR: A Unified framework for Deep Symbolic Regression



$\mathcal{P}_2 : (X, y)_2 \in \mathbb{R}^{n \times 2} \times \mathbb{R}^n$

① AIF identifies multiplicative separability, simplifying $\mathcal{P}$ into $\mathcal{P}_1$ and $\mathcal{P}_2$

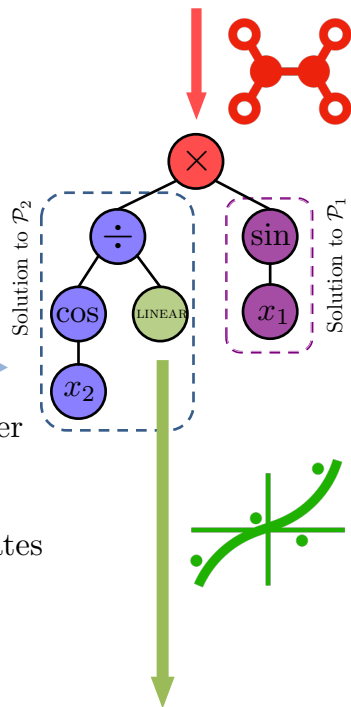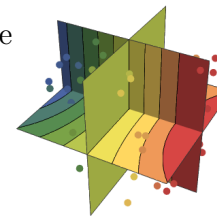③ DSR, conditioned on $\mathcal{P}_2$ data $(X, y)_2$ via LSPT, samples candidate solutions to $\mathcal{P}_2$

$(X, y)_2$ — encode → $\psi$ (Encoder) — condition → $\theta$ (Controller) — sample →

$(X, y)_2$ — encode → $\psi$ (Encoder) — condition → $\theta$ (Controller) — sample →

Solution to $\mathcal{P}_2$

Solution to $\mathcal{P}_1$

$\mathcal{P}_1 : (X, y)_1 \in \mathbb{R}^n \times \mathbb{R}^n$

② Solution to $\mathcal{P}_1$ is found by DSR

④ GP crossover operator recombines the candidates to form $\mathcal{P}_2$

⑤ Solve the solution to $\mathcal{P}_2$ for the LINEAR token, yielding $\bar{f}^{-1}(y)$

Solve for $\beta$ using sparse linear regression

$$\begin{bmatrix} | & | & | & | & | & | \\ 1 & x_2 & x_2^2 & x_2 x_3 & \ldots & x_2^3 \\ | & | & | & | & | & | \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = \begin{bmatrix} \bar{f}^{-1}(y_1) \\ \vdots \\ \bar{f}^{-1}(y_n) \end{bmatrix}$$
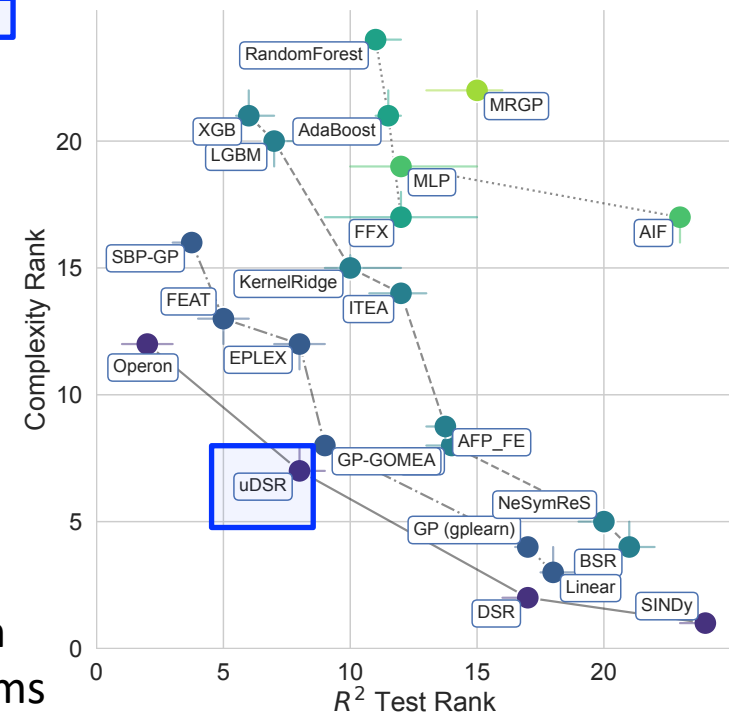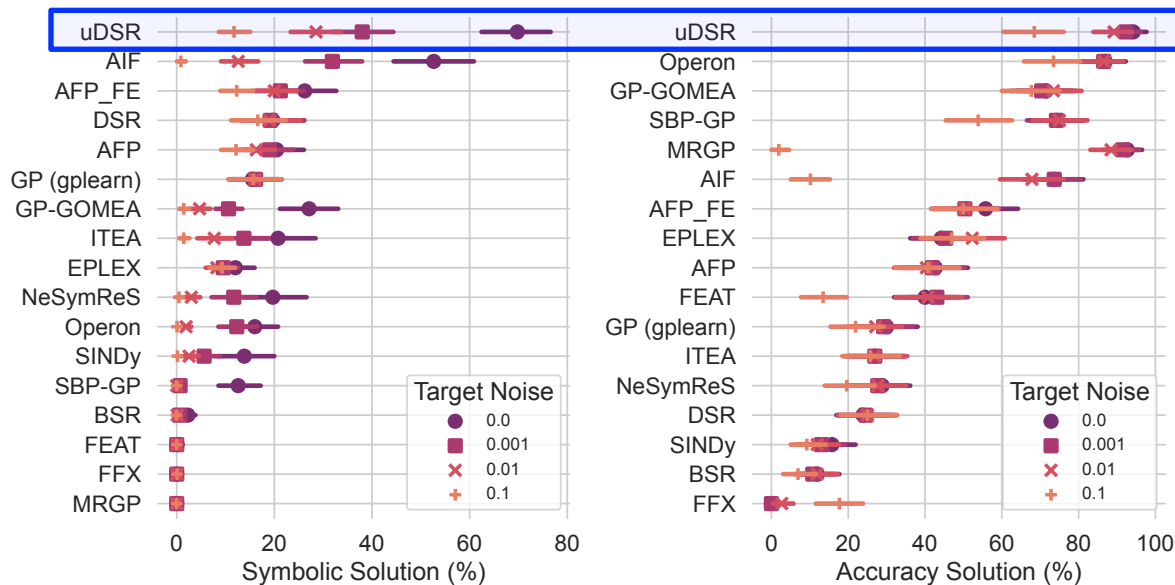
Final expression:

$$y = \sin(x_1) \times \frac{\cos(x_2)}{x_2^2 + 0.5 x_2 x_3 + 0.7}$$

The final LINEAR token is given by:

$$\text{LINEAR} = \boxed{x_2^2 + 0.5 x_2 x_3 + 0.7}$$

# Results on SRBench

- Benchmarking using the open-source pipeline SRBench (La Cava et al.,2021) (252 datasets from PMLB):



- uDSR **outperforms** all other 14 benchmarked methods in **symbolic** and **accuracy recovery** for ground-truth problems

- uDSR falls on the **Pareto frontier** (accuracy-complexity) on black-box SR problems