



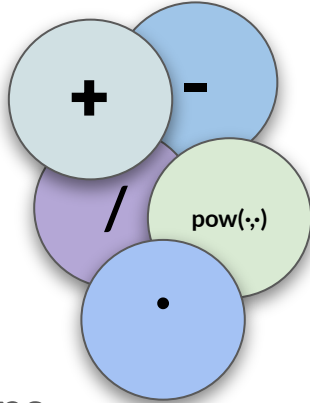
Complexity Aware Deep Symbolic Regression

Zachary Bastiani, Robert M. Kirby, Jacob Hochhalter, Shandian Zhe
University of Utah

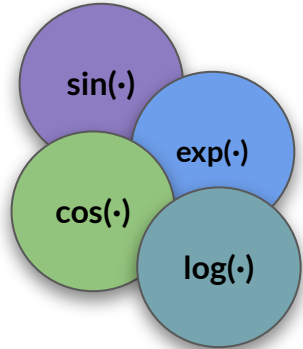


Symbolic Regression: Problem Statement

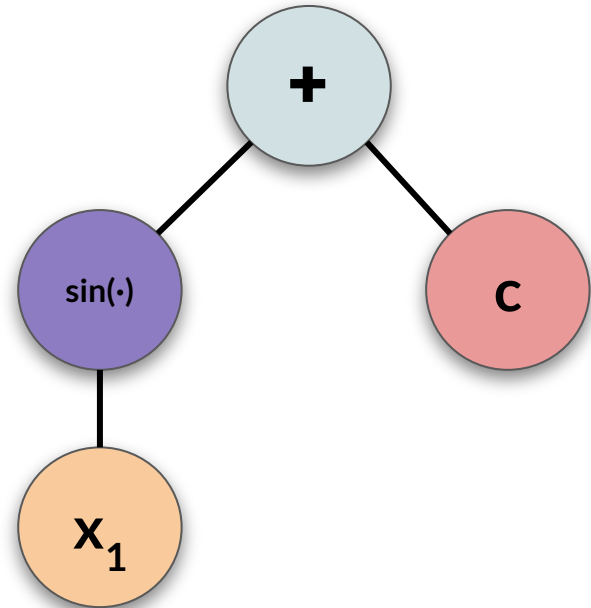
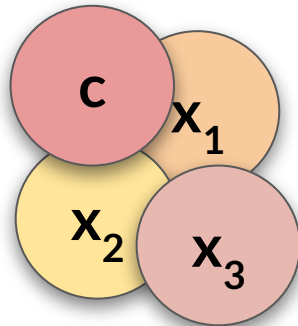
Binary Ops



Unary Ops



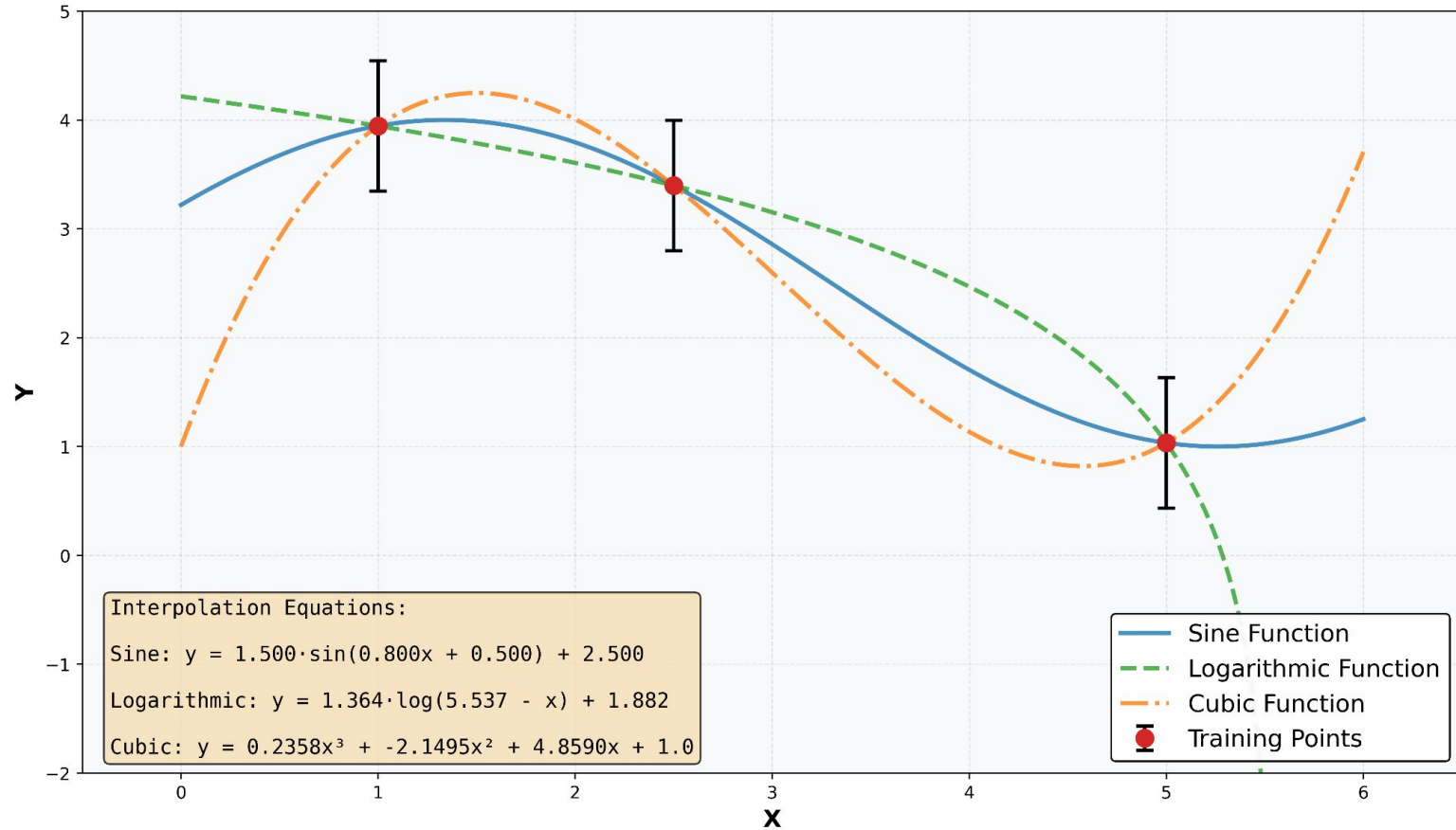
Variables



Deep Symbolic Regression: Reinforcement Learning

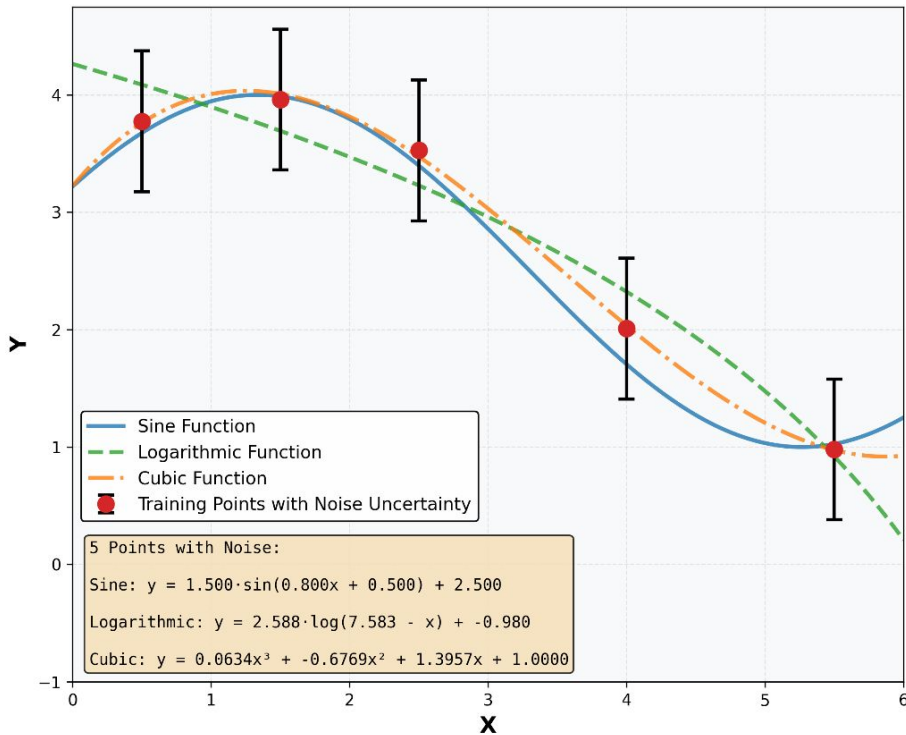


The Interpolation Problem

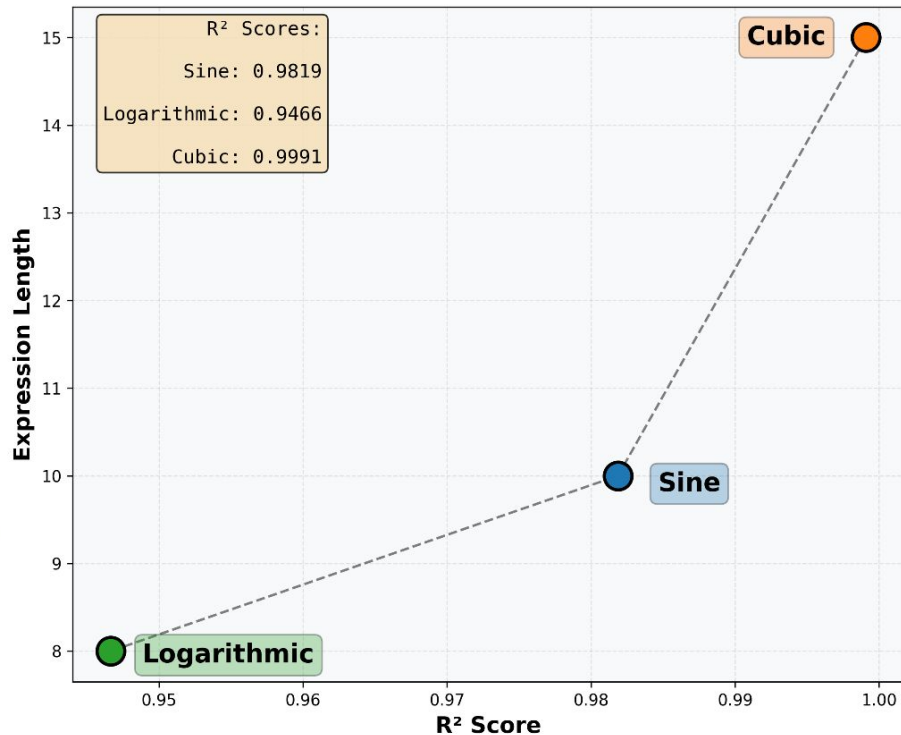


Pareto Frontier

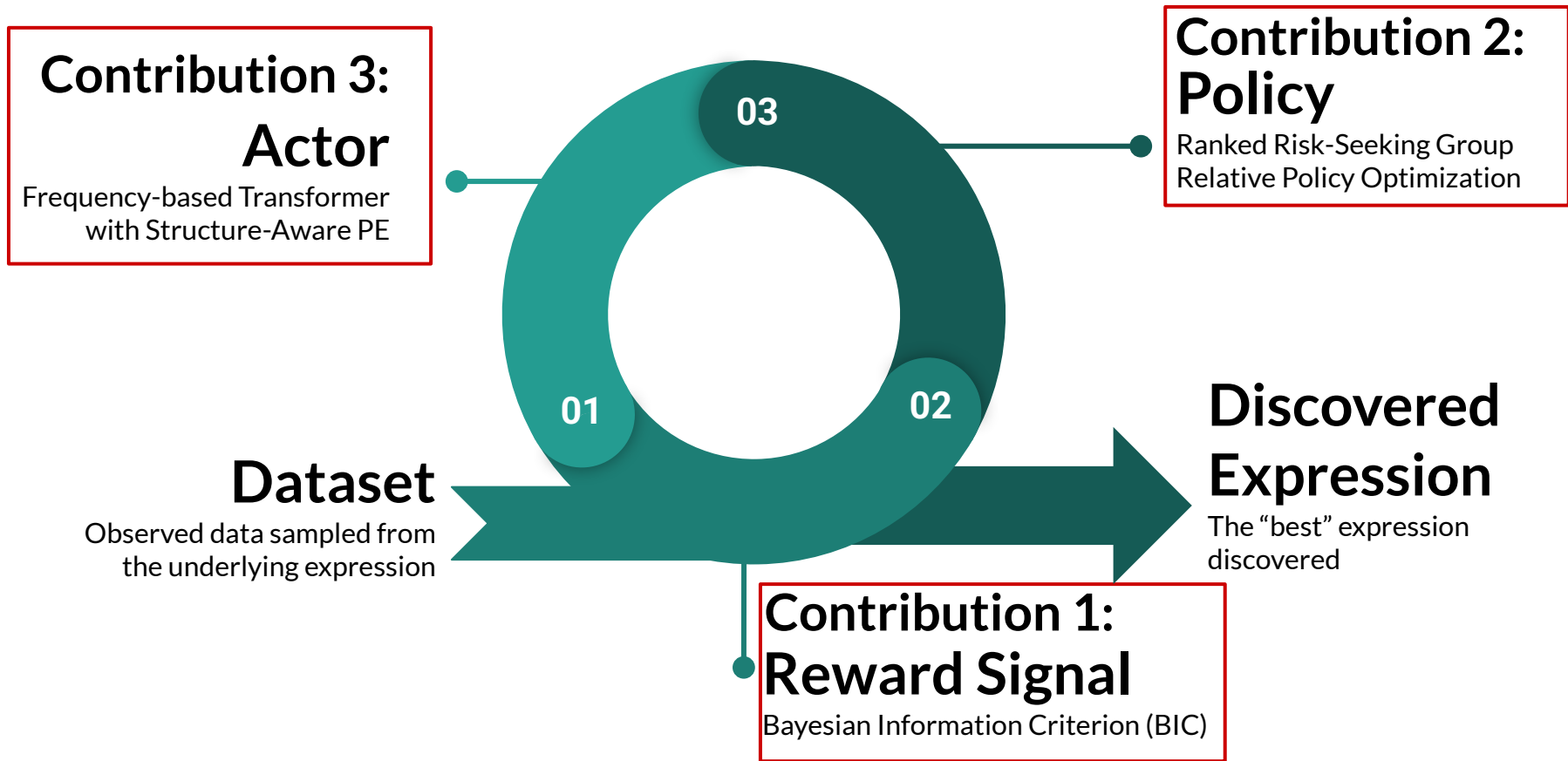
(A) Symbolic Regression on Noisy Data



(B) Pareto Frontier



Complexity Aware Deep Symbolic Regression (CADSR)



Contribution 1: Bayesian Information Criterion (BIC)

$$BIC = \underbrace{k}_{\text{Expression Length}} \log(\underbrace{S}_{\text{Dataset Size}}) + 2 \log\left(\underbrace{\prod_{i=1}^N p((y)_i | \tau, (\mathbf{x})_i)}_{\text{Expression Likelihood}}\right)$$

Expression	BIC Reward	Rank
Cubic	-9.3994	1
Sine	-2.4555	2
Logarithmic	-0.2785	3

Contribution 2: Group Relative Policy Optimization

$$\nabla \hat{J}(\theta, \alpha) = \frac{1}{\alpha B} \nabla_{\theta} \sum_{i=1}^B \sum_{j=1}^{|\tau^{(i)}|} \text{TR} \left[A_{ij} \frac{p_{\theta}(\tau^{(i)})}{p_{\theta_{\text{old}}}(\tau^{(i)})}, \epsilon \right] - \beta D_{KL} [p_{\theta}(\tau^{(i)}) || p_{\theta_{\text{ref}}}(\tau^{(i)})]$$

Trust Region Regularization

Ranked Reward:

$$A_{ij} = \lambda \cdot \text{ReLU} \left(1 - \frac{|\{\tau^{(k)} : R(\tau^{(k)}) > R(\tau^{(i)})\}|}{\alpha B / 100} \right)$$

Contribution 2: Ranked Risk-Seeking Advantage Matrix

Lemma 1: Given any continuous mapping, f , that can map unbounded reward function values to a bounded domain (e.g., $(-\infty, \infty) \rightarrow [0,1]$), suppose the reward function is continuous, there always exists a set of distinct rewards values that numerically create a tail barrier in the risk-seeking policy.

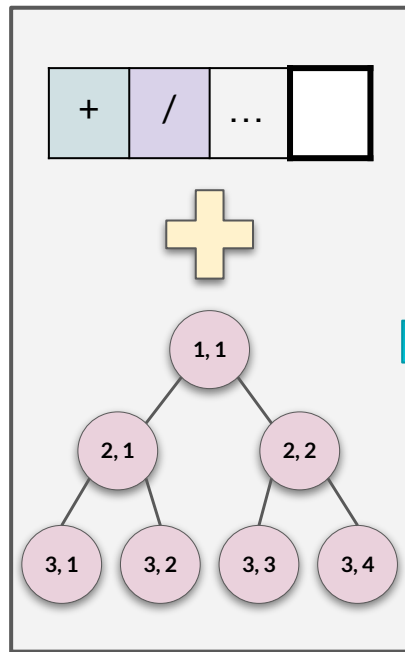
Lemma 2: By using the step function for reward mapping, the policy gradients with our BIC reward is unbiased and will not encounter any tail barrier.

Ranked Reward:

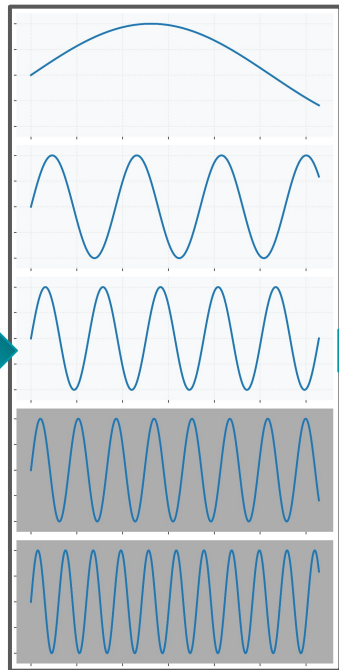
$$A_{ij} = \lambda \cdot \text{ReLU} \left(1 - \frac{|\{\tau^{(k)} : R(\tau^{(k)}) > R(\tau^{(i)})\}|}{\alpha B / 100} \right)$$

Contribution 3: DCT Layer and Structure Aware PE

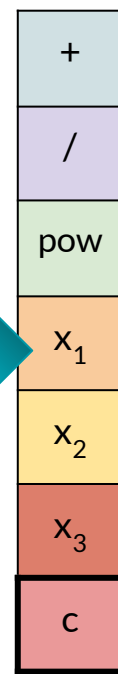
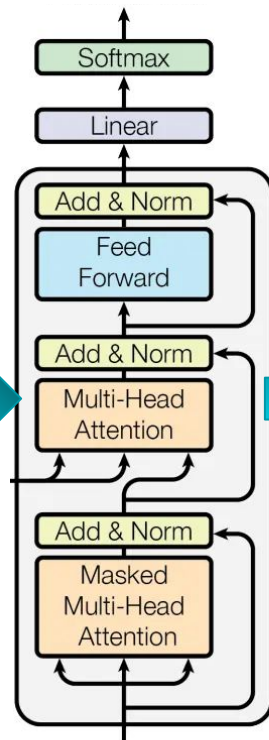
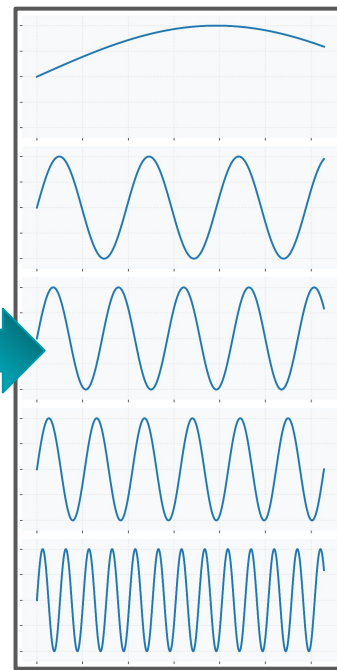
Positional Encoding



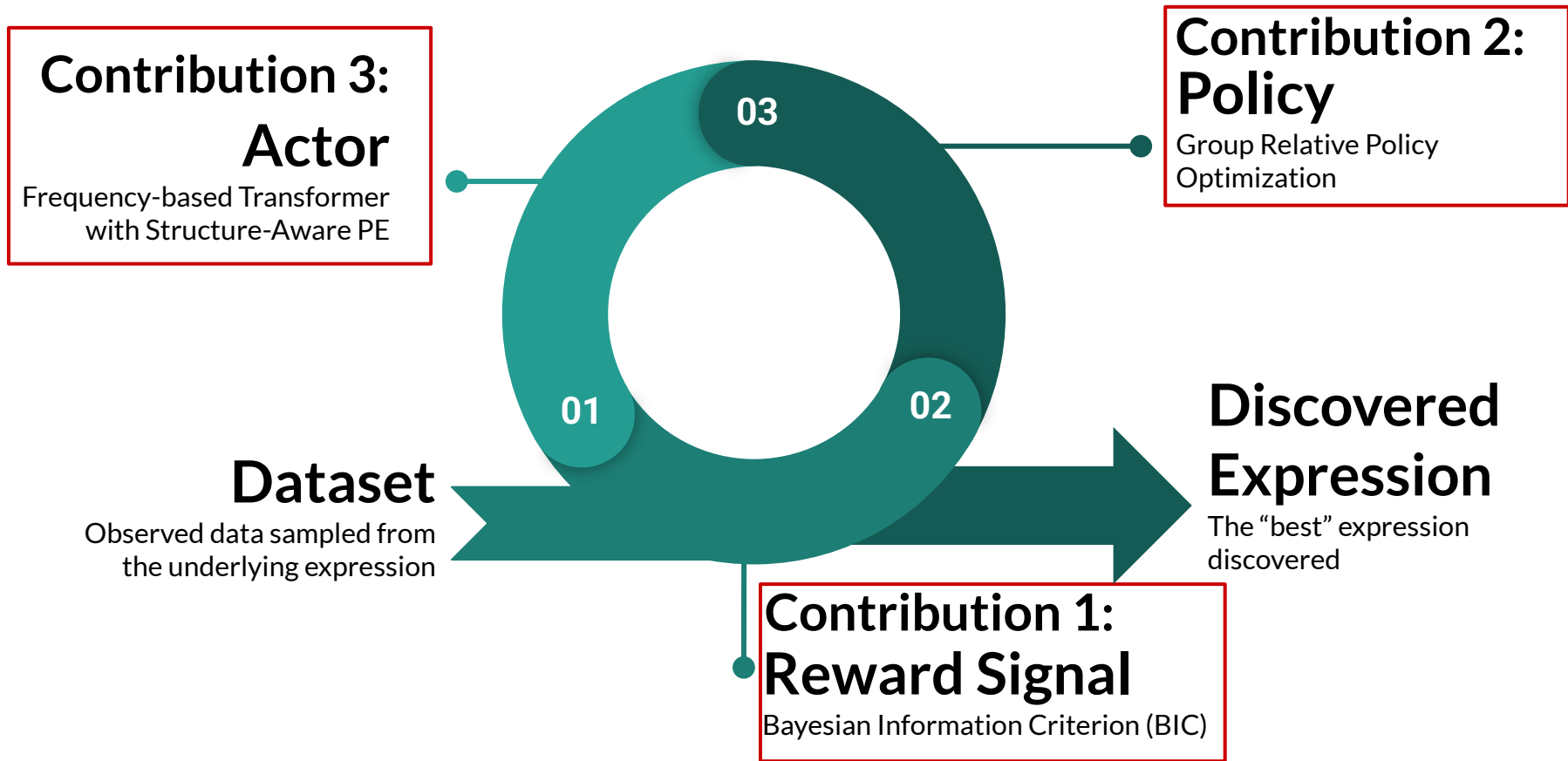
DCT Layer



I-DCT Layer



Complexity Aware Deep Symbolic Regression (CADSR)



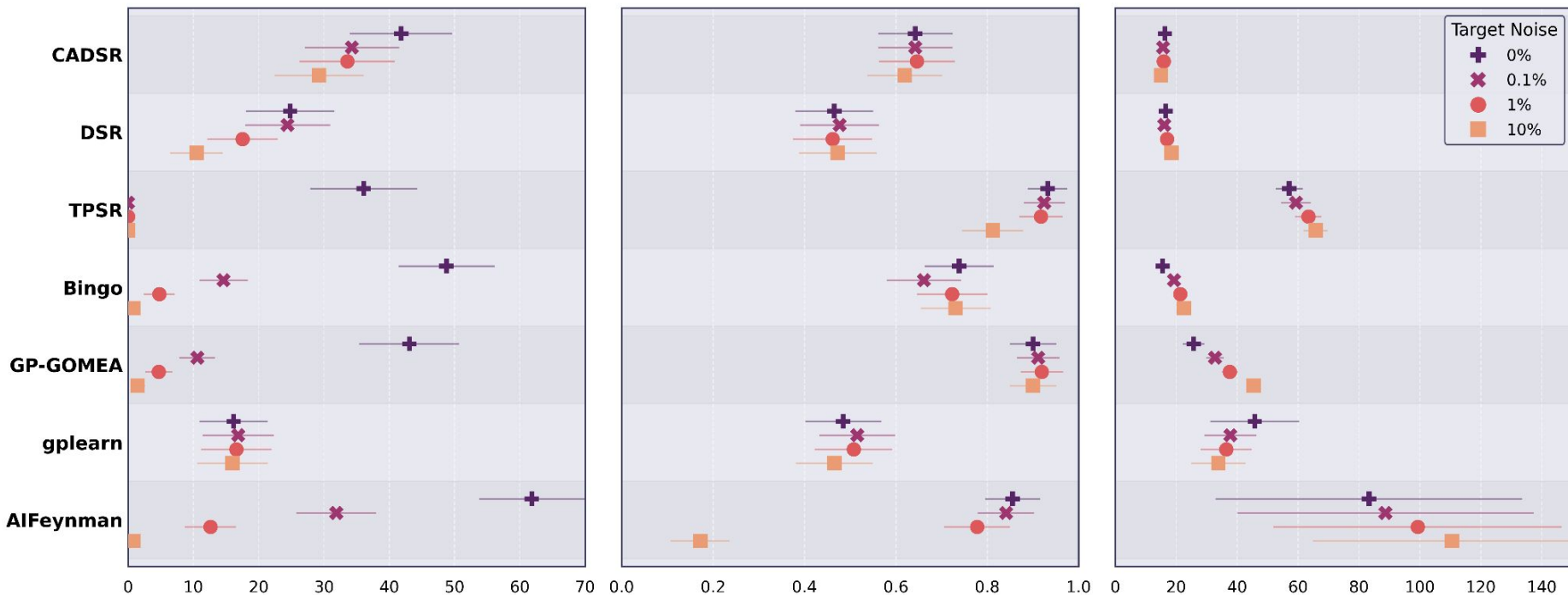
Results

Algorithm Comparison Across Noise Levels

Symbolic Solution Rate (%)

Accuracy Rate

Simplified Complexity



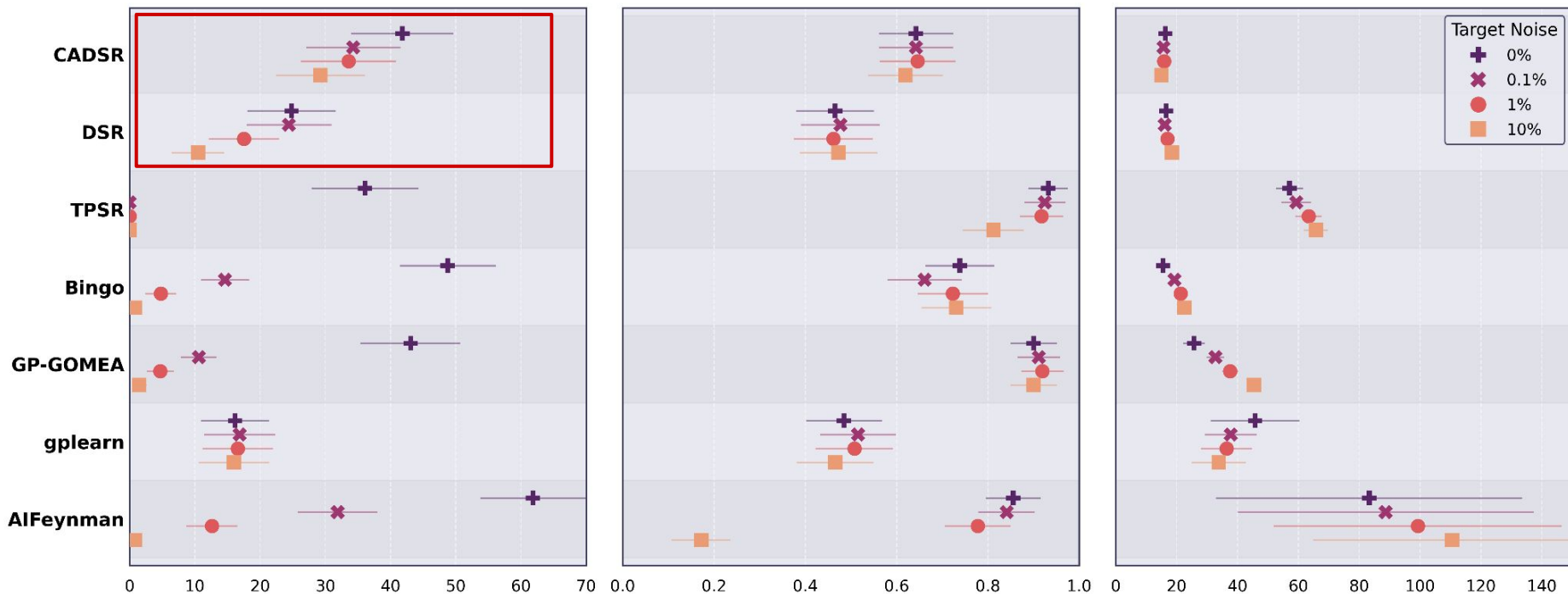
Results

Algorithm Comparison Across Noise Levels

Symbolic Solution Rate (%)

Accuracy Rate

Simplified Complexity



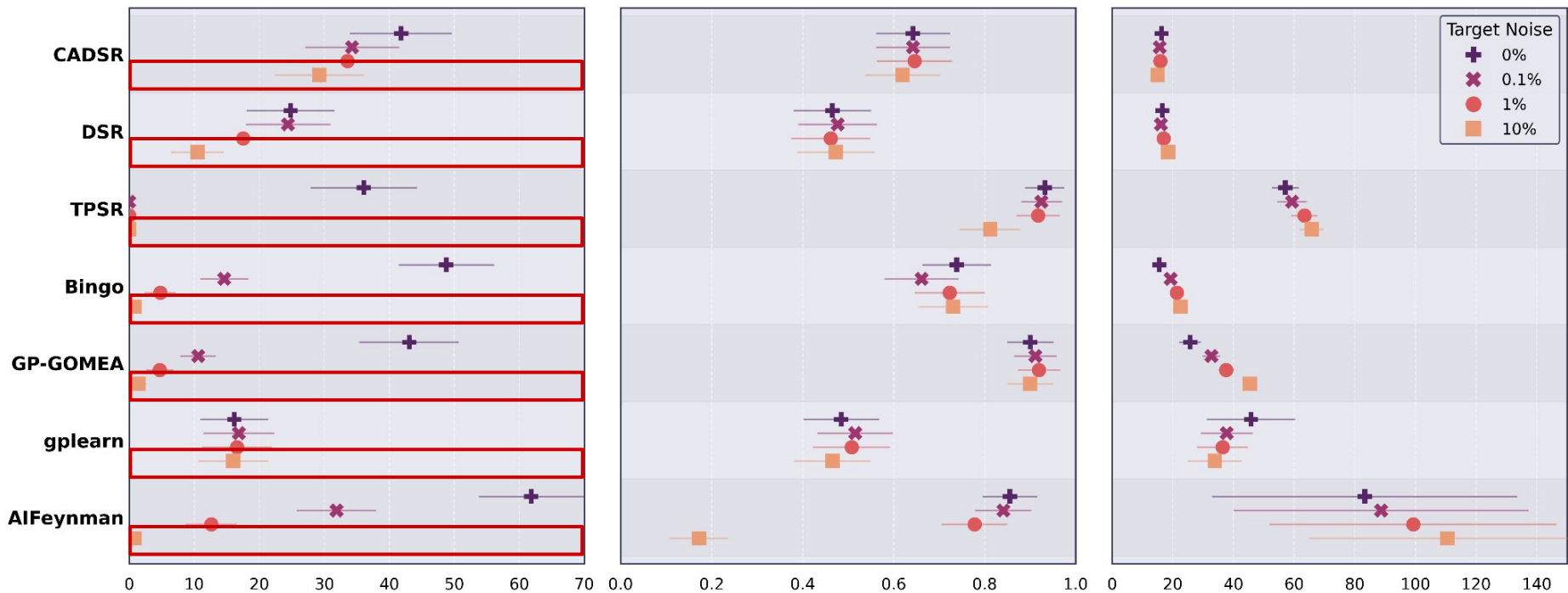
Results

Algorithm Comparison Across Noise Levels

Symbolic Solution Rate (%)

Accuracy Rate

Simplified Complexity



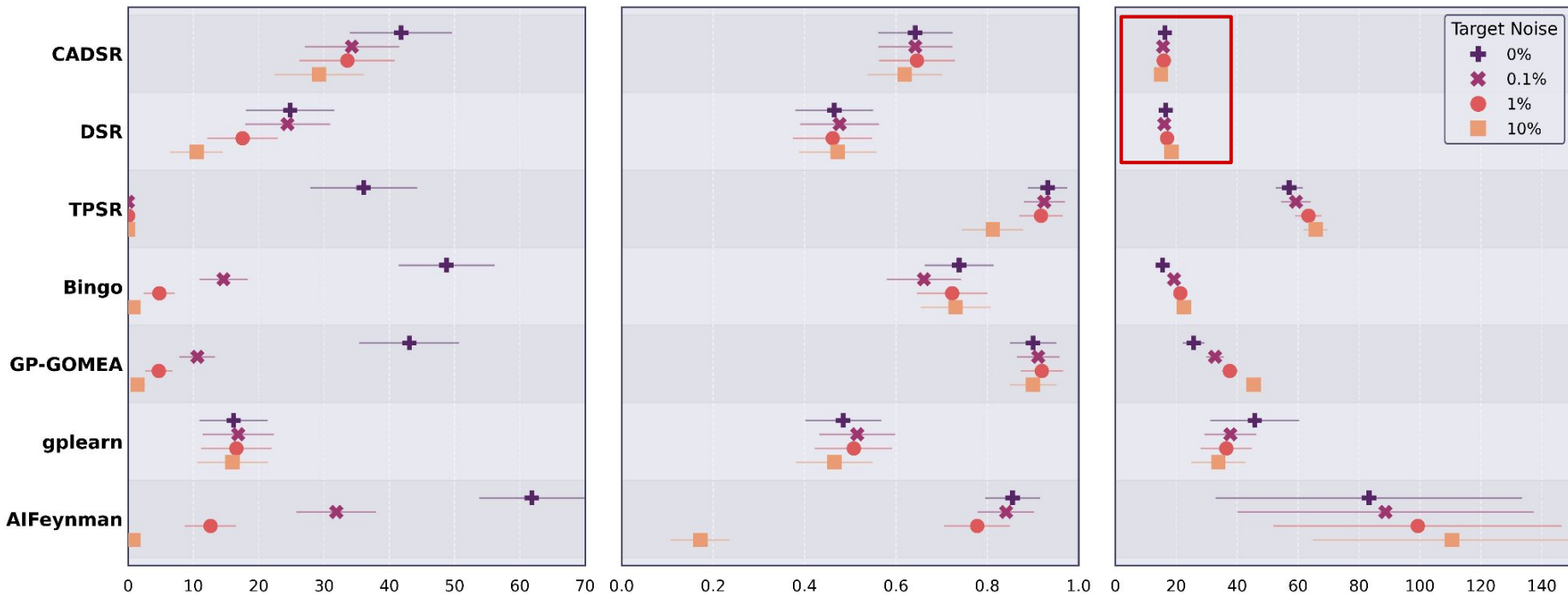
Results

Algorithm Comparison Across Noise Levels

Symbolic Solution Rate (%)

Accuracy Rate

Simplified Complexity



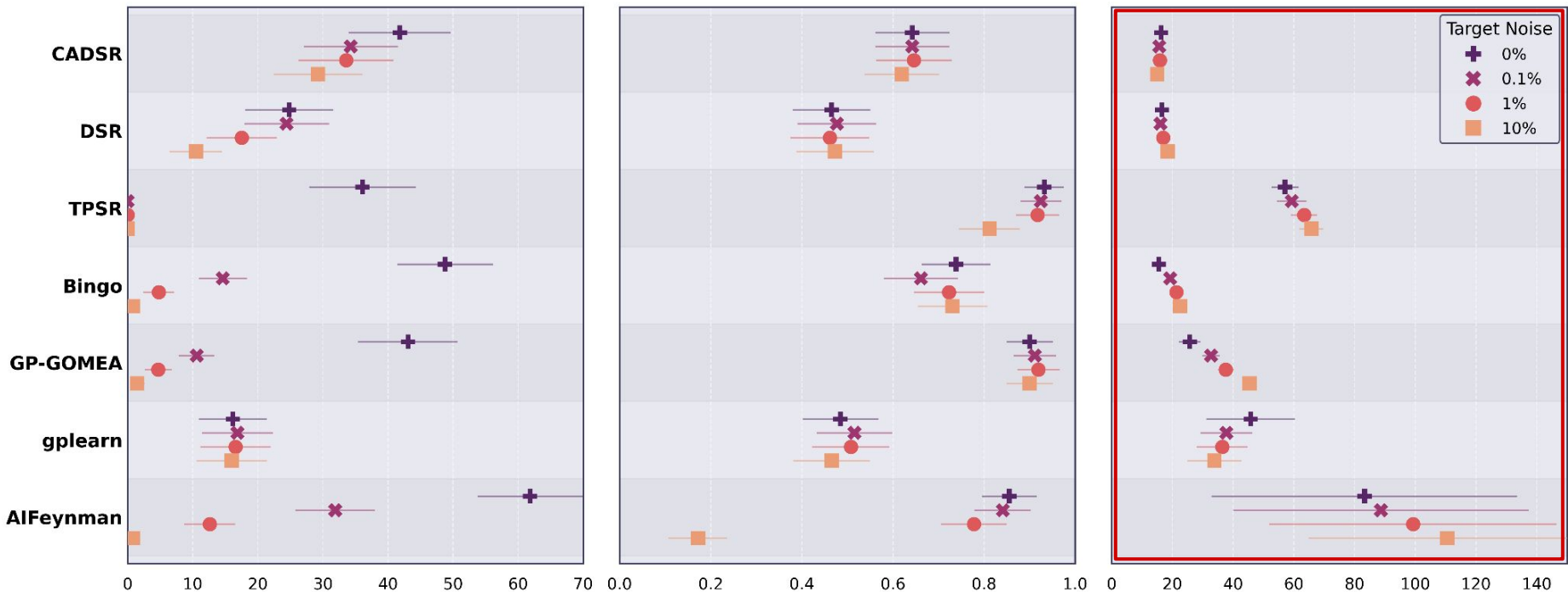
Results

Algorithm Comparison Across Noise Levels

Symbolic Solution Rate (%)

Accuracy Rate

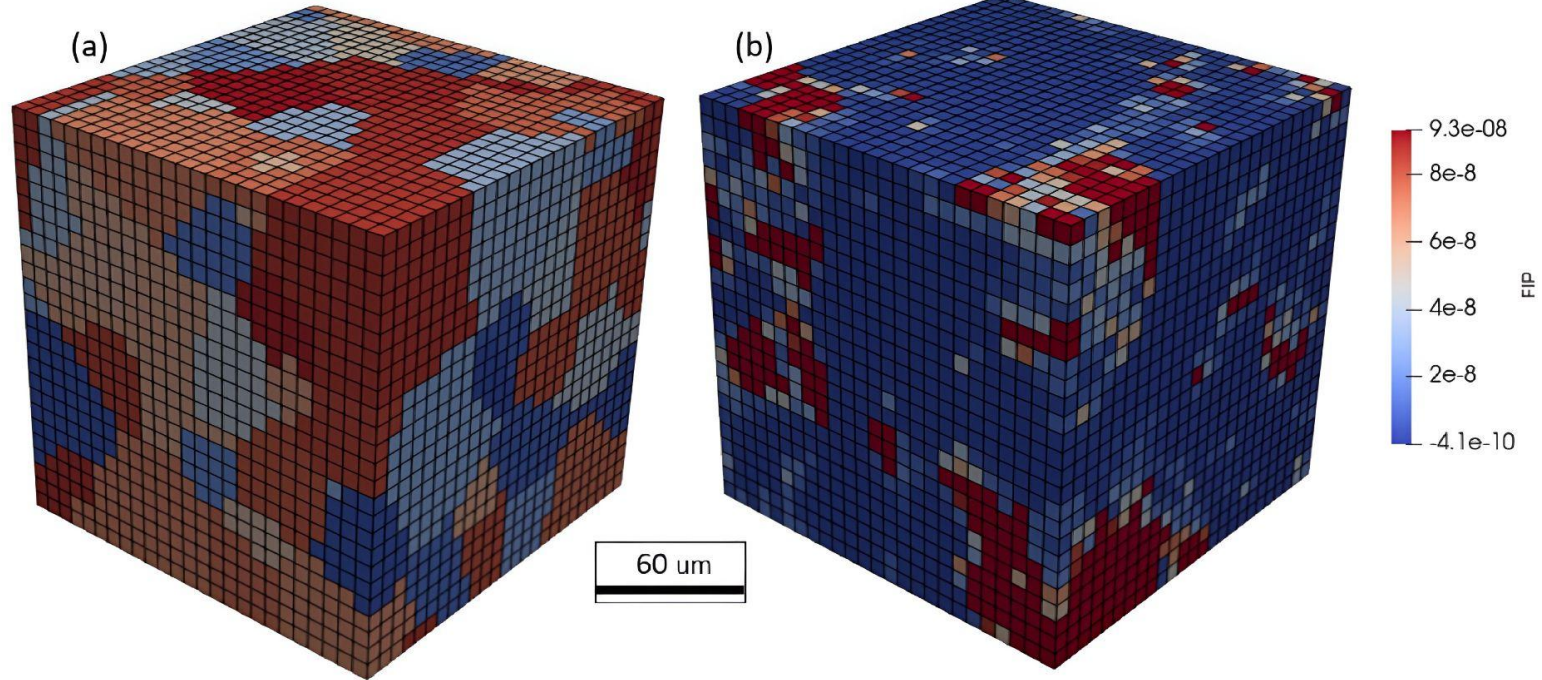
Simplified Complexity



Fracture Mechanics

Microstructure

FIP Values



Fracture Mechanics

Method	Train R ²	Test R ²	Runtime (min)	Complexity
CADSR	0.636	0.626	193.06	22.250
DSR	0.622	0.616	124.65	24.625
TPSR	0.043	0.046	11.12	40.750
GP-GOMEA	0.645	0.636	124.66	73.000
DySymNet	0.584	0.456	396.11	1696.875

CADSR Expression:

$$\text{FIP} = -44.7 + 66.1 \cdot S_{\max} - \frac{0.0442}{\text{PCG}} - \frac{36.4 \cdot S_{\text{var}} - 4.58}{S_{\text{avg}}}$$

Fracture Mechanics

Method	Train R ²	Test R ²	Runtime (min)	Complexity
CADSR	0.636	0.626	193.06	22.250
DSR	0.622	0.616	124.65	24.625
TPSR	0.043	0.046	11.12	40.750
GP-GOMEA	0.645	0.636	124.66	73.000
DySymNet	0.584	0.456	396.11	1696.875

CADSR Expression:

$$\text{FIP} = -44.7 + 66.1 \cdot S_{\max} - \frac{0.0442}{\text{PCG}} - \frac{36.4 \cdot S_{\text{var}} - 4.58}{S_{\text{avg}}}$$

Fracture Mechanics

Method	Train R ²	Test R ²	Runtime (min)	Complexity
CADSR	0.636	0.626	193.06	22.250
DSR	0.622	0.616	124.65	24.625
TPSR	0.043	0.046	11.12	40.750
GP-GOMEA	0.645	0.636	124.66	73.000
DySymNet	0.584	0.456	396.11	1696.875

CADSR Expression:

$$\text{FIP} = -44.7 + 66.1 \cdot S_{\max} - \frac{0.0442}{\text{PCG}} - \frac{36.4 \cdot S_{\text{var}} - 4.58}{S_{\text{avg}}}$$

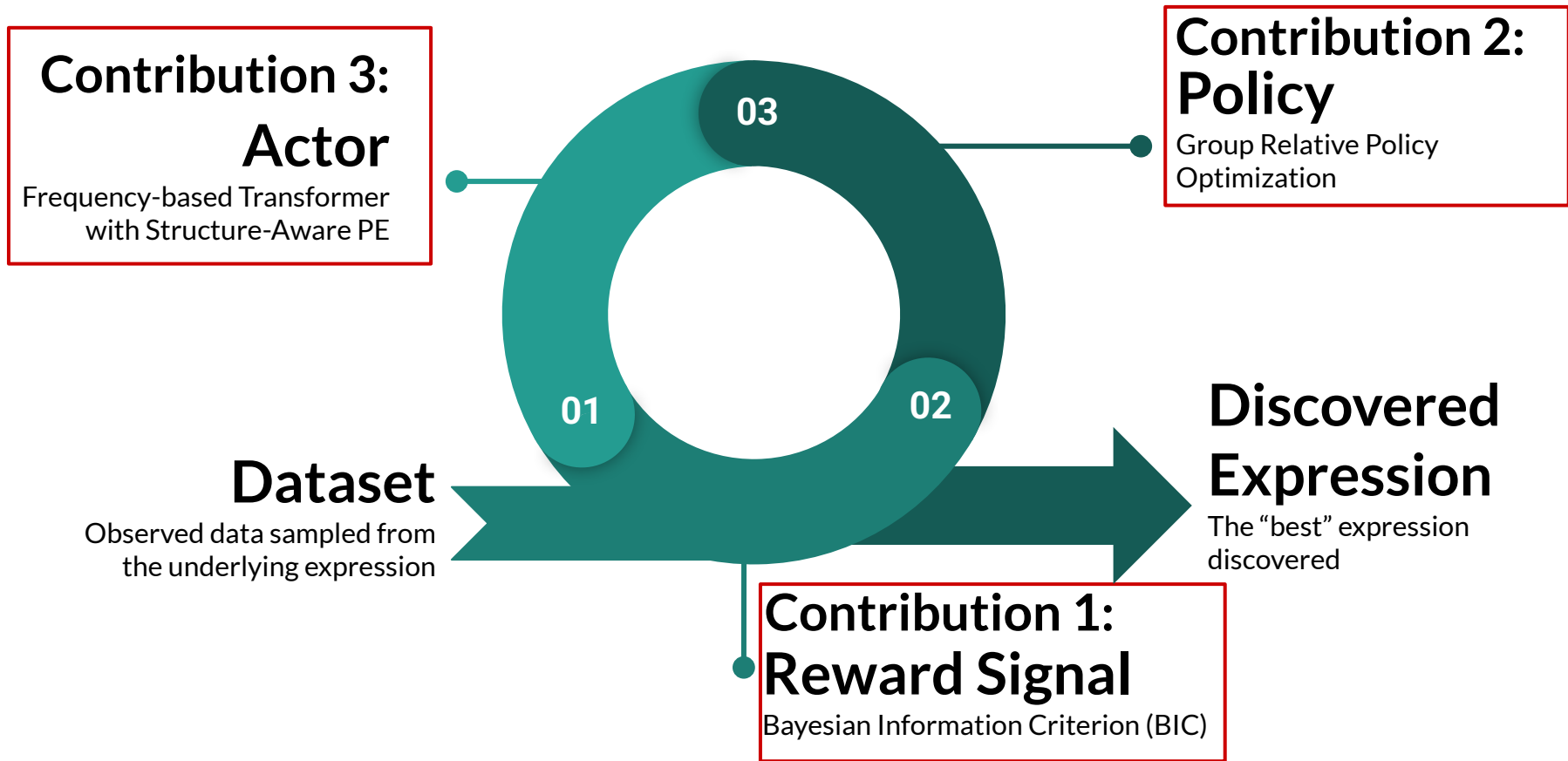
Fracture Mechanics

Method	Train R ²	Test R ²	Runtime (min)	Complexity
CADSR	0.636	0.626	193.06	22.250
DSR	0.622	0.616	124.65	24.625
TPSR	0.043	0.046	11.12	40.750
GP-GOMEA	0.645	0.636	124.66	73.000
DySymNet	0.584	0.456	396.11	1696.875

CADSR Expression:

$$\text{FIP} = -44.7 + 66.1 \cdot S_{\max} - \frac{0.0442}{\text{PCG}} - \frac{36.4 \cdot S_{\text{var}} - 4.58}{S_{\text{avg}}}$$

Complexity Aware Deep Symbolic Regression (CADSR)





Questions

Code, Results, and Dataset
are available at
<https://github.com/ZakBastiani/CADSR>

