

# Bingo: A Customizable Framework for Symbolic Regression with Genetic Programming

**David Randall<sup>1</sup>, Tyler Townsend<sup>2</sup>, Jacob Hochhalter<sup>1</sup>, Geoffrey Bomarito<sup>3</sup>**

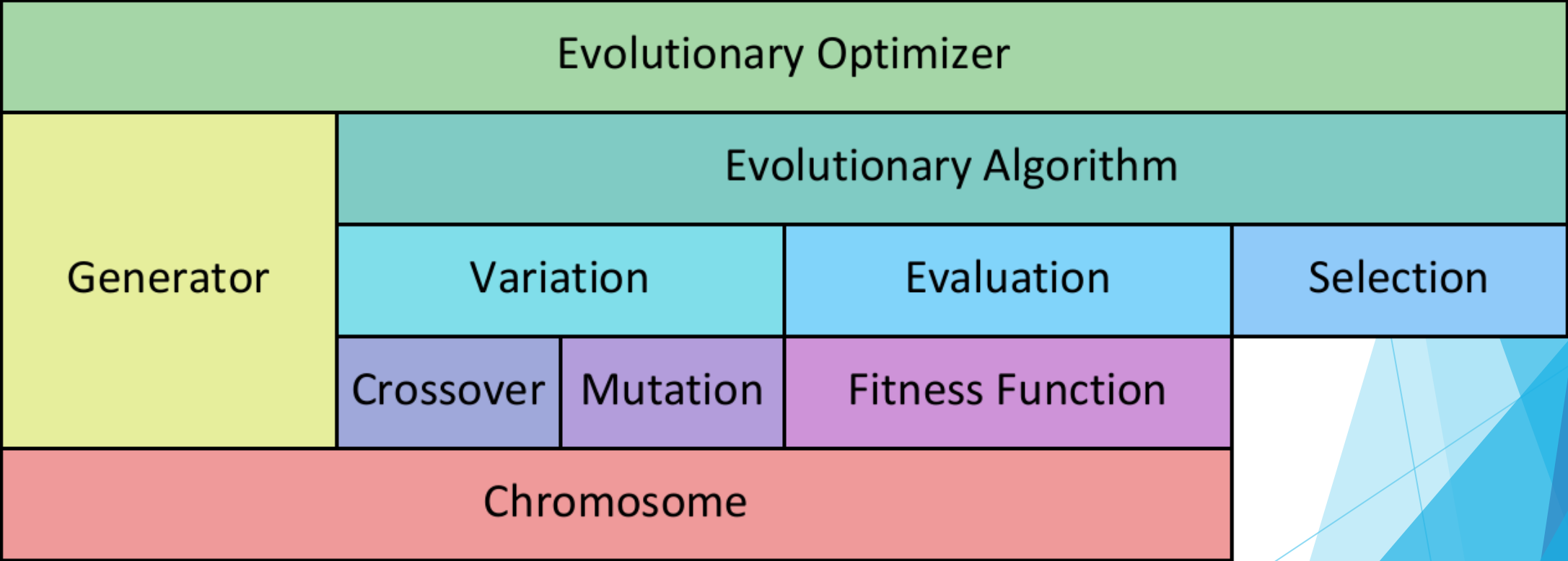
<sup>1</sup>University of Utah, <sup>2</sup>Microsoft (Views and opinions are personal and do not represent Microsoft), <sup>3</sup>NASA Langley Research Center

# What is Bingo?

- ▶ [Bingo](#) is an open-source framework for genetic programming for symbolic regression (GPSR) developed by NASA
- ▶ Made to be
  - ▶ Modular
  - ▶ Extendible
  - ▶ Efficient
- ▶ Why use Bingo over other GPSR packages?
  - ▶ High- and low-level interfaces for ease-of-use and customizability
  - ▶ Modularity makes it easy to compare and develop components of GPSR
  - ▶ Produces simple and accurate models

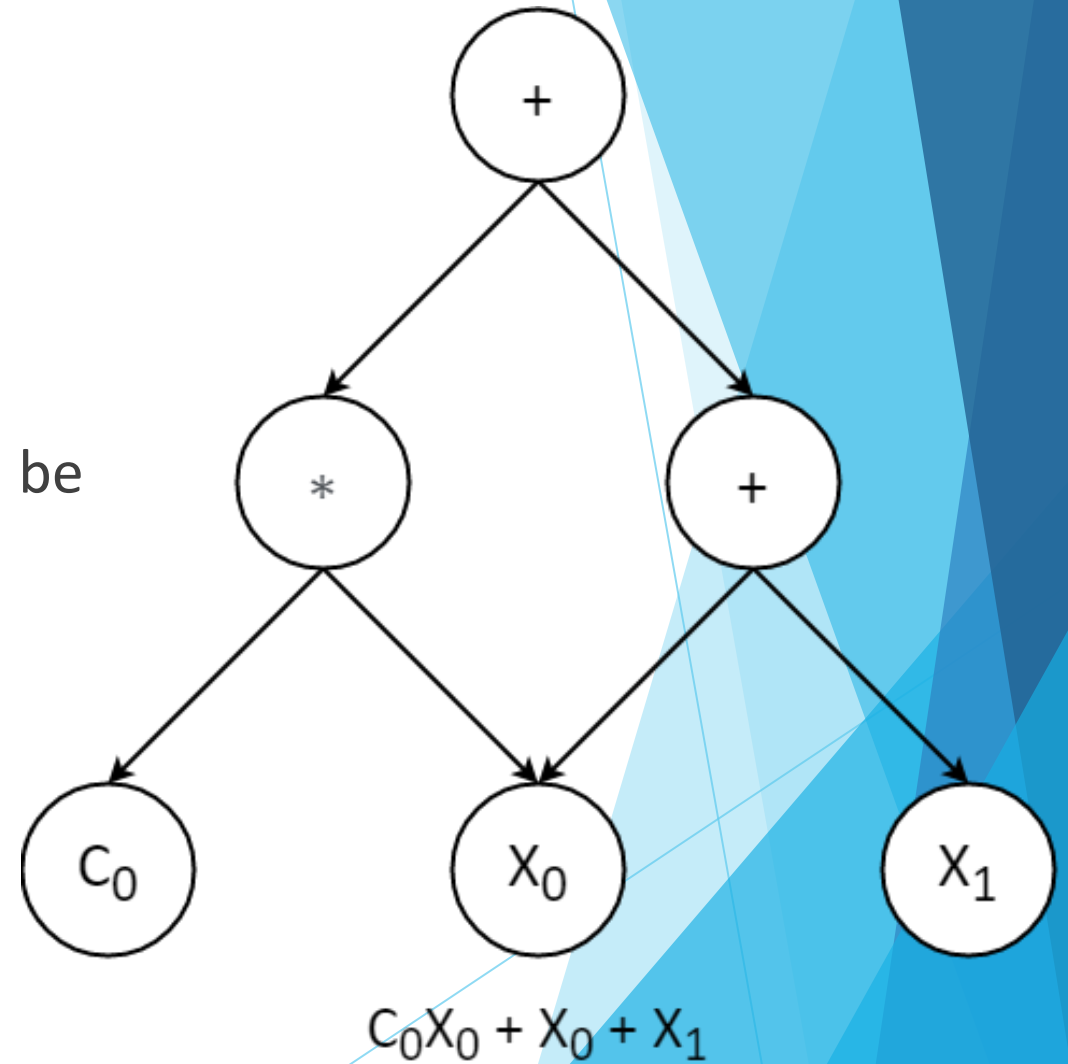


# Bingo's Structure

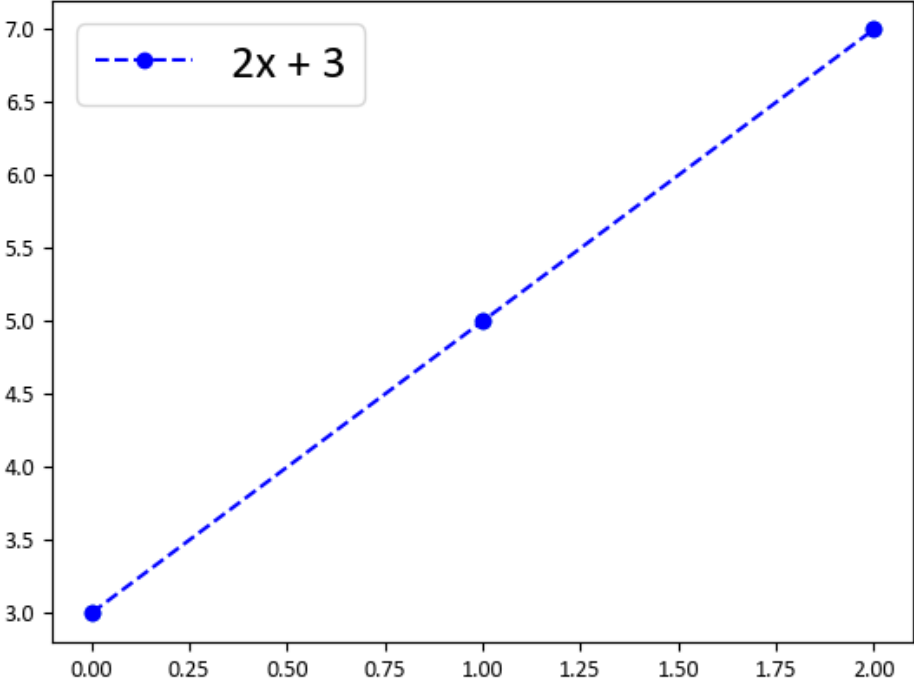
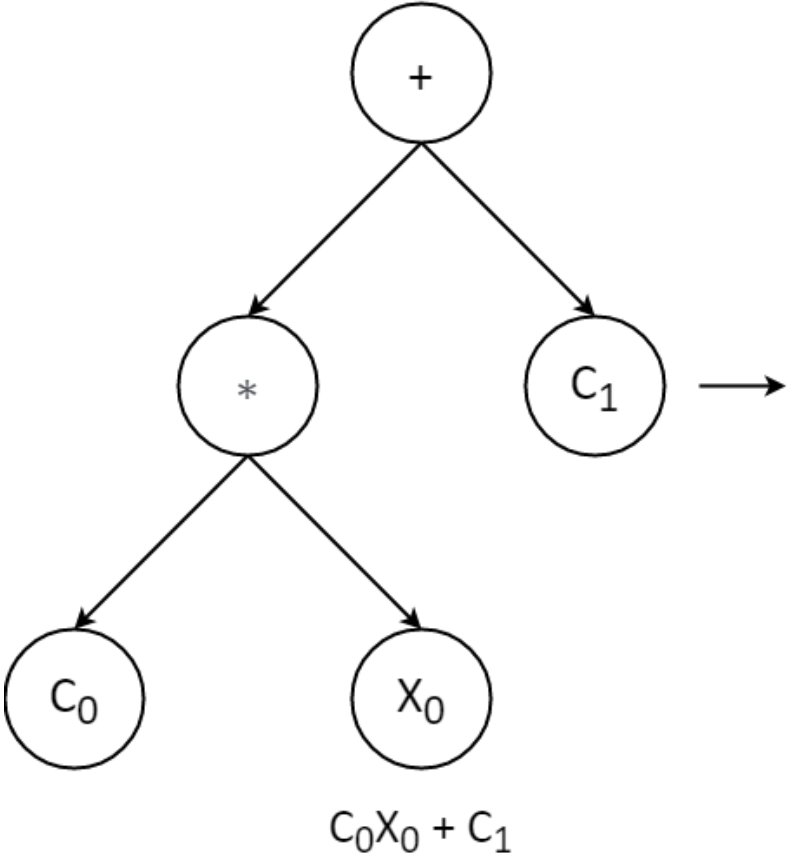


# Equations/Chromosomes

- ▶ Encoded as directed acyclic graphs (AGraphs)
  - ▶ Computational benefits over tree encodings [1]
- ▶ Have free-form constants/parameters that can be locally optimized
- ▶ Complexity
  - ▶ Measured as number of utilized nodes in the graph encoding



# Local Optimization of Parameters



$C_0 = 2$   
 $C_1 = 3$

# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
<b>0</b>	<b>constant</b>	<b>0</b>	<b>0</b>	<b><math>C_0</math></b>
1	variable	0	0	$X_0$
2	variable	1	1	$X_1$
3	*	0	1	$C_0X_0$
4	+	1	2	$X_0 + X_1$
5	sin	2	2	$\sin(X_1)$
6	+	3	4	$C_0X_0 + X_0 + X_1$

# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
<b>1</b>	<b>variable</b>	<b>0</b>	<b>0</b>	<b><math>X_0</math></b>
2	variable	1	1	$X_1$
3	*	0	1	$C_0X_0$
4	+	1	2	$X_0 + X_1$
5	sin	2	2	$\sin(X_1)$
6	+	3	4	$C_0X_0 + X_0 + X_1$

# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
1	variable	0	0	$X_0$
<b>2</b>	<b>variable</b>	<b>1</b>	<b>1</b>	<b><math>X_1</math></b>
3	*	0	1	$C_0X_0$
4	+	1	2	$X_0 + X_1$
5	sin	2	2	$\sin(X_1)$
6	+	3	4	$C_0X_0 + X_0 + X_1$



# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
1	variable	0	0	$X_0$
2	variable	1	1	$X_1$
<b>3</b>	<b>*</b>	<b>0</b>	<b>1</b>	<b><math>C_0X_0</math></b>
4	+	1	2	$X_0 + X_1$
5	sin	2	2	$\sin(X_1)$
6	+	3	4	$C_0X_0 + X_0 + X_1$

# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
1	variable	0	0	$X_0$
2	variable	1	1	$X_1$
3	*	0	1	$C_0X_0$
<b>4</b>	<b>+</b>	<b>1</b>	<b>2</b>	<b><math>X_0 + X_1</math></b>
5	sin	2	2	$\sin(X_1)$
6	+	3	4	$C_0X_0 + X_0 + X_1$

# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
1	variable	0	0	$X_0$
2	variable	1	1	$X_1$
3	*	0	1	$C_0X_0$
4	+	1	2	$X_0 + X_1$
<b>5</b>	<b>sin</b>	<b>2</b>	<b>2</b>	<b>sin(<math>X_1</math>)</b>
6	+	3	4	$C_0X_0 + X_0 + X_1$

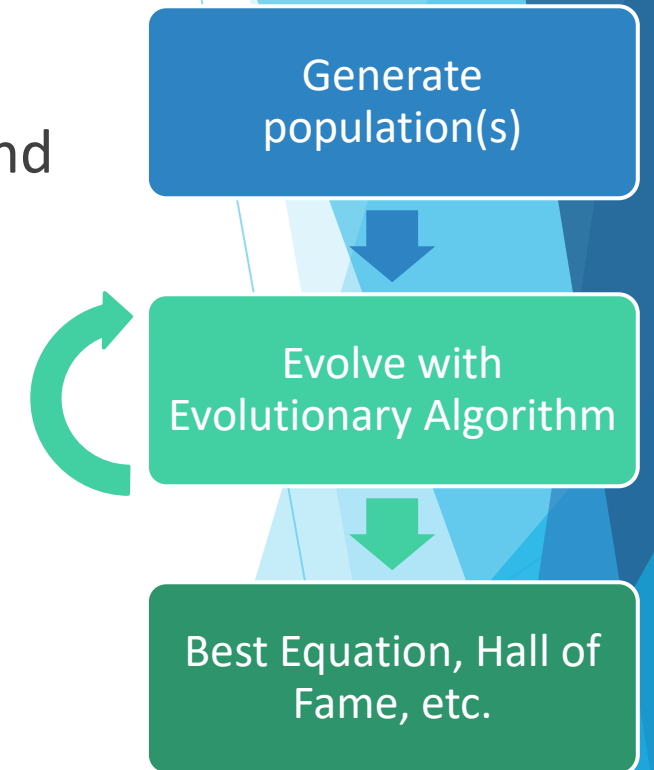
# Command Arrays

- ▶ Genome of equations
- ▶ Terminal nodes and operators
- ▶ Not all commands are utilized
  - ▶ Size of command array  $\neq$  complexity

i	node	parameter 1	parameter 2	expression
0	constant	0	0	$C_0$
1	variable	0	0	$X_0$
2	variable	1	1	$X_1$
3	*	0	1	$C_0X_0$
4	+	1	2	$X_0 + X_1$
5	sin	2	2	$\sin(X_1)$
<b>6</b>	<b>+</b>	<b>3</b>	<b>4</b>	<b><math>C_0X_0 + X_0 + X_1</math></b>

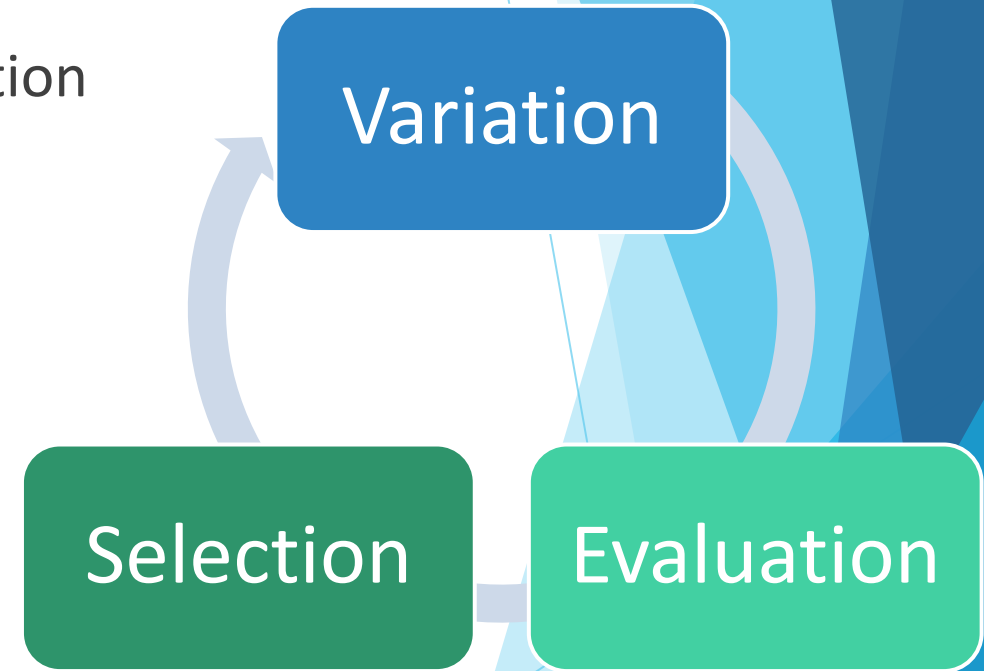
# Evolutionary Optimizer

- ▶ Responsible for generating an initial population (generator) and evolving that population (evolutionary algorithm)
- ▶ Potential generator customizations
  - ▶ Initial population with all unique equations
  - ▶ Seeding of initial population with parts of previously found equations
  - ▶ Filtering of random equations
- ▶ Potential genetic programming workflow customizations
  - ▶ Archipelago
  - ▶ Coevolution



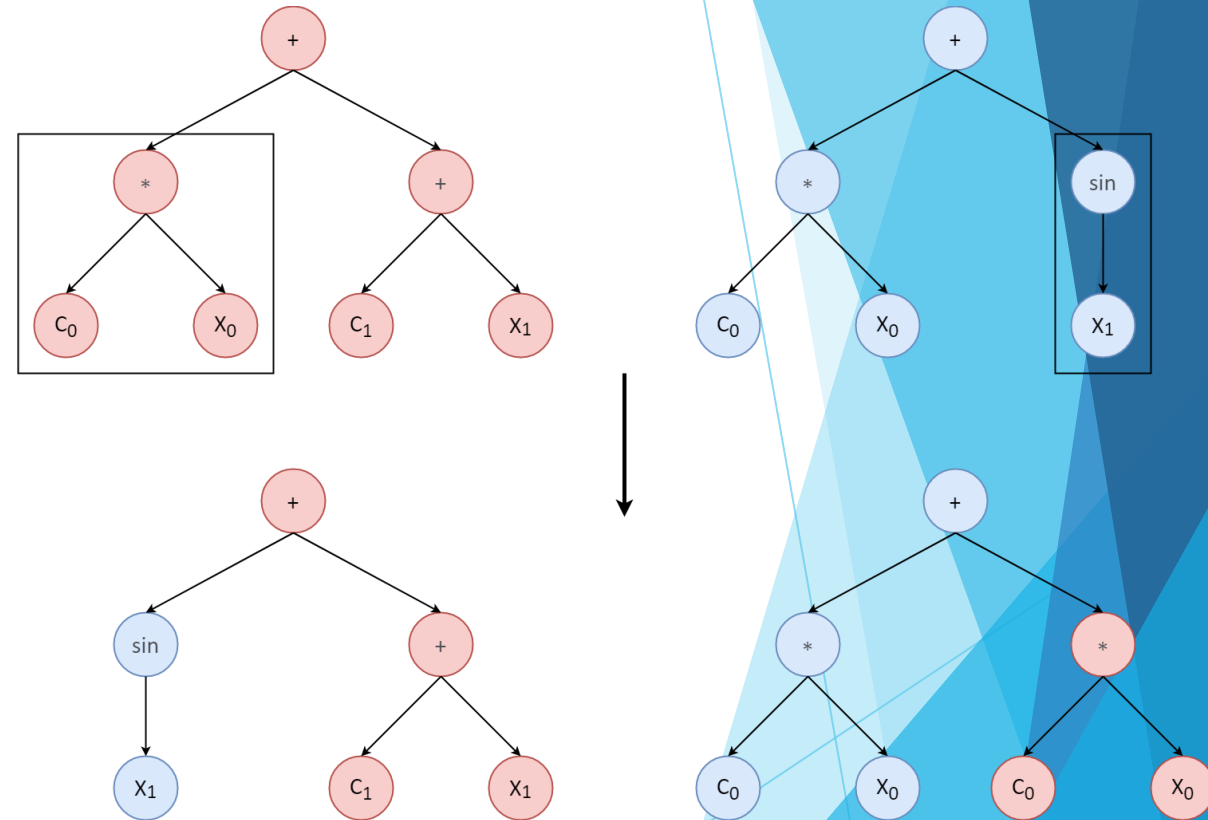
# Evolutionary Algorithm

- ▶ Responsible for performing evolution on a population
- ▶ Follows a traditional GPSR workflow by default
  - ▶ Variation
  - ▶ Evaluation
  - ▶ Selection
- ▶ Can introduce other evolutionary operations
  - ▶ Simplification



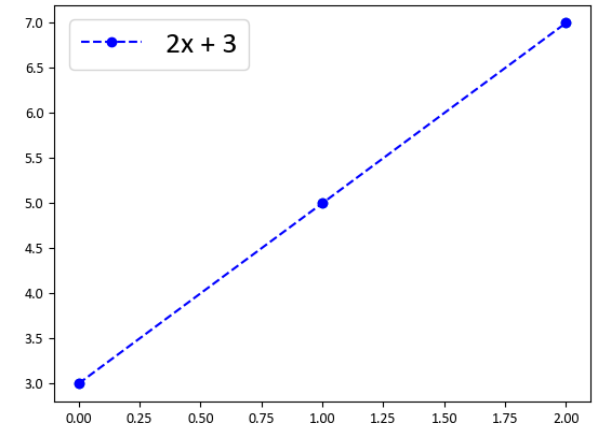
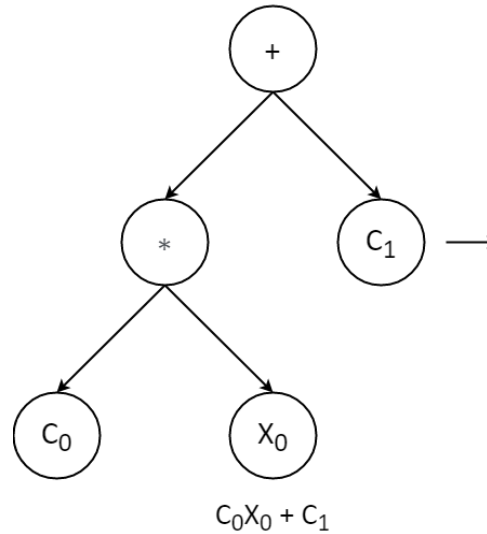
# Variation

- ▶ Takes a population and varies it to form new individuals
- ▶ Two main methods
  - ▶ Crossover and mutation
  - ▶ Crossover or mutation
- ▶ Potential customizations
  - ▶ Only select variations that are more fit than their original counterparts



# Evaluation

- ▶ Evaluates the fitness of individuals in a population
- ▶ Fitness functions
  - ▶ Explicit regression
  - ▶ Implicit regression
  - ▶ Continuous local optimization
- ▶ Potential customizations
  - ▶ Domain-specific constraints
  - ▶ Optimization with consideration of uncertainty

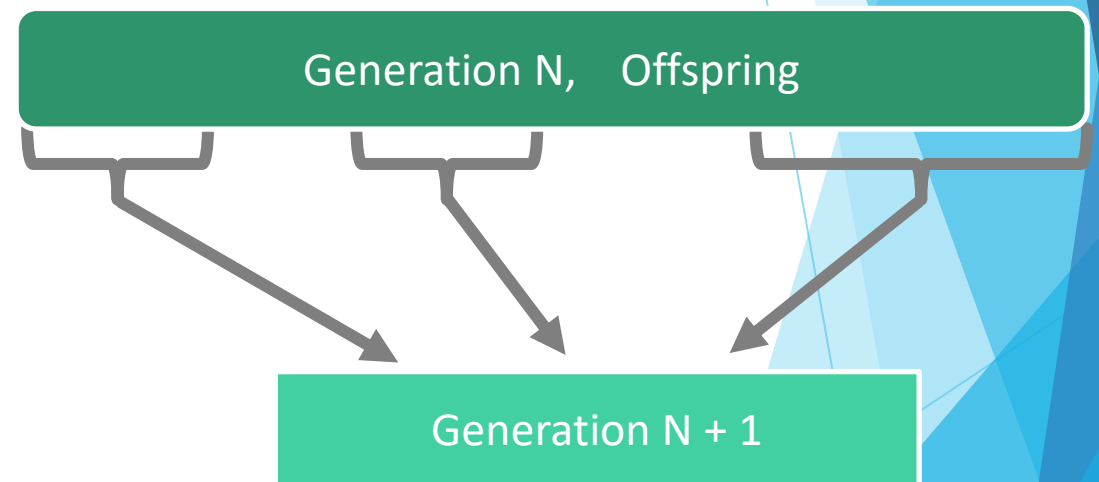


→  $C_0 = 2$   
 $C_1 = 3$



# Selection

- ▶ Selects among original population and its variants to form next generation's population
  - ▶ Deterministic crowding
  - ▶ Age-fitness Pareto selection [2]
  - ▶ Tournament
- ▶ Potential customizations
  - ▶ Probabilistic crowding



# Notable Customizations

- ▶ Fitness functions for mechanical engineering
  - ▶ Incorporation of automatic differentiation
- ▶ Tensorial GPSR
- ▶ Sequential Monte-Carlo and uncertainty quantification

$$\frac{\partial}{\partial X_0} \sin(X_0^2) = \cos(X_0^2) \cdot 2X_0$$

# What if I don't care about customization?

- ▶ “Out-of-the-box” scikit-learn<sup>1</sup> wrapper
- ▶ Easy configuration and simple interface for training
- ▶ Can use with other scikit-learn utilities
  - ▶ e.g., Cross validation classes, useful for hyperparameter tuning

```
reg = SymbolicRegressor(evolutionary_algorithm=AgeFitnessEA, ...)  
reg.fit(x_train, y_train)  
y_pred = reg.predict(x_test)
```

<sup>1</sup>This is not an endorsement by the National Aeronautics and Space Administration (NASA)

# Efficiency

- ▶ **Parallelism**
  - ▶ Distributed memory: parallel evolution of islands
  - ▶ Shared memory: parallel evaluation of individuals
- ▶ **Fitness predictors [3]**
  - ▶ Subsets of data used to predict the true fitness of individuals
  - ▶ Less computational effort required for evaluation
  - ▶ Similar ideas to support vectors and batch evaluation

# Efficiency (cont.)

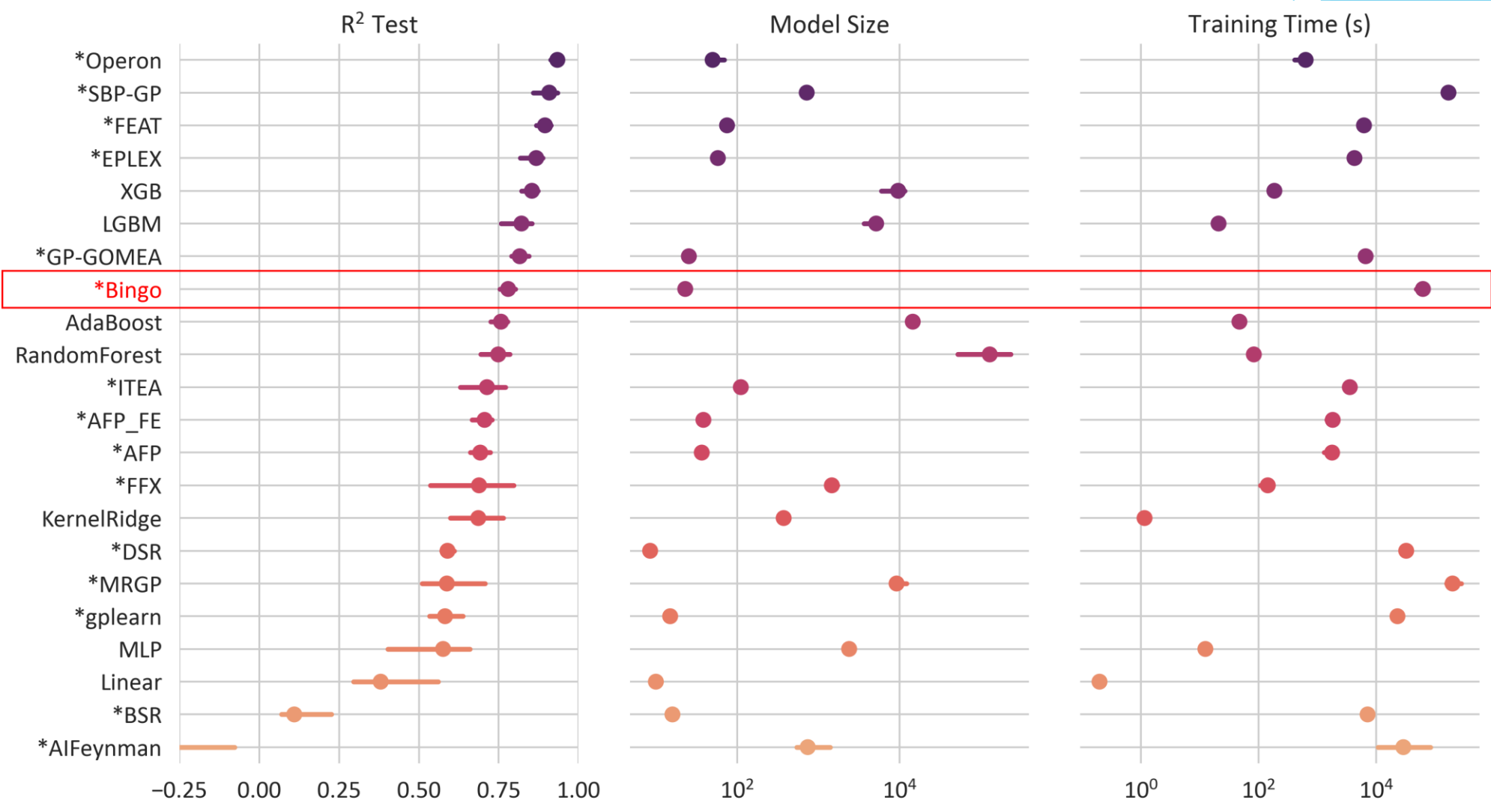
- ▶ Equation simplification
  - ▶ More complex equation -> more effort to evaluate
  - ▶ We use simplified versions of equations for evaluation
    - ▶ Doesn't modify the original genotype
- ▶ C++ backend
  - ▶ An optional backend that can be enabled so that some components are implemented in C++

# SRBench

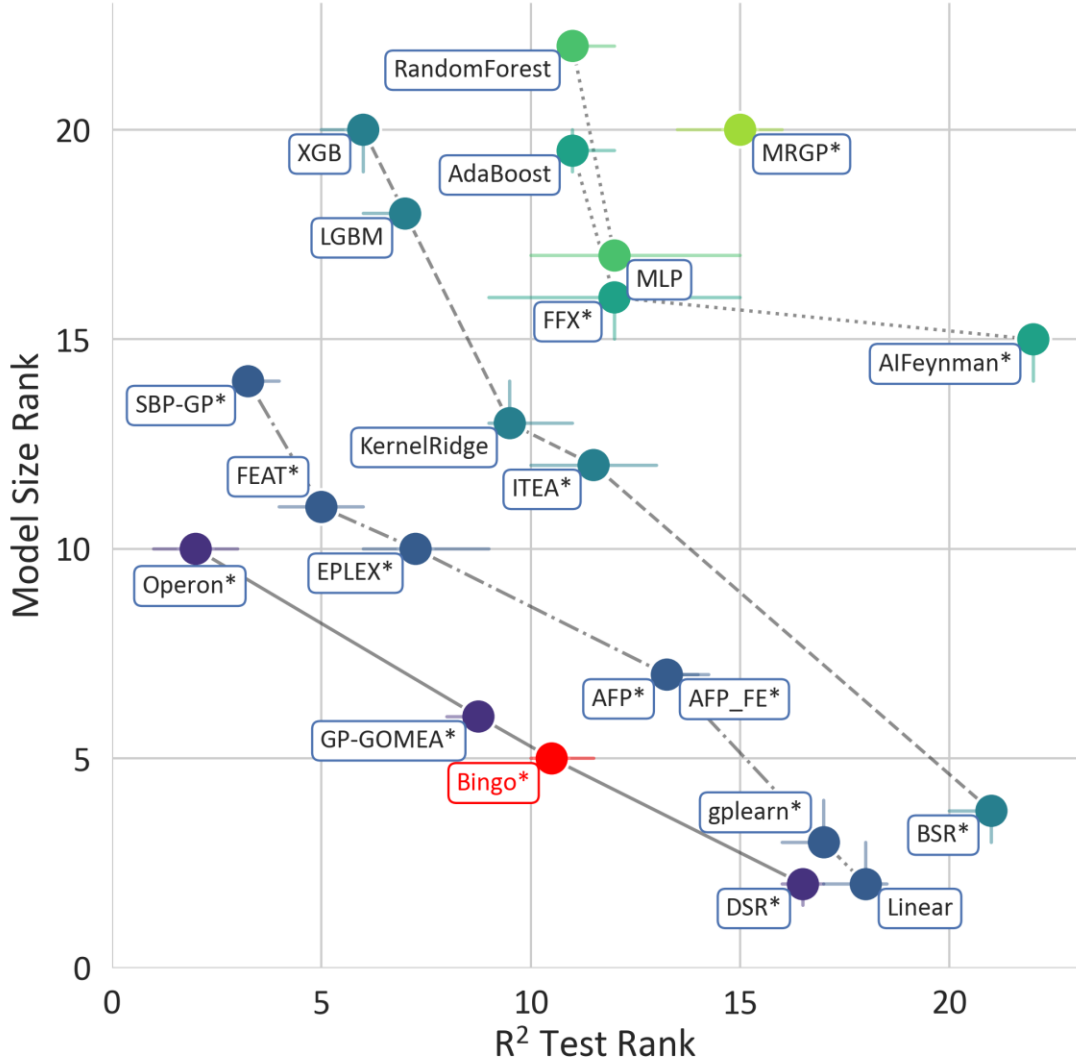
- ▶ SRBench<sup>2</sup> and Penn Machine Learning Benchmarks (PMLB)<sup>3</sup>
  - ▶ SRBench [4]
    - ▶ Open-source benchmark for SR and ML methods
    - ▶ Used for this workshop's competition
  - ▶ PMLB [5]
    - ▶ Mix of realistic and synthetic datasets
- ▶ Black-box problems
  - ▶ General regression
- ▶ Ground-truth problems
  - ▶ Recovering an underlying equation

<sup>2,3</sup>This is not an endorsement by the National Aeronautics and Space Administration (NASA)

# Results

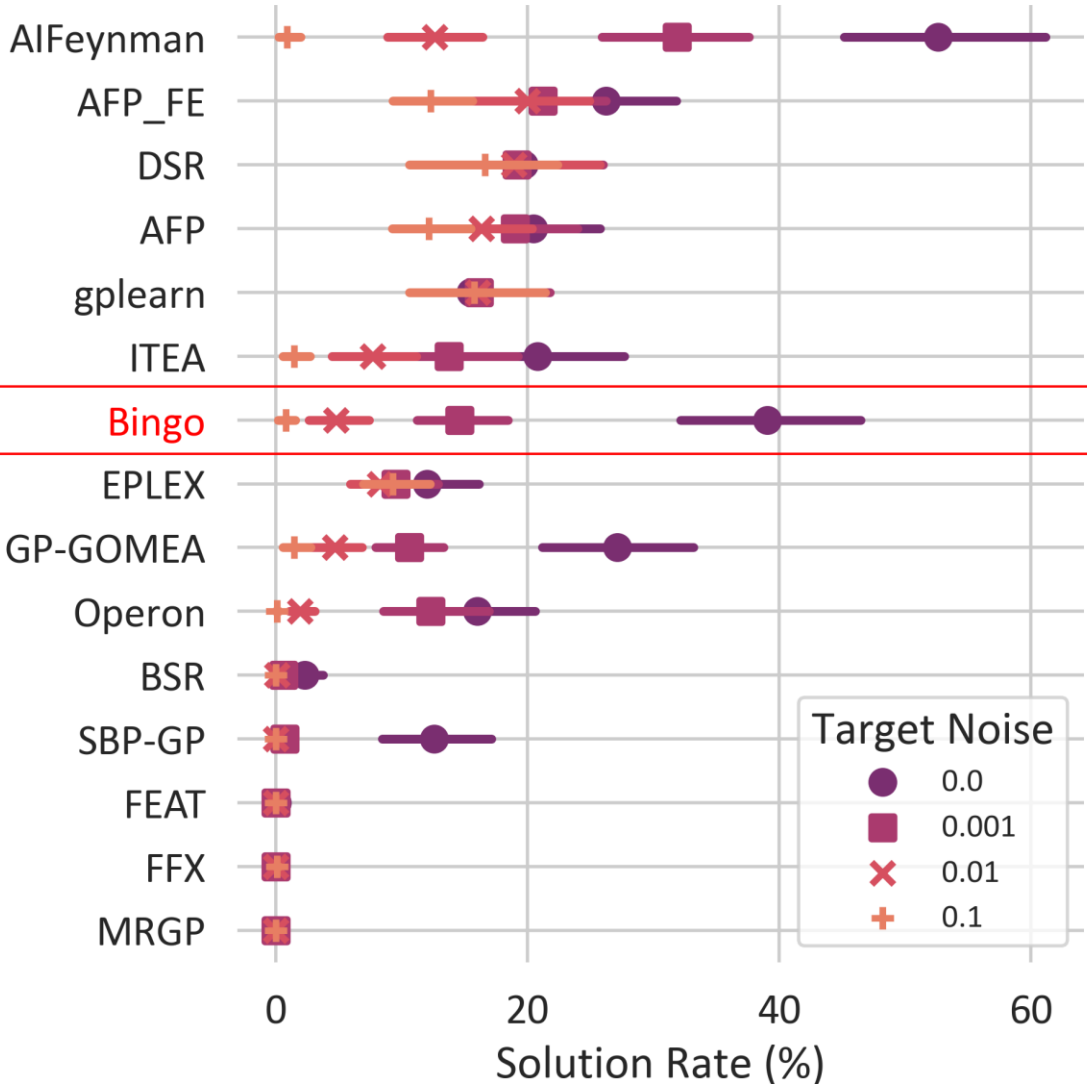


# Results (cont.)





# Results (cont.)



# Summary

- ▶ What is Bingo?
  - ▶ Flexible, modular framework for GPSR
  - ▶ Made to easily compare and create GPSR components
    - ▶ Sandbox
  - ▶ Can still be used for typical GPSR without modifications
    - ▶ SRBench
    - ▶ sklearn wrapper

# Future Work

- ▶ Faster training times
  - ▶ More parts implemented in C++
- ▶ Built-in integration with other tools
  - ▶ SymPy<sup>4</sup>
  - ▶ PyTorch<sup>5</sup>/TensorFlow<sup>6</sup>
- ▶ Better performance on general regression problems
  - ▶ Revisiting and analyzing performance on SRBench

<sup>4,5,6</sup>This is not an endorsement by the National Aeronautics and Space Administration (NASA)

# References

- ▶ [1] Michael Schmidt and Hod Lipson. 2007. Comparison of tree and graph encodings as function of problem complexity.
- ▶ [2] Michael D. Schmidt and Hod Lipson. 2010. Age-fitness pareto optimization.
- ▶ [3] Michael D. Schmidt and Hod Lipson. 2008. Coevolution of Fitness Predictors.
- ▶ [4] William La Cava, et al. 2021. Contemporary Symbolic Regression Methods and their Relative Performance.
- ▶ [5] Joseph D. Romano, et al. 2020. PMLB v1.0: an open source dataset collection for benchmarking machine learning methods.