

# CP-Analyses with Symbolic Regression

[H. Bahl, E. Fuchs, M. Menen, T. Plehn '25]

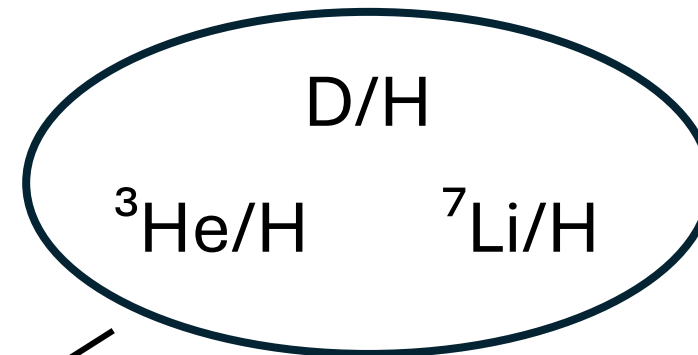
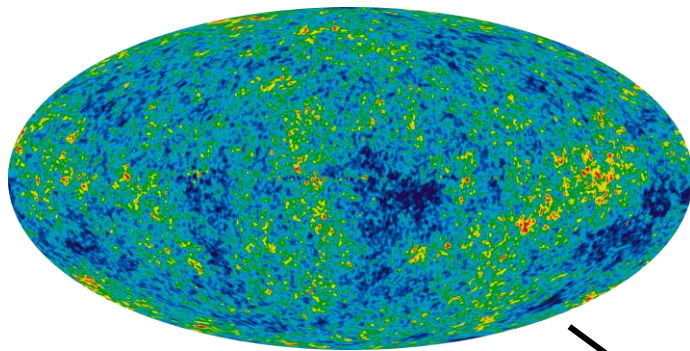
---

Quantum Universe Attract Workshop

Hamburg, 24.11.2025

Marco Menen



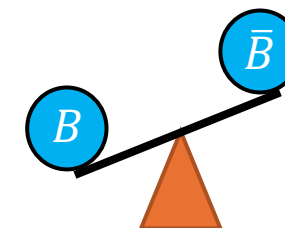


Baryon asymmetry  
of the Universe (BAU)

[Planck '18]

[S. Navas et al. '24]

$$\eta = \frac{n_B - n_{\bar{B}}}{s} \approx 8.7 \cdot 10^{-11}$$



Sakharov conditions  
for successful  
baryogenesis:

[A. D. Sakharov '67]

- ① Departure from thermal equilibrium
- ② CP (and C) violation
- ③ Baryon number violation

Not sufficient  
amounts in SM



BSM Higgs  
sectors

# The Higgs Characterization Model:

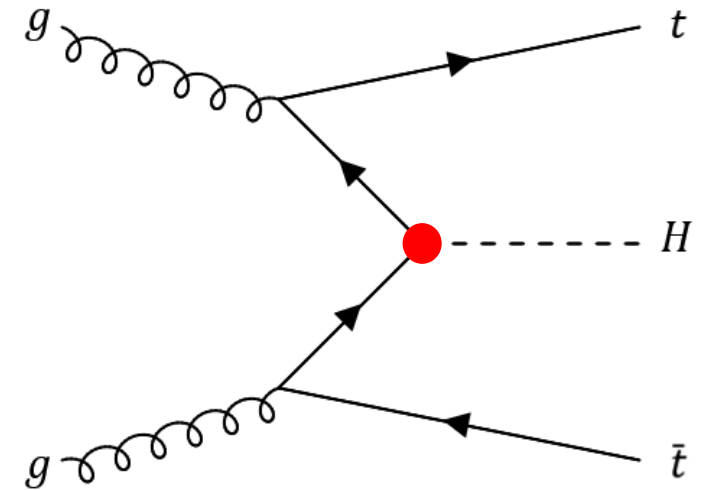
[P. Artoisenet et al. '13]

Modification of SM top-Yukawa coupling via CP-violating phase

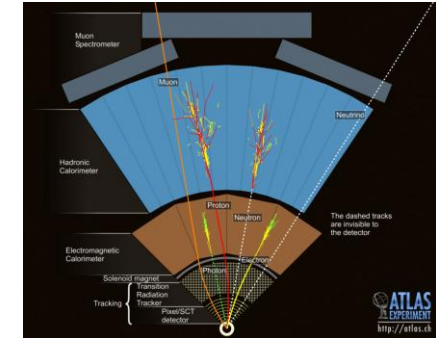
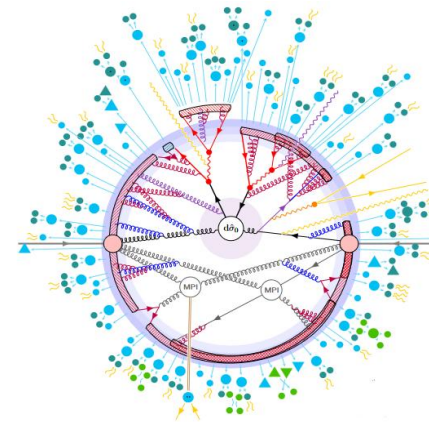
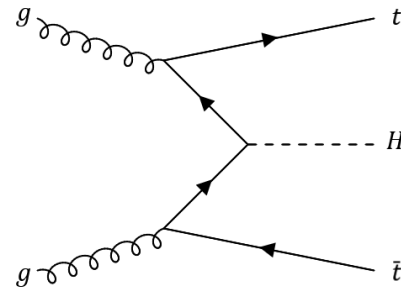
Total rate modifier  $g_t$ , CP-mixing angle  $\alpha_t$

$$\mathcal{L}_{t\bar{t}H}^{\text{HCM}} \supset \frac{y_t}{\sqrt{2}} t \, g_t (\cos \alpha_t + i\gamma_5 \sin \alpha_t) \bar{t} H$$

SM:  $g_t = 1 \quad \alpha_t = 0$



$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \chi_i y_{ij} \chi_j \phi + h.c. + |D_\mu \phi|^2 - V(\phi)$$



## Analysis

$$f(x_{\text{parton}}|x_{\text{reco}}) = ?$$

Symmetries imprinted

Fast evaluation

Easily reusable

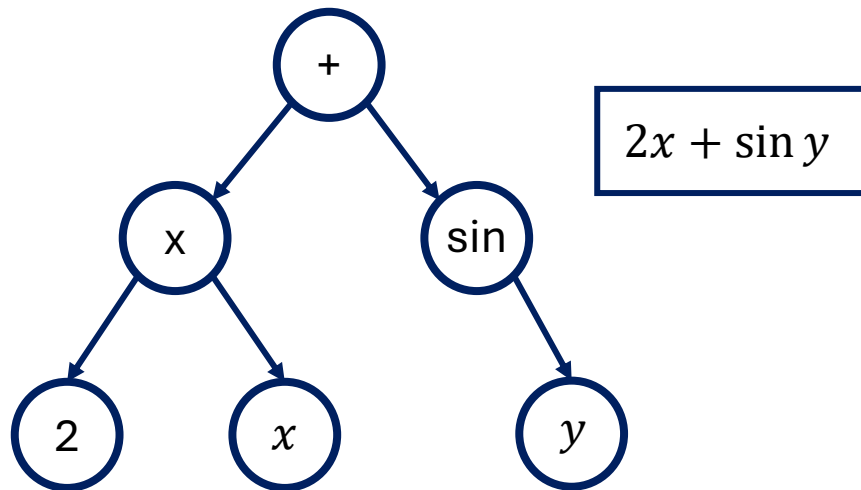
Maximal interpretability

# Symbolic Regression: Equation discovery from data

## PySR (Genetic algorithm)

[M. Cranmer '23]

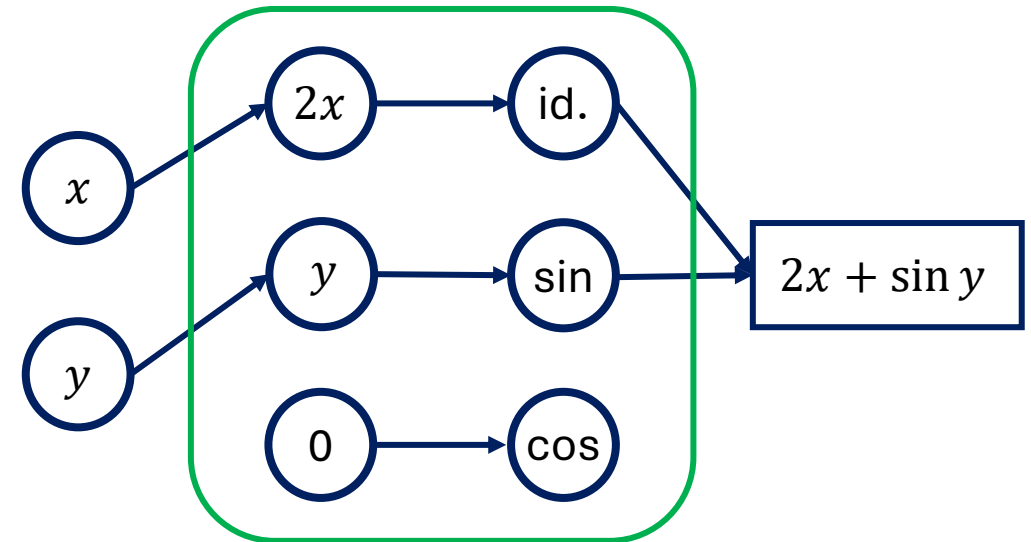
- Easy equation → complex equation
- Based on binary tree
- “Survival of the fittest”



## SymbolNet (Deep neural network)

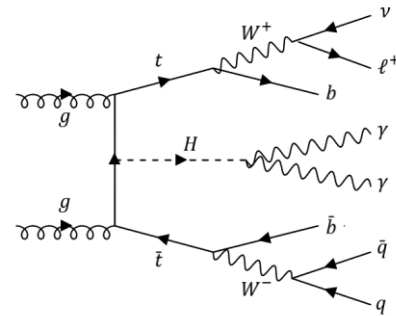
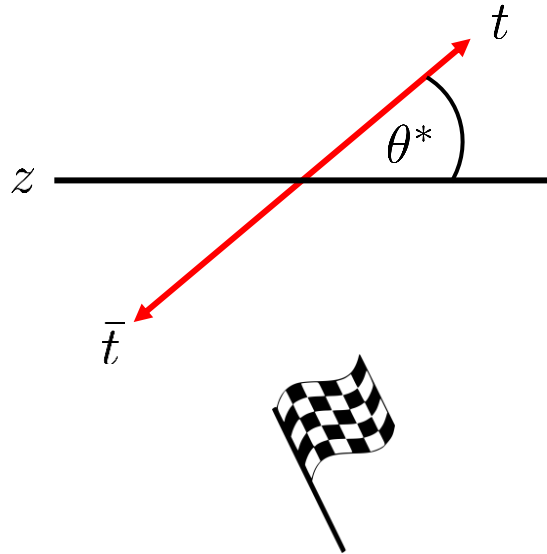
[Tsoi, Loncar, Dasu, Harris '24]

- Complex equation → easy equation
- Based on **symbolic layers**
- Backpropagation

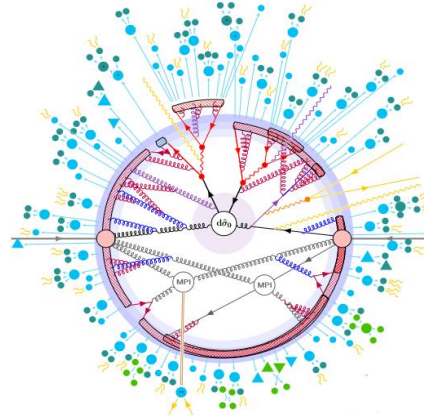


# Goal: Reconstruct Collins-Soper angle for $t\bar{t}H$

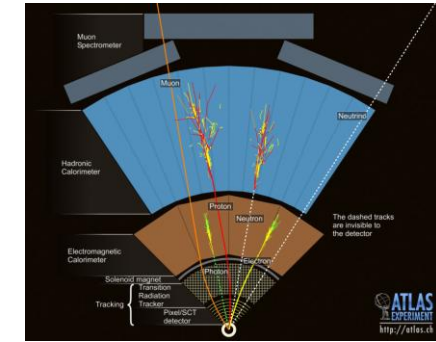
[J. C. Collins, D. E. Soper '77]



Scenario 1



• • •



Scenario 6

- Define benchmark scenarios with increasing complexity
- Scenario 1:  $t\bar{t}$ -frame variables, full neutrino information, no shower or detector effects
- Scenario 6: lab-frame variables, only  $E_T^{\text{miss}}$  information, an additional QCD jet, variables are smeared according to detector resolution

# Scenario 1:

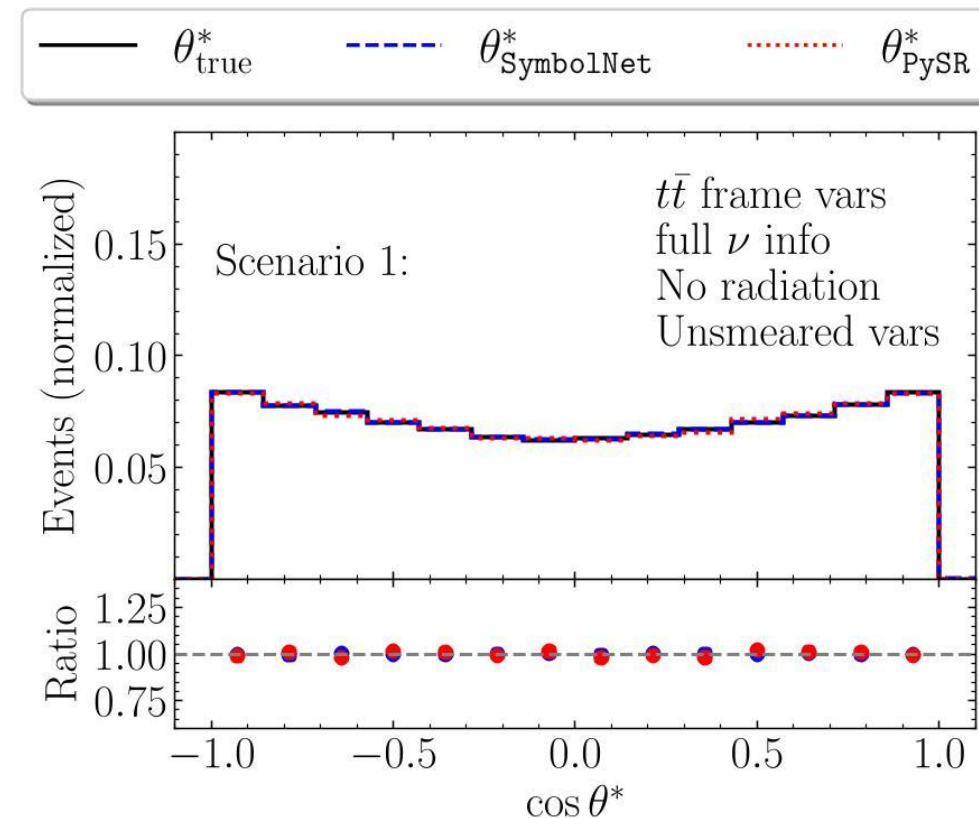
Equation for  $\cos \theta^*$  is known:

$$\cos \theta^* = \frac{\sum p_{z,i}^{t\bar{t}}}{\sum \|p_i^{t\bar{t}}\|}$$

Both methods reconstruct it well

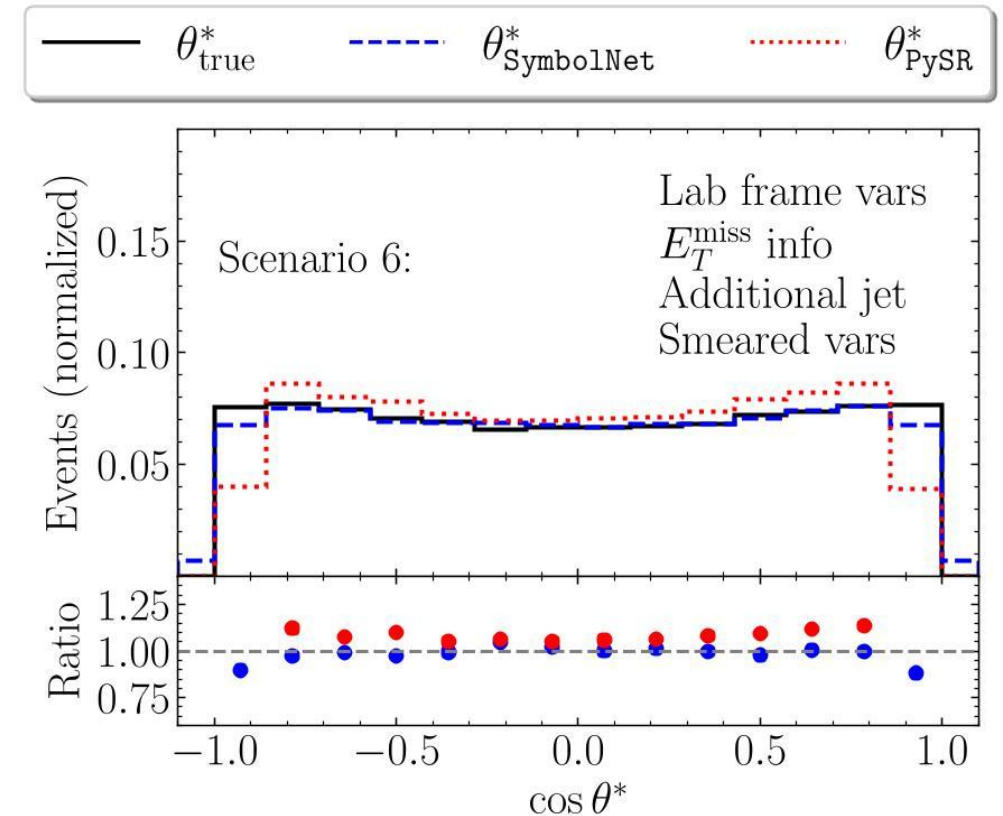
$$\cos \theta_{\text{PySR}}^* = \frac{p_{z,b} + p_{z,\bar{l}} + p_{z,\nu}}{\|p_b + p_{\bar{l}} + p_{\nu}\|_3}$$

$$\cos \theta_{\text{SymbolNet}}^* = \frac{\underbrace{p_{z,b} + p_{z,\bar{l}} + p_{z,\nu}}_{\text{Leptonically decaying top}} + \underbrace{p_{z,\bar{b}} + p_{z,q} + p_{z,\bar{q}}}_{\text{Hadronically decaying top}}}{\|p_b + p_{\bar{l}} + p_{\nu} + p_{z,\bar{b}} + p_{z,q} + p_{z,\bar{q}}\|_3}$$



## Scenario 6:

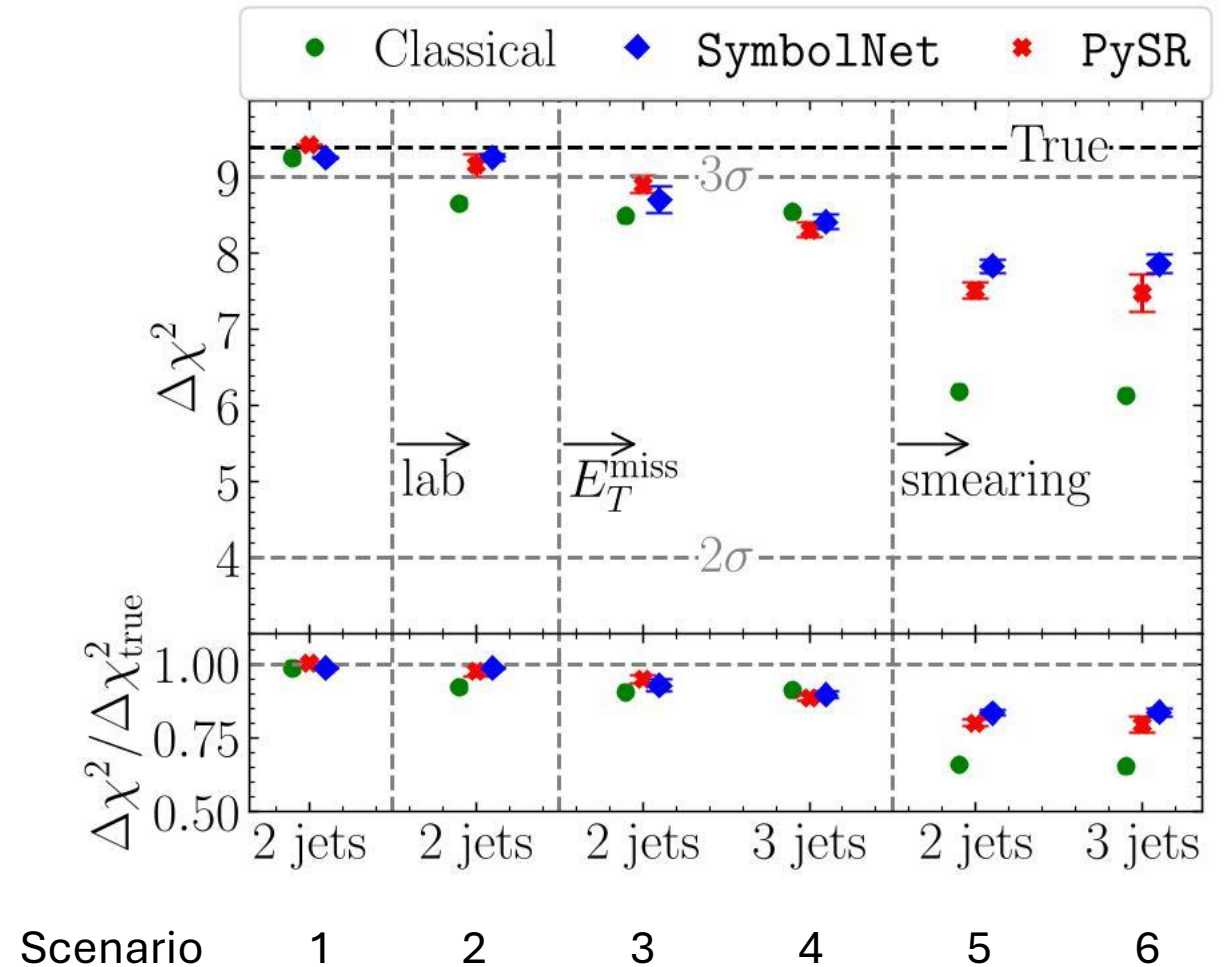
- Full information not available at detector level
- Center of the  $\cos \theta^*$  is reconstructed well by both algorithms
- SymbolNet performs slightly better than PySR
- But the SymbolNet equation is also much more complicated





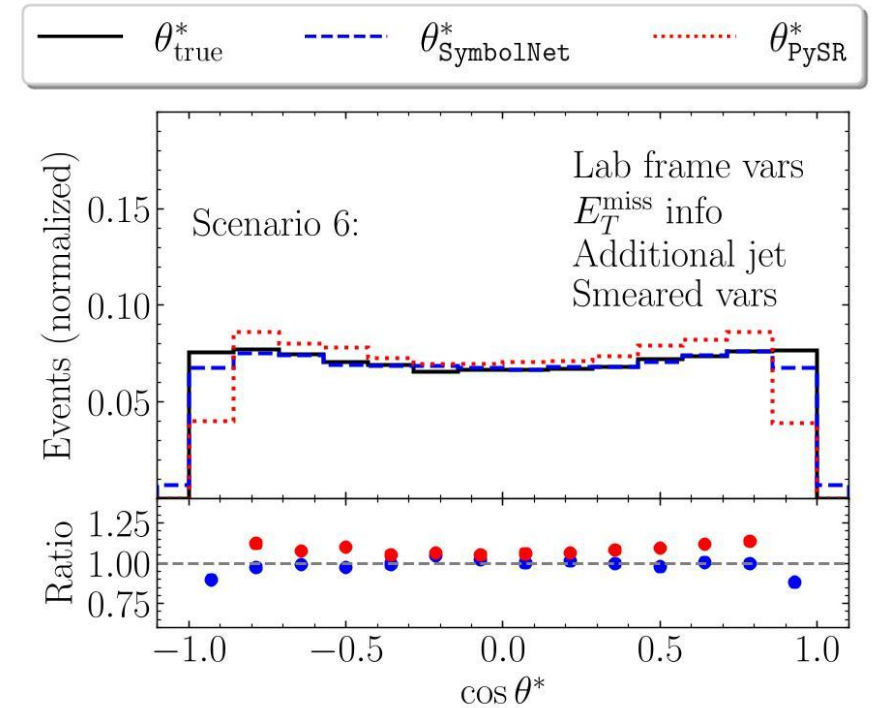
Significances to exclude  $\alpha_t = 45^\circ$  with a SM measurement at  $\mathcal{L} = 300\text{fb}^{-1}$

- „Classical“ method: Reconstruct top quark via mass requirements
- SR approaches outperform classical reconstruction
- SymbolNet performs slightly better in complex scenarios
- Around 80% of CP information is retained by SR



# Conclusions / Outlook

- Can find an equation to reconstruct a parton-level variable using detector-level information
- Retain around 80% of CP information
- SR algorithms outperform classical reconstruction method



⇒ Many potential applications, also outside of physics

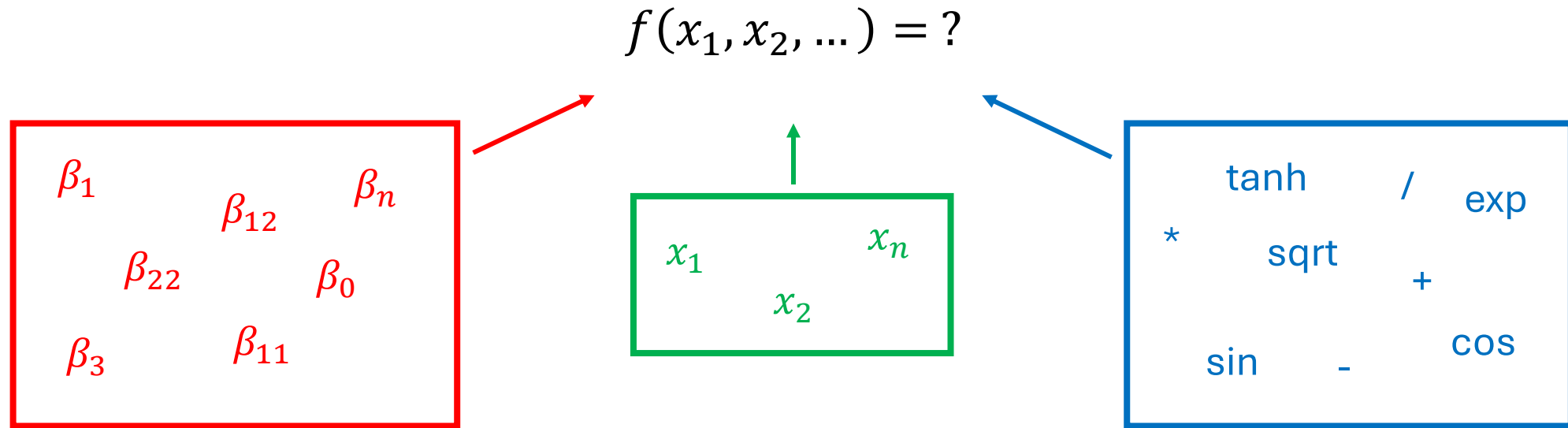
⇒ Further improvements might be possible from adding scientific knowledge

[Shojaee, Meidani, Gupta, Farimani, Reddy '24]

Backup

Default regression: Provide equation, estimate **parameters**

Symbolic regression: Estimate equation from **parameters**, **variables**, **operators**



General tradeoff: Goodness of fit  $\leftrightarrow$  complexity of equation

Very (NP) hard task  $\Rightarrow$  Tools needed

# PySR (Genetic algorithm):

[M. Cranmer '23]

Updates to populations via inner loop:

1. Select a tree to mutate determined via its fitness
2. Select a mutation to apply:

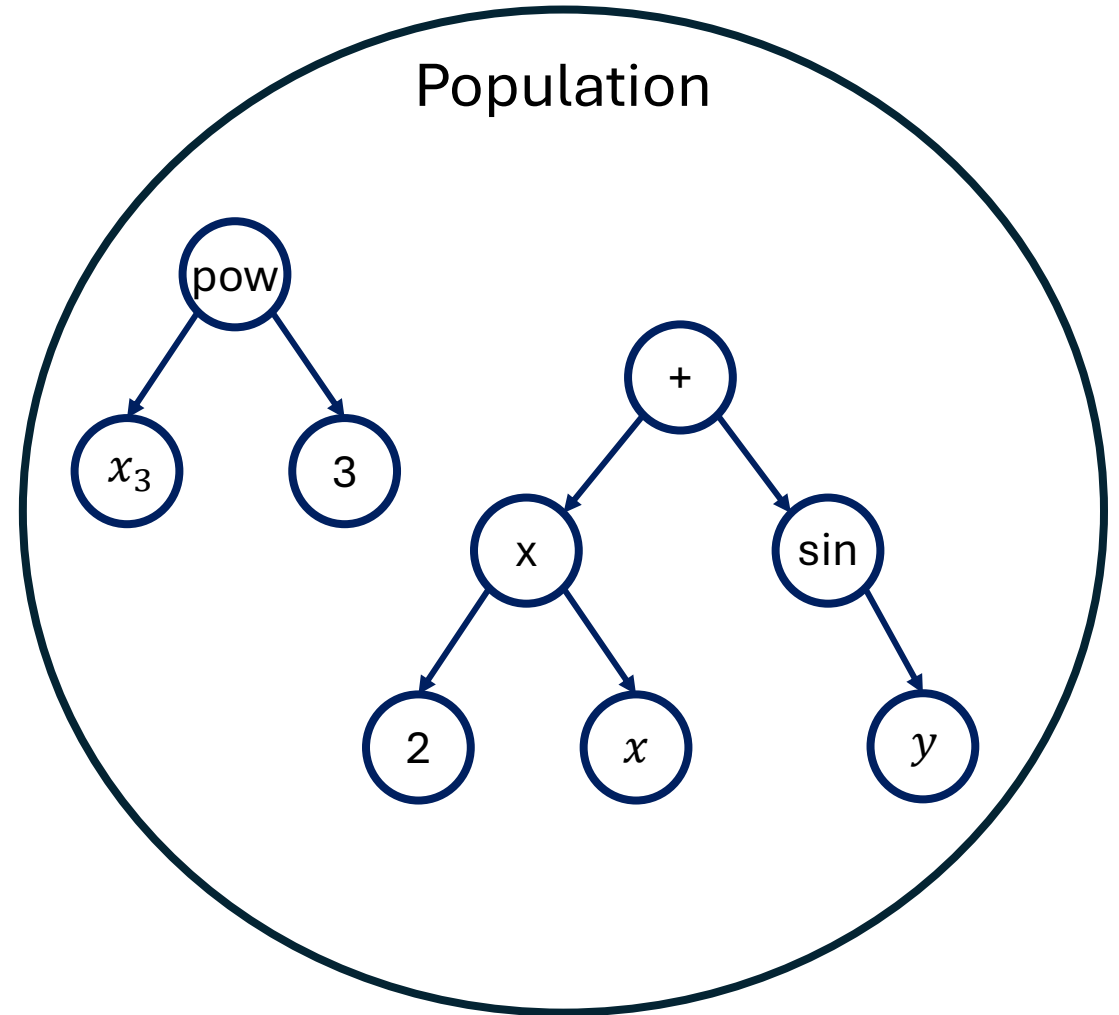
Operator mutation

Crossover

Simplification

Constant optimization

3. Replace the least fit member of the population by the mutated tree

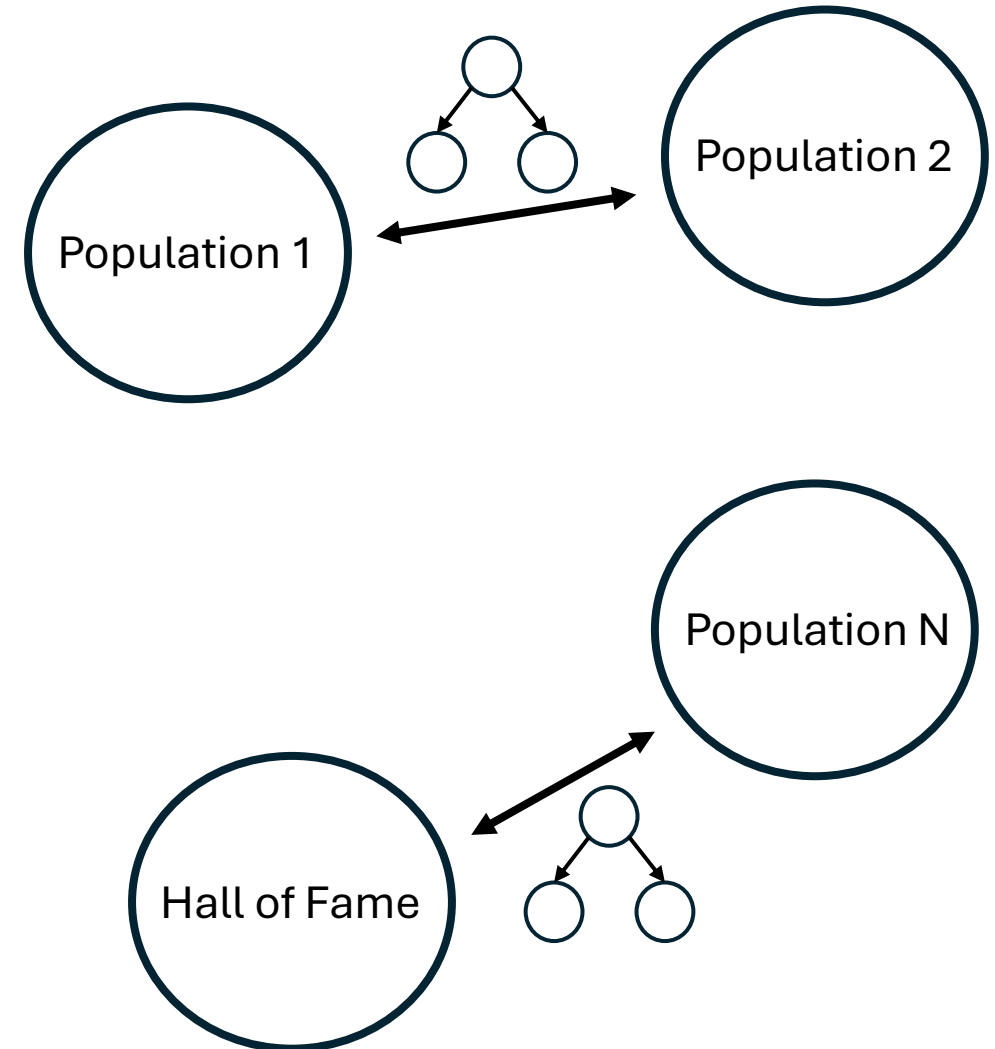


# PySR (Genetic algorithm):

[M. Cranmer '23]

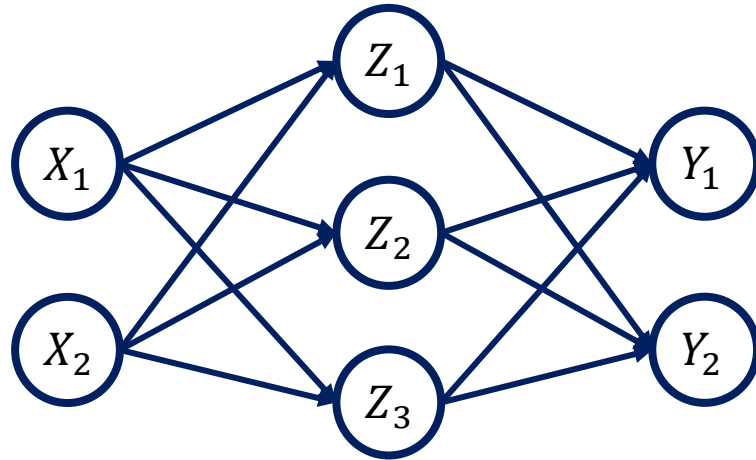
## Migration of trees via outer loop:

1. Multiple individual populations evolve independently
2. After a fixed number of mutations, trees migrate between populations
3. The trees with the highest fitness are stored in the hall of fame



# SymbolNet (DNN):

[Tsoi, Loncar, Dasu, Harris '24]

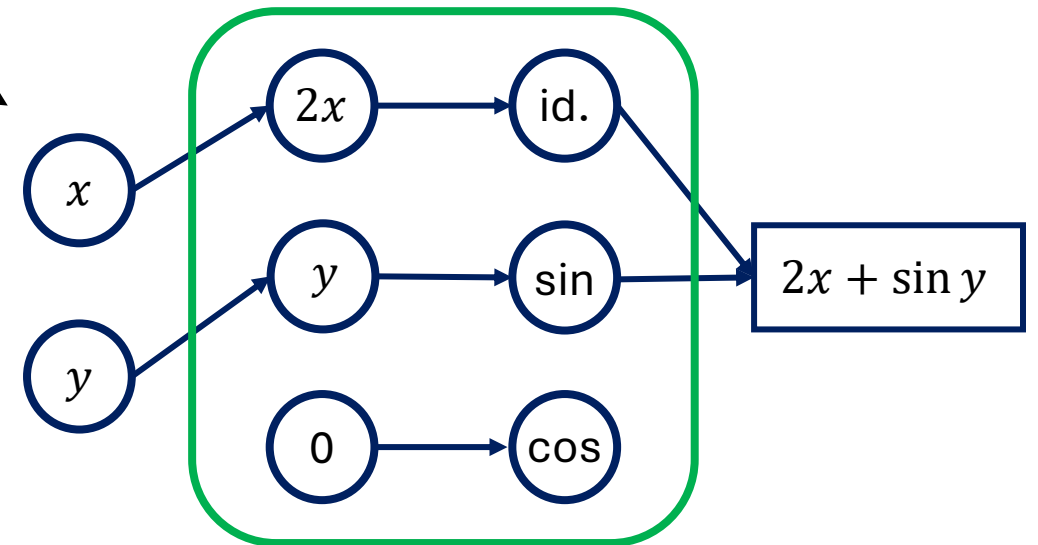


## Symbolic layers:

- Linear combination of input nodes
- Unary and binary operators are applied

Updates via backpropagation:

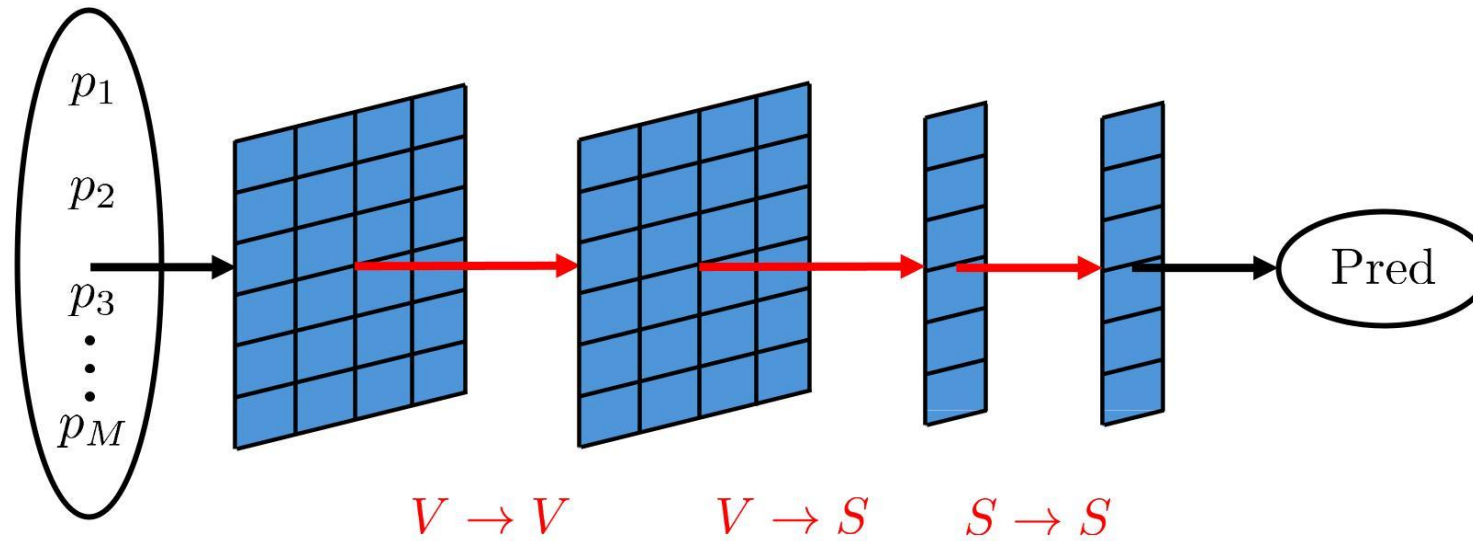
- Usual tuning of weights
- Removal of connections between nodes and operators
- Hard equation  $\rightarrow$  easy equation



# Vectorized SymbolNet:

[Bahl, Fuchs, Menen, Plehn '25]

Modification of SymbolNet to process 4-vectors:



$V \rightarrow V$ : Operators that preserve vector dimension (addition, Lorentz boost, ...)

$V \rightarrow S$ : Operators that collapse vector dimension (norm, scalar product, ...)

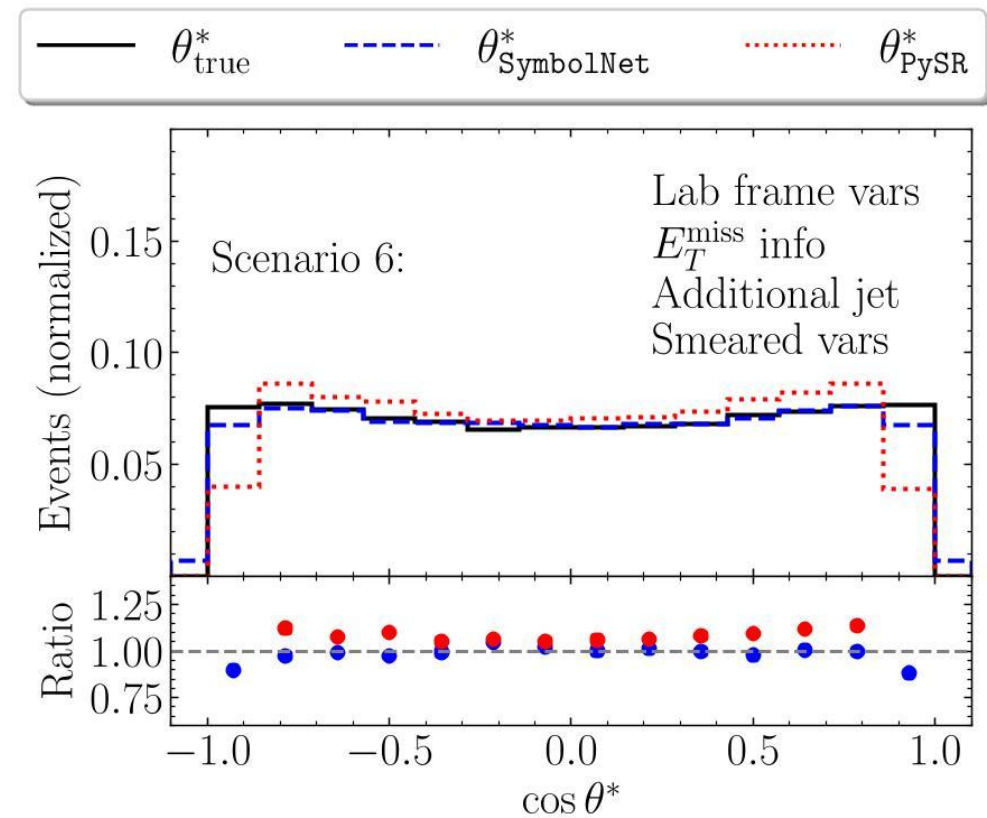
$S \rightarrow S$ : Default scalar operators



## Scenario 6:

$$\cos \theta_{\text{PySR}}^* = \sin \left[ \left( 1.114 p_{z,b} + 2.143 p_{z,\bar{l}} - 0.858 p_{z,\bar{b}} - 0.426 p_{z,q} - 1.088 p_{z,\bar{q}} \right) / \right. \\ \left. \left( E_T^{\text{miss}} E_q + E_b + E_{\bar{b}} + 1.85 E_{\bar{l}} + E_{\bar{q}} \right. \right. \\ \left. \left. - \left( p_{x,b} + p_{x,\bar{l}} \right) \left( p_{x,\bar{b}} + p_{x,q} + p_{x,\bar{q}} \right) - 0.863 - \frac{0.205 p_{z,q}}{\sqrt{E_q}} \right) \right]$$

$$\cos \theta_{\text{SymbolNet}}^* = 0.971 \left( \text{boost} \left[ -1.521 p_b - 3.066 p_{\bar{l}} + 1.421 p_{\bar{b}} + 2.218 p_q + 2.043 p_{\bar{q}} \right. \right. \\ \left. \left. - 2.845 p_b - 4.559 p_{\bar{l}} - 2.575 p_{\bar{b}} - 2.753 p_q - 2.676 p_{\bar{q}} \right] \right. \\ \left. - 0.256 \tanh \left( 0.646 p_{\bar{b}} - 0.634 p_{\bar{l}} + 0.648 p_q \right) \right) \Big|_z / \\ \left( \left\| \text{boost} \left[ -1.521 p_b - 3.066 p_{\bar{l}} + 1.421 p_{\bar{b}} + 2.218 p_q + 2.043 p_{\bar{q}} \right. \right. \right. \right. \\ \left. \left. \left. - 2.845 p_b - 4.559 p_{\bar{l}} - 2.575 p_{\bar{b}} - 2.753 p_q - 2.676 p_{\bar{q}} \right] \right\|_3 \right. \\ \left. + 0.774 \left( 0.071 \text{boost} \left[ -1.521 p_b - 3.066 p_{\bar{l}} + 1.421 p_{\bar{b}} + 2.218 p_q + 2.043 p_{\bar{q}} \right. \right. \right. \right. \\ \left. \left. \left. - 2.845 p_b - 4.559 p_{\bar{l}} - 2.575 p_{\bar{b}} - 2.753 p_q - 2.676 p_{\bar{q}} \right] \right. \right. \\ \left. \left. + \tanh \left( 0.646 p_{\bar{b}} - 0.634 p_{\bar{l}} + 0.648 p_q \right) \times \right. \right. \\ \left. \left. \left. \tanh \left( 0.646 p_{\bar{b}} - 0.634 p_{\bar{l}} + 0.648 p_q \right) \right) \right\|_3 - 0.223 \right)$$



Significances to exclude different  $\alpha_t$  with a SM measurement at  $\mathcal{L} = 300\text{fb}^{-1}$

- Total exclusion power heavily decreases towards  $\alpha_t = 0^\circ$
- Relative amount of CP information compared to parton-level drops for small  $\alpha_t$
- SR approaches outperform classical reconstruction for all  $\alpha_t$

Scenario 6:

