

# Symbolic Distillation of Neural Networks

## Executive Summary

Supervised by Miles Cranmer

Word Count: 998 words

Yihao Liu; CRSid:yl2063

June 2025

## Background and Significance

Deriving closed-form analytic expressions for complex physical systems has long been a formidable challenge: deep learning yields powerful “black-box” models with high predictive accuracy, but their internal representations provide little insight into underlying analytic laws. Modern Graph Neural Networks (GNNs) naturally mirror many-body interactions through inductive biases [1] such as permutation invariance [2], locality, and parameter sharing, providing a structured latent space suitable for interpretation. Concurrently, symbolic regression methods like PySR [3] apply evolutionary search to extract human-readable equations from data, balancing complexity against accuracy. We reproduce and extend the work [4], which integrates GNNs’ physics-aligned message representations with symbolic regression. We aim to bridge high-dimensional “black-box” prediction and interpretable deep learning in scientific research with a reproducible pipeline.

## Research Objectives and Methods

We first generate datasets by simulating four canonical 2D four-body systems, including spring, inverse-square-distance (gravity), Coulomb, and inverse-distance (fake gravity), each with disjoint trajectory splits for training and testing.

GNNs are well suited to these many-body dynamics because they treat particles as

graph nodes and pairwise interactions as edges, passing learned messages that aggregate local physical effects. In our implementation, both nodes and edges are modeled by MLPs: the node MLP embeds particle attributes (position, velocity, charge and mass), while the edge MLP computes interaction messages from neighboring nodes. The edge MLP corresponds to physical forces between particles, while the node MLP takes the sum of incoming messages together with each particle’s attributes and predicts its acceleration.

A theory from the original work [4] indicates that after model training the physical force information lie in a low-dimensional subspace, matching the true force law dimensionality, of the message activation of the edge model. In addition, the true force components form a linear transformation of the corresponding message components.

To investigate this compact representation, we apply and compare four variants of a one-step message-passing GNN:

- **Standard:** unrestricted message dimension.
- **Bottleneck:** forces messages into the same low dimension as the physical system.
- $\ell_1$  **regularization:** encourages sparsity in the message components.
- **KL regularization:** aligns messages to a standard Gaussian distribution.

For each variant, we extract the highest two message channels by variance and first apply linear regression against the ground-truth force components to measure how closely messages linearly align to physical laws. We then run symbolic regression, by using PySR, on those channels to recover analytic expressions.

Two extensions further explore the encoding mechanism: one tests acceleration versus raw forces, and the other examines whether individual channels specialize in distinct force directions.

## Main Results

**Linear Regression Analysis** We evaluate how well the selected message channels from each GNN variant linearly align with true force components across four 2D many-body systems. Table 1 summarizes the results:

There is an overall trend,

$$\text{Bottleneck} > \ell_1 > \text{KL} > \text{Standard},$$

System	Standard	Bottleneck	$\ell_1$	KL
Spring	0.343, 0.337	0.997, 0.998	0.833, 0.863	0.466, 0.487
Coulomb	0.011, 0.033	0.763, 0.162	0.443, 0.158	0.584, 0.350
Inverse-square	0.020, 0.043	0.668, 0.525	0.201, 0.206	0.163, 0.162
Inverse-distance	0.431, 0.425	0.409, 0.530	0.470, 0.508	0.391, 0.302

Table 1: Linear-combination fitting  $R^2$  values for the two highest-variance message channels.

confirming that stronger constraints on message dimensionality yield better linear alignment with the true forces.

We provide Figure 1 and Figure 2 to illustrate the linear transformation in the mass-spring system for the Standard and Bottleneck GNN. The Bottleneck variant shows much stronger linear relationship than the Standard variant.

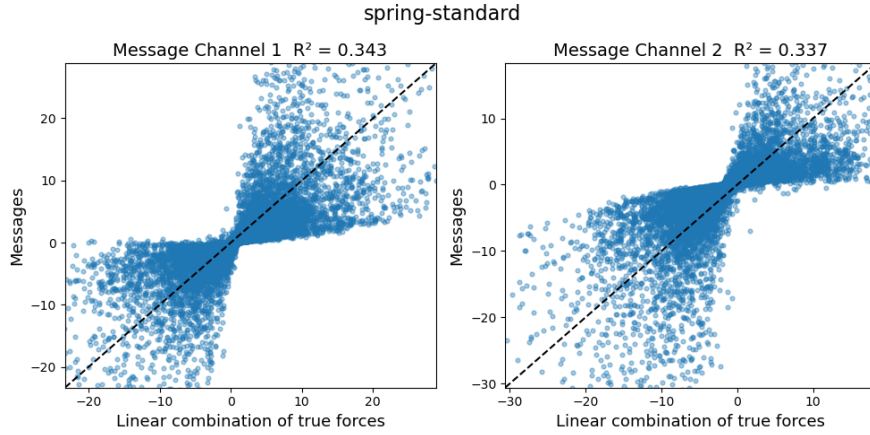


Figure 1: Spring system – Standard GNN variant. The dashed line shows  $y = x$ .

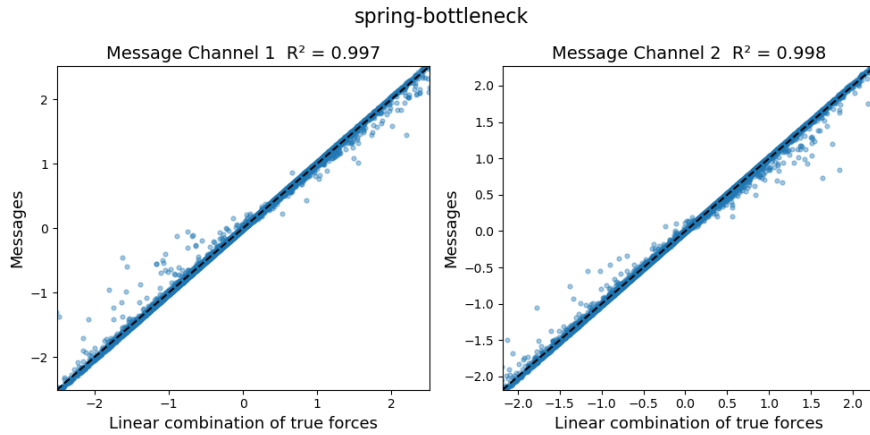


Figure 2: Spring system – Bottleneck GNN variant.

**Symbolic Regression Analysis** The symbolic regression outcomes, summarized in Table 2, largely echo the linear regression, except that the ordering of  $\ell_1$  versus KL differs.

System	Standard	Bottleneck	$\ell_1$	KL
Spring	×	✓	○★	★
Coulomb	×	✓	×	✓
Inverse-square	×	○★	×	○★
Inverse-distance	×	○	○	○

- ✓ perfect discovery matching the true law.
- ×
- partial discovery (kernel found, but missing factors and/or dimensions).
- ★ correct formula exists in the Pareto front but was not selected as “best” by PySR.
- ★ both partial and not top-ranked.

Table 2: Symbolic regression outcome markers for each system and GNN variant.

To illustrate the detailed results of using PySR to perform symbolic regression, we display the analytic equation recovered from the Spring-Bottleneck case:

$$\phi_{\text{edge}} = \left( \frac{1.2640961}{r} - 1.2788521 \right) (0.299832 \Delta x + \Delta y) \quad (1)$$

Here  $\Delta x = x_j - x_i$ ,  $\Delta y = y_j - y_i$ , and  $r = \sqrt{\Delta x^2 + \Delta y^2}$ . The small cross-term  $0.299832 \Delta y$  reflects a rotated basis in the two-dimensional space. By comparing with the standard Hooke’s law,

$$\mathbf{F}_{ij} = -k (r_{ij} - 1) \frac{\mathbf{r}_{ij}}{r_{ij}}, \quad \mathbf{r}_{ij} = (x_i - x_j, y_i - y_j), \quad r_{ij} = \|\mathbf{r}_{ij}\|. \quad (2)$$

$$F_{ij,x} = -k (r_{ij} - 1) \frac{x_i - x_j}{r_{ij}}, \quad F_{ij,y} = -k (r_{ij} - 1) \frac{y_i - y_j}{r_{ij}}, \quad (3)$$

we show that the near-perfect match confirms the rediscovery of the true Hooke’s law directly into the message channels.

Figure 3 illustrates the overall workflow of our pipeline, using the Spring-Bottleneck variant as a representative example.

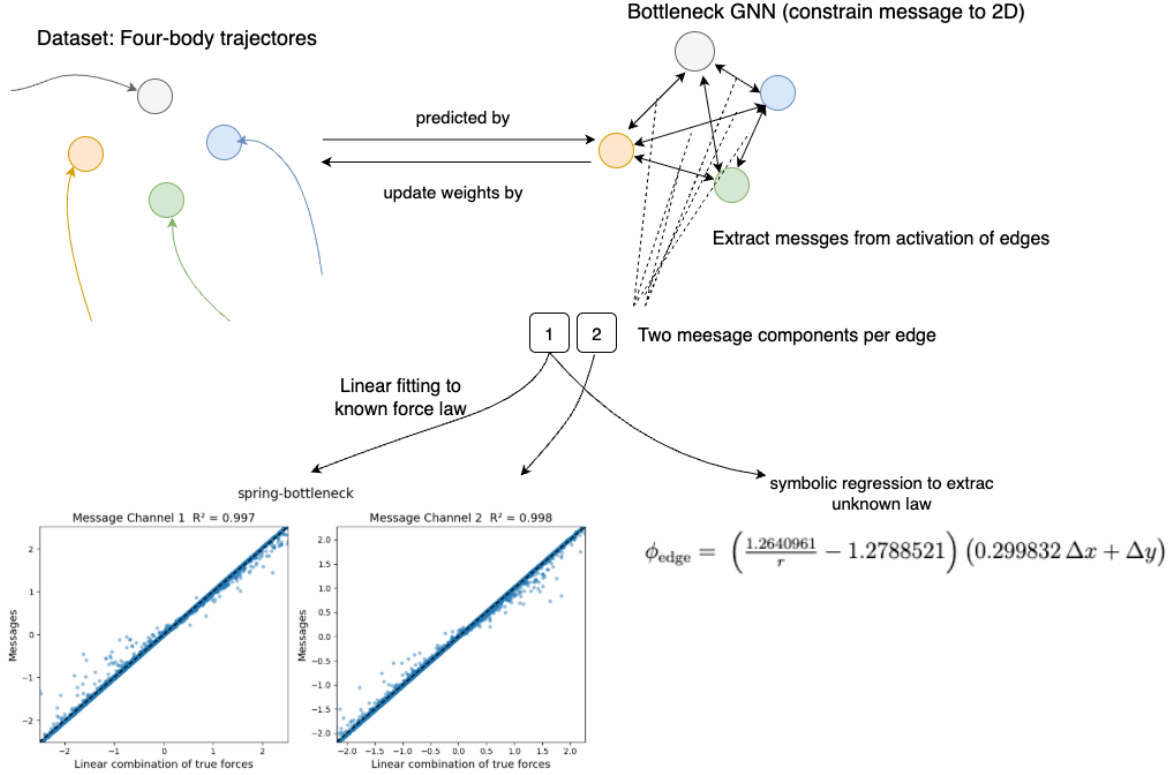


Figure 3: Overall Methodology (Spring-Bottleneck variant as a representative example).

## Extensions

**Extension i: Acceleration vs. Force Encoding** In the inverse-distance system, none of the GNN variants recovered the necessary source node mass. We hypothesize that the edge MLP messages (in this system only) align more naturally with the acceleration than with raw force. We re-run linear regression against ground-truth acceleration. The resulting  $R^2$  scores jumped dramatically as shown in Table 3, compared with Table 1. Crucially, this extension does not assert that GNNs always inherently learn acceleration over force, but demonstrates that PySR is reliable in extracting whatever physical law the model has learned.

Variant	$R^2$ (Channel 1)	$R^2$ (Channel 2)
Bottleneck	0.956	0.948
$\ell_1$	0.884	0.863
KL	0.571	0.529

Table 3: Linear regression  $R^2$  on acceleration (inverse-distance system)

**Extension ii: Channel-Wise Specialization** To assess whether message channels align uniquely to specific force components or mix multi-dimensional signals, we perform

separate linear fits of each channel to the  $x$ - and  $y$ -components of the true force. From several representative examples as shown in Table 4, we observe two patterns: some channels exhibit clear specialization, while others show mixed alignments. The encoding mechanism is therefore not fully transparent.

System	Channel	$R^2(F_x)$	$R^2(F_y)$
Spring – Bottleneck	<b>1</b>	<b>0.089</b>	<b>0.916</b>
	2	0.914	0.075
Spring – $\ell_1$ only x-dim discovered in SR	1	0.336	0.510
	2	0.835	0.023
Inverse-square – Bottleneck only x-dim discovered in SR	1	0.655	0.058
	2	0.035	0.519
Inverse-distance (acceleration) – Bottleneck	1	0.536	0.434
	2	0.620	0.315

Table 4: Component-wise channel regression  $R^2$  scores for selected cases.

Notably, for the highlighted Spring-Bottleneck case, even when a channel’s projection onto one component is weak, symbolic regression still recovers that subtle signal: the fitted analytic expression includes a small but nonzero coefficient for the weaker component (see Equation 1). This demonstrates PySR’s ability to capture low-amplitude contributions in the embedding.

## Conclusion

Our study validates that the combination of symbolic regression, inductive biases and compact message representation introduced in the original work [4] is both effective and reproducible across four canonical 2D four-body systems and four GNN variants. The consistency between linear-regression alignment and symbolic-regression success confirms that constraining messages reliably exposes true force laws. Furthermore, our two extensions deepen the analysis of constrained GNN encoding: while message channels inhabit low-dimensional subspaces matching physical laws, the exact internal mechanism remains only partially transparent. More importantly, these results underscore the robustness, sensitivity, and practical utility of PySR-based symbolic regression [3] for interpreting deep learning models for scientific research.

## Research Impact and Future Directions

Our work establishes a reliable pipeline for interpretable AI in scientific domains. This “white-box” approach empowers researchers to deploy deep learning for complex pre-

dictions while retaining analytic insight, and can be adopted for studies where interpretability is paramount.

For future researches, it will be valuable to study how GNN architectures internalize physical priors by varying system complexity and graph topology to map the limits of compact message representations. Additionally, a systematic study of the relationship between linear-regression alignment and symbolic-regression success could reveal predictive correlations.

## References

- [1] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu, “Relational inductive biases, deep learning, and graph networks,” 2018.
- [2] S. J. Prince, *Understanding Deep Learning*. The MIT Press, 2023.
- [3] M. Cranmer, “Interpretable machine learning for science with pysr and symbolicregression.jl,” 2023.
- [4] M. Cranmer, A. Sanchez-Gonzalez, P. Battaglia, R. Xu, K. Cranmer, D. Spergel, and S. Ho, “Discovering symbolic models from deep learning with inductive biases,” 2020.

## Appendix: Declaration of AI generation tools

Declaration of AI generation tools for the report part is made here:

Although some of the advise is not accepted, ChatGPT was used to help this report:

1. The format of the mathematical calculations was helped by it to make the process in a publication quality.
2. The format of plots, tables and itemized expression was helped to make things in a publication quality.
3. It suggested some alternative wording and other academic language issues throughout the whole executive summary.
4. It provided proofreading for some of text, including the description of the methodology, results and the overall logic flow among paragraphs.