

1 We thank the reviewers for giving positive and insightful evaluations of our paper. We will adjust the camera-ready
2 version to improve clarity of explanations and address all other comments. Specific responses are given below.

3 **Related work.** Thank you for introducing us to two new symbolic regression (SR) papers: Peterson (2020), and Sahoo
4 et al. (2018). We will discuss these in our paper. Note that our work is slightly orthogonal to these approaches: SR
5 algorithms like these and Eureqa, dCGP, and [34] are techniques to search an equation space with few input variables in
6 the raw symbolic regression setting. Our paper gives a way of extending any of these to high-dimensional problems, by
7 factorizing the problem into low-dimensional sub-problems corresponding to a neural network’s structure. Consider
8 attempting to fit the equivalent of both the edge model and node model, simultaneously, with any SR approach. In a
9 typical setup, the number of possible functional forms for these under a complexity upper bound of 10 tokens, and
10 only one latent variable, is approximately 10^9 equations each, but explored simultaneously, there are $(10^9)^2 = 10^{18}$
11 possibilities. However, by using our strategy to learn these models with a neural network first, and fit symbolic regression
12 to the edge and node model independently, we go from 10^{18} to 2×10^9 considered equations. This speedup grows for
13 multiple latents. Future work could replace Eureqa inside our framework with more sophisticated SR backends.

14 **Comparison of SR packages.** We will discuss a comparison of these SR packages, hyperparameters (fixed in Eureqa),
15 and why we chose Eureqa—we realize now that many at NeurIPS would find this useful. Furthermore, we will
16 demonstrate a comparison of our approach against a pure SR on one of our case studies. We have previously performed
17 this experiment and found that Eureqa performed very poorly without our framework.

18 **Udrescu & Tegmark (2019) [34].** This Eureqa alternative is optimized for rediscovering existing equations by, e.g.,
19 requiring as input the equation’s constants, which allows for dimensional analysis (a trick to recall equations in
20 physics). This approach does not seem applicable for discovering new equations so we chose Eureqa. Their description
21 “recursively break hard problems into simpler ones with fewer variables” refers to a common search strategy in SR to
22 determine feature importance (related to feature selection in decision tree learning), dissimilar to our framework.

23 **Other architectures.** Our procedure is not restricted to GNNs. However, the inductive bias of the architecture
24 determines the structure of the recovered equation. GNNs with sum-pool give the form $y_i = f(x_i, \sum_j g(x_i, x_j))$. We
25 focused on GNNs since this equation matches our problems, yet one can apply our framework to alternate inductive
26 biases. We will add a discussion on additional example architectures.

27 **Other domains.** Regarding tests for alternate domains, symbolic priors might work best for problems where we know
28 simple analytic equations already predict accurately: physics, chemistry, engineering, etc., whereas, e.g., biology and
29 behavioral sciences lack strong analytic models. Perhaps there exists a deep unknown reason for this, or maybe one
30 could use our framework to discover new equations in these domains. Regardless, we hope that our paper will stimulate
31 future research to apply our framework to new architectures and datasets.

32 **Lottery ticket hypothesis.** We thank the reviewers for making this connection to the bottleneck model results. We will
33 incorporate this in the discussion.

34 **Generalization on simulation data.** We did not test generalization performance on our simple n-body datasets because
35 our framework recovered the ground truth equation. Thus, since Newton’s law generalizes to any number of bodies,
36 this equation will give a loss of zero on any dataset with the same force law. However, in the dark matter experiment,
37 we also demonstrate better generalization (Section 4.3 - Symbolic generalization) than the trained neural network,
38 despite this simulation being entirely unlike a simple analytic equation. This simulation integrates the complex and
39 noisy (chaotic) dynamics of dark matter for a million CPU hours. Yet the analytic equation obtained via our framework,
40 describing a pattern in the output dark matter dataset, generalizes significantly better (error=0.0892) than the same
41 neural network it was extracted from (error=0.142). The question: “Where and why does a symbolic prior improve
42 generalization?” presents a very intriguing problem and motivates future research.

43 **Metrics.** We will update Table 1 to use R^2 as a metric instead of mean square error to make results more intuitive.

44 **Quantifying SR.** As suggested, we will create a table showing success/failure for each force law SR reconstruction.

45 **Table 2.** “Formula” represents the recovered analytic expression to predict overdensity. These formulas contain a
46 variable “ e_i ”, which represents a sum of some scalar function over nearby dark matter, given in the adjacent column.

47 **Generative models.** We have so far not experimented integrating generative models in our framework, but this sounds
48 like a very interesting approach, and we look forward to exploring this direction in future work.

49 **Broader impact.** We are also unaware of potential misuses, but will discuss any suggestions.

50 **Global models.** We studied GNNs with an edge model and node model (Newtonian, dark matter), a global model
51 (Hamiltonian), but not all three models together. One could approach such a GNN with the same technique: sparsify
52 latents, apply SR to each model, compose.