

---

# On the Expressiveness and Generalization of Hypergraph Neural Networks

---

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

## Abstract

1 This extended abstract describes a framework for analyzing the expressiveness,  
2 learning, and (structural) generalization of hypergraph neural networks (Hyper-  
3 GNNs). Specifically, we focus on how HyperGNNs can learn from finite datasets  
4 and generalize structurally to graph reasoning problems of arbitrary input sizes.  
5 Our first contribution is a fine-grained analysis of the expressiveness of Hyper-  
6 GNNs, that is, the set of functions that they can realize. Our result is a hierarchy  
7 of problems they can solve, defined in terms of various hyperparameters such as  
8 depths and edge arities. Next, we analyze the learning properties of these neural  
9 networks, especially focusing on how they can be trained on a finite set of small  
10 graphs and generalize to larger graphs, which we term structural generalization.  
11 Our theoretical results are further supported by the empirical results.  
12

## 13 1 Introduction

14 Reasoning over graph-structured data is an important task in many applications, including molecule  
15 analysis, social network modeling, and knowledge graph reasoning [1–3]. While we have seen  
16 great success of various relational neural networks, such as Graph Neural Networks [GNNs; 4] and  
17 Neural Logical Machines [NLM; 5] in a variety of applications [6–8], we do not yet have a full  
18 understanding of how different design parameters, such as the depth of the neural network, affects  
19 the expressiveness of these models, or how effectively these models generalize from limited data.

20 This paper analyzes the *expressiveness* and *generalization* of relational neural networks applied to  
21 *hypergraphs*, which are graphs with edges connecting more than two nodes. [Literature has shown](#)  
22 [that even when the inputs and outputs of models have only unary and binary relations, allowing](#)  
23 [intermediate hyperedge representations increases the expressiveness](#) [9, 10]. In this paper, we further  
24 formally show the “if and only if” conditions for the expressive power with respect to the edge  
25 arity. That is,  $k$ -ary hyper-graph neural networks are sufficient and necessary for realizing FOC- $k$ , a  
26 fragment of first-order logic [with counting quantification](#) which involves at most  $k$  variables. This  
27 is a helpful result because now we can determine whether a specific hypergraph neural network  
28 can solve a problem by understanding what form of logic formula can represent the solution to this  
29 problem. Next, we formally described the relationship between expressiveness and non-constant-  
30 depth networks. We state a conjecture about the “depth hierarchy,” and connect the potential proof of  
31 this conjecture to the distributed computing literature.

32 Furthermore, we prove, under certain assumptions, it is possible to train a hypergraph neural networks  
33 on a finite set of small graphs, and it will generalize to arbitrarily large graphs. This ability results  
34 from the weight-sharing nature of hypergraph neural networks. We hope our work can serve as a  
35 foundation for designing hypergraph neural networks: to solve a specific problem, what arity do you  
36 need? What depth do you need? Will my model have structural generalization (i.e., to larger graphs)?  
37 Our theoretical results are further supported by experiments, for empirical demonstrations.

## 38 2 Hypergraph Reasoning Problems and Hypergraph Neural Networks

39 A *hypergraph representation*  $G$  is a tuple  $(V, X)$ , where  $V$  is a set of entities (nodes), and  $X$  is  
40 a set of *hypergraph representation functions*. Specifically,  $X = \{X_0, X_1, X_2, \dots, X_k\}$ , where  
41  $X_j : (v_1, v_2, \dots, v_j) \rightarrow \mathcal{S}$  is a function mapping every tuple of  $j$  nodes to a value. We call  $j$  the  
42 *arity* of the hyperedge and  $k$  is the max arity of input hyperedges. The range  $\mathcal{S}$  can be any set of

43 discrete labels that describes relation type, or a scalar number (e.g., the length of an edge), or a vector.  
 44 We will use the arity 0 representation  $X_0(\emptyset) \rightarrow \mathcal{S}$  to represent any global properties of the graph.

45 A *graph reasoning function*  $f$  is a mapping from a hypergraph representation  $G = (V, X)$  to another  
 46 hyperedge representation function  $Y$  on  $V$ . As concrete examples, asking whether a graph is fully  
 47 connected is a graph classification problem, where the output  $Y = \{Y_0\}$  and  $Y_0(\emptyset) \rightarrow \mathcal{S}' = \{0, 1\}$  is  
 48 a global label; finding the set of disconnected subgraphs of size  $k$  is a  $k$ -ary hyperedge classification  
 49 problem, where the output  $Y = \{Y_k\}$  is a label for each  $k$ -ary hyperedges.

50 There are two main motivations and constructions of a neural network applied to graph reasoning  
 51 problems: message-passing-based and first-order-logic-inspired. Both approaches construct the  
 52 computation graph layer by layer. The input is the features of nodes and hyperedges, while the output  
 53 is the per-node or per-edge prediction of desired properties, depending on the task.

54 In a nutshell, within each layer, *message-passing-based* hypergraph neural networks, Higher-  
 55 Order GNNs [11], perform message passing between each hyperedge and its neighbours. Specif-  
 56 ically, we say the  $j$ -th neighbour set of a hyperedge  $u = (x_1, x_2, \dots, x_i)$  of arity  $i$  is  $N_j(u) =$   
 57  $\{(x_1, x_2, \dots, x_{j-1}, r, x_{j+1}, \dots, x_i)\}$ , where  $r \in V$ . Then, the all neighbours of node  $u$  is the union  
 58 of all  $N_j$ 's, where  $j = 1, 2, \dots, i$ .

59 On the other hand, first-order-logic-inspired hypergraph neural networks consider building neural  
 60 networks that can emulate first logic formulas. Neural Logic Machines [NLM; 5] are defined in  
 61 terms of a set of input hyperedges; each hyperedge of arity  $k$  is represented by a vector of (possibly  
 62 real) values obtained by applying all of the  $k$ -ary predicates in the domain to the tuple of vertices it  
 63 connects. Each layer in an NLM learns to apply a linear transformation with nonlinear activation and  
 64 quantification operators (analogous to the for all  $\forall$  and exists  $\exists$  quantifiers in first-order logic), on these  
 65 values. It is easy to prove, by construction, that given a sufficient number of layers and maximum  
 66 arity, NLMs can learn to realize any first-order-logic formula. For readers who are not familiar with  
 67 HO-GNNs [11] and NLMs [5], we include a mathematical summary of their computation graph in  
 68 Appendix B. Our analysis starts from the following theorem.

69 **Theorem 2.1.** HO-GNNs [11] are equivalent to NLMs in terms of expressiveness. Specifically, a  $B$ -  
 70 ary HO-GNN is equivalent to an NLM applied to  $B + 1$ -ary hyperedges. Proofs are in Appendix B.3.  
 71 Given Theorem 2.1, we can focus on just one single type of hypergraph neural network. Specifically,  
 72 we will focus on Neural Logic Machines [NLM; 5] because its architecture naturally aligns with  
 73 first-order logic formula structures, which will aid some of our analysis. An NLM is characterized  
 74 by hyperparameters  $D$  (depth), and  $B$  maximum arity. We are going to assume that  $B$  is a constant,  
 75 but  $D$  can be dependent on the size of the input graph. We will use  $\text{NLM}[D, B]$  to denote an NLM  
 76 family with depth  $D$  and max arity  $B$ . Other parameters such as the width of neural networks affects  
 77 the precise details of what functions can be realized, as it does in a regular neural network, but does  
 78 not affect the analyses in this extended abstract. Furthermore, we will be focusing on neural networks  
 79 with bounded precision, and briefly discuss how our results generalize to unbounded precision cases.  
 80

### 81 3 Expressiveness of Relational Neural Networks

82 We start from a formal definition of hypergraph neural network expressiveness.

83 **Definition 3.1** (Expressiveness). We say a model family  $\mathcal{M}_1$  is *at least expressive as*  $\mathcal{M}_2$ , written  
 84 as  $\mathcal{M}_1 \succcurlyeq \mathcal{M}_2$ , if for all  $M_2 \in \mathcal{M}_2$ , there exists  $M_1 \in \mathcal{M}_1$  such that  $M_1$  can realize  $M_2$ . A model  
 85 family  $\mathcal{M}_1$  is *more expressive than*  $\mathcal{M}_2$ , written as  $\mathcal{M}_1 \succ \mathcal{M}_2$ , if  $\mathcal{M}_1 \succcurlyeq \mathcal{M}_2$  and  $\exists M_1 \in \mathcal{M}_1$ ,  
 86  $\forall M_2 \in \mathcal{M}_2$ ,  $M_2$  can not realize  $M_1$ .

87 **Arity Hierarchy** We first aim to quantify how the maximum arity  $B$  of the network's representation  
 88 affects its expressiveness and find that, in short, even if the inputs and outputs of neural networks are  
 89 of low arity, the higher the maximum arity for intermediate layers, the more expressive the NLM is.

90 **Corollary 3.1** (Arity Hierarchy). For any maximum arity  $B$ , there exists a depth  $D^*$  such that:  
 91  $\forall D \geq D^*$ ,  $\text{NLM}[D, B + 1]$  is more expressive than  $\text{NLM}[D, B]$ . This theorem applies to both  
 92 fixed-precision\* and unbounded-precision networks.

93 *Proof sketch:* Our proof slightly extends the proof of Morris et al. [11]. First, the set of graphs distin-  
 94 guishable by  $\text{NLM}[D, B]$  is bounded by graphs distinguishable by a  $D$ -round order- $B$  Weisfeiler-  
 95 Leman test [12]. If models in  $\text{NLM}[D, B]$  cannot generate different outputs for two distinct

\*Fixed precision: the results of intermediate layers (tensors) are constant-sized (e.g.,  $W$  bits per entry).  
 Practical GNNs are all fixed-precision because real number types in modern computers have finite precision.

96 hypergraphs  $G_1$  and  $G_2$ , but there exists  $M \in \text{NLM}[D, B + 1]$  that *can* generate different outputs  
 97 for  $G_1$  and  $G_2$ , then we can construct a graph classification function  $f$  that  $\text{NLM}[D, B + 1]$  (with  
 98 some fixed precision) can realize but  $\text{NLM}[D, B]$  (even with unbounded precision) cannot.<sup>†</sup> The full  
 99 proof is described in Appendix C.1.

100 It is also important to quantify the minimum arity for realizing certain graph reasoning functions.

101 **Corollary 3.2** (FOL realization bounds). Let  $\text{FOC}_B$  denote a fragment of first order logic with at  
 102 most  $B$  variables, extended with counting quantifiers of the form  $\exists^{\geq n} \phi$ , which state that there are at  
 103 least  $n$  nodes satisfying formula  $\phi$  [13].

- 104 • (Upper Bound) Any function  $f$  in  $\text{FOC}_B$  can be realized by  $\text{NLM}[D, B]$  for some  $D$ .
- 105 • (Lower Bound) There exists a function  $f \in \text{FOC}_B$  such that for all  $D$ ,  $f$  cannot be realized by  
 106  $\text{NLM}[D, B - 1]$ .

107 *Proof:* The upper bound part of the claim has been proved by Barceló et al. [14] for  $B = 2$ . The  
 108 results generalize easily to arbitrary  $B$  because the counting quantifiers can be realized by sum  
 109 aggregation. The lower bound part can be proved by applying Section 5 of [13], in which they show  
 110 that  $\text{FOC}_B$  is equivalent to a  $(B - 1)$ -dimensional WL test in distinguishing non-isomorphic graphs.  
 111 Given that  $\text{NLM}[D, B - 1]$  is equivalent to the  $(B - 2)$ -dimensional WL test of graph isomorphism,  
 112 there must be an  $\text{FOL}_B$  formula that distinguishes two non-isomorphic graphs that  $\text{NLM}[D, B - 1]$   
 113 cannot. Hence,  $\text{FOL}_B$  cannot be realized by  $\text{NLM}[\cdot, B - 1]$ .

114 **Depth Hierarchy** We now study the dependence of the expressiveness of NLMs on depth  $D$ . **Neural**  
 115 **networks are generally defined to have a fixed depth, but allowing them to have a depth that is**  
 116 **dependent on the number of nodes  $n = |V|$  in the graph, in many cases, can substantially increase**  
 117 **their expressive power [15, see also Theorem 3.4 and Appendix C for examples].** In the following, we  
 118 define a *depth hierarchy* by analogy to the *time hierarchy* in computational complexity theory [16],  
 119 and we extend our notation to let  $\text{NLM}[O(f(n)), B]$  denote the class of adaptive-depth NLMs in  
 120 which the growth-rate of depth  $D$  is bounded by  $O(f(n))$ .

121 **Conjecture 3.3** (Depth hierarchy). For any maximum arity  $B$ , for any two functions  $f$  and  $g$ , if  
 122  $g(n) = o(f(n)/\log n)$ , that is,  $f$  grows logarithmically more quickly than  $g$ , then fixed-precision  
 123  $\text{NLM}[O(f(n)), B]$  is more expressive than fixed-precision  $\text{NLM}[O(g(n)), B]$ .

124 There is a closely related result for the *congested clique* model in distributed computing, where [17]  
 125 proved that  $\text{CLIQUE}(g(n)) \subsetneq \text{CLIQUE}(f(n))$  if  $g(n) = o(f(n))$ . This result does not have the  
 126  $\log n$  gap because the congested clique model allows  $\log n$  bits to transmit between nodes at each  
 127 iteration, while fixed-precision NLM allows only a constant number of bits. The reason why the result  
 128 on congested clique can not be applied to fixed-precision NLMs is that congested clique assumes  
 129 unbounded precision representation for each individual node.

130 However, Conjecture 3.3 is not true for NLMs with unbounded precision, because there is an upper  
 131 bound depth  $O(n^{B-1})$  for a model’s expressiveness power.<sup>‡</sup> That is, an unbounded-precision NLM  
 132 can not achieve stronger expressiveness by increasing its depth beyond  $O(n^{B-1})$ .

133 It is important to point out that, to realize a specific graph reasoning function, NLMs with different  
 134 maximum arity  $B$  may require different depth  $D$ . Fürer [18] provides a general construction for  
 135 problems that higher-dimensional NLMs can solve in asymptotically smaller depth than lower-  
 136 dimensional NLMs. In the following we give a concrete example for computing *S-T Connectivity- $k$* ,  
 137 which asks whether there is a path of nodes from  $S$  and  $T$  in a graph, with length  $\leq k$ .

138 **Theorem 3.4** (S-T Connectivity- $k$  with Different Max Arity). For any function  $f(k)$ , if  $f(k) = o(k)$ ,  
 139  $\text{NLM}[O(f(k)), 2]$  cannot realize S-T Connectivity- $k$ . That is, S-T Connectivity- $k$  requires depth  
 140 at least  $O(k)$  for a relational neural network with an maximum arity of  $B = 2$ . However, S-T  
 141 Connectivity- $k$  can be realized by  $\text{NLM}[O(\log k), 3]$ .

142 *Proof sketch.* For any integer  $k$ , we can construct a graph with two chains of length  $k$ , so that if we  
 143 mark two of the four ends as  $S$  or  $T$ , any  $\text{NLM}[k - 1, 2]$  cannot tell whether  $S$  and  $T$  are on the same  
 144 chain. The full proof is described in Appendix C.3.

145 There are many important graph reasoning tasks that do not have known depth lower bounds,  
 146 including all-pair connectivity and shortest distance [19, 20]. In Appendix C.3, we discuss the  
 147 concrete complexity bounds for a series of graph reasoning problems.

<sup>†</sup>Note that the arity hierarchy is applied to fixed-precision and unbounded-precision separately. For example,  
 $\text{NLM}[D, B]$  with unbounded precision is incomparable with  $\text{NLM}[D, B + 1]$  with fixed precision.

<sup>‡</sup>See appendix C.2 for a formal statement and the proof.

## 148 4 Learning and Generalization in Relational Neural Networks

149 Given our understanding of what functions can be realized by NLMs, we move on to the problems of  
 150 learning them: Can we effectively learn a NLMs to solve a desired task given a sufficient number of  
 151 input-output examples? In this paper, we show that applying *enumerative training* with examples  
 152 up to some fixed graph size can ensure that the trained neural network will generalize to all graphs  
 153 *larger* than those appearing in the training set.

154 A critical determinant of the generalization ability for NLMs is the aggregation function. Specifically,  
 155 Xu et al. [21] have shown that using *sum* as the aggregation function provides maximum expressive-  
 156 ness for graph neural networks. However, sum aggregation cannot be implemented in fixed-precision  
 157 models, because as the graph size  $n$  increases, the range of the sum aggregation also increases.

158 **Definition 4.1** (Fixed-precision aggregation function). An aggregation function is *fixed precision*  
 159 if it maps from any finite *set* of inputs with values drawn from *finite domains* to a *fixed finite* set  
 160 of possible output values; that is, the cardinality of the range of the function cannot grow with the  
 161 number of elements in the input set. Two useful fixed-precision aggregation functions are *max*, which  
 162 computes the dimension-wise maximum over the set of input values, and *fixed-precision mean*, which  
 163 approximates the dimension-wise mean to a fixed decimal place.

164 In order to focus on structural generalization in this section, we consider an *enumerative* training  
 165 paradigm. When the input hypergraph representation domain  $\mathcal{S}$  is a finite set, we can enumerate the  
 166 set  $\mathcal{G}_{\leq N}$  of all possible input hypergraph representations of size bounded by  $N$ . We first enumerate  
 167 all graph sizes  $n \leq N$ ; for each  $n$ , we enumerate all possible values assigned to the hyperedges  
 168 in the input. Given training size  $N$ , we enumerate all inputs in  $\mathcal{G}_{\leq N}$ , associate with each one the  
 169 corresponding ground-truth output representation, and train the model with these input-output pairs.

170 This has much stronger data requirements than the standard sampling-based training mechanisms in  
 171 machine learning. In practice, this can be approximated well when the input domain  $\mathcal{S}$  is small and  
 172 the input data distribution is approximately uniformly distributed. The enumerative learning setting  
 173 is studied by the *language identification in the limit* community [22], in which it is called *complete*  
 174 *presentation*. This is an interesting learning setting because even if the domain for each individual  
 175 hyperedge representation is finite, as the graph size can go arbitrarily large, the number of possible  
 176 inputs is enumerable but unbounded.

177 **Theorem 4.1** (Fixed-precision generalization under complete presentation). For any hypergraph  
 178 reasoning function  $f$ , if it can be realized by a fixed-precision relational neural network model  $\mathcal{M}$ ,  
 179 then there exists an integer  $N$ , such that if we train the model with complete presentation on all input  
 180 hypergraph representations with size smaller than  $N$ ,  $\mathcal{G}_{\leq N}$ , then for all  $M \in \mathcal{M}$ ,

$$\sum_{G \in \mathcal{G}_{\leq N}} 1[M(G) \neq f(G)] = 0 \implies \forall G \in \mathcal{G}_{\infty} : M(G) = f(G).$$

181 That is, as long as  $M$  fits all training examples, it will generalize to all possible hypergraphs in  $\mathcal{G}_{\infty}$ .

182 *Proof.* The key observation is that for any fixed vector representation length  $W$ , there are only a finite  
 183 number of distinctive models in a fixed-precision NLM family, *independent of the graph size  $n$* . Let  
 184  $W_b$  be the number of bits in each intermediate representation of a fixed-precision NLM. There are  
 185 at most  $(2^{W_b})^{2^{W_b}}$  different mappings from inputs to outputs. Hence, if  $N$  is sufficiently large to  
 186 enumerate all input hypergraphs, we can always identify the correct model in the hypothesis space.

187 Our results are related to the *algorithmic alignment* approach [23, 24]. In contrast to their Probably  
 188 Approximately Correct (PAC) Learning bounds for sample efficiency, our expressiveness results  
 189 directly quantifies whether a hypergraph neural network can be trained to realize a specific function.

## 190 5 Conclusion

191 In this extended abstract, we have shown the substantial increase of expressive power due to higher-  
 192 arity relations and increasing depth, and have characterized very powerful structural generalization  
 193 from training on small graphs to performance on larger ones. We further discuss the relationship  
 194 between these results and existing results in Appendix A. All theoretical results are further supported  
 195 by the empirical results, discussed in Appendix D. Although many questions remain open about the  
 196 overall generalization capacity of these models in continuous and noisy domains, we believe this  
 197 work has shed some light on their utility and potential for application in a variety of problems.

## References

- 198
- 199 [1] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural  
200 message passing for quantum chemistry. In *ICML*, 2017. 1
- 201 [2] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and  
202 Max Welling. Modeling relational data with graph convolutional networks. In *ESWC*, 2018.
- 203 [3] Qi Liu, Maximilian Nickel, and Douwe Kiela. Hyperbolic graph neural networks. *IEEE*  
204 *Transactions on Knowledge and Data Engineering*, 2017. 1
- 205 [4] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.  
206 The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008. 1
- 207 [5] Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural  
208 logic machines. In *ICLR*, 2019. 1, 2, 7, 8, 16
- 209 [6] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius  
210 Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan  
211 Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint*  
212 *arXiv:1806.01261*, 2018. 1
- 213 [7] Christian Merkwirth and Thomas Lengauer. Automatic generation of complementary descriptors  
214 with molecular graph networks. *Journal of chemical information and modeling*, 45(5):1159–  
215 1168, 2005.
- 216 [8] Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural  
217 execution of graph algorithms. In *ICLR*, 2020. 1
- 218 [9] Waiss Azizian and Marc Lelarge. Expressive power of invariant and equivariant graph neural  
219 networks. In *ICLR*, 2021. 1, 7
- 220 [10] Cristian Bodnar, Fabrizio Frasca, Nina Otter, Yuguang Wang, Pietro Lio, Guido F Montufar,  
221 and Michael Bronstein. Weisfeiler and lehman go cellular: Cw networks. In *NeurIPS*, 2021. 1
- 222 [11] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,  
223 Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural  
224 networks. In *AAAI*, 2019. 2, 7, 12
- 225 [12] AA Leman and B Weisfeiler. A reduction of a graph to a canonical form and an algebra arising  
226 during this reduction. *Nauchno-Tekhnicheskaya Informatsiya*, 2(9):12–16, 1968. 2
- 227 [13] Jin-Yi Cai, Martin Fürer, and Neil Immerman. An optimal lower bound on the number of  
228 variables for graph identification. *Combinatorica*, 12(4):389–410, 1992. 3, 12, 13
- 229 [14] Pablo Barceló, Egor V Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, and Juan Pablo Silva.  
230 The logical expressiveness of graph neural networks. In *ICLR*, 2020. 3, 7, 13
- 231 [15] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Uni-  
232 versal transformers. In *ICLR*, 2019. 3, 7
- 233 [16] Juris Hartmanis and Richard E Stearns. On the computational complexity of algorithms.  
234 *Transactions of the American Mathematical Society*, 117:285–306, 1965. 3
- 235 [17] Janne H Korhonen and Jukka Suomela. Towards a complexity theory for the congested clique.  
236 In *SPAA*, 2018. 3
- 237 [18] Martin Fürer. Weisfeiler-lehman refinement requires at least a linear number of iterations. In  
238 *ICALP*, 2001. 3
- 239 [19] Mauricio Karchmer and Avi Wigderson. Monotone circuits for connectivity require super-  
240 logarithmic depth. *SIAM J. Discrete Math.*, 3(2):255–265, 1990. 3, 7
- 241 [20] Shreyas Pai and Sriram V Pemmaraju. Connectivity lower bounds in broadcast congested clique.  
242 In *PODC*, 2019. 3
- 243 [21] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
244 networks? In *ICLR*, 2019. 4, 7
- 245 [22] E Mark Gold. Language identification in the limit. *Inf. Control.*, 10(5):447–474, 1967. 4
- 246 [23] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie  
247 Jegelka. What can neural networks reason about? In *ICLR*, 2020. 4, 7

- 248 [24] Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie  
 249 Jegelka. How neural networks extrapolate: From feedforward to graph neural networks. In  
 250 *ICLR*, 2021. 4, 7
- 251 [25] Yuan Li, Alexander Razborov, and Benjamin Rossman. On the ac<sup>0</sup> complexity of subgraph  
 252 isomorphism. *SIAM J. Comput.*, 46(3):936–971, 2017. 7
- 253 [26] Benjamin Rossman. *Average-case complexity of detecting cliques*. PhD thesis, Massachusetts  
 254 Institute of Technology, 2010. 7
- 255 [27] Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International  
 256 Conference on Learning Representations*, 2020. 7
- 257 [28] Davide Buffelli, Pietro Liò, and Fabio Vandin. Sizeshiftreg: a regularization method for  
 258 improving size-generalization in graph neural networks. In *NeurIPS*, 2022. 7
- 259 [29] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural  
 260 networks. In *AAAI*, 2019. 11
- 261 [30] Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha  
 262 Talukdar. Hypergcn: A new method of training graph convolutional networks on hypergraphs.  
 263 In *NeurIPS*, 2019.
- 264 [31] Song Bai, Feihu Zhang, and Philip HS Torr. Hypergraph convolution and hypergraph attention.  
 265 *Pattern Recognition*, 110:107637, 2021. 11
- 266 [32] Kaize Ding, Jianling Wang, Jundong Li, Dingcheng Li, and Huan Liu. Be more with less:  
 267 Hypergraph attention networks for inductive text classification. In *EMNLP*, 2020. 11
- 268 [33] Jing Huang and Jie Yang. Unignn: a unified framework for graph and hypergraph neural  
 269 networks. In *IJCAI*, 2021. 11
- 270 [34] Sandra Kiefer and Pascal Schweitzer. Upper bounds on the quantifier depth for graph differenti-  
 271 ation in first order logic. In *LICS*, 2016. 12
- 272 [35] Efim A Dinic. Algorithm for solution of a problem of maximum flow in networks with power  
 273 estimation. In *Soviet Math. Doklady*, volume 11, pages 1277–1280, 1970. 15

274

## Appendix

275 The appendix is organized as the following. In Appendix A, we discuss the related work. In  
 276 Appendix B, we provide a formalization of two types of hypergraph neural networks discussed in  
 277 the main paper, and proved their equivalence. In Appendix C, we prove the theorems for the arity  
 278 hierarchy and provide concrete examples for expressiveness analyses. Finally, in Appendix D, we  
 279 include additional experiment results to empirically illustrate the application of theorems discussed  
 280 in the paper.

### 281 A Related Work

282 Solving problems on graphs of arbitrary size is studied in many fields. NLMs can be viewed as circuit  
 283 families with constrained architecture. In distributed computation, the congested clique model can be  
 284 viewed as 2-arity NLMs, where nodes have identities as extra information. Common graph problems  
 285 including sub-structure detection[25, 26] and connectivity[19] are studied for lower bounds in terms  
 286 of depth, width and communication. This has been connected to GNNs for deriving expressiveness  
 287 bounds [27].

288 Studies have been conducted on the expressiveness of GNNs and their variants. Xu et al. [21] provide  
 289 an illuminating characterization of GNN expressiveness in terms of the WL graph isomorphism test.  
 290 [Azizian and Lelarge \[9\] analyze the expressiveness of higher-order Folklore GNNs by connecting them  
 291 with high-dimensional WL-tests. We have the similar results in the arity hierarchy.](#) Barceló et al. [14]  
 292 reviewed GNNs from the logical perspective and rigorously refined their logical expressiveness with  
 293 respect to fragments of first-order logic. Dong et al. [5] proposed Neural Logical Machines (NLMs)  
 294 to reason about higher-order relations, and showed that increasing order increases expressiveness. It is  
 295 also possible to gain expressiveness using unbounded computation time, as shown by the work of  
 296 Dehghani et al. [15] on dynamic halting in transformers.

297 It is interesting that GNNs may generalize to larger graphs. Xu et al. [23, 24] have studied the  
 298 notion of *algorithmic alignment* to quantify such structural generalization. Dong et al. [5] provided  
 299 empirical results showing that NLMs generalize to much larger graphs on certain tasks. [Buffelli  
 300 et al. \[28\] introduced a regularization technique to improve GNNs' generalization to larger graphs  
 301 and demonstrated its effectiveness empirically.](#) In Xu et al. [23], they analyzed and compared the  
 302 sample complexity of Graph Neural Networks. This is different from our notion of expressiveness  
 303 for realizing functions. In Xu et al. [24], they showed empirically on some problems (e.g., Max-  
 304 Degree, Shortest Path, and n-body problem) that algorithm alignment helps GNNs to extrapolate,  
 305 and theoretically proved the improvement by algorithm alignment on the Max-Degree problem. In  
 306 this extended abstract, instead of focusing on computing specific graph problems, we analyzed how  
 307 GNNs can extrapolate to larger graphs in a general case, based on the assumption of fixed precision  
 308 computation.

### 309 B Hypergraph Neural Networks

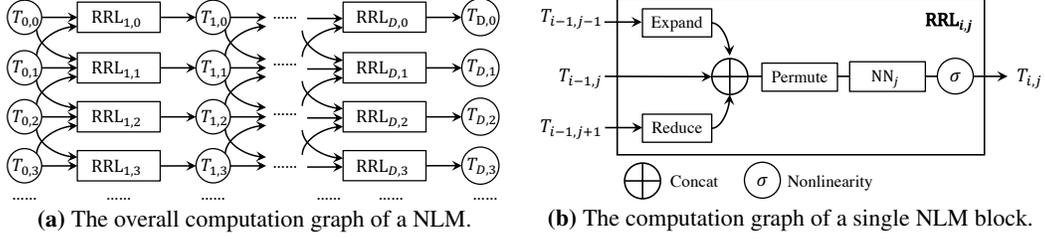
310 We now introduce two important hypergraph neural network implementations that can be trained to  
 311 solve graph reasoning problems: Higher-order Graph Neural Networks [HO-GNN; 11] and Neural  
 312 Logic Machines [NLM; 5]. They are equivalent to each other in terms of expressiveness. Showing  
 313 this equivalence allows us to focus the rest of the paper on analyzing a single model type, with the  
 314 understanding that the conclusions generalize to a broader class of hypergraph neural networks.

#### 315 B.1 Higher-order Graph Neural Networks

316 Higher-order Graph Neural Networks [HO-GNNs; 11] are Graph Neural Networks (GNNs) that  
 317 apply to hypergraphs. A GNN is usually defined based on two message passing operations.

- 318 • Edge update: the feature of each edge is updated by features of its ends.
- 319 • Node update: the feature of each node is updated by features of all edges adjacent to it.

320 However, computing only node-wise and edge-wise features does not handle higher-order relations,  
 321 such as triangles in the graph. In order to obtain more expressive power, GNNs have to extend to  
 322 hypergraphs of higher arity [11]. Specifically, HO-GNNs on  $B$ -ary hypergraph maintains features



**Figure 1:** The overall architecture of our Neural Logic Machines (NLMs). It follows the computation graph of NLM [5] and can be applied to hypergraphs.

323 for all  $B$ -tuple of nodes, and the neighborhood is extended to  $B$ -tuples accordingly: the feature of  
 324 tuple  $(v_1, v_2, \dots, v_B)$  is updated by the  $|V|$  element multiset (contain  $|V|$  elements for each  $u \in V$ )  
 325 of  $B$ -tuples of features

$$(H_i[u, v_2, \dots, v_B], H_{i-1}[v_1, u, v_2, \dots, v_B], \dots, H_{i-1}[v_1, \dots, v_{B-1}, u]) \quad (\text{B.1})$$

326 where  $H_{i-1}[\mathbf{v}]$  is the feature of tuple  $\mathbf{v}$  from the previous iteration.

327 We now introduce the formal definition of the high-dimensional message passing. We denote  $\mathbf{v}$   
 328 as a  $B$ -tuple of nodes  $(v_1, v_2, \dots, v_B)$ , and generalize the neighborhood to a higher dimension by  
 329 defining the neighborhood of  $\mathbf{v}$  as all node tuples that differ from  $\mathbf{v}$  at one position.

$$\text{Neighbors}(\mathbf{v}, u) = ((u, v_2, \dots, v_B), (v_1, u, v_3, \dots, v_B), \dots, (v_1, \dots, v_{B-1}, u)) \quad (\text{B.2})$$

$$N(\mathbf{v}) = \{\text{Neighbors}(\mathbf{v}, u) | u \in V\} \quad (\text{B.3})$$

330 Then message passing scheme naturally generalizes to high-dimensional features using the high-  
 331 dimensional neighborhood.

$$\text{Received}_i[\mathbf{v}] = \sum_u (\text{NN}_1 (H_{i-1}[\mathbf{v}]; \text{CONCAT}_{\mathbf{v}' \in \text{neighbors}(\mathbf{v}, u)} H_{i-1}[\mathbf{v}'])) \quad (\text{B.4})$$

## 332 B.2 Neural Logic Machines

333 A NLM is a multi-layer neural network that operates on hypergraph representations, in which  
 334 the hypergraph representation functions are represented as tensors. The input is a hypergraph  
 335 representation  $(V, X)$ . There are then several computational layers, each of which produces a  
 336 hypergraph representation with nodes  $V$  and a new set of representation functions. Specifically, a  
 337  $B$ -ary NLM produces hypergraph representation functions with arities from 0 up to a maximum  
 338 hyperedge arity of  $B$ . We let  $T_{i,j}$  denote the tensor representation for the output at layer  $i$  and arity  
 339  $j$ . Each entry in the tensor is a mapping from a set of node indices  $(v_1, v_2, \dots, v_j)$  to a vector in a  
 340 latent space  $\mathbb{R}^W$ . Thus,  $T_{i,j}$  is a tensor of  $j+1$  dimensions, with the first  $j$  dimensions corresponding  
 341 to  $j$ -tuple of nodes, and the last feature dimension. For convenience, we write  $h_{0,\cdot}$  for the input  
 342 hypergraph representation and  $h_{D,\cdot}$  for the output of the NLM.

343 Fig. 1a shows the overall architecture of NLMs. It has  $D \times B$  computation blocks, namely relational  
 344 reasoning layers (RRLs). Each block  $RRL_{i,j}$ , illustrated in Fig. 1b, takes the output from neighboring  
 345 arities in the previous layer,  $T_{i-1,j-1}$ ,  $T_{i-1,j}$  and  $T_{i-1,j+1}$ , and produces  $T_{i,j}$ . Below we show the  
 346 computation of each primitive operation in an RRL.

347 The *expand* operation takes tensor  $T_{i-1,j-1}$  (arity  $j-1$ ) and produces a new tensor  $T_{i-1,j-1}^E$  of arity  
 348  $j$ . The *reduce* operation takes tensor  $T_{i-1,j+1}$  (arity  $j+1$ ) and produces a new tensor  $T_{i-1,j+1}^R$  of  
 349 arity  $j+1$ . Mathematically,

$$\begin{aligned} T_{i-1,j-1}^E[v_1, v_2, \dots, v_j] &= T_{i-1,j-1}[v_1, v_2, \dots, v_{j-1}]; \\ T_{i-1,j+1}^R[v_1, v_2, \dots, v_j] &= \text{Agg}_{v_{j+1}} \{T_{i-1,j+1}[v_1, v_2, \dots, v_j, v_{j+1}]\}. \end{aligned}$$

350 Here,  $\text{Agg}$  is called the aggregation function of a NLM. For example, a sum aggregation function  
 351 takes the summation along the dimension  $j + 1$  of the tensor, and a max aggregation function takes  
 352 the max along that dimension.

353 The *concat* (concatenate) operation  $\oplus$  is applied at the “vector representation” dimension. The  
 354 *permute* operation generates a new tensor of the same arity, but it fuses the representations of  
 355 hyperedges that share the same set of entities but in different order, such as  $(v_1, v_2)$  and  $(v_2, v_1)$ .  
 356 Mathematically, for tensor  $X$  of arity  $j$ , if  $Y = \text{permute}(X)$  then

$$Y[v_1, v_2, \dots, v_j] = \text{Concat} \left\{ X[v_{\sigma_1}, v_{\sigma_2}, \dots, v_{\sigma_j}] \right\},$$

357 where  $\sigma \in S_j$  iterates over all permutations of  $\{1, 2, \dots, j\}$ .  $\text{NN}_j$  is a multi-layer perceptron (MLP)  
 358 applied to each entry in the tensor produced after permutation, with nonlinearity  $\sigma$  (e.g., ReLU).

359 It is important to note that we intentionally name the MLPs  $\text{NN}_j$  instead of  $\text{NN}_{i,j}$ . In generalized rela-  
 360 tional neural networks, for a given arity  $j$ , all MLPs across all layers  $i$  are shared. It is straightforward  
 361 to see that this “weight-shared” model can realize a “non-weight-shared” NLM that uses different  
 362 weights for MLPs at different layers when the number of layers is a constant. With a sufficiently  
 363 large length of the representation vector, we can simulate the computation of applying different  
 364 transformations by constructing block matrix weights. (A more formal proof is in Appendix B) The  
 365 advantage of this weight sharing is that the network can be easily extended to a “recurrent” model.  
 366 For example, we can apply the NLM for a number of layers that is a function of  $n$ , where  $n$  is the the  
 367 number of nodes in the input graph. Thus, we will use the term *layers* and *iterations* interchangeably.

368 Handling high-arity features and using deeper models usually increase the computational cost. In  
 369 appendix B.5, we show that the time and space complexity of NLM  $[D, B]$  is  $O(Dn^B)$ .

370 Note that even when hyperparameters such as the maximum arity and the number of iterations  
 371 are fixed, a NLM is still a *model family*  $\mathcal{M}$ : the weights for MLPs will be trained on some data.  
 372 Furthermore, each model  $M \in \mathcal{M}$  is a NLM with a specific set of MLP weights.

### 373 B.3 Expressiveness Equivalence of Relational Neural Networks

374 Since we are going to study both constant-depth and adaptive-depth graph neural networks, we first  
 375 prove the following lemma (for general multi-layer neural networks), which helps us simplify the  
 376 analysis.

377 **Lemma B.1.** A neural network with representation width  $W$  that has  $D$  different layers  
 378  $\text{NN}_1, \dots, \text{NN}_D$  can be realized by a neural network that applies a single layer  $\text{NN}'$  for  $D$  iter-  
 379 ations with width  $(D + 1)(W + 1)$ .

380 *Proof.* The representation for  $\text{NN}'$  can be partitioned into  $D + 1$  segments each of length  $W + 1$ .  
 381 Each segment consist of a “flag” element and a  $W$ -element representation, which are all 0 initially,  
 382 except for the first segment, where the flag is set to 1, and the representation is the input.

383  $\text{NN}'$  has the weights for all  $\text{NN}_1, \dots, \text{NN}_D$ , where weights  $\text{NN}_i$  are used to compute the representa-  
 384 tion in segment  $i + 1$  from the representation in segment  $i$ . Additionally, at each iteration, segment  
 385  $i + 1$  can only be computed if the flag in segment  $i$  is 1, in which case the flag of segment  $i + 1$   
 386 is set to 1. Clearly, after  $D$  iterations, the output of  $\text{NN}'_k$  should be the representation in segment  
 387  $D + 1$ .  $\square$

388 Due to Lemma B.1, we consider the neural networks that recurrently apply the same layer because  
 389 a) they are as expressive as those using layers of different weights, b) it is easier to analyze a single  
 390 neural network layer than  $D$  layers, and c) they naturally generalize to neural networks that runs for  
 391 adaptive number of iterations (e.g. GNNs that run  $O(\log n)$  iterations where  $n$  is the size of the input  
 392 graph).

393 We first describe a framework for quantifying if two hypergraph neural network models are equally  
 394 expressive on regression tasks (which is more general than classification problems). The frame-  
 395 work view the expressiveness from the perspective of computation. Specifically, we will prove the  
 396 expressiveness equivalence between models by showing that their computation can be aligned.

397 In complexity, we usually show a problem is at least as hard as the other one by showing a reduction  
 398 from the other problem to the problem. Similarly, on the expressiveness of NLMs, we can construct  
 399 reduction from model family  $\mathcal{A}$  to model family  $\mathcal{B}$  to show that  $\mathcal{B}$  can realize all computation that  $\mathcal{A}$   
 400 does, or even more. Formally, we have the following definition.

401 **Definition B.1** (Expressiveness reduction). For two model families  $\mathcal{A}$  and  $\mathcal{B}$ , we say  $\mathcal{A}$  can be  
 402 **reduced** to  $\mathcal{B}$  if and only if there is a function  $r : \mathcal{A} \rightarrow \mathcal{B}$  such that for each model instance  $A \in \mathcal{A}$ ,  
 403  $r(A) \in \mathcal{B}$  and  $A$  have the same outputs on all inputs. In this case, we say  $\mathcal{B}$  is at least as expressive  
 404 as  $\mathcal{A}$ .

405 **Definition B.2** (Expressiveness equivalence). For two model families  $\mathcal{A}$  and  $\mathcal{B}$ , if  $\mathcal{A}$  and  $\mathcal{B}$  can be  
 406 reduced to each other, then  $\mathcal{A}$  and  $\mathcal{B}$  are equally expressive. Note that this definition of expressiveness  
 407 equivalence generalizes to both classification and regression tasks.

408 **Equivalence between HO-GNNs and NLMs.** We will prove the equivalence between HO-GNNs  
 409 and NLMs by making reductions in both directions.

410 **Lemma B.2.** A  $B$ -ary HO-GNN with depth  $D$  can be realized by a NLM with maximum arity  $B + 1$   
 411 and depth  $2D$ .

412 *Proof.* We prove lemma B.2 by showing that one layer of GNNs on  $B$ -ary hypergraphs can be  
 413 realized by two NLM with maximum arity  $B + 1$ .

414 Firstly, a GNN layer maintain features of  $B$ -tuples, which are stored in correspondingly in an NLM  
 415 layer at dimension  $B$ . Then we will realize the message passing scheme using the NLM features of  
 416 dimension  $B$  and  $B + 1$  in two steps.

417 Recall the message passing scheme generalized to high dimensions (to distinguish, we use  $H$  for  
 418 HO-GNN features and  $T$  for NLM features.)

$$\text{Received}_i(\mathbf{v}) = \sum_u (\text{NN}_1 (H_{i-1,B}[\mathbf{v}]; \text{CONCAT}_{\mathbf{v}' \in \text{neighbors}(\mathbf{v}, u)} H_{i-1}[\mathbf{v}'])) \quad (\text{B.5})$$

419 At the first step, the Expand operation first raise the dimension to  $B + 1$  by expanding a non-related  
 420 variable  $u$  to the end, and the Permute operation can then swap  $u$  with each of the elements (or no  
 421 swap). Particularly,  $T_{i,B}[v_1, v_2, \dots, v_B]$  will be expand to

$$T_{i+1,B+1}[u, v_2, v_3, \dots, v_B, v_1], T_{i+1,B+1}[v_1, u, v_3, \dots, v_B, v_2], \dots, \\ T_{i+1,B+1}[v_1, v_2, \dots, v_{B-1}, u, v_B], \text{and } T_{i+1,B+1}[v_1, v_2, \dots, v_{B-1}, v_B, u]$$

422 Hence,  $T_{i+1,B+1}[v_1, v_2, v_3, \dots, v_B, u]$  receives the features from

$$T_{i,B}[v_1, v_2, \dots, v_B], T_{i,B}[u, v_2, v_3, \dots, v_B], T_{i,B}[v_1, u, v_3, \dots, v_B], \dots, T_{i,B}[v_1, v_2, \dots, v_{B-1}, u]$$

423 These features matches the input of  $\text{NN}_1$  in equation B.5, and in this layer  $\text{NN}_1$  can be applied to  
 424 compute things inside the summation.

425 Then at the second step, the last element is reduced to get what tuple  $\mathbf{v}$  should receive, so  $\mathbf{v}$  can be  
 426 updated. Since each HO-GNN layer can be realized by such two NLM layers, each  $B$ -ary HO-GNN  
 427 with depth  $D$  can be realized by a NLM of maximum arity  $(B + 1)$  and depth  $2D$ .  $\square$

428 To complete the proof we need to find a reduction from NLMs of maximum arity  $B + 1$  to  $B$ -ary  
 429 HO-GNNs. The key observation here is that the features of  $(B + 1)$ -tuples in NLMs can only be  
 430 expanded from sub-tuples, and the expansion and reduction involving  $(B + 1)$ -tuples can be simulated  
 431 by the message passing process.

432 **Lemma B.3.** The features of  $(B + 1)$ -tuples feature  $T_{i,B+1}[v_1, v_2, \dots, v_{B+1}]$  can be computed  
 433 from the following tuples

$$(T_{i,B}[v_2, v_3, \dots, v_{B+1}], T_{i,B}[v_1, v_3, \dots, v_{B+1}], \dots, T_{i,B}[v_1, v_2, \dots, v_B]).$$

434

435 *Proof.* Lemma B.3 is true because  $(B + 1)$ -dimensional representations can either be computed  
 436 from themselves at the previous iteration, or expanded from  $B$ -dimensional representations. Since  
 437 representations at all previous iterations  $j < i$  can be contained in  $T_{i,B}$ , it is sufficient to compute  
 438  $T_{i,B+1}[v_1, v_2, \dots, v_{B+1}]$  from all its  $B$ -ary sub-tuples.  $\square$

439 Then let's construct the HO-GNN for given NLM to show the existence of the reduction.

440 **Lemma B.4.** A NLM of maximum arity  $B + 1$  and depth  $D$  can be realized by a  $B$ -ary HO-GNN  
 441 with no more than  $D$  iterations.

442 *Proof.* We can realize the Expand and Reduce operation with only the  $B$ -dimensional features using  
 443 the broadcast message passing scheme. Note that Expand and Reduce between  $B$ -dimensional  
 444 features and  $(B + 1)$ -dimensional features in the NLM is a special case where claim B.3 is applied.

445 Let's start with Expand and Reduce operations between features of dimension  $B$  or lower. For the  
 446  $b$ -dimensional feature in the NLM, we keep  $n^b n^{B-b}$ § copies of it and store them the representation of  
 447 every  $B$ -tuple who has a sub-tuple¶ that is a permutation of the  $b$ -tuple. That is, for each  $B$ -tuple in  
 448 the  $B$ -ary HO-GNN, for its every sub-tuple of length  $b$ , we store  $b!$  representations corresponding to  
 449 every permutation of the  $b$ -tuple in the NLM. Keeping representation for all sub-tuple permutations  
 450 make it possible to realize the Permute operation. Also, it is easy to notice that Expand operation is  
 451 realized already, as all features with dimension lower than  $B$  are naturally expanded to  $B$  dimension  
 452 by filling in all possible combinations of the rest elements. Finally, the Reduce operation can be  
 453 realized using a broadcast casting message passing on certain position of the tuple.

454 Now let's move to the special case – the Expand and Reduce operation between features of dimensions  
 455  $B$  and  $B + 1$ . Claim B.3 suggests how the  $(B + 1)$ -dimensional features are stored in  $B$ -dimensional  
 456 representations in GNNs, and we now show how the Reduce can be realized by message passing.

457 We first bring in claim B.3 to the HO-GNN message passing, where we have  $\text{Received}_i[v]$  to be

$$\sum_u (\text{NN}_1 (T_{i-1,B}[v_2, v_3, \dots, v_B, u], T_{i-1,B}[v_1, v_3, \dots, v_B, u], \dots, T_{(i-1),B}[v_1, v_2, \dots, v_B]))$$

458 Note that the last term  $T_{i-1,B}[v_1, v_2, \dots, v_B]$  is contained in  $H_{i-1}(v)$  in equation B.5, and other  
 459 terms are contained in  $H_{i-1}(v')$  for  $v' \in \text{neighbors}(v, u)$ . Hence, equation B.5 is sufficient to  
 460 simulate the Reduce operation.  $\square$

461 **Theorem B.5.**  $B$ -ary HO-GNNs are equally expressive as NLMs with maximum arity  $B + 1$ .

462 *Proof.* This is a direct conclusion by combining Lemma B.2 and Lemma B.4.  $\square$

#### 463 B.4 Expressiveness of hypergraph convolution and attention

464 There exist other variants of hypergraph neural networks. In particular, hypergraph convolution[29–  
 465 31], attention[32] and message passing[33] focus on updating node features instead of tuple features  
 466 through hyperedges. These approaches can be viewed as instances of hypergraph neural networks,  
 467 and they have smaller time complexity because they do not model all high-arity tuples. However,  
 468 they are less expressive than the standard hypergraph neural networks with equal max arity.

469 These approaches can be formulated to two steps at each iteration. At the first step, each hyperedge is  
 470 updated by the features of nodes it connects.

$$h_{i,e} = \text{AGG}_{v \in e} f_{i-1,v} \tag{B.6}$$

471 At the second step, each node is updated by the features of hyperedges connecting it.

$$f_{i,v} = \text{AGG}_{v \in e} h_{i,e} \tag{B.7}$$

§  $n^k = n \times (n - 1) \times \dots \times (n - k + 1)$ .

¶ The sub-tuple does not have to be consecutive, but instead can be a any subset of the tuple that keeps the element order.

472 where  $f_{i,v}$  is the feature of node  $v$  at iteration  $i$ , and  $h_{i,v}$  is the aggregated message passing through  
 473 hyperedge  $e$  at iteration  $i + 1$ .

474 It is not hard to see that B.6 can be realized by  $B$  iterations of NLM layers with Expand operations  
 475 where  $B$  is the max arity of hyperedges. This can be done by expanding each node feature to every  
 476 high arity features that contain the node, and aggregate them at the tuple corresponding to each  
 477 hyperedge. Then, B.7 can also be realized by  $B$  iterations of NLM layers with Reduce operations, as  
 478 the tuple feature will finally be reduced to a single node contained in the tuple.

479 This approach has lower complexity compared to the GNNs we study applied on hyperedges, because  
 480 it only requires communication between nodes and hyperedges connecting to them, which takes  
 481  $O(|V| \cdot |E|)$  time at each iteration. Compared to them, NLMs takes  $O(|V|^B)$  time because NLMs  
 482 keep features of every tuple with max arity  $B$ , and allow communication from tuples to tuples instead  
 483 of between tuples and single nodes. An example is provided below that this approach can not solve  
 484 while NLMs can.

485 Consider a graph with 6 nodes and 6 edges forming two triangles  $(1, 2, 3)$  and  $(4, 5, 6)$ . Because of  
 486 the symmetry, the representation of each node should be identical throughout hypergraph message  
 487 passing rounds. Hence, it is impossible for these models to conclude that  $(1, 2, 3)$  is a triangle but  
 488  $(4, 2, 3)$  is not, based only on the node representations, because they are identical. In contrast, NLMs  
 489 with max arity 3 can solve them (as standard triangle detection problem in Table 1).

### 490 B.5 The Time and Space Complexity of NLMs

491 Handling high-arity features and using deeper models usually increase the computational cost in  
 492 terms of time and space. As an instance that use the architecture of RelNN, NLMs with depth  $D$   
 493 and max arity  $B$  takes  $O(Dn^B)$  time when applying to graphs with size  $n$ . This is because both  
 494 Expand and Reduce operation have linear time complexity with respect to the input size (which is  
 495  $O(n^B)$  at each iteration). If we need to record the computational history (which is typically the case  
 496 when training the network using back propagation), the space complexity is the same as the time  
 497 complexity.

498 GNNs applied to  $(B - 1)$ -ary hyperedges and depth  $D$  are equally expressive as RelNNs with depth  
 499  $O(D)$  and max arity  $B$ . Though up to  $(B - 1)$ -ary features are kept in their architecture, the broadcast  
 500 message passing scheme scale up the complexity by a factor of  $O(n)$ , so they also have time and  
 501 space complexity  $O(Dn^B)$ . Here the length of feature tensors  $W$  is treated as a constant.

## 502 C Arity and Depth Hierarchy: Proofs and Analysis

### 503 C.1 Proof of Theorem 3.1: Arity Hierarchy.

504 [11] have connected high-dimensional GNNs with high-dimensional WL tests. Specifically, they  
 505 showed that the  $B$ -ary HO-GNNs are equally expressive as  $B$ -dimensional WL test on graph  
 506 isomorphism test problem. In Theorem B.5 we proved that  $B$ -ary HO-GNNs are equivalent to NLM  
 507 of maximum arity  $B + 1$  in terms of expressiveness. Hence, NLM of maximum arity  $B + 1$  can  
 508 distinguish if two non-isomorphic graphs if and only if  $B$ -dimensional WL test can distinguish them.

509 However, Cai et al. [13] provided an construction that can generate a pair of non-isomorphic graphs  
 510 for every  $B$ , which can not be distinguished by  $(B - 1)$ -dimensional WL test but can be distinguished  
 511 by  $B$ -dimensional WL test. Let  $G_B^1$  and  $G_B^2$  be such a pair of graph.

512 Since NLM of maximum arity  $B + 1$  is equally expressive as  $B$ -ary HO-GNNs, there must be such a  
 513 NLM that classify  $G_B^1$  and  $G_B^2$  into different label. However, such NLM can not be realized by any  
 514 NLM of maximum arity  $B$  because they are proven to have identical outputs on  $G_B^1$  and  $G_B^2$ .

515 In the other direction, NLMs of maximum arity  $B + 1$  can directly realize NLMs of maximum arity  
 516  $B$ , which completes the proof.

### 517 C.2 Upper Depth Bound for Unbounded-Precision NLM.

518 The idea for proving an upper bound on depth is to connect NLMs to WL-test, and use the  $O(n^B)$   
 519 upper bound on number of iterations for  $B$ -dimensional test [34], and FOC formula is the key  
 520 connection.

521 For any fixed  $n$ ,  $B$ -dimensional WL test divide all graphs of size  $n$ ,  $\mathcal{G}_{=n}$ , into a set of equivalence  
 522 classes  $\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m\}$ , where two graphs belong to the same class if they can not be distinguished  
 523 by the WL test. We have shown that NLMs of maximum arity  $(B + 1)$  must have the same input for  
 524 all graphs in the same equivalence class. Thus, any NLM of maximum arity  $B + 1$  can be view as a  
 525 labeling over  $\mathcal{C}_1, \dots, \mathcal{C}_m$ .

526 Stated by Cai et al. [13],  $B$ -dimensional WL test are as powerful as  $\text{FOC}_{B+1}$  in differentiating graphs  
 527 graphs. Combined with the  $O(n^B)$  upper bound of WL test iterations, for each  $\mathcal{C}_i$ , there must be an  
 528  $\text{FOC}_{B+1}$  formula of quantifier depth  $O(n^B)$  that exactly recognize  $\mathcal{C}_i$  over  $\mathcal{G}_{=n}$ .

529 Finally, with unbounded precision, for any  $f(n)$ , NLM of maximum arity  $B + 1$  and depth  $f(n)$  can  
 530 compute all  $\text{FOC}_{B+1}$  formulas with quantifier depth  $f(n)$ . Note that there are finite number of such  
 531 formula because the supscript of counting quantifiers is bounded by  $n$ .

532 For any graph in some class  $\mathcal{C}_i$ , the class can be determined by evaluating these FOC formulas, and  
 533 then the label is determined. Therefore, any NLM of maximum arity  $B + 1$  can be realized by a  
 534 NLM of maximum arity  $B + 1$  and depth  $O(n^B)$ .

535 **C.3 Graph Problems**

$B = 4$	4-Clique Detection NLM[ $O(1), 4$ ]	4-Clique Count NLM[ $O(1), 4$ ]
$B = 3$	Triangle Detection NLM[ $O(1), 3$ ]	All-Pair Distance NLM[ $O(\log n), 3$ ]*
	Bipartiteness NLM[ $O(\log n), 3$ ]*	
	All-Pair Connectivity NLM[ $O(\log n), 3$ ]*	
	All-Pair Connectivity- $k$ NLM[ $O(\log k), 3$ ]*	
$B = 2$	$\text{FOC}_2$ Realization NLM[ $\cdot, 2$ ] [14]	S-T Distance NLM[ $O(n), 2$ ]
	3/4-Link Detection NLM[ $O(1), 2$ ]	Max Degree NLM[ $O(1), 2$ ]
	S-T Connectivity NLM[ $O(n), 2$ ]	Max Flow NLM[ $O(n^3), 2$ ]*
	S-T Connectivity- $k$ NLM[ $O(k), 2$ ]	
$B = 1$	Node Color Majority: NLM[ $O(1), 1$ ]	Count Red Nodes: NLM[ $O(1), 1$ ]
	<b>Classification Tasks</b>	<b>Regression Tasks</b>

**Table 1:** The minimum depth and arity of NLMs for solving graph classification and regression tasks. The \* symbol indicates that these are conjectured lower bounds.

536 We list a number of examples for graph classification and regression tasks, and we provide the  
 537 definitions and the current best known NLMs for learning these tasks from data. For some of the  
 538 problems, we will also show why they can not be solved by a simpler problems, or indicate them as  
 539 open problems.

540 **Node Color Majority.** Each node is assigned a color  $c \in \mathcal{C}$  where  $\mathcal{C}$  is a finite set of all colors. The  
 541 model needs to predict which color the most nodes have.

542 Using a single layer with sum aggregation, the model can count the number of nodes of color  $c$  for  
 543 each  $c \in \mathcal{C}$  on its global representation.

544 **Count Red Nodes.** Each node is assigned a color of red or blue. The model needs to count the  
 545 number of red nodes.

546 Similarly, using a single layer with sum aggregation, the model can count the number of red nodes on  
 547 its global representation.

548 **3-Link Detection.** Given an unweighted, undirected graph, the model needs to detect whether there  
 549 is a triple of nodes  $(a, b, c)$  such that  $a \neq c$  and  $(a, b)$  and  $(b, c)$  are edges.

550 This is equivalent to check whether there exists a node with degree at least 2. We can use a Reduction  
 551 operation with sum aggregation to compute the degree for each node, and then use a Reduction  
 552 operation with max aggregation to check whether the maximum degree of nodes is greater than or  
 553 equal to 2.

554 Note that this can not be done with 1 layer, because the edge information is necessary for the problem,  
 555 and they require at least 2 layers to be passed to the global representation.

556 **4-Link Detection.** Given an unweighted undirected graph, the model needs to detect whether there  
 557 is a 4-tuple of nodes  $(a, b, c, d)$  such that  $a \neq c, b \neq d$  and  $(a, b), (b, c), (c, d)$  are edges (note that a  
 558 triangle is also a 4-link).

559 This problem is equivalent to check whether there is an edge between two nodes with degrees  $\geq 2$ .  
 560 We can first reduce the edge information to compute the degree for each node, and then expand it  
 561 back to 2-dimensional representations, so we can check for each edge if the degrees of its ends are  
 562  $\geq 2$ . Then the results are reduced to the global representation with existential quantifier (realized by  
 563 max aggregation) in 2 layers.

564 **Triangle Detection.** Given a unweighted undirected graph, the model is asked to determine whether  
 565 there is a triangle in the graph i.e. a tuple  $(a, b, c)$  so that  $(a, b), (b, c), (c, a)$  are all edges.

566 This problem can be solved by NLM [4,3]: we first expand the edge to 3-dimensional representations,  
 567 and determine for each 3-tuple if they form a triangle. The results of 3-tuples require 3 layers to be  
 568 passed to the global representation.

569 We can prove that Triangle Detection indeed requires breadth at least 3. Let  $k$ -regular graphs be  
 570 graphs where each node has degree  $k$ . Consider two  $k$ -regular graphs both with  $n$  nodes, so that  
 571 exactly one of them contains a triangle<sup>†</sup>. However, NLMs of breadth 2 has been proven not to be  
 572 stronger than WL test on distinguish graphs, and thus can not distinguish these two graphs (WL test  
 573 can not distinguish any two  $k$ -regular graphs with equal size).

574 **4-Clique Detection and Counting.** Given an undirected graph, check existence of, or count the  
 575 number of tuples  $(a, b, c, d)$  so that there are edges between every pair of nodes in the tuple.

576 This problem can be easily solved by a NLM with breadth 4 that first expand the edge information to  
 577 the 4-dimensional representations, and for each tuple determine whether its is a 4-clique. Then the  
 578 information of all 4-tuples are reduced 4 times to the global representation (sum aggregation can be  
 579 used for counting those).

580 Though we did not find explicit counter-example construction on detecting 4-cliques with NLMs of  
 581 breadth 3, we suggest that this problem can not be solved with NLMs with 3 or lower breadth.

582 **Connectivity.** The connectivity problems are defined on unweighted undirected graphs. S-T con-  
 583 nectivity problems provides two nodes  $S$  and  $T$  (labeled with specific colors), and the model needs  
 584 to predict if they are connected by some edges. All pair connectivity problem require the model to  
 585 answer for every pair of nodes. Connectivity- $k$  problems have an additional requirement that the  
 586 distance between the pair of nodes can not exceed  $k$ .

587 S-T connectivity- $k$  can be solved by a NLM of breadth 2 with  $k$  iterations. Assume  $S$  is colored with  
 588 color  $c$ , at every iteration, every node with color  $c$  will spread the color to its neighbors. Then, after  $k$   
 589 iterations, it is sufficient to check whether  $T$  has the color  $c$ .

590 With NLMs of breadth 3, we can use  $O(\log k)$  matrix multiplications to solve connectivity- $k$  between  
 591 every pair of nodes. Since the matrix multiplication can naturally be realized by NLMs of breadth 3  
 592 with two layers. All-pair connectivity problems can all be solved with  $O(\log k)$  layers.

593 **Theorem C.1** (S-T connectivity- $k$  with NLM). S-T connectivity- $k$  can not be solved by a NLM of  
 594 maximum arity within  $o(k)$  iterations.

595 *Proof.* We construct two graphs each has  $2k$  nodes  $u_1, \dots, u_k, v_1, \dots, v_k$ . In both graph, there are  
 596 edges  $(u_i, u_{i+1})$  and  $(v_i, v_{i+1})$  for  $1 \leq i \leq k - 1$  i.e. there are two links of length  $k$ . We then set  
 597  $S = u_1, T = u_n$  and  $S = u_1, T = v_n$  the the two graphs.

598 We will analysis GNNs as NLMs are proved to be equivalent to them by scaling the depth by a  
 599 constant factor. Now consider the node refinement process where each node  $x$  is refined by the  
 600 multiset of labels of  $x$ 's neighbors and the multiset of labels of  $x$ 's non-neighbors.

601 Let  $C_j^{(i)}(x)$  be the label of  $x$  in graph  $j$  after  $i$  iterations, at the beginning, WLOG, we have

$$C_1^{(0)}(u_1) = 1, C_1^{(0)}(u_n) = 2C_1^{(0)}(u_1) = 1, C_1^{(0)}(v_n) = 2$$

<sup>†</sup>Such construction is common. One example is  $k = 2, n = 6$ , and the graph may consist of two separated triangles or one hexagon

602 and all other nodes are labeled as 0.

603 Then we can prove by induction: after  $i \leq \frac{k}{2} - 1$  iterations, for  $1 \leq t \leq i + 1$  we have

$$604 \quad C_1^{(u_t)} = C_2^{(i)}(u_t), C_1^{(v_t)} = C_2^{(i)}(v_t)$$

$$C_1^{(u_{k-t+1})} = C_2^{(i)}(v_{k-t+1}), C_1^{(v_{k-t+1})} = C_2^{(i)}(u_{k-t+1})$$

605 and for  $i + 2 \leq t \leq k - i - 1$  we have

$$C_1^{(u_t)} = C_2^{(i)}(u_t), C_1^{(v_t)} = C_2^{(i)}(v_t)$$

606 This is true because before  $\frac{k}{2}$  iterations are run, the multiset of all node labels are identical for the  
 607 two graphs (say  $S^{(i)}$ ). Hence each node  $x$  is actually refined by its neighbors and  $S^{(i)}$  where  $S^{(i)}$   
 608 is the same for all nodes. Hence, before running  $\frac{k}{2}$  iterations when the message between  $S$  and  $T$   
 609 finally meets in the first graph, GNN can not distinguish the two graphs, and thus can not solve the  
 610 connectivity with distance  $k - 1$ .  $\square$

611 **Max Degree.** The max degree problem gives a graph and ask the model to output the maximum  
 612 degree of its nodes.

613 Like we mentioned in 3-link detection, one layer for computing the degree for each node, and another  
 614 layer for taking the max operation over nodes should be sufficient.

615 **Max Flow.** The Max Flow problem gives a directional graph with capacities on edges, and indicate  
 616 two nodes  $S$  and  $T$ . The models is then asked to compute the amount of max-flow from  $S$  to  $T$ .

617 Notice that the Breadth First Search (BFS) component in Dinic’s algorithm[35] can be implemented  
 618 on NLMs as they does not require node identities (all new-visited nodes can augment to their non-  
 619 visited neighbors in parallel). Since the BFS runs for  $O(n)$  iteration, and the Dinic’s algorithm runs  
 620 BFS  $O(n^2)$  times, the max-flow can be solved by NLMs with in  $O(n^3)$  iterations.

621 **Distance.** Given a graph with weighted edges, compute the length of the shortest between specified  
 622 node pair (S-T Distance) or all node pairs (All-pair Distance).

623 Similar to Connectivity problems, but Distance problems now additionally record the minimum  
 624 distance from  $S$  (for S-T) or between every node pairs (for All-pair), which can be updated using min  
 625 operator (using Min-plus matrix multiplication for All-pair case).

## 626 D Experiments

627 We now study how our theoretical results on model expressiveness and learning apply to relational  
 628 neural networks trained with gradient descent on practically meaningful problems. We begin by  
 629 describing two synthetic benchmarks: graph substructure detection and relational reasoning.

630 In the graph substructure detection dataset, there are several tasks of predicting whether there input  
 631 graph contained a sub-graph with specific structure. The tasks are: *3-link* (length-3 path), *4-link*,  
 632 *triangle*, and *4-clique*. These are important graph properties with many potential applications.

633 The relational reasoning dataset is composed of two family-relationship prediction tasks and two  
 634 connectivity-prediction tasks. They are all *binary edge classification* tasks. In the family-relationship  
 635 prediction task, the input contains the *mother* and *father* relationships, and the task is to predict the  
 636 *grandparent* and *uncle* relationships between all pairs of entities. In the connectivity-prediction tasks,  
 637 the input is the edges in an undirected graph and the task is to predict, for all pairs of nodes, whether  
 638 they are connected with a path of length  $\leq 4$  (*connectivity-4*) and whether they are connected with a  
 639 path of arbitrary length (*connectivity*). The data generation for all datasets is included in Appendix D.

### 640 D.1 Experiment Setup

641 For all problems, we have 800 training samples, 100 validation samples, and 300 test samples for  
 642 each different  $n$  we are testing the models on.

643 We then provide the details on how we synthesize the data. For most of the problems, we generate  
 644 the graph by randomly selecting from all potential edges i.e. the Erdős–Rényi model. We sample the

645 number of edges around  $n$ ,  $2n$ ,  $n \log n$  and  $n^2/2$ . For all problems, with 50% probability the graph  
 646 will first be divided into 2, 3, 4 or 5 parts with equal number of components, where we use the first  
 647 generated component to fill the edges for rest of the components. Some random edges are added  
 648 afterwards. This make the data contain more isomorphic sub-graphs, which we found challenging  
 649 empirically.

650 **Substructure Detection.** To generate a graph that does not contain a certain substructure, we  
 651 randomly add edges when reaching a maximal graph not containing the substructure or reaching the  
 652 edge limit. For generating a graph that does contain the certain substructure, we first generate one  
 653 that does not contain, and then randomly replace present edges with missing edges until we detect  
 654 the substructure in the graph. This aim to change the label from “No” to “Yes” while minimizing  
 655 the change to the overall graph properties, and we found that data generated using edge replacing is  
 656 much more difficult for neural networks compared to random generated graphs from scratch.

657 **Family Tree.** We generate the family trees using the algorithm modified from [5]. We add people  
 658 to the family one by one. When a person is added, with probability  $p$  we will try to find a single  
 659 woman and a single man, get them married and let the new children be their child, and otherwise the  
 660 new person is introduced as a non-related person. Every new person is marked as single and set the  
 661 gender with a coin flip.

662 We adjust  $p$  based on the ratio of single population:  $p = 0.7$  when more than 40% of the population  
 663 are single, and  $p = 0.3$  when less than 20% of the population are single, and  $p = 0.5$  otherwise.

664 **Connectivity.** For connectivity problems, we use the similar generation method as the substructure  
 665 detection. We sample the query pairs so that the labels are balanced.

## 666 D.2 Model Implementation Details

667 For all models, we use a hidden dimension 128 except for 3-dimensional HO-GNN and 4-dimensional  
 668 NLM where we use hidden dimension 64.

669 All model have 4 layers that each has its own parameters, except for connectivity where we use  
 670 the recurrent models that apply the second layer  $k$  times, where  $k$  is sampled from integers in  
 671  $[2 \log n, 3 \log n]$ . The depths are proven to be sufficient for solving these problems (unless the model  
 672 itself can not solve).

673 All models are trained for 100 epochs using adam optimizer with learning rate  $3 \times 10^{-4}$  decaying at  
 674 epoch 50 and 80.

675 We have varied the depth, the hidden dimension, and the activation function of different models.  
 676 We select sufficient hidden dimension and depth for every model and problem (i.e., we stop when  
 677 increasing depth or hidden dimension doesn’t increase the accuracy). We tried linear, ReLU, and  
 678 Sigmoid activation functions, and ReLU performed the best overall combinations of models and  
 679 tasks.

## 680 D.3 Results

681 Our main results on all datasets are shown in Table 2 and Table 3. We empirically compare relational  
 682 neural networks with different maximum arity  $B$ , different model architecture (GNN and NLM), and  
 683 different aggregation functions (max and sum). All models use sigmoidal activation for all MLPs.  
 684 For each task on both datasets we train on a set of small graphs ( $n = 10$ ) and test the trained model  
 685 on both small graphs and large graphs ( $n = 10$  and  $n = 30$ ). We summarize the findings below.

686 **Expressiveness.** We have seen a theoretical equal expressiveness between GNNs and NLMs applied  
 687 to hypergraphs. That is, a GNN applied to  $B$ -ary hyperedges is equivalent to a  $(B + 1)$ -ary NLM.  
 688 Table 2 and 3 further suggest their similar performance on tasks when trained with gradient descent.

689 Formally, triangle detection requires NLMs with at least  $B = 3$  to solve. Thus, we see that all  
 690 NLMs with arity  $B = 2$  fail on this task, but models with  $B = 3$  perform well. Formally, 4-clique  
 691 is realizable by NLMs with maximum arity  $B = 4$ , but we failed to reliably train models to reach  
 692 perfect accuracy on this problem. It is not yet clear what the cause of this behavior is.

693 **Structural generalization.** We discussed the structural generalization properties of NLMs in Sec-  
 694 tion 4, in a learning setting based on fixed-precision networks and enumerative training. This setting

Model	Agg.	3-link		4-link		triangle		4-clique	
		$n = 10$	$n = 30$						
1-ary GNN	Max	70.0 $\pm$ 0.0	82.7 $\pm$ 0.0	92.0 $\pm$ 0.0	91.7 $\pm$ 0.0	73.7 $\pm$ 3.2	50.2 $\pm$ 1.8	55.3 $\pm$ 4.0	46.2 $\pm$ 1.3
	Sum	100.0 $\pm$ 0.0	89.4 $\pm$ 0.4	100.0 $\pm$ 0.0	86.1 $\pm$ 1.2	77.7 $\pm$ 8.5	48.6 $\pm$ 1.6	53.7 $\pm$ 0.6	55.2 $\pm$ 0.8
2-ary NLM	Max	65.3 $\pm$ 0.6	54.0 $\pm$ 0.6	93.0 $\pm$ 0.0	95.7 $\pm$ 0.0	51.0 $\pm$ 1.7	49.2 $\pm$ 0.4	55.0 $\pm$ 0.0	45.7 $\pm$ 0.0
	Sum	100.0 $\pm$ 0.0	88.3 $\pm$ 0.0	100.0 $\pm$ 0.0	67.4 $\pm$ 16.4	82.0 $\pm$ 2.6	48.3 $\pm$ 0.0	53.0 $\pm$ 0.0	54.4 $\pm$ 1.5
2-ary GNN	Max	78.7 $\pm$ 0.6	76.0 $\pm$ 17.3	97.7 $\pm$ 4.0	98.6 $\pm$ 2.5	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	55.0 $\pm$ 0.0	45.7 $\pm$ 0.0
	Sum	100.0 $\pm$ 0.0	51.2 $\pm$ 7.9	100.0 $\pm$ 0.0	45.7 $\pm$ 7.6	100.0 $\pm$ 0.0	49.2 $\pm$ 1.0	61.0 $\pm$ 5.6	54.3 $\pm$ 0.0
3-ary NLM	Max	100.0 $\pm$ 0.0	59.0 $\pm$ 6.9	45.9 $\pm$ 0.4					
	Sum	100.0 $\pm$ 0.0	87.6 $\pm$ 11.0	100.0 $\pm$ 0.0	65.4 $\pm$ 14.3	100.0 $\pm$ 0.0	80.6 $\pm$ 8.8	73.7 $\pm$ 13.8	53.3 $\pm$ 8.8
3-ary GNN	Max	79.0 $\pm$ 0.0	86.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	84.0 $\pm$ 0.0	93.3 $\pm$ 0.0
	Sum	100.0 $\pm$ 0.0	84.1 $\pm$ 18.6	100.0 $\pm$ 0.0	61.1 $\pm$ 15.0	100.0 $\pm$ 0.0	95.1 $\pm$ 7.3	80.5 $\pm$ 0.7	66.2 $\pm$ 19.6
4-ary NLM	Max	100.0 $\pm$ 0.0	82.0 $\pm$ 1.7	93.1 $\pm$ 0.2					
	Sum	100.0 $\pm$ 0.0	59.1 $\pm$ 5.3	100.0 $\pm$ 0.0	67.7 $\pm$ 24.1	100.0 $\pm$ 0.0	82.1 $\pm$ 12.8	84.0 $\pm$ 0.0	67.0 $\pm$ 18.9

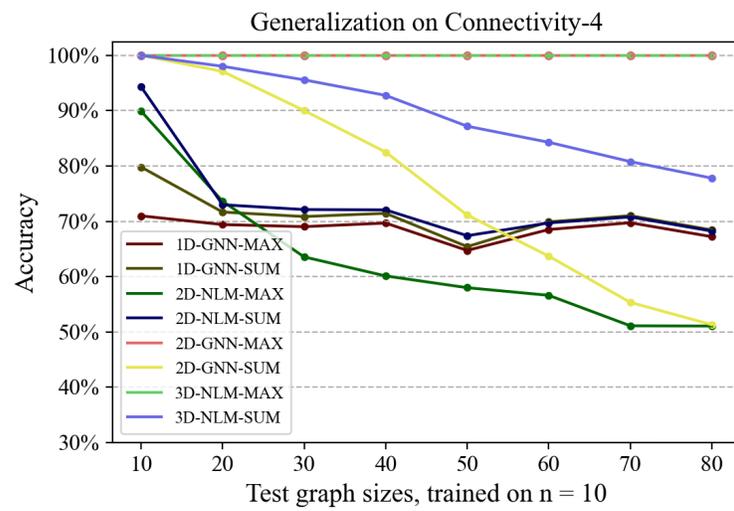
**Table 2:** Overall accuracy on relational reasoning problems. All models are trained on  $n = 10$ , and tested on  $n = 30$ . The standard error of all values are computed based on three random seeds.

Model	Agg.	grand parent		uncle		connectivity-4**		connectivity	
		$n = 20$	$n = 80$	$n = 20$	$n = 80$	$n = 10$	$n = 80$	$n = 10$	$n = 80$
1-ary GNN	Max	84.0 $\pm$ 0.3	64.8 $\pm$ 0.0	93.6 $\pm$ 0.3	66.1 $\pm$ 0.0	72.6 $\pm$ 3.6	67.5 $\pm$ 0.5	85.6 $\pm$ 0.3	75.1 $\pm$ 1.9
	Sum	84.7 $\pm$ 0.1	64.4 $\pm$ 0.0	94.3 $\pm$ 0.2	66.2 $\pm$ 0.0	79.6 $\pm$ 0.1	68.3 $\pm$ 0.1	87.1 $\pm$ 0.3	75.0 $\pm$ 0.2
2-ary NLM	Max	82.3 $\pm$ 0.5	65.6 $\pm$ 0.1	93.1 $\pm$ 0.0	66.6 $\pm$ 0.0	91.2 $\pm$ 0.2	51.0 $\pm$ 0.6	88.9 $\pm$ 2.6	67.1 $\pm$ 4.8
	Sum	82.9 $\pm$ 0.1	64.6 $\pm$ 0.1	93.4 $\pm$ 0.0	66.7 $\pm$ 0.2	96.0 $\pm$ 0.4	68.3 $\pm$ 0.5	84.0 $\pm$ 0.0	71.9 $\pm$ 0.0
2-ary GNN	Max	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	84.0 $\pm$ 0.0	71.9 $\pm$ 0.0				
	Sum	100.0 $\pm$ 0.0	35.7 $\pm$ 0.0	100.0 $\pm$ 0.0	33.9 $\pm$ 0.0	100.0 $\pm$ 0.0	51.3 $\pm$ 5.3	84.0 $\pm$ 0.0	71.9 $\pm$ 0.0
3-ary NLM	Max	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.1				
	Sum	100.0 $\pm$ 0.0	35.7 $\pm$ 0.0	100.0 $\pm$ 0.0	50.8 $\pm$ 29.4	100.0 $\pm$ 0.0	77.8 $\pm$ 11.8	100.0 $\pm$ 0.0	88.2 $\pm$ 8.0
3-ary NLM <sub>HE</sub>	Max	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	100.0 $\pm$ 0.0	N/A	N/A	N/A	N/A
	Sum	100.0 $\pm$ 0.0	35.7 $\pm$ 0.0	100.0 $\pm$ 0.0	33.8 $\pm$ 29.4	N/A	N/A	N/A	N/A

**Table 3:** Overall accuracy on relational reasoning problems. Models for family-relationship prediction are trained on  $n = 20$ , while models for connectivity problems are trained on  $n = 10$ . All model are tested on  $n = 80$ . The standard error of all values are computed based on three random seeds. The 3-ary NLMs marked with “HE” have hyperedges in inputs, where each family is represented by a 3-ary hyperedge instead of two parent-child edges, and the results are similar to binary edges.

695 can be *approximated* by training NLMs with max aggregation and sigmoidal activation on sufficient  
 696 data.

697 We run a case study on the problem connectivity-4 about how the generalization performance changes  
 698 when the test graph size gradually becomes larger. Figure 2 show how these models generalize  
 699 to gradually larger graphs with size increasing from 10 to 80. From the curves we can see that  
 700 only models with sufficient expressiveness can get 100% accuracy on the same size graphs, and  
 701 among them the models using max aggregation generalize to larger graphs with no performance  
 702 drop. 2-ary GNN and 3-ary NLM that use max aggregation have sufficient expressiveness and better  
 703 generalization property. They achieve 100% accuracy on the original graph size and generalize  
 704 perfectly to larger graphs.



**Figure 2:** How the performance of models drop when generalizing to larger graphs on the problem connectivity-4 (trained on graphs with size 10).