
SHINE: SubHypergraph Inductive Neural nEtnetwork

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Hypergraph neural networks can model multi-way connections that are beyond
2 pairwise associations among nodes of the graphs. Multi-way connections are
3 common in many real-world applications and, in particular, genetic medicine. In
4 particular, genetic pathways or broadly speaking gene sets encode relationships
5 among multiple genes that collectively drive a molecular function, which can be
6 naturally modeled as hyperedges connecting all involved nodes (e.g., genes). Thus,
7 hypergraph-guided embedding can capture functional relations in learned represen-
8 tations. Existing hypergraph neural network models often focus on node-level or
9 graph-level inference. There is an unmet need in learning powerful representations
10 of subgraphs of hypergraphs in real-world applications. For example, a cancer pa-
11 tient can be viewed as a subgraph of genes harboring mutations in the patient, while
12 all the genes are connected by hyperedges that correspond to pathways representing
13 specific molecular functions. To achieve accurate inductive subgraph prediction,
14 we propose SubHypergraph Inductive Neural nEtnetwork (SHINE). SHINE uses
15 informative genetic pathways that encode molecular functions as hyperedges to
16 connect genes as nodes. SHINE jointly optimizes the objectives of end-to-end
17 subgraph classification and hypergraph nodes' similarity regularization. SHINE
18 simultaneously learns representations for both genes and pathways using strongly
19 dual attention message passing. These learned representations are then aggregated
20 via a subgraph attention layer to derive a subgraph representation, which is used in
21 training a multilayer perceptron for subgraph inferencing. We evaluated SHINE
22 against a wide array of state-of-the-art hypergraph neural networks, XGBoost,
23 NMF and polygenic risk score models, using diverse large scale NGS and curated
24 datasets. SHINE outperformed all comparison models significantly, and yielded
25 interpretable models with functional insights on disease molecular mechanisms.

26 1 Introduction

27 Hypergraph neural networks have recently emerged as a series of successful methods to model
28 multi-way connections that are beyond pairwise associations among nodes of the graphs. Multi-way
29 connections are common in many real-world applications and, in particular, genetic medicine. From
30 genetic medicine's perspective, pathways or broadly speaking gene sets encode the relationship
31 among multiple genes that collectively correspond to a molecular function [1], which account for the
32 mechanisms of pathogenesis more intuitively and accurately than individual genes. Genetic pathways
33 or gene sets encode functional relations among multiple genes, which can be naturally modeled as
34 hyperedges connecting all involved nodes (e.g., genes). Thus, hypergraph-guided embedding can
35 capture functional relations in learned representations.

36 Existing hypergraph neural network models often adopt the semi-supervised learning (SSL) paradigm
37 to assign labels to initially unlabeled nodes in a hypergraph [2–4]. Other methods have focused on
38 learning graph representations [5, 6]. Node-level and graph-level representations give either local or

39 overarching views of the graphs, i.e., at the two extremes of hypergraph topological structures. There
40 is an unmet need in learning powerful representations of subgraphs in hypergraphs. Such capabilities
41 are important in genetic medicine. For example, cancer patients can be viewed as subgraphs of genes
42 that harbor mutations, while all the genes are connected by hyperedges that correspond to pathways
43 or gene sets representing specific molecular functions. Powerful subgraph representations will lead
44 to the capability to more accurately account for the subject’s pathophysiology. For regular graphs
45 where edges connect node pairs, several subgraph representation learning algorithms were proposed,
46 including methods that can use the learned representations to make predictions for subgraphs with
47 fixed sizes [7] or varying sizes [8]. There are currently few if any work on inductive inference
48 for varying-sized subhypergraphs. In this work, we propose a new framework named SHINE:
49 SubHypergraph Inductive Neural nEtnetwork. We will share our source code upon publication. Our
50 contributions are as follows:

- 51 • To the best of our knowledge, SHINE is the first model to effectively learn subgraph represen-
52 tations for hypergraphs, use the learned representations (for seen subgraphs) and inductively
53 infer representations (for unseen subgraphs) for downstream subgraph predictions.
- 54 • Novel applications in the field of genetic medicine on Next Generation Sequencing (NGS)
55 datasets across diverse diseases show significant performance improvements by SHINE over
56 a wide array of state-of-the-art baselines.
- 57 • In addition to learning and inductively inferring subgraph representations, SHINE simul-
58 taneously learns the representations of nodes and hyperedges. This brings interpretation
59 advantages, allowing assessing pathways (hyperedges) correlations and reasoning multiple
60 molecular functions mutually interacting and collectively contributing to the disease onset.

61 2 Related Work

62 **Graph Neural Networks.** Graph representation learning maps graphs or their components to
63 vector representations and has attracted growing attention over the past decade. Recently, graph
64 neural networks (GNNs), which can learn a distributed representation for a graph or a node in a graph,
65 are widely applied to a variety of areas including social network analysis [9, 10], molecular structure
66 inference [11, 12] and text mining [13, 14]. GNN recursively updates the representation of a node
67 in a graph by aggregating the feature vectors of its neighbors and itself, e.g. [15]. The graph-level
68 representations can then be obtained through set pooling (e.g., [16]) or graph coarsening (e.g., [17]) to
69 aggregate the node representations in the graph. The reader is referred to a comprehensive book [18]
70 on the topic of graph neural networks.

71 **Hypergraph neural network.** Hypergraph neural networks [5, 3, 2, 6] have become a popular
72 approach for learning on multiway relations from data. Early work on hypergraph learning, e.g.,
73 [19], formulated hypergraph message passing using spectral theory of hypergraphs. This formulation
74 and its variants [2, 3, 20–23] essentially adopted clique expansion to extend graph convolutional
75 network (GCN) for hypergraph learning. Others methods applied attention mechanism to aggregate
76 the information across the hypergraph [5, 6] or directly learned node representations to preserve the
77 proximity of nodes sharing a hyperedge or having similar neighborhoods [24]. In both formulations,
78 messages were passed to the node of interest from its immediate neighbors, and added layers allow
79 propagation of messages to a farther neighborhood.

80 **Subgraph representation learning and prediction.** Recent studies on subgraph embeddings and
81 prediction starts with learning representations of small subgraphs. For example, [7] encoded small
82 fixed-sized subgraphs for subgraph evolution prediction. SubGNN [8] learned representations for
83 varying-sized subgraphs through neighborhood, position and structure channels using random patches
84 distributed throughout the graph. [25] and [26] learned subgraph representations by pooling across
85 local structures to aid the predictions of the entire graphs.

86 The intersection of hypergraph neural network and subgraph representation learning is currently
87 underexplored. While the above methods focus on either hypergraph learning or subgraph learning,
88 none of the methods consider subgraph prediction for hypergraphs. Technically, subgraphs can be
89 viewed as a hyperedge and studies on link prediction could predict the existence of a hyperedge [27].
90 However, few if any such studies addressed the problems of differentiating the classes of the subgraphs,
91 which is especially important in genetic medicine where subgraphs and hyperedges have different

Table 1: Common notations used throughout the paper.

Symbol	Definition	Symbol	Definition
\mathcal{H}	An undirected hypergraph	$ \mathcal{E} $	Number of hypergraph hyperedges
\mathcal{N}	Set of hypergraph nodes	$ \mathcal{N} $	Number of hypergraph nodes
\mathcal{E}	Set of hypergraph hyperedges	d	Hidden layer size
\mathbf{H}	Hyperedge incidence matrix	\circ	Operation composition
n	Number of subgraphs	$*$	Element-wise multiplication

92 real-world meanings. For example, a hyperedge corresponds to a genetic pathway from curated
 93 knowledge and a subgraph corresponds to a patient with mutated genes as nodes.

94 To sum, there is a major unmet need regarding varying-sized subgraph inference for hypergraphs,
 95 and even more so in the inductive learning setting. Our proposed framework SHINE provides an
 96 end-to-end framework that operates on hypergraphs and performs inductive subgraph inferencing.

97 3 Methods

98 We first outline the workflow of SHINE, refer to Table 1 for symbols used throughout the paper. This
 99 study considers both genetic pathways and genetic variants. For genetic variants, we first filter the
 100 called variants from exome sequencing data and keep those variants that pass multiple quality control
 101 steps. We then calculate aggregated mutation rate at the gene level. For pathways (gene sets), we
 102 select the gene sets from curated databases and remove redundancy. We then construct the hypergraph
 103 using gene as nodes and pathways as hyperedges. We develop a strongly dual attention message
 104 passing algorithm to propagate information between nodes and hyperedges and across layers. We
 105 develop a weighted subgraph attention mechanism to learn the subgraph representation by integrating
 106 representations of hypergraph nodes (e.g., genes harboring mutations). We next explain each step.

107 3.1 Collecting Genetic Pathways

108 To construct the graph and to establish comparison models, we use the pathways with mutated
 109 genes as an approximation of the disrupted molecular functions. We use the Molecular Signatures
 110 Database (MSigDB) [1] to obtain a comprehensive collection of curated genetic pathways. We
 111 focus on MSigDB’s curated gene sets collection, which contains human gene sets that are canonical
 112 representations of a biological process compiled by domain experts. Gene sets in this collection are
 113 curated from various sources, including the biomedical literature and pathway databases such as
 114 Reactome [28] and KEGG [29]. Some pathways may overlap with others and have been filtered by
 115 MSigDB to remove interset redundancy. There are 21,587 genes in MSigDB Pathways. Genes with
 116 unknown functions are not included in the pathways and not used for classification, as our focus here
 117 is on interpretable modeling through inferencing with known molecular functions. Adding genes with
 118 unknown functions to study their impact will be our future work. With each pathway as a hyperedge,
 119 the incidence matrix entry \mathbf{H}_{ij} is 1 if the hyperedge j contains the node i , and 0 otherwise.

120 3.2 Hypergraph Learning

121 We first review the hypergraph analysis theory. Different from a simple graph, a hyperedge in a
 122 hypergraph connects two or more vertices. A hypergraph is defined as $\mathcal{H} = (\mathcal{N}, \mathcal{E})$, which includes a
 123 set of nodes \mathcal{N} , a set of hyperedges \mathcal{E} . In the case of genetic medicine, we model genes as hypergraph
 124 nodes, i.e., $\mathcal{N} = \{g_1, \dots, g_{|\mathcal{N}|}\}$, and pathways as hyperedges, i.e., $\mathcal{E} = \{p_1, \dots, p_{|\mathcal{E}|}\}$, where $|\mathcal{N}|, |\mathcal{E}|$
 125 are sizes of the nodes and hyperedges respectively. The hypergraph’s topological structure can be
 126 denoted by an $|\mathcal{N}| \times |\mathcal{E}|$ incidence matrix \mathbf{H} , whose entries are defined as

$$\mathbf{H}_{ij} = \begin{cases} 1 & \text{if node } g_i \in \text{hyperedge } p_j \\ 0 & \text{if node } g_i \notin \text{hyperedge } p_j \end{cases} \quad (1)$$

127 Generally speaking, each node in the hypergraph may be accompanied by a d -dimensional node
 128 feature/embedding matrix $\mathbf{N} \in R^{|\mathcal{N}| \times d}$, where each row corresponds to a node’s feature/embedding.
 129 The hypergraph with its topological structure and node features can be represented succinctly as
 130 $\mathcal{H} = (\mathbf{H}, \mathbf{N})$. For our experiments, we choose to use one-hot initialization of node embeddings, i.e.,

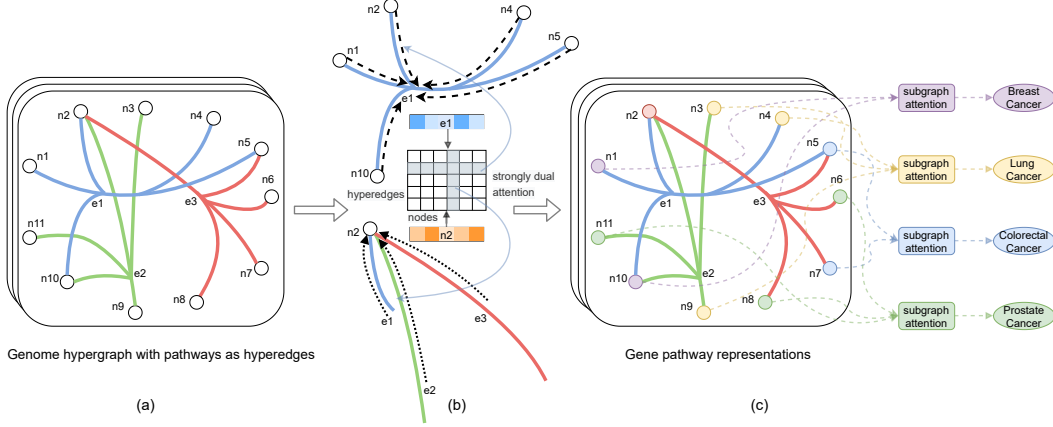


Figure 1: SHINE’s strongly dual attention mechanism for message passing for the genome hypergraph, and its use of subgraph attention to integrate gene nodes in the feature learning for patients.

131 $\mathbf{N} = \mathbf{I}$, and let SHINE learn gene embeddings through message passing. This choice also increases
 132 the difficulty for the overall learning algorithm, thus puts SHINE under stress test.

133 Fig. 1 (a) shows a schematic of the constructed genome hypergraph with nodes denoted by circles
 134 and hyperedges denoted by colored arcs. While a pathway can contain multiple genes, a gene can
 135 also contribute to multiple pathways. That is, we can have multiple hyperedges incident on the same
 136 node (gene), as can be seen in Fig. 1 (a) nodes n_2, n_5, n_{10} .

137 3.3 Strongly Dual Attention Message Passing

138 **Hyperedge attention over nodes.** For a hyperedge edge $p_j \in \mathcal{E}$, in order to update its hidden
 139 representation at layer k , we aggregate the information from its incident nodes using the following
 140 attention mechanism. We first calculate the hyperedge attention over nodes as in

$$a_E(p_j, g_i) = \frac{\exp(\mathbf{c}^T \mathbf{s}(p_j, g_i))}{\sum_{g_{i'} \in p_j} \exp(\mathbf{c}^T \mathbf{s}(p_j, g_{i'}))} \quad (2)$$

141 where \mathbf{c} is a trainable context vector and the attention ready state $\mathbf{s}(p_j, g_i)$ for a hyperedge-node pair
 142 (p_j, g_i) is calculated by mapping and combining nodes’ and hyperedges’ states from the $(k - 1)$ th
 143 layer as in

$$\mathbf{s}(p_j, g_i) = \text{LeakyReLU}\left(\left(\mathbf{W}_N \mathbf{h}_N^{k-1}(g_i) + \mathbf{b}_N\right) * \left(\mathbf{W}_E \mathbf{h}_E^{k-1}(p_j) + \mathbf{b}_E\right)\right) \quad (3)$$

144 where $*$ denotes element-wise product, \mathbf{W}_N and \mathbf{b}_N (\mathbf{W}_E and \mathbf{b}_E) are the transformation weights
 145 and bias of the nodes (the hyperedges) for the attention ready state. This is motivated by the observa-
 146 tion that different nodes (genes) contribute differently to the hyperedges (pathways), thus we need
 147 proper attentions across the nodes to up- or down-weight their contributions when aggregating their
 148 representations to compute the representation of the hyperedge. Once we have the hyperedge attention
 149 over nodes, we calculate the hyperedge’s representation in layer k from the nodes’ representations in
 150 layer $k - 1$ (equation 4) where σ is the nonlinearity layer (ReLU in our experiment).

$$\mathbf{h}_E^k(p_j) = \sigma\left(\sum_{g_i \in p_j} a_E(p_j, g_i) \mathbf{h}_N^{k-1}(g_i)\right) \quad (4)$$

151 **Node attention over hyperedges.** For a node $g_i \in \mathcal{N}$, in order to update its hidden representation
 152 at layer k , we aggregate the information from its incident hyperedges using the following attention
 153 mechanism. We first calculate the node attention over hyperedges as in

$$a_N(g_i, p_j) = \frac{\exp(\mathbf{c}^T \mathbf{s}(p_j, g_i))}{\sum_{p_{j'} \ni g_i} \exp(\mathbf{c}^T \mathbf{s}(p_{j'}, g_i))} \quad (5)$$

154 where \mathbf{c} is the same trainable context vector as used in hyperedge attention calculation and the attention
 155 ready state $\mathbf{s}(p_j, g_i)$ for hyperedge-node pair (p_j, g_i) is calculated as in equation 3. This allows us to

156 have the node’s attention over the hyperedges, i.e., we can weight hyperedges’ contributions when
 157 aggregating their representations to compute the representation of the node. We calculate the node’s
 158 representation in layer k from the hyperedges’ representations in layer $k - 1$ as in

$$\mathbf{h}_N^k(g_i) = \sigma \left(\sum_{p_j \ni g_i} a_N(g_i, p_j) \mathbf{h}_E^{k-1}(p_j) \right) \quad (6)$$

159 Note that different from HyperGAT, here the calculation of hyperedge’s and node’s attentions share
 160 the same underlying dual-attention matrix as shown in Fig. 1 (b), which is essentially unstandardized
 161 covariance matrix. Such parameter sharing across hyperedges and nodes allows us to cross-regulate
 162 the learning of their mutual attentions to prevent overfitting. This difference not only allows for
 163 simplification of the model, but also is more consistent with the notion of duality. The dual \mathcal{H}^* of
 164 the hypergraph \mathcal{H} is a hypergraph with \mathcal{H} ’s vertices and edges interchanged, and we should have
 165 $(\mathcal{H}^*)^* = \mathcal{H}$. It is easily provable that the dual-attentions for $(\mathcal{H}^*)^*$ is the same as those for \mathcal{H} .
 166 Such a self-dual statement is generally not true for the attentions proposed in HyperGAT due to
 167 their unsymmetrical way of calculating the node-level and the edge-level attentions, despite that the
 168 HyperGAT attention was termed as “dual” attention. For this reason, we term our attention message
 169 passing scheme as strongly dual attention message passing.

170 **Hypergraph regularization.** One important intuition about graph and hypergraph convolutional
 171 network is that the learned representations for nodes with similar context of (hyper)edges should
 172 be similar. In the case of a simple graph $G = (V, E)$, this is to minimize the summed distance
 173 $\sum_{(u,v) \in E} \|h_u - h_v\|^2$ or its weighted variants. Instead of using it as an explicit regularizer, graph or
 174 hypergraph convolutional networks leverage an appropriately defined graph or hypergraph Laplacian.
 175 As noted in [19], the hypergraph Laplacian is $\Delta = \mathbf{I} - \Theta$ where \mathbf{I} is the identity matrix and Θ is
 176 defined as (let \mathbf{W} be a diagonal matrix with diagonal entries as hyperedge weights)

$$\Theta = \mathbf{D}_v^{-1/2} \mathbf{H} \mathbf{W} \mathbf{D}_e^{-1} \mathbf{H}^T \mathbf{D}_v^{-1/2} \quad (7)$$

177 Here, different from hypergraph convolutional networks, we use explicit regularization on the
 178 similarity of representations of nodes with similar hyperedge context. Let \mathbf{X} be the matrix of the
 179 learned nodes’ representations, where row $\mathbf{X}_i = \mathbf{h}_N^K(g_i)$ and the K th layer is the last hypergraph
 180 message passing layer. We can define the regularizer as

$$\mathcal{L}_{reg} = \sum_{i,j} ((\mathbf{X}_i \mathbf{X}_i^T - 2\mathbf{X}_i \mathbf{X}_j^T + \mathbf{X}_j \mathbf{X}_j^T) * \Theta_{ij}) \quad (8)$$

181 Intuitively, the more hyperedges are incident on the node pair i, j , the more we should penalize their
 182 representational differences. On the other hand, the regularizer down-weights the penalization if
 183 a hyperedge connects many nodes or if a node has many incident hyperedges, indicating lack of
 184 specificity for hyperedges and nodes, respectively.

185 3.4 Weighted Subgraph Attention

186 The multiple layers of strongly dual attention message passing allow learning the nodes’ and hyper-
 187 edges’ representations. However, the instance for the classification algorithm is a subgraph (e.g., a
 188 patient, who has mutations in multiple genes (nodes)). From the hypergraph perspective, a patient
 189 $j (1 \leq j \leq n)$ can be considered as a subhypergraph \mathcal{G}_j whose nodes (genes) have mutations in j and
 190 are a subset of those of \mathcal{H} . This is shown in Fig. 1 (c) where different node colors in the hypergraph
 191 denote different patients. In order to calculate the subgraph’s representation from its component
 192 nodes’ representations at the K th layer, we use the following weighted subgraph attention (WSA)
 193 mechanism. In fact, none of the previous hypergraph methods support subgraph inferring, and we
 194 had to add our WSA module to those models for subgraph inferring as well. We first compute the
 195 subgraph attention over nodes (e.g., g_i ’s) as in

$$a(\mathcal{G}_j, g_i) = \frac{\exp(\mathbf{M}_{j_i} \mathbf{b}^T \mathbf{h}_N^K(g_i))}{\sum_{g_{i'} \in \mathcal{G}_j} \exp(\mathbf{M}_{j_{i'}} \mathbf{b}^T \mathbf{h}_N^K(g_{i'}))} \quad (9)$$

196 where \mathbf{b} is a trainable context vector, \mathbf{M} is the mutation rate feature matrix with each row corre-
 197 sponding to a patient and each column corresponding to a gene. Thus, equation 9 is a mutation rate
 198 weighted subgraph attention mechanism. This choice conforms to the intuition that the rate of a
 199 mutation is more informative than a categorical indicator of the mutation’s occurrence. With these

200 subgraph level attentions, we compute the patient (subgraph) representation from the K th layer’s
 201 gene representations as in

$$\mathbf{h}(\mathcal{G}_j) = \sigma \left(\sum_{g_i \in \mathcal{G}_j} a(\mathcal{G}_j, g_i) \mathbf{h}_N^K(g_i) \right) \quad (10)$$

202 We then stacked the learned patient representations to form the new patient feature matrix, as in

$$\mathbf{S} = [\mathbf{h}(\mathcal{G}_1)^T \mid \mathbf{h}(\mathcal{G}_2)^T \mid \dots \mid \mathbf{h}(\mathcal{G}_n)^T]^T \quad (11)$$

203 where each row is a patient (subgraph) embedding.

204 3.5 Inductive Classification on Subgraphs

205 Let the learned feature matrix be \mathbf{S} and feed it into a softmax classifier

$$\mathbf{Z} = \text{softmax}(\mathbf{W}^{(1)} (\text{ReLU} \circ \text{FC})^{(2)}(\mathbf{S}) + \mathbf{W}^{(0)}) \quad (12)$$

206 where (2) in the superscript indicates two MLP layers (FC =Fully Connected layer). The loss function
 207 is defined as the cross-entropy error over all subjects in all classes as in

$$\mathcal{L} = - \sum_{j \in \mathcal{Y}_D} \sum_{f=1}^F \mathbf{Y}_{jf} \ln \mathbf{Z}_{jf} + \mathcal{L}_{reg} \quad (13)$$

208 where \mathcal{Y}_D is the training set of subjects that have labels and F is the dimension of the output labels,
 209 which is equal to the number of classes. \mathbf{Y} is the label indicator matrix. Note that the subgraph
 210 attention layer allows us to compute any patient’s representation, which effectively eliminates the
 211 need for access to test set patient features during training, making the model inductive. Existing
 212 models such as HyperGCN and HGNN are transductive, we cascade our subgraph attention layer
 213 on top of these models and make them inductive to serve as our comparison models. We implement
 214 SHINE on PyTorch, and run it on NVIDIA V100 GPUs. We train SHINE for up to 6000 epochs
 215 using Adam [30] and stop training if the validation loss does not decrease for 10 consecutive epochs.

216 4 Experiments

217 We conducted experiments on real-world datasets in genetic medicine. Both datasets have more than
 218 20 different classes, indicating significant complexity of the prediction tasks. These datasets are
 219 different in nature, e.g., curated from literature vs. obtained directly from high-throughput sequencing,
 220 and multi-class vs. multi-class multi-label. Our experiments are motivated by the fact that massive
 221 genomic data call for novel methods and present unique technical challenges, in this case inductive
 222 subgraph inferencing on hypergraph. The summary statistics for each dataset is shown in Table 2,
 223 and the description of each dataset is as follows. Most pathways have small to medium sizes, see
 224 Table 2 pathway sizes IQR. In fact, even for 95 percentile, the pathway size is just over 200. On the
 225 other hand, we observed that the larger the pathway (hyperedge), the more subgraph it is incident
 226 on, and the less attention our model will give it as a discriminative feature. The DisGeNet and the
 227 TCGA-MC3 datasets are publicly available and this study is approved by Institutional Review Board.

228 4.1 Disease Type Prediction with DisGeNet Data

229 In this experiment, we have used the DisGeNet dataset [31] that is a collection of mutated genes
 230 involved in human diseases compiled from expert curated repositories, GWAS catalogs, animal
 231 models and the scientific literature. In the following text, we abuse terminology to use “gene” to
 232 really mean “variants in the gene”. We model genes as hypergraph nodes and diseases as hyperedges.
 233 Each disease is labeled with one or more of 22 MeSH codes, and the task is a multi-class multi-label
 234 classification problem. We used 6:2:2 train:validation:test partition, and the split distribution is shown
 235 in the appendix. The DisGeNet dataset has 6226 pathways and 9133 genes involved in 8383 diseases.

236 4.2 Cancer Type Prediction with NGS Somatic Mutations Data

237 In this experiment, we have used the consensus somatic mutations for TCGA subjects produced by
 238 the Multi-Center Mutation Calling in Multiple Cancers (MC3) project [32]. Aiming to enable robust

Table 2: Real-world hypergraph datasets for subgraph inference used in our work. For hyperedge sizes, shown are medians and interquartile ranges. From the distribution, it is clear that the hyperedge size has a positive skewness.

Dataset	# hypernodes	# hyperedges	Hyperedge size	# classes	# subgraphs
DisGeNet	9133	6226	25 (12 - 57)	22	8383
TCGA-MC3	18059	6229	33 (15 - 77)	25	9012

239 cross-tumor-type analyses, the MC3 approach applied an ensemble of 7 mutation-calling algorithms
 240 and assigned a PASS identifier to a mutation that was called by 2 or more variant callers out of the total
 241 7 callers [33]. The MC3 approach accounted for variance and batch effects introduced by the rapid
 242 advancement of DNA extraction, hybridization-capture, and sequencing over time. Following this
 243 approach, we restricted our analysis to PASS calls in order to maintain sample sizes and uniformity
 244 in mutation calling. Each subject is labeled with one of 25 cancer types, and the task is a multi-class
 245 classification problem. We used 6:2:2 train:validation:test partition, stratified by cancer types, and the
 246 split distribution is shown in the appendix. The TCGA-MC3 dataset has 6229 pathways and 18059
 247 genes involved in 9012 subjects in total.

248 4.3 Baselines

249 We compared SHINE with the following state-of-the-art baselines, we use validation datasets to tune
 250 parameters and hyperparameters, please see the appendix for details.

- 251 • Hypergraph neural networks (**HGNN**) [2] uses clique expansion to transform the hypergraph
 252 to graph and computes the graph Laplacian. HGNN then uses Chebyshev approximation to
 253 derive a simplified hypergraph convolution operation.
- 254 • **HyperGCN** [3] constructed Hypergraph Laplacian by reducing the hypergraph to graph and
 255 representing a hyperedge by a selected pairwise simple edge connecting two most unlike
 256 nodes, and adds the remaining nodes in the hyperedge as mediators.
- 257 • Hypergraph attention networks (**HyperGAT**) [6] learns node representations by aggregating
 258 information from nodes to edges and vice versa. Different from SHINE, HyperGAT’s
 259 alternating attention mechanism is not strictly dual attention. Moreover, there was no regu-
 260 larization that nodes with similar context of hyperedges should have similar representations.
- 261 • Polygenic risk score (**PRS**) is currently a standard practice widely used in genetic medicine
 262 and calculates liability to a disease according to their genotype profile using regression.
 263 Ridge regression is a frequent choice for PRS, as adopted in our experiments [34].
- 264 • Non-negative matrix factorization (**NMF**) is a dimensionality reduction tool with much
 265 success in discovering low-dimensional structure from high-dimensional omics data and
 266 enabling inference of complex biological processes [35].
- 267 • **XGBoost** XGBoost is a scalable end-to-end tree boosting system and a state-of-the-art
 268 machine learning method [36] that frequently achieves the top results on many machine
 269 learning challenges.

270 To assess whether performance changes are due to added information (e.g., pathway information)
 271 and/or better utilization of the added information, we run PRS, NMF, and XGBoost in the following
 272 three settings: gene features only, pathway features only, and both gene and pathway features.

273 5 Results

274 The held-out test set micro-averaged F1 scores (micro-F1) for our proposed method SHINE and
 275 all the other comparison models are in Table 3. Comparing all the models, we can see that SHINE
 276 clearly outperforms a comprehensive array of state-of-the-art baselines in various configurations, with
 277 non-overlapping standard deviation intervals. PRS is indeed a competitive baseline, as can be seen
 278 from its close performance compared with XGBoost that frequently topped many machine learning
 279 challenges’ leaderboards. Previously state-of-the-art hypergraph neural network models (HyperGCN,
 280 HGNN, HyperGAT) do not always outperform PRS and XGBoost (e.g., on the TCGA-MC3 dataset).
 281 On the other hand, pathway as features do improve performance, whether alone or jointly with genes.

Table 3: Held-out test set micro-F1 on real-world datasets. Standard deviations are provided from runs with 10 random seeds. SHINE significantly outperforms all the state-of-the-art comparison models. PRS: polygenic risk score. NMF: non-negative matrix factorization. Best model in bold.

Model <i>Metrics</i>	Feature	DisGeNet Dataset Test Micro F1	TCGA-MC3 Dataset Test Micro F1
PRS	gene	0.6303	0.4981
PRS	pathway	0.6461	0.5047
PRS	gene+pathway	0.6512	0.5042
XGBoost	gene	0.6259 ± 0.0012	0.4927 ± 0.0058
XGBoost	pathway	0.6467 ± 0.0035	0.4936 ± 0.0092
XGBoost	gene+pathway	0.6486 ± 0.0036	0.5117 ± 0.0084
NMF	gene	0.6167 ± 0.0040	0.4181 ± 0.0125
NMF	pathway	0.5867 ± 0.0039	0.4842 ± 0.0057
NMF	gene+pathway	0.5847 ± 0.0045	0.4839 ± 0.0032
HyperGCN	gene+pathway	0.6638 ± 0.0028	0.4384 ± 0.0095
HGNN	gene+pathway	0.6809 ± 0.0027	0.4504 ± 0.0042
HyperGAT	gene+pathway	0.6495 ± 0.0050	0.4721 ± 0.0032
SHINE	gene+pathway	0.6955 ± 0.0034	0.5319 ± 0.0049

282 This comparison shows that genetic pathway information is useful to disease type classification,
283 consistent to the intuition that pathways encode molecular functional mechanisms that underlie the
284 disease etiology. However, properly utilizing such information is non-trivial, as evidenced by the
285 difficulty to outperform PRS and XGBoost models by NMF models and hypergraph models including
286 HyperGCN, HGNN and HyperGAT. Given that difficulty, SHINE still attained the best performance
287 on each dataset. Intuitively speaking, HyperGCN and HGNN focus on similarity regularization: hy-
288 pergraph nodes with similar context of hyperedges should have similar representations. HyperGAT’s
289 attention mechanism gears more towards minimizing the classification loss. SHINE, in an attempt
290 to balance the similarity regularization with the end-to-end classification task via strongly dual
291 attention mechanism, achieved better trade-off between the two objectives and effectively integrated
292 the functional pathway’s (hyperedge) information with individual gene’s (node) information. Also
293 note that in general, we have some performance drop when moving from the DisGeNet dataset to
294 the TCGA-MC3 dataset, likely due to the fact that the former uses genetic features from curated
295 literatures and the latter is from high throughput sequencing intended for data-driven discovery.
296 Both complex classification tasks (>20 classes) are uniquely challenging because diseases may
297 have overlapping disrupted molecular functions (genetic pathways, hyperedges), especially for the
298 TCGA-MC3 experiment that is distinguishing subcategories of similar diseases as they are loosely
299 all cancers. In addition, both tasks exhibit class distributional shift between the train and the test
300 datasets, as shown in Appendix Tables 1 and 2, and have been designed to require inductive inference
301 on subgraphs with highly variable hyperedge sizes. Strong performance of SHINE on these tasks thus
302 suggests that our model can leverage its relational inductive biases for more robust generalization.

303 **Model interpretation.** SHINE simultaneously learns the representations of nodes and hyperedges,
304 which are then used to learn and inductively infer subgraph representations. This brings model
305 interpretation advantages as it allows assessing pathways (hyperedges) correlations and reasoning
306 multiple molecular functions mutually interacting and collectively contributing to the disease onset.
307 We identify the top pathways that are enriched in different cancers using the attention weights
308 learned for SHINE, as shown in Table 4. From the table, we see that many of the listed pathways
309 reflect innate key events in the development of individual or multiple types of cancers, consistent
310 with genetic and medical knowledge from wet lab (e.g., TNF/Stress Related Signaling [37]). We
311 showcase interpretations for breast cancer and lung cancer here, and refer the reader to the appendix
312 for full interpretation of Table 4. For breast cancer, TNF α is not only closely involved in its onset,
313 progression and in metastasis formation, but also linked to therapy resistance [37]. Regarding the
314 4-1BB pathway, studies have suggested HER2/4-1BB bispecific molecule as a candidate of alternative
315 therapeutic strategy to patients in HER2-positive breast cancer [38]. VIP/PACAP and their receptors
316 have prominent roles in transactivation of the Epidermal growth factor (EGF) family and growth
317 effects in breast cancer [39]. For lung cancer, the ErbB3 receptor recycling controlled by neuroregulin
318 receptor degradation protein-1 is linked to lung cancer and small inhibitory RNA (siRNA) to ErbB3

Table 4: Top enriched genetic pathways associated with different cancer risks. The text color indicates the source database for pathways that MSigDB integrated: [BioCarta](#), Reactome, [WikiPathways](#), [Pathway Interaction Database](#), [KEGG](#).

BRCA	LUAD	LGG	HNSC
Stress pathway	PTK6 stabilizes HIF1 α	Citrate cycle TCA cycle	Apoptotic factor response
4-1BB pathway	ErbB3 pathway	Cytosine methylation	Programmed cell death
VIP pathway	Hypertrophic cardiomyopathy	TCA cycle and deficiency of pyruvate dehydrogenase	MECP2 regulates neuronal receptors and channels
CD40 pathway	Diseases of metabolism	Glutathione metabolism	FRA pathway
TOLL pathway	TFAP2 regulates growth factors transcription	Digestion of dietary carbohydrate	Caspase activation via extrinsic apoptotic signalling

319 shows promise as a therapeutic approach to treatment of lung adenocarcinoma [40]. Lung cancer is
 320 also modulated by multiple miRNAs interacting with the TFAP2 family [41].

321 6 Discussion, Limitation and Future Work

322 In addition to being significantly more accurate and interpretable, SHINE uses inductive subgraph
 323 inferencing that works well with minibatch, and scales well to large scale problems, as showcased
 324 by real-world experiments. On the other hand, our work comes with limitations. We assumed that
 325 the hyperedges are known in advance. However, in reality as our domain knowledge increases and
 326 evolves, we need to account for unknown hyperedges and, better, simultaneously discover novel
 327 hyperedges from data while predicting disease classes. Such a task has important clinical utilities in
 328 genetic medicine to discover new genetic pathways that may underlie disease etiology, and will be
 329 our future work. Another line of future work is to derive a hypergraph coarsening model on top of
 330 SHINE. SHINE currently has flat hypergraph layout and does not learn hierarchical representations
 331 of hypergraphs. In genetic medicine, often a group of hotspot mutations occur in a limited set of
 332 genes, displaying locality patterns. A flexible hypergraph coarsening model that can effectively
 333 learn hierarchical network structure out of the hypergraphs can shed light on the organizations of the
 334 hyperedges (e.g., pathways representing synergistic molecular functions). From the application point
 335 of view, detecting tumor subtypes is often interesting, and we expect to apply our method to such
 336 detections when large shared dataset will become available. To certain extent, the TCGA labels we
 337 used reflect subtypes of organ-specific primary tumors, e.g., LUAD vs. LUSC in lung cancer, KIRC
 338 vs. KIRP in kidney cancer. On the other hand, identifying drivers genes and pathways for cancer
 339 types (as experimented in this paper) continue to be biologically important, as shown in [33].

340 7 Conclusions

341 We proposed a novel framework termed SubHypergraph Inductive Neural nEtnetwork (SHINE) for
 342 inductive subgraph inferencing on hypergraphs, designed for jointly optimizing the objectives of
 343 end-to-end subgraph classification and similarity regularization for representations of hypergraph
 344 nodes with similar context of hyperedges. We showed that SHINE improved the performance
 345 (micro-F1) of the learned model for disease type prediction for complex (>20 classes) genetic
 346 medicine datasets of different characteristics and under different settings (e.g., multi-class and/or
 347 multi-label). Genetic pathways directly correspond to molecular mechanisms and functions, which
 348 are more informative than individual genes and are represented as hyperedges in SHINE. The novel
 349 formulation of disease classification as a subgraph inferencing problem allows a hypergraph neural
 350 network to link correlated pathways, i.e., interacting molecular mechanisms, to disease etiology. This
 351 leads to better performance with added interpretability. We compared SHINE with a wide array
 352 of state-of-the-art hypergraph neural networks, XGBoost, NMF, and PRS models with different
 353 configurations of genes and pathways as features. SHINE consistently outperformed all state-of-the-
 354 art baselines significantly in each of the disease classification and cancer classification tasks. Feature
 355 analysis of the learned pathway groups that are automatically identified by SHINE in a data-driven
 356 fashion offered significant clinical insights about multiple molecular mechanisms that interact and
 357 are associated with disease types and status.

References

- 358
- 359 [1] Arthur Liberzon, Chet Birger, Helga Thorvaldsdóttir, Mahmoud Ghandi, Jill P Mesirov, and
360 Pablo Tamayo. The molecular signatures database hallmark gene set collection. *Cell systems*, 1
361 (6):417–425, 2015.
- 362 [2] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural
363 networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages
364 3558–3565, 2019.
- 365 [3] Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha
366 Talukdar. Hypergcn: A new method of training graph convolutional networks on hypergraphs.
367 *arXiv preprint arXiv:1809.02589*, 2018.
- 368 [4] Naganand Yadati. Neural message passing for multi-relational ordered and recursive hyper-
369 graphs. *Advances in Neural Information Processing Systems*, 33, 2020.
- 370 [5] Ruochi Zhang, Yuesong Zou, and Jian Ma. Hyper-sagnn: a self-attention based graph neural
371 network for hypergraphs. In *International Conference on Learning Representations*, 2020. URL
372 <https://openreview.net/forum?id=ryeHuJBtPH>.
- 373 [6] Kaize Ding, Jianling Wang, Jundong Li, Dingcheng Li, and Huan Liu. Be more with less: Hy-
374 pergraph attention networks for inductive text classification. *arXiv preprint arXiv:2011.00387*,
375 2020.
- 376 [7] Changping Meng, S Chandra Mouli, Bruno Ribeiro, and Jennifer Neville. Subgraph pattern neu-
377 ral networks for high-order graph evolution prediction. In *Proceedings of the AAAI Conference*
378 *on Artificial Intelligence*, volume 32, 2018.
- 379 [8] Emily Alsentzer, Samuel G Finlayson, Michelle M Li, and Marinka Zitnik. Subgraph neural
380 networks. *arXiv preprint arXiv:2006.10538*, 2020.
- 381 [9] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
382 graphs. In *NIPS*, pages 1024–1034, 2017.
- 383 [10] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure
384 Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *ACM*
385 *SIGKDD*, pages 974–983. ACM, 2018.
- 386 [11] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel,
387 Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning
388 molecular fingerprints. In *NeurIPS*, pages 2224–2232, 2015.
- 389 [12] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
390 message passing for quantum chemistry. In *ICML*, pages 1263–1272. JMLR. org, 2017.
- 391 [13] Liang Yao, Chengsheng Mao, and Yuan Luo. Graph convolutional networks for text classifica-
392 tion. In *AAAI*, 2019.
- 393 [14] Hao Peng, Jianxin Li, Yu He, Yaopeng Liu, Mengjiao Bao, Lihong Wang, Yangqiu Song,
394 and Qiang Yang. Large-scale hierarchical text classification with recursively regularized deep
395 graph-cnn. In *WWW*, pages 1063–1072, 2018.
- 396 [15] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
397 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 398 [16] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for
399 sets. *arXiv preprint arXiv:1511.06391*, 2015.
- 400 [17] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L Hamilton, and Jure
401 Leskovec. Hierarchical graph representation learning with differentiable pooling. *arXiv preprint*
402 *arXiv:1806.08804*, 2018.
- 403 [18] William L Hamilton. Graph representation learning. *Synthesis Lectures on Artificial Intelligence*
404 *and Machine Learning*, 14(3):1–159, 2020.

- 405 [19] Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. Learning with hypergraphs: Clus-
 406 tering, classification, and embedding. *Advances in neural information processing systems*, 19:
 407 1601–1608, 2006.
- 408 [20] Taisong Jin, Liujuan Cao, Baochang Zhang, Xiaoshuai Sun, Cheng Deng, and Rongrong Ji.
 409 Hypergraph induced convolutional manifold networks. In *IJCAI*, pages 2670–2676, 2019.
- 410 [21] Jianwen Jiang, Yuxuan Wei, Yifan Feng, Jingxuan Cao, and Yue Gao. Dynamic hypergraph
 411 neural networks. In *IJCAI*, pages 2635–2641, 2019.
- 412 [22] Sai Nageswar Satchidanand, Harini Ananthapadmanaban, and Balaraman Ravindran. Extended
 413 discriminative random walk: A hypergraph approach to multi-view multi-relational transductive
 414 learning. In *IJCAI*, pages 3791–3797, 2015.
- 415 [23] Fuli Feng, Xiangnan He, Yiqun Liu, Liqiang Nie, and Tat-Seng Chua. Learning on partial-order
 416 hypergraphs. In *Proceedings of the 2018 World Wide Web Conference*, pages 1523–1532, 2018.
- 417 [24] Ke Tu, Peng Cui, Xiao Wang, Fei Wang, and Wenwu Zhu. Structural deep embedding for
 418 hyper-networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32,
 419 2018.
- 420 [25] Kexin Huang and Marinka Zitnik. Graph meta learning via local subgraphs. *Advances in Neural
 421 Information Processing Systems*, 33, 2020.
- 422 [26] Qingyun Sun, Jianxin Li, Hao Peng, Jia Wu, Yuanxing Ning, Phillip S Yu, and Lifang He. Sugar:
 423 Subgraph neural network with reinforcement pooling and self-supervised mutual information
 424 mechanism. *arXiv preprint arXiv:2101.08170*, 2021.
- 425 [27] Muhan Zhang, Zhicheng Cui, Shali Jiang, and Yixin Chen. Beyond link prediction: Predicting
 426 hyperlinks in adjacency space. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 427 volume 32, 2018.
- 428 [28] David Croft, Gavin O’Kelly, Guanming Wu, Robin Haw, Marc Gillespie, Lisa Matthews,
 429 Michael Caudy, Phani Garapati, Gopal Gopinath, Bijay Jassal, et al. Reactome: a database of
 430 reactions, pathways and biological processes. *Nucleic acids research*, 39(suppl_1):D691–D697,
 431 2010.
- 432 [29] Minoru Kanehisa, Miho Furumichi, Mao Tanabe, Yoko Sato, and Kanae Morishima. Kegg:
 433 new perspectives on genomes, pathways, diseases and drugs. *Nucleic acids research*, 45(D1):
 434 D353–D361, 2017.
- 435 [30] DP Kingma and JL Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- 436 [31] Janet Piñero, Àlex Bravo, Núria Queralt-Rosinach, Alba Gutiérrez-Sacristán, Jordi Deu-Pons,
 437 Emilio Centeno, Javier García-García, Ferran Sanz, and Laura I. Furlong. DisGeNET: a
 438 comprehensive platform integrating information on human disease-associated genes and variants.
 439 *Nucleic Acids Research*, 45(D1):D833–D839, 10 2016. ISSN 0305-1048. doi: 10.1093/nar/
 440 gkw943. URL <https://doi.org/10.1093/nar/gkw943>.
- 441 [32] Kyle Ellrott, Matthew H Bailey, Gordon Saksena, Kyle R Covington, Cyriac Kandath, Chip
 442 Stewart, Julian Hess, Singer Ma, Kami E Chiotti, Michael McLellan, et al. Scalable open
 443 science approach for mutation calling of tumor exomes using multiple genomic pipelines. *Cell
 444 systems*, 6(3):271–281, 2018.
- 445 [33] Matthew H Bailey, Collin Tokheim, Eduard Porta-Pardo, Sohini Sengupta, Denis Bertrand,
 446 Amila Weerasinghe, Antonio Colaprico, Michael C Wendl, Jaegil Kim, Brendan Reardon, et al.
 447 Comprehensive characterization of cancer driver genes and mutations. *Cell*, 173(2):371–385,
 448 2018.
- 449 [34] Shing Wan Choi, Timothy Shin-Heng Mak, and Paul F O’Reilly. Tutorial: a guide to performing
 450 polygenic risk score analyses. *Nature Protocols*, 15(9):2759–2772, 2020.

- 451 [35] Genevieve L Stein-O'Brien, Raman Arora, Aedin C Culhane, Alexander V Favorov, Lana X
452 Garmire, Casey S Greene, Loyal A Goff, Yifeng Li, Aloune Ngom, Michael F Ochs, et al. Enter
453 the matrix: factorization uncovers knowledge from omics. *Trends in Genetics*, 34(10):790–805,
454 2018.
- 455 [36] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of*
456 *the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages
457 785–794, 2016.
- 458 [37] María Florencia Mercogliano, Sofia Bruni, Patricia V Elizalde, and Roxana Schillaci. Tumor
459 necrosis factor α blockade: an opportunity to tackle breast cancer. *Frontiers in oncology*, 10,
460 2020.
- 461 [38] Marlon J Hinner, Rachida Siham Bel Aiba, Thomas J Jaquin, Sven Berger, Manuela Carola Dürr,
462 Corinna Schlosser, Andrea Allersdorfer, Alexander Wiedenmann, Gabriele Matschiner, Julia
463 Schüller, et al. Tumor-localized costimulatory t-cell engagement by the 4-1bb/her2 bispecific
464 antibody-anticalin fusion prs-343. *Clinical Cancer Research*, 25(19):5878–5889, 2019.
- 465 [39] Terry W Moody, Bernardo Nuche-Berenguer, and Robert T Jensen. Vip/pacap, and their
466 receptors and cancer. *Current opinion in endocrinology, diabetes, and obesity*, 23(1):38, 2016.
- 467 [40] G Sithanandam and LM Anderson. The erbb3 receptor in cancer and cancer gene therapy.
468 *Cancer gene therapy*, 15(7):413–448, 2008.
- 469 [41] Damian Kołat, Żaneta Kałuzińska, Andrzej K Bednarek, and Elżbieta Płuciennik. The biological
470 characteristics of transcription factors ap-2 α and ap-2 γ and their importance in various types of
471 cancers. *Bioscience reports*, 39(3), 2019.

472 Checklist

- 473 1. For all authors...
- 474 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
475 contributions and scope? [Yes]
- 476 (b) Did you describe the limitations of your work? [Yes]
- 477 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 478 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
479 them? [Yes]
- 480 2. If you are including theoretical results...
- 481 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 482 (b) Did you include complete proofs of all theoretical results? [N/A]
- 483 3. If you ran experiments...
- 484 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
485 mental results (either in the appendixal material or as a URL)? [Yes]
- 486 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
487 were chosen)? [Yes]
- 488 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
489 ments multiple times)? [Yes]
- 490 (d) Did you include the total amount of compute and the type of resources used (e.g., type
491 of GPUs, internal cluster, or cloud provider)? [Yes]
- 492 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 493 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 494 (b) Did you mention the license of the assets? [Yes]
- 495 (c) Did you include any new assets either in the appendixal material or as a URL? [Yes]
- 496 (d) Did you discuss whether and how consent was obtained from people whose data you're
497 using/curating? [Yes] We discuss the IRB approval process in Section 4.

- 498 (e) Did you discuss whether the data you are using/curating contains personally identifiable
499 information or offensive content? [N/A]
- 500 5. If you used crowdsourcing or conducted research with human subjects...
- 501 (a) Did you include the full text of instructions given to participants and screenshots, if
502 applicable? [N/A]
- 503 (b) Did you describe any potential participant risks, with links to Institutional Review
504 Board (IRB) approvals, if applicable? [N/A]
- 505 (c) Did you include the estimated hourly wage paid to participants and the total amount
506 spent on participant compensation? [N/A]