
Towards characterizing the value of edge embeddings in Graph Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph neural networks (GNNs) are the dominant approach to solving machine
2 learning problems defined over graphs. Despite much theoretical and empirical
3 work in recent years, our understanding of finer-grained aspects of architectural
4 design for GNNs remains impoverished. In this paper, we consider the benefits
5 of architectures that maintain and update edge embeddings. On the theoretical
6 front, under a suitable computational abstraction for a layer in the model, as well
7 as memory constraints on the embeddings, we show that there are natural tasks
8 on graphical models for which architectures leveraging edge embeddings can be
9 much shallower. Our techniques are inspired by results on time-space tradeoffs in
10 theoretical computer science. Empirically, we show architectures that maintain edge
11 embeddings almost always improve on their node-based counterparts—frequently
12 significantly so in topologies that have “hub” nodes.

13 1 Introduction

14 Graph neural networks (GNNs) have emerged as the dominant approach for solving machine learning
15 tasks on graphs. Over the span of the last decade, many different architectures have been proposed,
16 both in order to improve different notions of efficiency, and to improve performance on a variety
17 of benchmarks. Nevertheless, theoretical and empirical understanding of the impact of different
18 architectural design choices remains elusive.

19 One previous line of work (Xu et al., 2018) has focused on characterizing the representational
20 limitations stemming from the *symmetry-preserving* properties of GNNs when the node features are
21 not informative (also called “anonymous GNNs”) — in particular, relating GNNs to the Weisfeiler-
22 Lehman graph isomorphism test (Leman & Weisfeiler, 1968). Another line of work (Oono & Suzuki,
23 2019) focuses on the potential pitfalls of the *(over)smoothing effect* of deep GNN architectures, with
24 particular choices of weights and non-linearities, in an effort to explain the difficulties of training
25 deep GNN models. Yet another (Black et al., 2023) focuses on training difficulties akin to vanishing
26 introduced by “*bottlenecks*” in the graph topology.

27 In this paper, we focus on the benefits of maintaining and updating *edge embeddings* over the course
28 of the computation of the GNN. A typical way to parametrize a layer l of a GNN (Xu et al., 2018) is
29 to maintain, for each node v in the graph, a node embedding $h_v^{(l)}$, which is calculated as

$$a_v^{(l+1)} = \text{AGGREGATE}\left(h_u^{(l)} : u \in N_G(v)\right) \quad h_v^{(l+1)} = \text{COMBINE}\left(a_v^{(l+1)}, h_v^{(l)}\right) \quad (1)$$

30 where $N_G(v)$ denotes the neighborhood of vertex v . These updates can be viewed as implementing a
31 (trained) message-passing algorithm, in which nodes pass messages to their neighbors, which are
32 then aggregated and combined with the current state (i.e., embedding) of a node. The initial node
33 embeddings $h_v^{(0)}$ are frequently part of the task specification (e.g., a vector of fixed features that can

34 be associated with each node). When this is not the case, they can be set to fixed values (e.g., the
35 all-ones vector) or random values.

36 But a more expressive way to parametrize a layer of computation is to maintain, for each *edge* e , an
37 edge embedding $h_e^{(l)}$ which is calculated as:

$$a_e^{(l+1)} = \text{AGGREGATE}\left(h_a^{(l)} : a \in M_G(e)\right) \quad h_e^{(l+1)} = \text{COMBINE}\left(a_e^{(l+1)}, h_e^{(l)}\right) \quad (2)$$

38 where $M_G(e)$ denotes the “neighborhood” of edge e : that is, all edges a that share a vertex with e ¹.

39 This paradigm is at least as expressive as (1): we can simulate a layer of (1) by designating the
40 embedding of an edge to be the concatenation of the node embeddings of its endpoints, and noticing
41 that $M_G(e)$ includes all the neighbors of both endpoints of e . In particular, if a task has natural initial
42 node embeddings, then their concatenations along edges can be used as initial edge embeddings.
43 Additionally, there may be tasks where initial features are most naturally associated with edges (e.g.,
44 attributes of the relationship between two nodes) — or the final predictions of the network are most
45 naturally associated with edges (e.g., in link prediction, where we want to decide which potential
46 links are true links).

47 GNNs that fall in the general paradigm of (2) have been used for various applications – including link
48 prediction (Cai et al., 2021; Liang & Pu, 2023) as well as reasoning about relations between objects
49 (Battaglia et al., 2016), molecular property prediction (Gilmer et al., 2017; Choudhary & DeCost,
50 2021), and detecting clusters of communities in graphs (Chen et al., 2017) – with robust empirical
51 benefits. These approaches instantiate the edge-based paradigm in a plethora of ways. However, it is
52 difficult to disentangle to what degree performance improvements come from added information from
53 domain-specific initial edge embeddings, versus properties of the particular architectural choices for
54 the aggregation functions in (2), versus inherent benefits of the edge-based paradigm itself (whether
55 representational, or via improved training dynamics).

56 We focus on *theoretically and empirically* quantifying the added *representational* benefit from
57 maintaining edge embeddings. Viewing the GNN as a computational model, we can think of
58 the intermediate embeddings as a “scratch pad”. Since we maintain more information per layer
59 compared to the node-based paradigm (1), we might intuitively hope to be able to use a shallower
60 edge embedding model. However, formally proving depth lower bounds both for general neural
61 networks (Telgarsky, 2016) and for specific architectures (Sanford et al., 2024b,a) frequently requires
62 non-trivial theoretical insights – as is the case for our question of interest. In this paper, we show that:

- 63 • *Theoretically*, for certain graph topologies, edge embeddings can have substantial *representational*
64 benefits in terms of the depth of the model, when the amount of memory (i.e., total bit complexity)
65 per node or edge embedding is bounded. Our results illuminate some subtleties of using particular
66 lenses to understand design aspects of GNNs: for instance, we prove that taking memory into
67 account reveals depth separations that the classical lens of invariance (Xu et al., 2018) alone cannot.
- 68 • *Empirically*, when given the same input information, edge-based models almost always lead to
69 performance improvements compared to their node-based counterpart — and often by a large
70 margin if the graph topology includes “hub” nodes with high degree.

71 2 Overview of results

72 2.1 Representational benefits from maintaining edge embeddings.

73 Our theoretical results elucidate the representational benefits of maintaining edge embeddings. More
74 precisely, we show that there are natural tasks on graphs that can be solved by a *shallow* model
75 maintaining constant-size edge embeddings, but can only be solved by a model maintaining constant-
76 size node embeddings if it is much *deeper*.

77 To reason about the impact of depth on the representational power of edge-embedding-based and
78 node-embedding-based architectures, we introduce two *local computation models*. In the node-
79 embedding case, we assume each node of the graph G supports a processor that maintains a state
80 with a *fixed amount of memory*. In one round of computation, each node receives messages from the

¹The graph is assumed to be undirected, as is most common in the GNN literature.

81 adjacent nodes, which are aggregated by the node into a new state. In this abstraction, we think of the
 82 memory of the processor as the total bits of information each embedding can retain, and we think
 83 of one round of the protocol as corresponding to one layer of a GNN. The edge-embedding case is
 84 formalized in a similar fashion, except that the processors are placed on the edges of the graph, and
 85 two edge processors are “adjacent” if the edges share a vertex in common. In both cases, the input is
 86 distributed across the edges of the graph, and is only locally accessible.

87 With this setup in mind, our first result focuses on *probabilistic inference* on graphs, specifically, the
 88 task of maximum a-posteriori (MAP) estimation in a pairwise graphical model on a graph $G = (V, E)$.
 89 For this task, given edge attributes describing the pairwise interactions $\phi_{\{a,b\}}$, the goal is to compute
 90 $\arg \max_{x \in \{0,1\}^V} p_\phi(x)$, where $p_\phi(x) \propto \exp(\sum_{\{a,b\} \in E} \phi_{\{a,b\}}(x_a, x_b))$.

91 **Theorem (Informal).** *Consider the task of using a GNN to calculate MAP (maximum a-posteriori)*
 92 *values in a pairwise graphical model, in which the pairwise interactions are given as input embeddings*
 93 *to a node-embedding or edge-embedding architecture. Then, there exists a graph with $O(n)$ vertices*
 94 *and edges, such that:*

- 95 • Any node message-passing protocol with T rounds and $O(1)$ bits of memory per node processor
 96 requires $T = \Omega(\sqrt{n})$.
- 97 • There is an edge message-passing protocol with $O(1)$ rounds and $O(1)$ bits of memory.

98 The proof techniques are of standalone interest: the lower bound on node message-passing protocols is
 99 inspired by tracking the “flow of information” in the graph, reminiscent of graph pebbling techniques
 100 used to prove time-space tradeoffs in theoretical computer science (Grigor’ev, 1976; Abrahamson,
 101 1991). The formal result is Theorem 1, and the proof sketch is included in Section 5.

102 **The view from symmetry.** Above, we are not imposing any *symmetry constraints* – that is,
 103 invariance of the computation at a node or edge to its identity and the identities of its neighbors.
 104 Indeed, the edge message-passing protocol constructed above is highly non-symmetric. However,
 105 we show there is a (different, but also natural) task where *even symmetric* edge message-passing
 106 protocols achieve a better depth/memory tradeoff than node message-passing protocols. We state the
 107 informal result below; the formal result is Theorem 4.

108 **Theorem (Informal).** *Let n be a positive integer. There is a graph G with $O(n)$ vertices and $O(n)$*
 109 *edges, and a computational task on G , such that:*

- 110 • Any node message-passing protocol with T rounds and $O(1)$ bits of memory per node processor
 111 requires $T = \Omega(\sqrt{n})$ to solve this task.
- 112 • There is a symmetric edge message-passing protocol that solves this task with $O(1)$ rounds and
 113 $O(1)$ bits of memory.

114 **Importance of the memory lens.** The memory constraints are crucial for the results above. Without
 115 memory constraints, we can show that the node message-passing architecture can simulate the edge
 116 message-passing architecture, while only increasing the depth by 1 (Proposition 3). Moreover, the
 117 *symmetric* node message-passing architecture can simulate the *symmetric* edge message-passing
 118 architecture, again while only increasing the depth by 1. We state the informal result below; the
 119 formal result is Theorem 5.

120 **Theorem (Informal).** *For any graph G , any symmetric edge message-passing protocol on G with T*
 121 *rounds can be represented by a symmetric node message-passing protocol with $T + 1$ rounds.*

122 We note that unlike prior work that focuses on understanding the representational power of GNN
 123 architectures under symmetry constraints (Xu et al., 2018) — which requires that the initial node
 124 features are the same for all nodes — our simulation theorem above holds for arbitrary choices of
 125 initial node features.

126 We view this as evidence that many fine-grained properties of architectural design for GNNs cannot
 127 be adjudicated by solely considering them through the lens of symmetries of the network.

128 2.2 Empirical benefits of edge-based architectures.

129 The theory, while only characterizing representational power, suggests that architectures that main-
130 tain edge embeddings should have strictly better performance compared to their node embedding
131 counterparts. We verify this in both real-life benchmarks and natural synthetic sandboxes.

132 First, we consider several popular GNN benchmarks (inspired by both predicting molecular properties,
133 and image-like data), and show that equalizing for all other aspects of the architecture (e.g., depth,
134 dimensionality of the embeddings) — the accuracy the edge-based architectures achieve is at least
135 as good as their node-based counterparts. Note, the goal of these experiments is *not* to propose
136 a new architecture — there are already a variety of (very computationally efficient) GNNs that in
137 some manner maintain edge embeddings. The goal is to confirm that — all other things being equal
138 — the representational advantages of edge-based architectures do not introduce additional training
139 difficulties. Details are included in Section 8.1.

140 Next, we consider two synthetic settings to stress test the performance of edge-based architectures.
141 Inspired by the graph topology that provides a theoretical separation between edge and node-based
142 protocols (Theorem 1 and Theorem 4), we consider graphs in which there is a hub node, and tasks
143 that are “naturally” solved by an edge-based architecture. Precisely, we consider a star graph, in
144 which the labels on the leaves are generated by a “planted” edge-based architecture with randomly
145 chosen weights. The node-based architecture, on the other hand, has to pass messages between the
146 leaves indirectly through the center of the star. Empirically, we indeed observe that the performance
147 of edge-based architectures is significantly better. Details are included in Section 8.2

148 Finally, again inspired by the theoretical setting in Theorem 1, we consider probabilistic inference
149 on *tree graphs* — precisely, learning a GNN that calculates node marginals for an Ising model, a
150 pairwise graphical model in which the pairwise interactions are just the product of the end points. An
151 added motivation for this setting is the fact that belief propagation — a natural algorithm to calculate
152 the marginals — can be written as an edge-based message-passing algorithm. Again, empirically we
153 see that edge-based architectures perform at least as well as node-based architectures. This advantage
154 is maintained even if we consider “directed” versions of both architectures, in which case embeddings
155 are maintained to be sent along each direction of the edge, and the message for the outgoing direction
156 of an edge depends only on the embeddings corresponding to the incoming directions of the edges.
157 Details are included in Section 8.3.

158 3 Related Works

159 **The symmetry lens on GNNs:** The most extensive theoretical work on GNNs has concerned itself
160 with the representational power of different GNN architectures, while trying to preserve equivariance
161 (to permuting the neighbors) of each layer. (Xu et al., 2018) connected the expressive power of such
162 architectures to the Weisfeiler-Lehman (WL) test for graph isomorphism. Subsequent works (Maron
163 et al., 2019; Zhao et al., 2021) focused on strengthening the representational power of the standard
164 GNN architectures from the perspective of symmetries—more precisely, to simulate the k -WL test,
165 which for k as large as the size of the graph becomes as powerful as testing graph isomorphism. Our
166 work suggests that this perspective may be insufficient to fully understand the representational power
167 of different architectures.

168 **GNNs as a computational machine:** Two recent papers (Loukas, 2019, 2020) considered properties
169 of GNNs when viewed as “local computation” machines, in which a layer of computation allows a
170 node to aggregate the current values of the neighbors (in an arbitrary fashion, without necessarily
171 considering symmetries). Using reductions from the CONGEST model, they provide lower bounds
172 on width and depth for the standard node-embedding based architecture. However, they do not
173 consider architectures with edge embeddings, which is a focus of our work.

174 **Communication complexity methods to prove representational separations:** Tools from dis-
175 tributed computation and communication complexity have recently been applied not only to under-
176 stand the representational power of GNNs (Loukas, 2019, 2020), but also the representational power
177 of other architectures like transformers (Sanford et al., 2024b,a). In particular, (Sanford et al., 2024a)
178 draws a connection between number of rounds for a MPC (Massively Parallel Computation) protocol,
179 and the depth of attention-based architectures.

180 **GNNs for inference and graphical models:** The paper (Xu & Zou, 2023) considers the approx-
 181 imation power of GNNs for calculating marginals for pairwise graphical models, if the family of
 182 potentials satisfies strong symmetries. They do not consider the role of edge embeddings or memory.

183 4 Setup

184 **Notation.** We will denote the graph associated with the GNN as $G = (V, E)$, denoting the vertex
 185 set as V and the edge set as E . The graph induces adjacency relations on both edges and nodes,
 186 namely for $v, v' \in V$ and $e, e' \in E$, we have: $v \sim v'$ if $\{v, v'\} \in E$; $v \sim e$ if $e = \{u, v\}$ for some
 187 $u \in V$; and $e \sim e'$ if e, e' share at least one vertex. For all graphs considered in this paper, we
 188 assume that $\{v, v\} \in E$ for all $v \in V$, so that adjacency is reflexive. We then define adjacency
 189 functions $N_G : V \cup E \rightarrow V$ and $M_G : V \cup E \rightarrow E$ as $N_G(a) := \{v \in V : a \sim v\}$ and
 190 $M_G(a) := \{e \in E : a \sim e\}$.

191 **Local memory-constrained computation.** In order to reason about the required depth with dif-
 192 ferent architectures, we will define a mathematical abstraction for one layer of computation in the
 193 GNN. We will define two models for local computation, one for each of the edge-embedding and
 194 node-embedding architecture. Unlike much prior work on GNNs and distributed computation, we
 195 will also have *memory* constraints — more precisely, we will constrain the bit complexity of the node
 196 and edge embeddings being maintained.

197 In both models, there is an underlying graph $G = (V, E)$, and the goal is to compute a function
 198 $g : \Phi^E \rightarrow \{0, 1\}^V$, where Φ is the fixed-size *input alphabet*, via several rounds of message-passing
 199 on the graph G . This domain of g is Φ^E because in *both* models, the inputs are given on the edges of
 200 the graph — the node model will just be unable to store any *additional* information on the edges. As
 201 we will see in Section 5, this is a natural setup for probabilistic inference on graphs.

202 In both models, a protocol is parametrized by the number of rounds T required, and the amount
 203 of memory B required per local processor. For notational convenience, for $B \in \mathbb{N}$ we define
 204 $\mathcal{X}_B := \{0, 1\}^B$, i.e. the length- B binary strings. Recall that $N_G(v), M_G(v)$ denote the sets of
 205 vertices and edges adjacent to vertex v in graph G , respectively.

206 **Definition 1** (Node message-passing protocol). Let $T, B \in \mathbb{N}$ and let $G = (V, E)$ be a graph. A
 207 *node message-passing protocol* P on graph G with T rounds and B bits of memory is a collection of
 208 functions $(f_{t,v})_{t \in [T], v \in V}$ where $f_{t,v} : \mathcal{X}_B^{N_G(v)} \times \Phi^{M_G(v)} \rightarrow \mathcal{X}_B$ for all t, v . For an *input* $I \in \Phi^E$,
 209 the *computation* of P at a round $t \in [T]$ is the map $P_t(\cdot; I) : V \rightarrow \mathcal{X}_B$ defined inductively by
 210 $P_t(v; I) := f_{t,v}((P_{t-1}(v'; I))_{v' \in N_G(v)}, (I(e))_{e \in M_G(v)})$ where $P_0 \equiv 0$. We say that P *computes* a
 211 function $g : \Phi^E \rightarrow \{0, 1\}^V$ on inputs $\mathcal{I} \subseteq \Phi^E$ if $P_T(v; I)_1 = g(I)_v$ for all $v \in V$ and all $I \in \mathcal{I}$.

212 In words, the value computed by vertex v at round t is some function of the previous values stored
 213 at the neighbors $v' \in N_G(v)$, as well as possibly the problem inputs on the edges adjacent to v (i.e.
 214 $(I(e))_{e \in M_G(v)}$). Note that $P_t(v; I)$ may indeed depend on $P_{t-1}(v; I)$, due to our convention that
 215 $v \in N_G(v)$. We can define the edge message-passing protocol analogously:

216 **Definition 2** (Edge message-passing protocol). Let $T, B \in \mathbb{N}$ and let $G = (V, E)$ be a graph. An
 217 *edge message-passing protocol* P on graph G with T rounds and B bits of memory is a collection
 218 of functions $(f_{t,e})_{t \in [T], e \in E}$ where $f_{t,e} : \mathcal{X}_B^{M_G(e)} \times \Phi \rightarrow \mathcal{X}_B$ for all t, e , together with a collection
 219 of functions $(\tilde{f}_v)_{v \in [V]}$ where $\tilde{f}_v : \mathcal{X}_B^{M_G(v)} \rightarrow \{0, 1\}$. For an *input* $I \in \Phi^E$, the *computation*
 220 of P at a timestep $t \in [T]$ is the map $P_t(\cdot; I) : E \rightarrow \mathcal{X}_B$ defined inductively by: $P_t(e; I) :=$
 221 $f_{t,e}((P_{t-1}(e'; I))_{e' \in M_G(e)}, I(e))$ where $P_0 \equiv 0$. We say that P *computes* a function $g : \Phi^E \rightarrow$
 222 $\{0, 1\}^V$ on inputs $\mathcal{I} \subseteq \Phi^E$ if $\tilde{f}_v((P_T(e; I))_{e \in M_G(v)}) = g(I)_v$ for all $v \in V$ and all $I \in \mathcal{I}$.

223 **Remark 3** (Relation to distributed computation literature). These models are very related to clas-
 224 sical models in distributed computation like LOCAL (Linial, 1992) and CONGEST (Peleg, 2000).
 225 However, the latter models ignore memory constraints, so we cannot usefully port lower and upper
 226 bounds from this literature.

227 **Remark 4** (Computational efficiency). In the definitions above, we allow the update rules $f_{t,v}, f_{t,e}$
 228 to be arbitrary functions. In particular, a priori they may not be efficiently computable. However, our
 229 results showing a function can be implemented by an edge message-passing protocol (Theorem 1,
 230 Part 2 and Theorem 4, Part 2) in fact use simple functions (computable in linear time in the size

231 of the neighborhood), implying the protocol can be implemented in parallel (with one processor
 232 per node/edge respectively) with parallel time complexity $O(TB \cdot \max_v |M_G(v)|)$. On the other
 233 hand, for the results showing a function cannot be implemented by a node message-passing protocol
 234 (Theorem 1, Part 1 and Theorem 4, Part 1), we prove an impossibility result for a *stronger* model
 235 (one in which the computational complexity of $f_{t,v}$ is unrestricted) — which makes our results only
 236 *stronger*.

237 **Symmetry-constrained protocols.** Typically, GNNs are architecturally constrained to respect the
 238 symmetries of the underlying graph. Below we formalize the most natural notion of symmetry in our
 239 models of computation. Note, our abstraction of a round in the message-passing protocol generalizes
 240 the notion of a layer in a graph neural network—and the abstraction defined below correspondingly
 241 generalizes the standard definition of permutation equivariance (Xu et al., 2018). We use the notation
 242 $\{\!\!\{\}$ to denote a multiset.

243 **Definition 5** (Symmetric node message-passing protocol). A node message-passing protocol $P =$
 244 $(f_{t,v})_{t \in [T], v \in V}$ on graph $G = (V, E)$ is *symmetric* if there are functions $(f_t^{\text{sym}})_{t \in [T]}$ so that for every
 245 $t \in [T]$ and $v \in V$, the function $f_{t,v}$ can be written as:

$$f_{t,v}((c(v'))_{v' \in N_G(v)}, (I(e))_{e \in M_G(v)}) = f_t^{\text{sym}}(c(v), \{\!\!\{(c(v'), I(\{v, v'\})) : v' \in N_G(v)\}\!\!\}).$$

246 **Definition 6** (Symmetric edge message-passing protocol). An edge message-passing protocol $P =$
 247 $((f_{t,e})_{t \in [T], e \in E}, (\tilde{f}_v)_{v \in V})$ on graph $G = (V, E)$ is *symmetric* if there are functions $(f_t^{\text{sym}})_{t \in [T]}$ and
 248 \tilde{f}^{sym} so that for every $t \in [T]$ and $e = \{u, v\} \in E$, the function $f_{t,e}$ can be written as:

$$f_{t,e}((c(e'))_{e' \in M_G(e)}, I(e)) = f_t^{\text{sym}}(I(e), c(e), \{\!\!\{\!\!\{c(\{u, v'\}) : v' \in N_G(u)\}\!\!\}, \{\!\!\{c(\{u', v\}) : u' \in N_G(v)\}\!\!\}\!\!\}),$$

249 and for every $v \in V$, \tilde{f}_v can be written as $\tilde{f}_v((c(e))_{e \in M_G(v)}) = \tilde{f}^{\text{sym}}(\{\!\!\{c(e) : e \in M_G(v)\}\!\!\})$.

250 5 Depth separation between edge and node message passing protocols under 251 memory constraints

252 We will consider a common task in probabilistic inference on a *pairwise graphical model*: calculating
 253 the MAP (maximum a-posterior) configuration.

254 **Definition 7** (Pairwise graphical model). For any graph $G = (V, E)$, the *pairwise graphical model*
 255 on G with potential functions $\phi_{\{a,b\}} : \{0, 1\}^2 \rightarrow \mathbb{R}$ is the distribution $p_\phi \in \Delta(\{0, 1\}^V)$ defined as
 256 $p_\phi(x) \propto \exp(-\sum_{\{a,b\} \in E} \phi_{\{a,b\}}(x_a, x_b))$.

257 **Definition 8** (MAP evaluation). Let $\Phi \subseteq \{\phi : \{0, 1\}^2 \rightarrow \mathbb{R}\}$ be a finite set of potential functions.
 258 A *MAP (maximum a-posteriori) evaluator* for G (with potential function class Φ) is any function
 259 $g : \Phi^E \rightarrow \{0, 1\}^V$ that satisfies $g(\phi) \in \arg \max_{x \in \{0, 1\}^V} p_\phi(x)$ for all $\phi \in \Phi^E$.

260 With this setup in mind, we will show that there exists a pairwise graphical model, and a local function
 261 class Φ , such that an edge message passing protocol can implement MAP evaluation with a constant
 262 number of rounds and a constant amount of memory, while any node message protocol with T rounds
 263 and B bits of memory requires $TB = \Omega(\sqrt{|V|})$. Precisely, we show:

264 **Theorem 1** (Main, separation between node and edge message-passing protocols). *Fix $n \in \mathbb{N}$. There*
 265 *is a graph G with $O(n)$ vertices and $O(n)$ edges, and a function class Φ of size $O(1)$, so that:*

- 266 1. *Let g be any MAP evaluator for G with potential function class Φ . Any node message-passing*
 267 *protocol on G with T rounds and B bits of memory that computes g requires $TB \geq \sqrt{n} - 1$.*
- 268 2. *There is an edge message-passing protocol $(f_{t,e})_{t,e}$ on G with $O(1)$ rounds and $O(1)$ bits of*
 269 *memory that computes a MAP evaluator for G with potential function class Φ . Additionally, for*
 270 *all t, e , the update rule $f_{t,e}$ can be evaluated in $O(|M_G(e)|)$ time.*

271 We provide a proof sketch of the main techniques here, and relegate the full proofs to Appendix A.
 272 The graph G that exhibits the claimed separation is a disjoint union of \sqrt{n} path graphs, with an
 273 additional “hub vertex” that is connected to all other vertices in the graph (Fig. 1). The intuition for
 274 the separation is that MAP estimation requires information to disseminate from one end of each path
 275 to the other, and the hub node is a bottleneck for node message-passing but not edge message-passing.
 276 We expand upon both aspects of this intuition below.

277 **Lower bound for node message-passing protocols:** Our main technical lemma for the first half
 278 of the theorem is Lemma 2. It gives a generic framework for lower bounding the complexity of any
 279 node message-passing protocol that computes some function g , by exhibiting a set of nodes $S \subset V$
 280 where computing g requires large “information flow” from distant nodes. More precisely, for any
 281 fixed set of “bottleneck nodes” K , consider the radius- T neighborhood of S when K is removed
 282 from the graph. In any T -round protocol, input data from outside this neighborhood can only reach
 283 S by passing through K . But the total number of bits of information computed by K throughout
 284 the protocol is only $TB|K|$. This gives a bound on the number of values achievable by g on S . We
 285 formalize this argument below (proof in Appendix A):

286 **Lemma 2.** *Let $G = (V, E)$ be a graph. Let P be a node message-passing protocol on G with T
 287 rounds and B bits of memory, which computes a function $g : \Phi^E \rightarrow \{0, 1\}^V$. Pick any disjoint sets
 288 $K, S \subseteq V$. Define $H := G[\bar{K}], F := M_G(N_H^{T-1}(S))$.*

289 *Then: $TB \geq \frac{1}{|K|} \log \max_{I_F \in \Phi^F} \left| \left\{ g_S(I_F, I_{\bar{F}}) : I_{\bar{F}} \in \Phi^{\bar{F}} \right\} \right|$.*

290 **Remark 9.** The proof technique is inspired by and related to classic techniques (specifically, Grig-
 291 oriev’s method) for proving time-space tradeoffs for restricted models of computation like branching
 292 programs ((Grigor’ev, 1976), see Chapter 10 in Savage (1998) for a survey). There, one defines
 293 the “flow” of a function, which quantifies the existence of subsets of coordinates, such that setting
 294 them to some value, and varying the remaining variables results in many possible outputs. In our
 295 case, the choice of subsets is inherently tied to the topology of the graph G . Our technique is also
 296 inspired by and closely related to the “light cone” technique for proving round lower bounds in the
 297 LOCAL computation model (Linial, 1992). However, our technique takes advantage of bottlenecks
 298 in the graph to prove stronger lower bounds (which would be impossible in the LOCAL model where
 299 memory constraints are ignored).

300 The proof of Part 1 of Theorem 1 now follows from an application of Lemma 2 with a particular
 301 choice of K and S . Specifically, we choose K to be the “hub” node (i.e. $K = \{0\}$) and S to be the
 302 set of left endpoints of each path. To show that any MAP evaluator has large information flow to
 303 S (in the quantitative sense of Lemma 2), it suffices to observe that in a pairwise graphical model
 304 on G where a different external field is applied to the right endpoint of each path, and all pairwise
 305 interactions along paths are positive, the MAP estimate on each vertex in S must match the external
 306 field on the corresponding right endpoint.

307 **Upper bound for edge message-passing protocols:** The key observation for constructing a
 308 constant-round edge message-passing protocol for MAP estimation on G is that all of the input data
 309 can be collected on the edges adjacent to the hub vertex. At this point, every such edge has access to
 310 all of the input data, and hence can evaluate the function. If G were an arbitrary graph, this final step
 311 would potentially be NP-hard. However, since the induced subgraph after removing the hub vertex
 312 is a disjoint union of paths, in fact there is a linear-time dynamic programming algorithm for MAP
 313 estimation on G (Lemma 6). This completes the proof overview for Theorem 1.

314 The separation discussed above crucially relies on the existence of a high-degree vertex in G . When
 315 the maximum degree of G is bounded by some parameter Δ , it turns out that any edge message-
 316 passing protocol can be simulated by a node message-passing protocol with roughly the same number
 317 of rounds and only a Δ factor more memory per processor. The idea is for each node to simulate the
 318 computation that would have been performed (in the edge message-passing protocol) on the adjacent
 319 edges. The following proposition formalizes this idea (proof in Appendix A):

320 **Proposition 3.** *Let $T, B \geq 1$. Let $G = (V, E)$ be a graph with maximum degree Δ . Let P be an
 321 edge message-passing protocol on G with T rounds and B bits of memory. Then there is a node
 322 message-passing protocol P' on G that computes P with $T + 1$ rounds and $O(\Delta B)$ bits of memory.*

323 6 Depth separation under memory and symmetry constraints

324 One drawback of the separation in the previous section is that the constructed edge protocol was
 325 highly non-symmetric, whereas in practice GNN protocols are typically architecturally constrained to
 326 respect the symmetries of the underlying graph. In this section we prove that there is a separation
 327 between the memory/round trade-offs for node and edge message-passing protocols even under
 328 additional symmetry constraints.

329 **Theorem 4.** Let $n \in \mathbb{N}$. There is a graph $G = (V, E)$ with $O(n)$ vertices and $O(n)$ edges, and a
 330 function $g : \{0, 1\}^E \rightarrow \{0, 1\}^V$, so that:

- 331 1. Any node message-passing protocol on G with T rounds and B bits of memory that computes g
 332 requires $TB \geq \Omega(\sqrt{n})$.
- 333 2. There is a symmetric edge message-passing protocol on G with $O(1)$ rounds and $O(\log n)$ bits of
 334 memory that computes g .

335 For intuition, we sketch the proof of a relaxed version of the theorem where the input alphabet is $[n]$.
 336 It is conceptually straightforward to adapt the construction to binary alphabet (essentially, by adding
 337 new vertices and using a unary encoding). We defer the full proof to Appendix B.

338 Let $G = (V, E)$ be a star graph with root node 0 and leaves $\{1, \dots, n\}$. We define a function
 339 $g : [n]^E \rightarrow \{0, 1\}^V$ by $g(I)_v = 1$ if and only if there is some edge $e \neq \{0, v\}$ such that $I(e) =$
 340 $I(\{0, v\})$, i.e. the input on edge $\{0, v\}$ equals the input on some other edge. Since g is defined to
 341 be equivariant to relabelling the edges, and all edges are incident to each other, it is straightforward
 342 to see that there is a symmetric one-round edge message-passing protocol that computes g with
 343 $O(\log n)$ memory (in contrast, the edge message-passing protocol constructed in Section 5 was not
 344 symmetric, as it required that the edges incident to the high-degree vertex were labelled by which path
 345 they belonged to). However, there is no low-memory, low-round *node* message-passing algorithm.
 346 Informally, this is because vertex 0 is an information bottleneck, and $\Omega(n)$ bits of information need
 347 to pass through it. Similar to in Section 5, this intuition can be made formal using Lemma 2.

348 7 Symmetry alone provides no separation

349 In the previous sections we saw that examining *memory constraints* yields a separation between
 350 different GNN architectures (whether or not we take symmetry into consideration). In this section,
 351 we consider what happens if we solely consider *symmetry constraints* (that is, constraints imposed by
 352 requiring that the computation in a round of the protocol is invariant to permutations of the order of
 353 the neighbors). This viewpoint was initiated by [Xu et al. \(2018\)](#), who showed that when the initial
 354 node features are uninformative (that is, the same for each node), a standard GNN necessarily outputs
 355 the same answer for two graphs that are 1-Weisfeiler-Lehman equivalent (that is, graphs that cannot
 356 be distinguished by the Weisfeiler-Lehman test, even though they may not be isomorphic).

357 To be precise, we revisit the representational power of symmetric GNN architectures in the setting
 358 where the input features may be distinct and informative. We show that *if we remove the memory*
 359 *constraints* from Section 5, but *impose permutation invariance* for the computation in each round,
 360 any function that is computable by a T -layer edge message-passing protocol can be computed
 361 by a $(T + 1)$ -layer node message-passing protocol. Note that this statement is incomparable to
 362 Proposition 3 because we impose constraints on symmetry, but remove constraints on memory.

363 **Theorem 5** (No separation under symmetry constraints). Let $T \geq 1$. Let P be a symmetric edge
 364 message-passing protocol (Definition 6) on graph $G = (V, E)$ with T rounds. Then there is a
 365 $(T + 1)$ -round symmetric node message-passing protocol (Definition 5) P' on G that computes the
 366 same function as P .

367 **Remark 10.** Theorem 5 and its proof are closely related to the fact that the 1-Weisfeiler-Lehman
 368 test is equivalent to the 2-Weisfeiler-Lehman test, which was reintroduced in the context of higher-
 369 order GNNs ([Huang & Villar, 2021](#)). However, the k -Weisfeiler-Lehman test only characterizes the
 370 representational power of k -GNNs with uninformative input features (i.e. that are identical for all
 371 nodes). Theorem 5 shows that even with arbitrary input features on the edges, the computation of
 372 a GNN with edge embeddings and symmetric updates can be simulated by a GNN with only node
 373 embeddings, without losing symmetry.

374 To prove Theorem 5, note that it suffices to simulate the protocol P for which the update rules
 375 $f^{\text{sym}}, \tilde{f}^{\text{sym}}$ in Definition 6 are identity functions on the appropriate domains. In order to simulate
 376 P , we construct a symmetric node message-passing protocol P' for which the computation at time
 377 $t + 1$ and node v on input I is the multiset of features computed by P at time t at edges adjacent to v :
 378 $Q_t(v; I) := \{P_t(e; I) : e \in M_G(v)\}$. This is possible since the computation of P at time t and edge
 379 $e = (u, v)$ is $P_t(e; I) = (I(e), P_{t-1}(e; I), \{Q_{t-1}(u; I), Q_{t-1}(v; I)\})$. The node message-passing

380 protocol is tracking $Q_{t-1}(\cdot; I)$; moreover, it can recursively compute $P_{t-1}(e; I)$ using the same
 381 formula. See Appendix C for the formal proof.

382 8 Empirical benefits of edge-based architectures

383 In this section we demonstrate that the representational advantages the theory suggests are borne out
 384 by experimental evaluations, both on real-life benchmarks and two natural synthetic tasks we provide.
 385 Note that all the experiments were done on a machine with 8 Nvidia A6000 GPUs.

386 8.1 Performance on common benchmarks

387 First we compare the performance of the most basic GNN architecture (Graph Convolutional Network,
 388 Kipf & Welling (2016)) with node vs. edge embeddings. In the notation of (1)–(2), the AGGREGATE
 389 and COMBINE operations are integrated as a transformation that looks like (3) or (4):²

$$h_v^{(l+1)} = h_v^{(l)} + \sigma(W^{(l)} \text{MEAN}(h_w^{(l)} : w \in N_G(v) \setminus \{v\})) \quad (3)$$

$$h_e^{(l+1)} = h_e^{(l)} + \sigma(W^{(l)} \text{MEAN}(h_f^{(l)} : f \in M_G(e) \setminus \{e\})) \quad (4)$$

390 for trained matrices $W^{(l)}$ and a choice of non-linearity σ . The only difference between these
 391 architectures is that in the latter case, the message passing happens over the *line graph* of the original
 392 graph (i.e. the neighborhood of an edge is given by the other edges that share a vertex with it) —
 393 thus, this can be viewed as an ablation experiment in which the only salient difference is the type of
 394 embeddings being maintained. To also equalize the information in the input embeddings, we only
 395 use the node embeddings in the benchmarks we consider: for the edge-based architecture (2), we
 396 initialize the edge embeddings by the concatenation of the node embeddings of the endpoints.

397 In Table 1, we show that *this single change* (without any other architectural modifications) uniformly
 398 results in the edge-based architecture at least matching the performance of the node-based architecture,
 399 sometimes improving upon it. *Note, the purpose of this table is not to advocate a new GNN*
 400 *architecture*³— but to confirm that the increased representational power of the edge-based architecture
 401 indicated by the theory also translates to improved performance when the model is trained. For each
 402 benchmark, we follow the best performing training configuration as in (Dwivedi et al., 2023).

Model	ZINC	MNIST	CIFAR-10	Peptides-Func	Peptides-Struct
	MAE (↓)	ACCURACY (↑)	ACCURACY (↑)	AP (↑)	MAE (↓)
GCN	0.3430 ± 0.034	95.29 ± 0.163	55.71 ± 0.381	0.6816 ± 0.007	0.2453 ± 0.0001
Edge-GCN (Ours)	0.3297 ± 0.011	94.37 ± 0.065	57.44 ± 0.387	0.6867 ± 0.004	0.2437 ± 0.0005

Table 1: Comparison of node-based (3) and edge-based (4) GCN architectures across various graph benchmarks. The performance of the edge-based architecture robustly matches or improves the node-based architecture.

403 8.2 A synthetic task for topologies with node bottlenecks

404 The topologies of the graphs in Theorem 1 and Theorem 4 both involve a “hub” node, which is
 405 connected to all other nodes in the graph. Intuitively, in node-embedding architectures, such nodes
 406 have to mediate messages between many pairs of other nodes, which is difficult when the node is
 407 constrained by memory. To empirically stress test this intuition, we produce a synthetic dataset
 408 and train a GNN to solve a regression task on a graph with a fixed *star-graph* topology—a simpler
 409 topology than the constructions in Theorem 1 and Theorem 4—but capturing the core aspect of
 410 both. A star graph is a graph with a center node v_0 , a set of n leaf nodes $\{v_i\}_{i \in [n]}$, and edge set
 411 $\{\{v_0, v_i\}_{i \in [n]}\}$. A training point in the dataset is a list $(x_i, y_i)_{i=1}^n$ where x_i is the *input feature* and
 412 y_i is the *label* for leaf node v_i .

413 The input features are in \mathbb{R}^{10} , and sampled from a standard Gaussian. The labels y_i are produced
 414 as outputs of a *planted* edge-based architecture. Namely, for a standard edge-based GCN as in (4),
 415 we randomly choose values for the matrices $\{W_i\}_{i \in [k]}$ for some number of layers k , and set the
 416 labels to be the output of this edge-based GCN, when the initial edge features to the GCN are set
 417 as $h_{\{v_0, v_i\}}^{(0)} := x_i$, i.e. the input feature x_i at the corresponding leaf i . In Table 2, we show the

²This is the “residual” parametrization, which we use in experiments unless otherwise stated.

³In particular, the edge-based architecture is often much more computationally costly to evaluate.

418 performance of edge-based and node-based architectures on this dataset, varying the number of leaves
 419 n in the star graph and the depth k of the planted edge-based model. In each case, the numbers
 420 indicate RMSE of the best-performing edge-based and node-based architecture, sweeping over depths
 421 up to 10 ($2\times$ the planted model), widths $\in \{16, 32, 64\}$, and a range of learning rates.

422 Since the planted edge-based model satisfies both *invariance* constraints (by design of the GCN
 423 architecture) and *memory* constraints (since the planted model maintains 10-dimensional embeddings),
 424 we view these results as empirical corroboration of Theorem 4—and even for simpler topologies than
 425 the proof construction.

Number of Leaves	Depth of Planted Model (RMSE)					
	5		3		1	
	Edge	Node	Edge	Node	Edge	Node
64	0.004	0.3790	0.011	0.3596	0.008	0.3752
32	0.003	0.3664	0.005	0.3626	0.003	0.3614
16	0.007	0.3336	0.002	0.2100	0.002	0.2847

Table 2: Performance (in RMSE \downarrow) of edge-based and node-based architectures on a star-graph topology. The first number is the performance of the best edge-based model, and the second is the best node-based model, across a range of depths up to 10 ($2\times$ the planted model), widths $\in \{16, 32, 64\}$, and a range of learning rates.

426 8.3 A synthetic task for inference in Ising models

427 Finally, motivated by the probabilistic inference setting in Theorem 1, we consider a synthetic
 428 sandbox of using GNNs to predict the values of marginals in an Ising model (Ising, 1924; Onsager,
 429 1944) – a natural type of pairwise graphical model where each node takes a value in $\{\pm 1\}$, and
 430 each edge potential is a weighted product of the edge endpoint values. Concretely, the probability
 431 distribution of an Ising model over graph $G = (V, E)$ has the form: $\forall x \in \{\pm 1\}^n : p_{J,h}(x) \propto$
 432 $\exp\left(\sum_{\{i,j\} \in E} J_{\{i,j\}} x_i x_j + \sum_{i \in V} h_i x_i\right)$.

433 Similar to in Section 8.2, we construct a training set where the graph G and edge potentials
 434 stay fixed (precisely, $J_{i,j} = 1$ for all $\{i, j\} \in E$). A training data-point consists of a vector of node
 435 potentials $\{h_i\}_{i \in [n]}$, and labels $\{\mathbb{E}[x_i]\}_{i \in [n]}$ consisting of the marginals from the resulting Ising
 436 model $p_{J,h}$. The node potentials are sampled from a standard Gaussian distribution.

437 There is a natural connection between GNNs and calculating marginals: a classical way to calculate
 438 $\{\mathbb{E}[x_i]\}$ when G is a *tree* is to iterate a message passing algorithm called *belief propagation* (5),
 439 in which for each edge $\{i, j\}$ and direction $i \rightarrow j$, a message $\nu_{i \rightarrow j}^{(t+1)}$ is calculated that depends on
 440 messages $\{\nu_{k \rightarrow i}^{(t)}\}_{\{k,i\} \in E}$. The belief-propagation updates (5) naturally fit the general edge-message
 441 passing paradigm from (2). In fact, they fit even more closely a “directed” version of the paradigm, in
 442 which each edge $\{i, j\}$ maintains two embeddings $h_{i \rightarrow j}, h_{j \rightarrow i}$, such that the embedding for direction
 443 $h_{i \rightarrow j}$ depends on the embeddings $\{h_{k \rightarrow i}\}_{\{k,i\} \in E}$ — and it is possible to derive a similar “directed”
 444 node-based architecture (See Appendix E.2). For both the undirected and directed version of the
 445 architecture, we see that maintaining edge embeddings gives robust benefits over maintaining node
 446 embeddings—for a variety of tree topologies including complete binary trees, path graphs, and
 447 uniformly randomly sampled trees of a fixed size. More details are included in Appendix E.

448 9 Conclusions and future work

449 Graph neural networks are the best-performing machine learning method for many tasks over graphs.
 450 There is a wide variety of GNN architectures, which frequently make opaque design choices and
 451 whose causal influence on the final performance is difficult to understand and estimate. In this paper,
 452 we focused on understanding the impact of maintaining edge embeddings on the representational
 453 power, as well as the subtleties of considering constraints like memory and invariance. One significant
 454 downside of maintaining edge embeddings is the *computational* overhead on dense graphs. Hence, a
 455 fruitful direction for future research would be to explore more computationally efficient variants of
 456 edge-based architectures that preserve their representational power and performance.

457 **References**

- 458 Karl Abrahamson. Time-space tradeoffs for algebraic problems on general sequential machines.
459 *Journal of Computer and System Sciences*, 43(2):269–289, 1991.
- 460 Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks
461 for learning about objects, relations and physics. *Advances in neural information processing*
462 *systems*, 29, 2016.
- 463 Mitchell Black, Zhengchao Wan, Amir Nayyeri, and Yusu Wang. Understanding oversquashing in
464 gnns through the lens of effective resistance. In *International Conference on Machine Learning*,
465 pp. 2528–2547. PMLR, 2023.
- 466 Lei Cai, Jundong Li, Jie Wang, and Shuiwang Ji. Line graph neural networks for link prediction.
467 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5103–5113, 2021.
- 468 Zhengdao Chen, Xiang Li, and Joan Bruna. Supervised community detection with line graph neural
469 networks. *arXiv preprint arXiv:1705.08415*, 2017.
- 470 Kamal Choudhary and Brian DeCost. Atomistic line graph neural network for improved materials
471 property predictions. *npj Computational Materials*, 7(1):185, 2021.
- 472 Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and
473 Xavier Bresson. Benchmarking graph neural networks. *Journal of Machine Learning Research*, 24
474 (43):1–48, 2023.
- 475 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
476 message passing for quantum chemistry. In *International conference on machine learning*, pp.
477 1263–1272. PMLR, 2017.
- 478 Dmitrii Yur’evich Grigor’ev. Application of separability and independence notions for proving lower
479 bounds of circuit complexity. *Zapiski Nauchnykh Seminarov POMI*, 60:38–48, 1976.
- 480 Johan Håstad and Avi Wigderson. The randomized communication complexity of set disjointness.
481 *Theory of Computing*, 3(1):211–219, 2007.
- 482 Ningyuan Teresa Huang and Soledad Villar. A short tutorial on the weisfeiler-lehman test and its
483 variants. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal*
484 *Processing (ICASSP)*, pp. 8533–8537. IEEE, 2021.
- 485 Ernst Ising. *Beitrag zur theorie des ferro-und paramagnetismus*. PhD thesis, Grefe & Tiedemann
486 Hamburg, Germany, 1924.
- 487 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
488 *arXiv preprint arXiv:1609.02907*, 2016.
- 489 AA Leman and Boris Weisfeiler. A reduction of a graph to a canonical form and an algebra arising
490 during this reduction. *Nauchno-Tekhnicheskaya Informatsiya*, 2(9):12–16, 1968.
- 491 Jinbi Liang and Cunlai Pu. Line graph neural networks for link weight prediction. *arXiv preprint*
492 *arXiv:2309.15728*, 2023.
- 493 Nathan Linial. Locality in distributed graph algorithms. *SIAM Journal on computing*, 21(1):193–201,
494 1992.
- 495 Andreas Loukas. What graph neural networks cannot learn: depth vs width. *arXiv preprint*
496 *arXiv:1907.03199*, 2019.
- 497 Andreas Loukas. How hard is to distinguish graphs with graph neural networks? *Advances in neural*
498 *information processing systems*, 33:3465–3476, 2020.
- 499 Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph
500 networks. *Advances in neural information processing systems*, 32, 2019.
- 501 Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University
502 Press, 2009.

- 503 Lars Onsager. Crystal statistics. i. a two-dimensional model with an order-disorder transition. *Physical*
504 *Review*, 65(3-4):117, 1944.
- 505 Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node
506 classification. *arXiv preprint arXiv:1905.10947*, 2019.
- 507 David Peleg. *Distributed computing: a locality-sensitive approach*. SIAM, 2000.
- 508 Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Transformers, parallel computation, and logarithmic
509 depth. *arXiv preprint arXiv:2402.09268*, 2024a.
- 510 Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations of
511 transformers. *Advances in Neural Information Processing Systems*, 36, 2024b.
- 512 John E Savage. *Models of computation*, volume 136. Addison-Wesley Reading, 1998.
- 513 Matus Telgarsky. Benefits of depth in neural networks. In *Conference on learning theory*, pp.
514 1517–1539. PMLR, 2016.
- 515 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
516 networks? *arXiv preprint arXiv:1810.00826*, 2018.
- 517 Tuo Xu and Lei Zou. Rethinking and extending the probabilistic inference capacity of gnns. In *The*
518 *Twelfth International Conference on Learning Representations*, 2023.
- 519 Lingxiao Zhao, Wei Jin, Leman Akoglu, and Neil Shah. From stars to subgraphs: Uplifting any gnn
520 with local structure awareness. *arXiv preprint arXiv:2110.03753*, 2021.

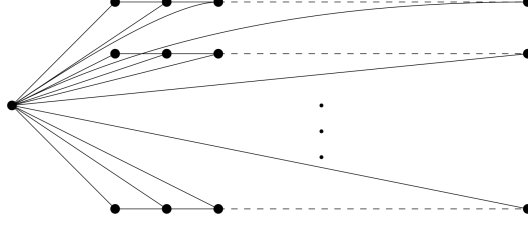


Figure 1: The graph G for which Theorem 1 exhibits a separation between edge message-passing and node message-passing. The graph consists of \sqrt{n} paths of length \sqrt{n} , as well as a single “hub vertex” connected to all other vertices.

521 Appendix

522 A Omitted Proofs from Section 5

523 In this section we give the omitted proofs from Section 5. In particular we give the formal proof of
 524 Theorem 1, which states that there is a depth separation between edge message-passing protocols and
 525 node message-passing protocols for a natural MAP estimation problem on the underlying graph G .
 526 We also remark that a quantitatively stronger (and in fact tight) separation is possible if one considers
 527 general tasks rather than MAP estimation tasks – see Appendix D.

528 *Proof of Lemma 2.* First, we argue by induction that for each $t \in [T]$ and $v \in V \setminus K$, $P_t(v; I)$
 529 is determined by $I_{M_G(N_H^{t-1}(v))}$ and $(P_\ell(k; I))_{\ell \in [t], k \in K}$. Indeed, by definition, $P_1(v; I)$ is deter-
 530 mined by $I_{M_G^1(v)}$ for any $v \in V \setminus K$. For any $t > 1$ and $v \in V \setminus K$, $P_t(v; I)$ is determined
 531 by $(P_{t-1}(v'; I))_{v' \in N_G(v)}$ and $(I(e))_{e \in M_G(v)}$. Note that $N_G(v) \subseteq N_H(v) \cup K$. Thus, using the
 532 induction hypothesis for each $v' \in N_H(v)$, we get that $(P_{t-1}(v'; I))_{v' \in N_G(v)}$ is determined by
 533 $\bigcup_{v' \in N_H(v)} I_{M_G(N_H^{t-2}(v'))}$ and $(P_\ell(k; I))_{\ell \in [t], k \in K}$. So $P_t(v; I)$ is determined by $I_{M_G(N_H^{t-1}(v))}$ and
 534 $(P_\ell(k; I))_{\ell \in [t], k \in K}$, completing the induction.

535 Since P computes g and $S \subseteq V \setminus K$, we get that $g_S(I)$ is determined by $I_{M_G(N_H^{T-1}(S))} = I_F$ and
 536 $(P_\ell(k; I))_{\ell \in [T], k \in K}$. Thus, for any fixed $I_F \in \Phi_F$, we have

$$\left| \left\{ g_S(I_F, I_{\bar{F}}) : I_{\bar{F}} \in \Phi^{\bar{F}} \right\} \right| \leq \left| \left\{ (P_\ell(k; (I_F, I_{\bar{F}})))_{\ell \in [T], k \in K} : I_{\bar{F}} \in \Phi^{\bar{F}} \right\} \right| \leq |\mathcal{X}_B|^{T|K|} = 2^{TB|K|}.$$

537 The lemma follows. \square

538 *Proof of Theorem 1.* Let G be the graph on vertex set $V := \{0\} \cup [\sqrt{n}] \times [\sqrt{n}]$ with edge set defined
 539 below (see also Fig. 1):

$$E := \{ \{0, (i, j)\} : i, j \in [\sqrt{n}] \} \cup \{ \{(i, j), (i+1, j)\} : 2 \leq i \leq \sqrt{n}, 1 \leq j \leq \sqrt{n} \}.$$

540 Define

$$\Phi := \{ (x_a, x_b) \mapsto \mathbb{1}[x_a \neq x_b], (x_a, x_b) \mapsto \mathbb{1}[x_a \neq 1 \vee x_b \neq 1], (x_a, x_b) \mapsto \mathbb{1}[x_a \neq 0 \vee x_b \neq 0], (x_a, x_b) \mapsto 0 \}.$$

541 First, let $g : \Phi^E \rightarrow \{0, 1\}^V$ be any MAP evaluator for G with potential function class Φ , and consider
 542 any node message-passing protocol on G with T rounds and B bits of memory that computes g . Let
 543 $K = \{0\}$ and $S = \{(1, j) : j \in [\sqrt{n}]\}$. Suppose that $T \leq \sqrt{n} - 2$. Let $F := M_G(N_H^{T-1}(S))$ and
 544 note that $\{(\sqrt{n}-1, j), (\sqrt{n}, j)\} \notin F$ for all $j \in [\sqrt{n}]$. Let $I_F : F \rightarrow \Phi$ be the mapping that assigns
 545 the function $(x_a, x_b) \mapsto 0$ to each edge $\{0, (i, j)\} \in F$ and $(x_a, x_b) \mapsto \mathbb{1}[x_a \neq x_b]$ to each edge
 546 $\{(i, j), (i+1, j)\} \in F$. We claim that

$$\left| \left\{ g_S(I_F, I_{\bar{F}}) : I_{\bar{F}} \in \Phi^{\bar{F}} \right\} \right| \geq 2^{\sqrt{n}}.$$

547 Indeed, for any string $y \in \{0, 1\}^{\sqrt{n}}$, consider the mapping $I_{\bar{F}} : \bar{F} \rightarrow \Phi$ that assigns the function
 548 $(x_a, x_b) \mapsto \mathbb{1}[x_a \neq y_j \vee x_b \neq y_j]$ to each edge $\{(\sqrt{n}-1, j), (\sqrt{n}, j)\} \in F$, assigns $(x_a, x_b) \mapsto 0$

549 to each edge $\{0, (i, j)\} \in E \setminus F$, and assigns $(x_a, x_b) \mapsto \mathbb{1}[x_a \neq x_b]$ to all remaining edges in $E \setminus F$.
 550 Then every minimizer of

$$\min_{x \in \{0,1\}^V} \sum_{\{a,b\} \in E} I_{\{a,b\}}(x_a, x_b)$$

551 satisfies $x_{(1,j)} = \dots = x_{(\sqrt{n},j)} = y_j$ for all $j \in [\sqrt{n}]$. Hence, $g_S(I_F, I_{\bar{F}}) = y$. Since y was chosen
 552 arbitrarily, this proves the claim. But now Lemma 2 implies that $TB \geq \sqrt{n}$.

553 We now construct an edge message-passing protocol P on G with $T = 3$ and $B = 4$. We (arbitrarily)
 554 identify Φ with $\{0, 1\}^2$. For all $i, j \in \sqrt{n}$, define

$$\begin{aligned} f_{1, \{(i,j), (i+1,j)\}}(x, y) &:= y && \text{if } i < \sqrt{n} \\ f_{2, \{0, (i,j)\}}(x, y) &:= (x_{\{(i,j), (i+1,j)\}}, x_{\{0, (i,j)\}}) && \text{if } i < \sqrt{n} \\ f_{3, \{0, (i,j)\}}(x, y) &:= (g_0(J(x)), g_{(i,j)}(J(x))) \end{aligned}$$

555 where the second line is well-defined since edge $\{0, (i, j)\}$ is adjacent to both itself and edge
 556 $\{(i, j), (i+1, j)\}$; and in the third line the function is computing g_0 and $g_{(i,j)}$ on the input $J(x) \in \Phi^E$
 557 defined as

$$J(x)_e := \begin{cases} (x_{\{0, (k, \ell)\}})_{1:2} & \text{if } e = \{(k, \ell), (k+1, \ell)\} \\ (x_{\{0, (k, \ell)\}})_{3:4} & \text{if } e = \{0, (k, \ell)\} \end{cases},$$

558 where we use the notation $v_{a:b}$ for a vector v and indices $a, b \in \mathbb{N}$ to denote $(v_a, v_{a+1}, \dots, v_b)$. Note
 559 that $J(x)$ is a well-defined function of x for every edge $\{0, (i, j)\}$, because $\{0, (i, j)\} \sim \{0, (k, \ell)\}$
 560 for all $i, j, k, \ell \in [n]$. Finally, define all other functions $f_{t,e}$ to compute the all-zero function, and
 561 define

$$\tilde{f}_v(x) := \begin{cases} (x_{\{0, (1,1)\}})_{1:2} & \text{if } v = 0 \\ (x_{\{0, v\}})_{3:4} & \text{otherwise} \end{cases}.$$

562 This function is well-defined since $v = 0$ is adjacent to edge $\{0, (1, 1)\}$ and any vertex $v \in V \setminus \{0\}$
 563 is adjacent to edge $\{0, v\}$.

564 Fix any $I \in \Phi^E$. From the definition, it's clear that $P_2(\{0, (i, j)\}; I) = (I_{\{(i,j), (i+1,j)\}}, I_{\{0, (i,j)\}})$
 565 for all I and $(i, j) \in [\sqrt{n} - 1] \times [\sqrt{n}]$. Hence $J((P_2(e'; I))_{e' \in M_G(e)})_e = I$ for all edges e of the
 566 form $(0, \{(i, j)\})$, and so $P_3(\{0, (i, j)\}; I) = (g_0(I), g_{(i,j)}(I))$ for all $(i, j) \in [\sqrt{n}] \times [\sqrt{n}]$. This
 567 means that $\tilde{f}_v((P_3(e; I))_{e \in M_G(v)}) = g(I)_v$ for all $v \in V$, so the protocol indeed computes g .

568 It remains to argue about the computational complexity of the updates $f_{t,e}$. It's clear that for all
 569 $e \in E$ and $t \in \{1, 2\}$, the function $f_{t,e}$ can be evaluated in input-linear time. The only case that
 570 requires proof is when $t = 3$ and $e = \{0, (i, j)\}$ for some $i, j \in \sqrt{n}$. In this case $|M_G(e)| = \Theta(n)$,
 571 so it suffices to give an algorithm for evaluating the function $g : \Phi^E \rightarrow \{0, 1\}^V$ on an explicit input
 572 J in $O(n)$ time. This can be accomplished via dynamic programming (Lemma 6). \square

573 **Lemma 6.** Fix $n \in \mathbb{N}$. Let G, Φ be as defined in Theorem 1. Then there is an $O(n)$ -time algorithm
 574 that computes a MAP evaluator for G with potential function class Φ .

575 *Proof.* Fix any $J \in \Phi^E$. As preliminary notation, for each $c, c_0 \in \{0, 1\}$ and $i, j \in \sqrt{n}$, let
 576 $V(i, j) := \{0\} \cup \{(k, j) : 1 \leq k \leq i\}$, and let $E(i, j)$ be the edge set of the induced subgraph
 577 $G[V(i, j)]$. Let

$$\begin{aligned} \hat{x}_{i,j}(c, c_0; J) &:= \arg \min_{\substack{x \in \{0,1\}^{V(i,j)} \\ x_0 = c_0 \wedge x_{(i,j)} = c}} \sum_{(a,b) \in E(i,j)} J_{\{a,b\}}(x_a, x_b), \\ \hat{C}_{i,j}(c, c_0; J) &:= \min_{\substack{x \in \{0,1\}^{V(i,j)} \\ x_0 = c_0 \wedge x_{(i,j)} = c}} \sum_{(a,b) \in E(i,j)} J_{\{a,b\}}(x_a, x_b). \end{aligned}$$

578 For each $j \in [\sqrt{n}]$, let

$$\hat{x}_j(c_0; J) := \hat{x}_{\sqrt{n}, j} \left(\left(\arg \min_{c \in \{0,1\}} \hat{C}_{\sqrt{n}, j}(c, c_0; J) \right), c_0; J \right).$$

579 Finally, let $\hat{x}(c_0; J) \in \{0, 1\}^V$ be the vector which takes value c_0 on vertex 0, and value $\hat{x}_j(c_0; J)_i$
580 on vertex (i, j) for all $i, j \in \sqrt{n}$. Let

$$\hat{x}(J) := \arg \max_{c_0 \in \{0, 1\}} p_J(\hat{x}(c_0; J)).$$

581 We claim that $\hat{x}(J)$ is a maximizer of $p_J(x)$. Indeed, for any fixed $c_0 \in \{0, 1\}$, $\hat{x}(c_0; J)$ is a
582 maximizer of $p_J(x)$ subject to $x_0 = c_0$, because under this constraint the maximization problem
583 decomposes into \sqrt{n} independent maximization problems, one for each path in G , which by definition
584 are solved by $\hat{x}_1(c_0; J), \dots, \hat{x}_{\sqrt{n}}(c_0; J)$.

585 Moreover, it's straightforward to see that for any fixed j , $\hat{C}_j(c_0; J)$ can be computed in $O(\sqrt{n})$
586 time by dynamic programming. Indeed for any i, j , $\hat{C}_{i,j}(c, c_0; J)$ can be computed in $O(1)$ time
587 from $\hat{C}_{i-1,j}(0, c_0; J)$ and $\hat{C}_{i-1,j}(1, c_0; J)$ as well as $J_{\{0, (i,j)\}}$ and $J_{\{(i-1,j), (i,j)\}}$. Once the values
588 $\hat{C}_{i,j}(c, c_0; J)$ have been computed for all $i \in [\sqrt{n}]$ and $c \in \{0, 1\}$, the vector $\hat{x}_j(c_0; J)$ can be
589 computed in $O(\sqrt{n})$ time via a reverse scan over $i = \sqrt{n}, \dots, 1$. It follows that $\hat{x}(J)$ can be
590 computed in $O(n)$ time. \square

591 *Proof of Proposition 3.* We claim that there is a node message-passing protocol P' on G with $T + 1$
592 rounds that at each time $t \in [T + 1]$ has computed

$$P'_t(v; I) = (P'_{t-1}(e; I))_{e \in M_G(v)}.$$

593 We argue inductively. Since $P_0 \equiv 0$, it's clear that this can be achieved for $t = 1$. Fix any $t > 1$
594 and suppose that $P'_{t-1}(u; I) = (P'_{t-2}(e; I))_{e \in M_G(u)}$ for all $u \in V$ and inputs I . For each $v \in V$, we
595 define a function $f'_{t,v}$ by

$$f'_{t,v}((c(v'))_{v' \in N_G(v)}, (I(e))_{e \in M_G(v)})_{e^*} := f_{t-1, e^*}((c(v)_e)_{e \in M_G(v)}, (c(v^*)_e)_{e \in M_G(v^*)}, I(e^*))$$

596 for each $e^* = (v, v^*) \in M_G(v)$. Then by definition and the inductive hypothesis, we have

$$\begin{aligned} P'_t(v; I)_{e^*} &= f'_{t,v}((P'_{t-1}(v'; I))_{v' \in N_G(v)}, (I(e))_{e \in M_G(v)})_{e^*} \\ &= f_{t-1, e^*}((P'_{t-1}(v; I)_e)_{e \in M_G(v)}, (P'_{t-1}(v^*; I)_e)_{e \in M_G(v^*)}, I(e^*)) \\ &= f_{t-1, e^*}((P_{t-2}(e; I))_{e \in M_G(v)}, (P_{t-2}(e; I)_e)_{e \in M_G(v^*)}, I(e^*)) \\ &= P_{t-1}(e^*; I) \end{aligned}$$

597 for any edge $e^* = (v, v^*) \in E$, since $M_G(e) = M_G(v) \cup M_G(v^*)$. This completes the induction
598 and shows that $P'_{T+1}(v; I) = (P_T(e; I))_{e \in M_G(v)}$ for all v, I . Replacing $f'_{T+1, v}$ by $\tilde{f}_{T, v} \circ f'_{T+1, v}$
599 completes the proof. \square

600 B Omitted Proofs from Section 6

601 In this section we provide a formal proof of Theorem 4. For notational convenience, define $m = \lfloor \sqrt{n} \rfloor$.
602 We define a graph $G = (V, E)$ that is a perfect n -ary tree of depth two. Formally, the graph G has
603 vertex set $V = \{0\} \cup [m] \cup ([m] \times [m])$. Vertex 0 is adjacent to each $i \in [m]$, and each $i \in [m]$ is
604 additionally adjacent to (i, j) for all $j \in [m]$. We define a function $g : \{0, 1\}^E \rightarrow \{0, 1\}^V$ as follows.
605 On input $I \in \{0, 1\}^E$, for each edge $e \in E$, define the *input summation* at e to be

$$C(I)_e := \sum_{e' \in M_G(e)} I(e').$$

606 Intuitively, one may think of $C(I)_e$ as simulating the input on e in the “large alphabet” construction
607 described in Section 6. Next, define

$$\begin{aligned} g(I)_{(u,j)} &:= 0. \\ g(I)_u &:= \mathbb{1}[\#\{e \in M_G(\{0, u\}) : C(I)_e = C(I)_{\{0, u\}}\} > m + 1]. \\ g(I)_0 &:= \mathbb{1}[\exists u \in [m] : g(I)_u = 1]. \end{aligned}$$

608 In words, $g(I)_u$ is the indicator for the event that, among the $2m + 1$ edges adjacent to $\{0, u\}$
609 (which include $\{0, u\}$ itself), more than $m + 1$ edges have the same input summation as $\{0, u\}$.

610 At a high level, this definition of g was designed to satisfy three criteria. First, $g(I)_u$ depends on
611 the input values on other branches of the tree: in particular, if $I_{\{0,v\}} = 0$ for all $v \in [n]$, then
612 $C(I)_e = C(I)_{\{0,u\}}$ for all edges e in the subtree of u , so $g(I)_u$ exactly measures the event that there
613 is *at least one* edge e outside the subtree of u for which $C(I)_e = C(I)_{\{0,u\}}$. Second, there is no
614 concise “summary” of I such that $g(I)_u$ can be determined from this summary in conjunction with
615 the inputs on the subtree of u . Third, $g(I)$ is equivariant to re-labelings of the tree.

616 The first two criteria, together with the fact that the root vertex 0 is an “information bottleneck” for
617 G , can be used to show that any node message-passing algorithm that computes g on G requires
618 either large memory or many rounds. The third criterion enables construction of a symmetric edge
619 message-passing protocol for g . The arguments are formalized in the claims below.

620 **Claim 7.** *For graph G and function g as defined above, any node message-passing protocol on G
621 that computes g with T rounds and B bits of memory requires $TB \geq \Omega(m)$.*

622 *Proof.* Consider any input $I \in \{0, 1\}^E$ with $I(\{0, u\}) = 0$ for all $u \in [m]$. Then for any $u, j \in [m]$,
623 we have

$$C(I)_{\{u,(u,j)\}} = C(I)_{\{0,u\}} = \sum_{i=1}^m I(\{u, (u, i)\}).$$

624 Thus $g(I)_u = 1$ if and only if there exists some $v \in [m] \setminus \{u\}$ with $C(I)_{\{0,u\}} = C(I)_{\{0,v\}}$, or
625 equivalently $\sum_{i=1}^m I(\{u, (u, i)\}) = \sum_{i=1}^m I(\{v, (v, i)\})$.

626 Fix T, B and suppose that P is a node message-passing protocol on G that computes g with T rounds
627 and B bits of memory. Define sets of vertices $K := \{0\}$ and $S := \{1, \dots, m/2\}$. Let $H := G[\overline{K}]$
628 and $F := M_G(N_H^{T-1}(S))$. Then for any T , we have that

$$F = \{\{0, u\} : 1 \leq u \leq m/2\} \cup \{\{u, (u, j)\} : 1 \leq u \leq m/2, 1 \leq j \leq m\}.$$

629 Define a vector $I_F \in \Phi^F$ by

$$\begin{aligned} I_{\{0,u\}} &= 0 \text{ for } 1 \leq u \leq m/2 \\ I_{\{u,(u,j)\}} &= \mathbb{1}[j \leq u] \text{ for } 1 \leq u \leq m/2, 1 \leq j \leq m. \end{aligned}$$

630 Now fix any $x \in \{0, 1\}^S$. We claim that there is some $I_{\overline{F}} \in \Phi^{\overline{F}}$ such that $g_S(I_F, I_{\overline{F}}) = x$. Indeed,
631 let us define $I_{\overline{F}}$ by:

$$\begin{aligned} I_{\{0,v\}} &= 0 \text{ for } m/2 < v \leq m \\ I_{\{v,(v,j)\}} &= x_{v-m/2} \mathbb{1}[j \leq v - m/2] \text{ for } m/2 < v \leq m, 1 \leq j \leq m. \end{aligned}$$

632 Then $C(I)_{\{0,u\}} = u$ for all $1 \leq u \leq m/2$, and $C(I)_{\{0,v\}} = (v-m/2)x_{v-m/2}$ for all $m/2 < v \leq m$.
633 It follows that for any $1 \leq u \leq n/2$, $x_u = 1$ if and only if there exists some $v \in [m] \setminus u$ with
634 $C(I)_{\{0,u\}} = C(I)_{\{0,v\}}$, and hence $x_u = g(I)_u$. We conclude that

$$\left| \left\{ g_S(I_F, I_{\overline{F}}) : I_{\overline{F}} \in \Phi^{\overline{F}} \right\} \right| \geq 2^{m/2}.$$

635 Applying Lemma 2 we conclude that $TB \geq \Omega(m)$ as claimed. \square

636 **Claim 8.** *For graph G and function g as defined above, there is a symmetric edge message-passing
637 protocol on G that computes g with $O(1)$ rounds and $O(\log m)$ bits of memory.*

638 *Proof.* In the first round, each edge processor reads its input value. In the second round, each edge
639 processor sums the values computed by all neighboring edges (including itself). In the third round,
640 each edge processor computes the indicator for the event that strictly more than $m + 1$ neighboring
641 edges (including itself) have the same value as itself. In the final aggregation round, the output of a
642 vertex is the indicator for the event that any neighbor has value 1.

643 By construction, the value computed by any edge e after the second round is exactly $C(I)_e$. Thus,
644 after the third round, the value computed by any edge $\{0, u\}$ is exactly $g(I)_u$. Moreover, the value
645 computed by any edge $\{u, (u, j)\}$ is 0 after the third round, since such edges only have $m + 1$
646 neighbors. It follows by construction of the final aggregation step that the protocol computes g . \square

647 *Proof of Theorem 4.* Immediate from Claims 7 and 8. \square

648 **C Omitted Proofs from Section 7**

649 *Proof of Theorem 5.* Without loss of generality, we may assume that the functions $(f_t^{\text{sym}})_{t \in [T]}$ and
 650 \tilde{f}^{sym} are all the identity function (on the appropriate domains). The reason is that any symmetric
 651 edge message-passing protocol \tilde{P} on T rounds may be simulated by running P and then applying a
 652 universal function (depending only on \tilde{P}) to each node's output value – see Lemma 9.

653 We argue by induction that for each $t \in [T]$, there is a $(t+1)$ -round symmetric node message-passing
 654 protocol that, on any input I , computes the function $Q_t(u; I) := \{\{P_t(e; I) : e \in M_G(u)\}\}$ for every
 655 node $u \in V$. Consider $t = 1$. For any $e = (u, v) \in E$, we have by symmetry and the initial
 656 assumption that

$$P_1(e; I) = (I(e), 0, \{\{0 : v' \in N_G(u)\}, \{0 : u' \in N_G(v)\}\}).$$

657 We define a two-round node message-passing protocol on G where the first update at node u computes

$$P'_1(u; I) = \{\{I(\{u, v\}) : v \in N_G(u)\}\}$$

658 and the second update at node u computes

$$\begin{aligned} (P'_1(u; I), \{\{(P'_1(v; I), I(\{u, v\})) : v \in N_G(u)\}\}) &\mapsto \{\{(I(\{u, v\}), 0, |N_G(u)|, |P'_1(v; I)|) : v \in N_G(u)\}\} \\ &\mapsto \{\{(I(\{u, v\}), 0, \{|N_G(u)|, |P'_1(v; I)|\}) : v \in N_G(u)\}\} \\ &= \{\{P_1(\{u, v\}; I) : v \in N_G(u)\}\} =: P'_2(u; I) \end{aligned}$$

659 since $|P'_1(v; I)| = |N_G(v)|$. By construction, this protocol is symmetric, which proves the induction
 660 for step $t = 1$.

661 Now pick any $t > 1$. For any $e = \{u, v\} \in E$, we have

$$P_t(e; I) = (I(e), P_{t-1}(e; I), \{\{Q_{t-1}(u; I), Q_{t-1}(v; I)\}\})$$

662 By the induction hypothesis, there is a t -round symmetric node message-passing protocol P' that, at
 663 node v on input I , computes

$$P'_t(v; I) = \{\{P_{t-1}(\{v, v'\}; I) : v' \in N_G(v)\}\} = Q_{t-1}(v; I).$$

664 Note that since $P_{t-1}(e; I)$ is an element of the tuple $P_t(e; I)$, for each $1 \leq s \leq t-1$ there is a fixed
 665 function γ_s such that $\gamma_s(Q_{t-1}(v; I)) = Q_s(v; I)$ for all v, I . Using this fact, we extend P' to $t+1$
 666 rounds, defining the update at round $t+1$ and node u as follows:

$$\begin{aligned} (P'_t(u; I), \{\{(P'_t(v; I), I(\{u, v\})) : v \in N_G(u)\}\}) & \\ = (Q_{t-1}(u; I), \{\{(Q_{t-1}(v; I), I(\{u, v\})) : v \in N_G(u)\}\}) & \\ \mapsto (Q_{1:t-1}(u; I), \{\{(Q_{1:t-1}(v; I), I(\{u, v\})) : v \in N_G(u)\}\}) & \\ \mapsto (Q_{1:t-1}(u; I), \{\{(Q_{1:t-1}(v; I), I(\{u, v\})) : v \in N_G(u)\}\}) & \\ = \{\{(I(\{u, v\}), \{Q_{1:t-1}(u; I), Q_{1:t-1}(v; I)\}) : v \in N_G(u)\}\} & \\ \mapsto \{\{(I(\{u, v\}), P_{t-1}(\{u, v\}; I), \{Q_{t-1}(u; I), Q_{t-1}(v; I)\}) : v \in N_G(u)\}\} =: P'_{t+1}(u; I) & \end{aligned}$$

667 where $Q_{1:t-1}(u; I)$ refers to the tuple $(Q_1(u; I), \dots, Q_{t-1}(u; I))$. The first map is well-defined due
 668 to the existence of the functions $\gamma_1, \dots, \gamma_{t-1}$, and the final map is well-defined because the definition
 669 of $P_{t-1}(\{u, v\}; I)$ can be iteratively unpacked, and it is ultimately a function of

$$(I(\{u, v\}), \{\{Q_{1:t-1}(u; I), Q_{1:t-1}(v; I)\}\}).$$

670 This shows that P' computes $Q_t(v; I)$ at node u on input I . By construction, P' is symmetric. This
 671 completes the induction. Since $Q_T(u; I)$ is precisely the output of P at node u on input I (after
 672 the node aggregation step), this shows that P can be simulated by a $(T+1)$ -round symmetric node
 673 message-passing protocol on G . \square

674 **Lemma 9.** Let $T \geq 1$, and let $P = ((f_{t,e})_{t \in [T], e \in E}, (\tilde{f}_v)_{v \in V})$ be a symmetric edge message-passing
 675 protocol on $G = (V, E)$ with T rounds. Consider the T -round edge message-passing protocol
 676 $P^\circ = ((f_{t,e}^\circ)_{t \in [T], e \in E}, (\tilde{f}_v^\circ)_{v \in V})$ where for all t, e ,

$$f_{t,e}^\circ((c(e'))_{e' \in M_G(e)}, I(e)) := (I(e), c(e), \{\{c(\{u, v'\}) : v' \in N_G(u)\}\}, \{\{c(\{u', v\}) : u' \in N_G(v)\}\}),$$

677 and for every $v \in V$,

$$\tilde{f}_v^\circ((c(e))_{e \in M_G(v)}) := \{\{c(e) : e \in M_G(v)\}\}.$$

678 Then there is a function h such that $\tilde{f}_v((P_T(e; I))_{e \in M_G(v)}) = h(\tilde{f}_v^\circ((P_T^\circ(e; I))_{e \in M_G(v)}))$ for all
 679 v, I .

680 *Proof.* We prove by induction that for each $t \in \{0, \dots, T\}$ there is a function h_t such that $P_t(e; I) =$
681 $h_t(P_t^\circ(e; I))$ for all e, I . For $t = 0$ this is immediate from the convention that $P_0 \equiv P_0^\circ \equiv 0$. Fix any
682 $t \in \{1, \dots, T\}$. Since P is symmetric, there is a function f_t^{sym} so that for all $e = (u, v) \in E$ and
683 inputs I ,

$$\begin{aligned} P_t(e; I) &= f_t^{\text{sym}}(I(e), P_{t-1}(e; I), \{\!\{P_{t-1}(\{u, v'\}; I) : v' \sim u\}\!\}, \{\!\{P_{t-1}(\{u', v\}; I) : u' \sim v\}\!\}) \\ &= f_t^{\text{sym}}(I(e), h_{t-1}(P_{t-1}^\circ(e; I)), \{\!\{h_{t-1}(P_{t-1}^\circ(\{u, v'\}; I)) : v' \sim u\}\!\}, \{\!\{h_{t-1}(P_{t-1}^\circ(\{u', v\}; I)) : u' \sim v\}\!\}) \end{aligned}$$

684 which is indeed a well-defined function (independent of e, I) of

$$P_t^\circ(e; I) = (I(e), P_{t-1}^\circ(e; I), \{\!\{P_{t-1}^\circ(\{u, v'\}; I) : v' \sim u\}\!\}, \{\!\{P_{t-1}^\circ(\{u', v\}; I) : u' \sim v\}\!\}).$$

685 This completes the induction. Finally, since P is symmetric, there is a function \tilde{f}^{sym} such that
686 $\tilde{f}_v((P_T(e; I))_{e \in M_G(v)}) = \tilde{f}^{\text{sym}}(\{\!\{P_T(e; I) : e \in M_G(v)\}\!\})$ for all v, I . Hence we can write

$$\begin{aligned} \tilde{f}_v((P_T(e; I))_{e \in M_G(v)}) &= \tilde{f}^{\text{sym}}(\{\!\{P_T(e; I) : e \in M_G(v)\}\!\}) \\ &= \tilde{f}^{\text{sym}}(\{\!\{h_T(P_T^\circ(e; I)) : e \in M_G(v)\}\!\}) \end{aligned}$$

687 which is a well-defined function (independent of v, I) of $\{\!\{P_T^\circ(e; I) : e \in M_G(v)\}\!\}$ as needed. \square

688 D A quantitatively tight depth/memory separation

689 For each $n \in \mathbb{N}$, let $K_n := ([n], E_n)$ be the complete graph on $[n]$. In this section we show that there
690 is a function that can be computed by an edge message-passing protocol on K_n with constant rounds
691 and constant memory per processor, but for which any node message-passing protocol with T rounds
692 and B bits of memory requires $TB \geq \Omega(n)$. We remark that this separation is quantitatively tight
693 due to Proposition 3, although it is possible that a larger (e.g. even super-polynomial in n) depth
694 separation may be possible if the node message-passing protocol is restricted to constant memory per
695 processor.

696 At a technical level, the lower bound proceeds via a reduction from the *set disjointness problem* in
697 communication complexity, similar to the lower bounds in [Loukas \(2019\)](#).

698 **Definition 11.** Fix $m \in \mathbb{N}$. The set disjointness function $\text{DISJ}_m : \{0, 1\}^m \times \{0, 1\}^m \rightarrow \{0, 1\}$ is
699 defined as

$$\text{DISJ}_m(A, B) := \mathbb{1}[\forall i \in [m] : A_i B_i = 0].$$

700 The following fact is well-known; see e.g. discussion in [Håstad & Wigderson \(2007\)](#).

701 **Lemma 10.** *In the two-party deterministic communication model, the deterministic communication*
702 *complexity of DISJ_m is at least m .*

703 The main result of this section is the following:

704 **Theorem 11.** *Fix any even $n \in \mathbb{N}$. Define $g : \{0, 1\}^{E_n} \rightarrow \{0, 1\}^n$ by*

$$g(I)_v := \mathbb{1}[\exists \{i, j\} \in E_n : i, j \leq n/2 \wedge I(\{i, j\}) = I(\{n+1-i, n+1-j\}) = 1]$$

705 *for all $I \in \{0, 1\}^{E_n}$ and $v \in [n]$. Then the following properties hold:*

- 706 • *Any node message-passing protocol on K_n with T rounds and B bits of memory that*
707 *computes g requires $TB \geq \Omega(n)$*
- 708 • *There is an edge message-passing protocol on K_n with $O(1)$ rounds and $O(1)$ bits of*
709 *memory that computes g .*

710 *Proof.* Let $m := \binom{n/2}{2}$. Let $P = (f_{t,v})_{t,v}$ be a node message-passing protocol on K_n that computes
711 g with T rounds and B bits of memory. We design a two-party communication protocol for DISJ_m
712 as follows. Suppose that Alice holds input $X \in \{0, 1\}^m$ and Bob holds input $Y \in \{0, 1\}^m$. Let us
713 index the edges $\{i, j\} \in E_n$ with $i, j \leq n/2$ by $[m]$, and similarly index the edges $\{i, j\} \in E_n$ with
714 $i, j > n/2$ by $[m]$, in such a way that edge $\{i, j\}$ has the same index as edge $\{n+1-i, n+1-j\}$.
715 Let $I \in \{0, 1\}^{E_n}$ be defined by

$$I(\{i, j\}) := \begin{cases} X_{\{i, j\}} & \text{if } i, j \leq n/2 \\ Y_{\{i, j\}} & \text{if } i, j > n/2. \\ 0 & \text{otherwise} \end{cases}$$

716 Initially, Alice computes $\hat{P}_0(v) := 0$ for all $v \in \{1, \dots, n/2\}$, and Bob computes $\hat{P}_0(v) := 0$ for all
717 $v \in \{n/2 + 1, \dots, n\}$. The communication protocol then proceeds in T rounds. At round $t \in [T]$,
718 Alice sends $(\hat{P}_{t-1}(v))_{1 \leq v \leq n/2}$ to Bob, and Bob sends $(\hat{P}_{t-1}(v))_{n/2+1 \leq v \leq n}$ to Alice. Alice then
719 computes

$$\hat{P}_t(v) := f_{t,v}((\hat{P}_{t-1}(v'))_{v' \in [n]}, (I(e))_{e \in M_{K_n}(v)})$$

720 for each $1 \leq v \leq n/2$, and Bob computes the same for each $n/2 < v \leq n$. Note that for any $i \leq n/2$
721 and edge $e \in M_{K_n}(i)$, Alice can compute $I(e)$. Similarly, for any $i > n/2$ and edge $e \in M_{K_n}(i)$,
722 Bob can compute $I(e)$. Thus, this computation is well-defined. After round T , Alice and Bob output
723 $1 - \hat{P}_T(1)$ and $1 - \hat{P}_T(n)$ respectively.

724 This defines a communication protocol. Since $\hat{P}_t(v) \in \{0, 1\}^B$ for each $v \in [n]$ and $t \in [T]$, the total
725 number of bits communicated is at most nBT . Moreover, by induction it's clear that Alice and Bob
726 output $1 - P_T(1; I)$ and $1 - P_T(n; I)$ respectively. By assumption that P computes g and the fact
727 that $g(I)_v = 1 - \text{DISJ}_m(X, Y)$ for all $v \in [n]$, we have that $1 - P_T(1; I) = 1 - P_T(n; I) = 0$ if
728 $\text{DISJ}_m(I) = 0$, and $1 - P_T(1; I) = 1 - P_T(n; I) = 1$ if $\text{DISJ}_m(I) = 1$. Thus, this communication
729 protocol computes DISJ_m . By Lemma 10, it follows that $nBT \geq m = \Omega(n^2)$, so $BT = \Omega(n)$ as
730 claimed.

731 Next, we exhibit an edge message-passing protocol on K_n that computes g with six rounds and one
732 bit of memory. For $1 \leq t \leq 6$ and $e \in E_n$, define $f_{t,e} : \{0, 1\}^{M_G(e)} \times \{0, 1\} \rightarrow \{0, 1\}$ as follows:

$$\begin{aligned} f_{1,\{i,j\}}(x, y) &:= y \\ f_{2,\{i,j\}}(x, y) &:= x_{\{n+1-i,j\}} \\ f_{3,\{i,j\}}(x, y) &:= x_{\{i,n+1-j\}} \\ f_{4,\{i,j\}}(x, y) &:= \mathbb{1}[y = x_{\{i,j\}} \wedge i, j \leq n/2] \\ f_{5,\{i,j\}}(x, y) &:= \mathbb{1}[\exists k \in [n] : x_{\{i,k\}} = 1] \\ f_{6,\{i,j\}}(x, y) &:= \mathbb{1}[\exists k \in [n] : x_{\{i,k\}} = 1]. \end{aligned}$$

733 Also define $\tilde{f}_v : \{0, 1\}^{M_G(v)} \rightarrow \{0, 1\}$ for each $v \in [n]$ by $\tilde{f}_v(x) := x_{\{x,1\}}$. It can be checked that
734 the computation of P at timestep $t = 6$ is

$$P_6(\{i, j\}; I) := \mathbb{1}[\exists k, \ell \in [n/2] : I(\{k, \ell\}) = I(\{n+1-k, n+1-\ell\})] = g(I).$$

735 From the definition of \tilde{f} , it follows that P computes g . □

736 E Further details on synthetic task over Ising models

737 E.1 Background on belief propagation

738 A classical way to calculate the marginals $\{\mathbb{E}[x_i]\}$ of an Ising model, when the associated graph is a
739 tree, is to iterate the message passing algorithm:

$$\nu_{i \rightarrow j}^{(t+1)} = \tanh \left(h_i + \sum_{k \in \partial_i \setminus j} \tanh^{-1} \left(\tanh(J_{ik}) \nu_{k \rightarrow i}^{(t)} \right) \right) \quad (5)$$

740 When the graph is a tree, it is a classical result ((Mezard & Montanari, 2009), Theorem 14.1) that the
741 above message-passing algorithm converges to values ν^* that yield the correct marginals, namely:

$$\mathbb{E}[x_i] = \tanh \left(h_i + \sum_{k \in \partial_i} \tanh^{-1} \left(\tanh(J_{ik}) \nu_{k \rightarrow i}^* \right) \right).$$

742 The reason the updates converge to the correct values on a tree topology is that they implicitly
743 simulate a dynamic program. Namely, we can write down a recursive formula for the marginal of
744 node i which depends on sums spanning each of the subtrees of the neighbors of i (i.e., for each
745 neighbor j , the subgraph containing j that we would get if we removed edge $\{i, j\}$).

746 If we root the tree at an arbitrary node r , we can see that after completing a round of message
747 passing from the leaves to the root, and another from the root to the leaves, each subtree of i will be
748 (inductively) calculated correctly.

749 Moreover, even though the updates (5) are written over edges, the dynamic programming view makes
 750 it clear an equivalent message-passing scheme can be written down where states are maintained over
 751 the *nodes* in the graph. Namely, for each node v , we can maintain two values $h_{v,\text{down}}$ and $h_{v,\text{up}}$,
 752 which correspond to the values that will be used when v sends a message upwards (towards the root)
 753 or downwards (away from the root). Then, for appropriately defined functions F, G (depending on
 754 the potentials J and h), one can “simulate” the updates in (5):

$$h_{v,\text{up}}^{(t+1)} \leftarrow F \left(\{h_{w,\text{up}}^{(t)} : w \in v \cup \text{Children}(v)\} \right) \quad (6)$$

$$h_{v,\text{down}}^{(t+1)} \leftarrow G \left(h_{\text{Parent}(v),\text{down}}^{(t)}, \{h_{w,\text{up}}^{(t)}\}_{w \in \text{Children}(v)} \right) \quad (7)$$

755 Intuitively, $h_{v,\text{up}}$ captures the effective external field induced by the subtree rooted at v on
 756 $\text{Parent}(v)$. After the upward messages propagate, the root r can compute its correct marginal.
 757 Once $h_{\text{Parent}(v),\text{down}}^{(t)}$ is the correct marginal for $\text{Parent}(v)$ at some step, $h_{v,\text{down}}^{(t)}$ will be the correct
 758 marginal for v at all subsequent steps.

759 E.2 GCN-based architectures to calculate marginals

760 The belief-propagation updates (5) naturally fit the general edge-message passing paradigm from
 761 (2). In fact, they fit even more closely a “directed” version of the paradigm, in which each edge
 762 $\{i, j\}$ maintains two embeddings $h_{i \rightarrow j}, h_{j \rightarrow i}$, such that the embedding for direction $h_{i \rightarrow j}$ depends
 763 on the embeddings $\{h_{k \rightarrow i}\}_{\{k, i\} \in E}$. With this modification to the standard edge GCN architecture
 764 Eq. (4), it is straightforward to implement (5) with one layer, using a particular choice of activation
 765 functions and weight matrices W (since, in particular, in our dataset all edge potentials $J_{i,j}$ are set
 766 to 1). Similarly, with a directed version of the node GCN architecture Eq. (3), where each node
 767 maintains an “up” embedding as well as a “down” embedding, it is straightforward to implement the
 768 “node-based” dynamic programming solution (6)-(7).

769 We call the architectures that do not maintain directionality Node-U and Edge-U (depending on
 770 whether they use a node-based or edge-based GCN). We call the “directed” architectures Node-D and
 771 Edge-D respectively. Since there are only initial node features (input as node potentials $\{h_i\}_{i \in V}$), for
 772 the edge based architectures we initialize the edge features as a concatenation of the node features of
 773 the endpoints of the edge. The results we report for each architecture are the best over a sweep of
 774 depth $\in \{5, 10, 15, 20, 25, 30\}$ and width $\in \{10, 32, 64\}$.

775 E.3 Edge-based models improve over node-based models

776 In Figure 2 we show the results for several tree topologies: a complete binary tree (of size 31), a
 777 path graph (of size 30), and uniformly randomly chosen trees of size 30 (the results in Figure 2 are
 778 averaged over 3 samples of tree). The architectures in the legend (Node-U, Edge-U, Node-D, Edge-D)
 779 are based on a standard GCN, and detailed in Section E.2

780 We can see that for both the undirected and directed versions, adding edge embeddings improves
 781 performance. The improved performance of all directed versions compared to their undirected
 782 counterpart is not very surprising: the standard, undirected GCN architecture treats all neighbors
 783 symmetrically — hence, the directed versions can more easily simulate something akin to the belief
 784 propagation updates (5) as well as the node-based dynamic programming (6)-(7).

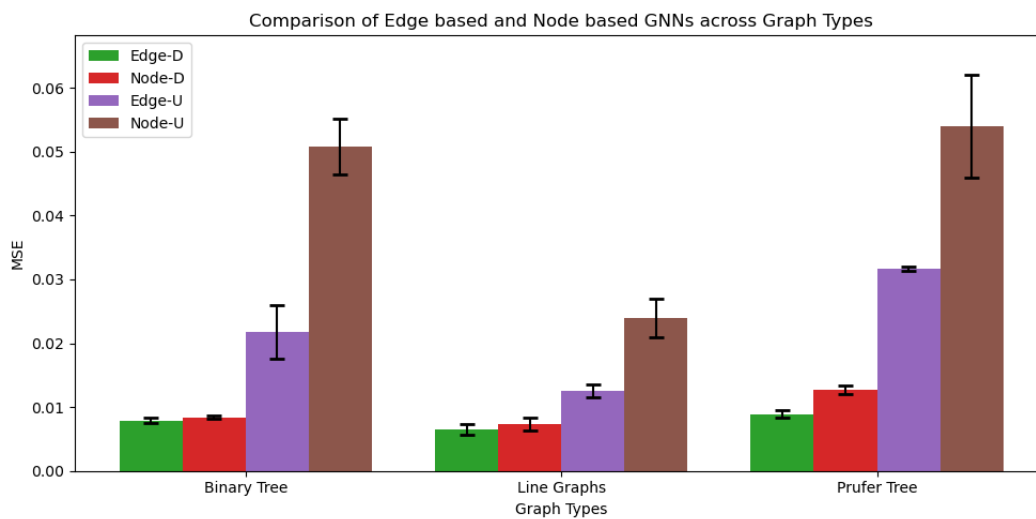


Figure 2: Comparison of four architectures for calculating node marginals in an Ising model. The architectures considered are node-embedding (3) and edge-embedding (4) versions of a GCN (correspondingly labeled Node-U and Edge-U), as well as their “directed” counterparts, as described in Section E.2, correspondingly labeled Node-D and Edge-D. The x-axis groups results according to the topology of the graph, the y-axis is MSE (lower is better). The mean and variances are reported over 3 runs for the best choice of depth and width over the sweep described in Section E.2.