
Pruning Edges and Gradients to Learn Hypergraphs from Larger Sets

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

This paper aims for set-to-hypergraph prediction, where the goal is to infer the set of relations for a given set of entities. This is a common abstraction for applications in particle physics, biological systems, and combinatorial optimization. We address two common scaling problems encountered in set-to-hypergraph tasks that limit the size of the input set: the exponentially growing number of hyperedges and the run-time complexity, both leading to higher memory requirements. We make three contributions. First, we propose to predict and supervise the *positive* edges only, which changes the asymptotic memory scaling from exponential to linear. Second, we introduce a training method that encourages iterative refinement of the predicted hypergraph, which allows us to skip iterations in the backward pass for improved efficiency and constant memory usage. Third, we combine both contributions in a single set-to-hypergraph model that enables us to address problems with larger input set sizes. We provide ablations for our main technical contributions and show that our model outperforms prior state-of-the-art, especially for larger sets.

1 Introduction

Inferring the relational structure for a given set of entities is a common abstraction for many applications, including vertex reconstruction in particle physics [1, 2], inferring higher-order interactions in biological and social systems [3, 4] or combinatorial optimization problems, such as finding the convex hull or Delaunay triangulation [2, 5]. The wide spectrum of different application areas underlines the expressivity of this abstract task, which is known in machine learning as set-to-hypergraph prediction [2]. Here, the hypergraph generalizes the pairwise relations in a graph to multi-way relations, a.k.a. hyperedges. In biological systems, multi-way relationships (hyperedges) among genes and proteins are relevant for protein complexes and metabolic reactions [6]. A subgroup in a social network can be understood as a hyperedge that connects subgroup members [7] and in images interacting objects can be modeled by scene graphs [8] which is useful for counting objects [9]. We distinguish the set-to-hypergraph task from the related, but different, task of link prediction that aims to discover the missing edges in a graph, given the set of vertices *and a subset of the edges* [10]. For the set-to-hypergraph problem considered in this paper, we start with a set of nodes without any edges.

A common approach to set-to-hypergraph problems is to decide for *every* edge, whether it exists or not [2]. For a set of n nodes, the number of all possible hyperedges grows in $\mathcal{O}(2^n)$, which already becomes intractable for moderately sized n . This is the *scaling* problem of set-to-hypergraph prediction that we will address in this paper. Combinatorial optimization challenges, like set-to-hypergraph prediction, introduce the additional problem of *complexity*. For example, convex hull finding in d dimensions has a run time complexity of $\mathcal{O}(n \log(n) + n^{\lfloor \frac{d}{2} \rfloor})$ [11]. This means that larger input sets require more computation regardless of the quality of the hypergraph prediction algorithm. Indeed, it was observed in [2] that for larger set sizes performance was worse. In this paper, we aim to address the scaling and complexity problems in order to predict hypergraphs from larger sets.

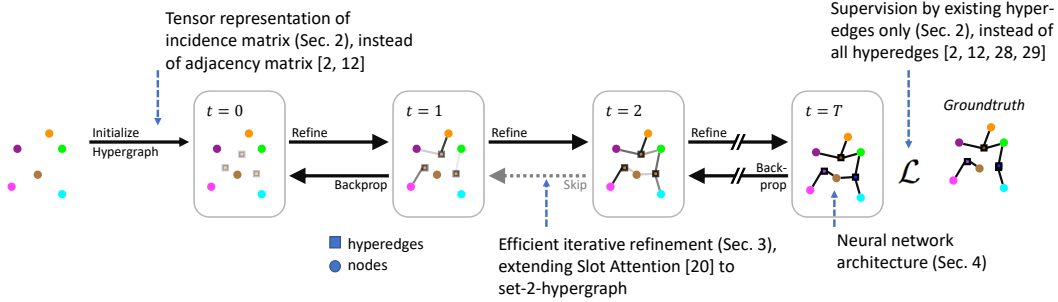


Figure 1: Outline of the paper and our contributions.

41 We make three contributions in this paper. **1.** We propose a scalable set-to-hypergraph framework that
 42 represents the hypergraph structure with a pruned incidence matrix — the incidence matrix without
 43 the edges (rows) that have no incident nodes (Section 2). We prove that computing the loss over the
 44 full incidence matrix is equivalent solely using the pruned version, thus improving the asymptotic
 45 memory requirements from $\mathcal{O}(2^n)$ to $\mathcal{O}(mn)$, that is linear in the input set size. **2.** To address the
 46 complexity problem, we introduce in Section 3 a training method that encourages iterative refinement
 47 of the predicted hypergraph with memory requirements scaling constant in the number of refinement
 48 steps. This addresses the need for more compute by complex problems in a scalable manner. Third,
 49 we combine in Section 4 the efficient representation from the first contribution with the requirements
 50 of the scalable training method from the second contribution in a recurrent model that performs
 51 iterative refinement on a pruned incidence matrix. Our model handles different input set sizes and
 52 varying numbers of edges while respecting the permutation symmetry of both. Figure 1 illustrates
 53 the three contributions, the neural network, and the organization of the paper. In our experiments in
 54 Section 5, we provide an in-depth ablation of each of our technical contributions and compare our
 55 model against prior work on common set-to-hypergraph benchmarks.

56 **Preliminary.** Hypergraphs generalize normal graphs by replacing the normal edges that connect
 57 exactly two nodes with hyperedges that connect an arbitrary number of nodes. Since we only consider
 58 the general version, we shorten hyperedges to edges in the remainder of the paper. Given a set of
 59 nodes and their corresponding input feature vectors X , in the set-to-hypergraph task we want to learn
 60 a neural network f that predicts for all possible edges whether the target hypergraph contains that
 61 edge. Two edges are equivalent if and only if they are incident to the same nodes. For example,
 62 the input set could consist of objects in an image and the edges would represent their relations.
 63 Next, we provide an overview of the Set2Graph neural network [2]. We focus on it because most
 64 previous networks follow a similar structure. In Set2Graph, f consists of a collection of functions
 65 $f := (F^1, F^2, \dots, F^k)$, where F^k maps a set of nodes to a set of k -edges (edges with k incident
 66 nodes). All F^k are composed of three steps: a set-to-set model maps the input set to a latent set,
 67 a broadcasting step forms all possible k -tuples from the latent set elements, and a final graph-to-
 68 graph model that predicts for each k -edge whether it exists or not. The output of every F^k is then
 69 represented by an adjacency tensor, a tensor with k dimensions each of length n . Serviansky et al. [2]
 70 show that this can approximate any continuous set-to- k -edge function, and by extension, the family
 71 of F^k functions can approximate any continuous set-to-hypergraph function. Since the asymptotic
 72 memory scaling of F^k is in $\mathcal{O}(n^k)$, modeling k -edges beyond $k > 2$ already becomes intractable in
 73 many settings and one has to apply heuristics to recover higher-order edges from pairwise edges [2].

74 2 Scaling by pruning the non-existing edges

75 In this section, we propose a solution for the memory scaling problem encountered in set-to-
 76 hypergraph tasks. The goal is to learn a model $f(X) = \mathcal{H}$ that maps a set X of input vectors
 77 to the hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E}, I)$ that consists of latent node and edge features and the incidence
 78 matrix. In what follows, we describe each constituent of \mathcal{H} and how the loss is computed.

79 **Nodes.** Each input element $x \in X$ gets an “identity” as a node in the hypergraph, meaning if a
 80 subset of the nodes is connected by an edge then there exists a relation between the corresponding
 81 input elements. We differentiate between the input elements and the nodes of the hypergraph, as we

82 expect the latter to be represented as latent vectors $\mathbf{v} \in \mathcal{V} = \{\mathbf{v}_i \in \mathbb{R}^{d_{\mathcal{V}}} | \mathbf{x}_i \in X\}$. The importance
83 of this becomes clear in the discussion on the training objective later on.

84 **Edges.** The set of all possible edges can be expressed using the power set $\mathcal{P}(\mathcal{V}) \setminus \{\emptyset\}$, that is the
85 set of all subsets of \mathcal{V} minus the empty set. Different from the situation with the nodes, we do not
86 know which edge will be part of the hypergraph, since this is what we want to predict. Listing out
87 all $2^{|\mathcal{V}|-1}$ possible edges and deciding for each edge whether it exists or not, becomes intractable
88 very quickly. Furthermore, we observe that in many cases the number of existing edges is much
89 smaller than the total number of possible edges. We leverage this observation by keeping only a fixed
90 number of edges m that is sufficient to cover all (or most) hypergraphs for a given task. Thus, we
91 improve the memory requirement from $\mathcal{O}(2^{|\mathcal{V}|}d_{\mathcal{E}})$ to $\mathcal{O}(md_{\mathcal{E}})$, where $d_{\mathcal{E}}$ is the vector size of the
92 edge representations $e \in \mathcal{E}$. All possible edges that do not have an edge representation in \mathcal{E} are
93 *implicitly* predicted to not exist. After specifying the training objective, we provide a more formal
94 argument on why pruning all but m edges from \mathcal{E} is sound.

95 **Incidence matrix.** Since both the nodes and edges are represented by latent vectors, we require
96 an additional component for specifying the connections. Different from previous approaches, we
97 use the incidence matrix over adjacency tensors [2, 12]. The two differ in that incidence describes
98 whether an edge is connected to a node, while adjacency describes whether an edge between a subset
99 of nodes exists. The rows of the incidence matrix \mathbf{I} correspond to the edges and the columns to the
100 nodes. Thus, an entry $\mathbf{I}_{i,j} \in [0, 1]$ represents the probability of the i -th edge being incident to the
101 j -th node. Theoretically, we can express any hypergraph in both representations, but pruning edges
102 becomes especially simple in the incidence matrix, where we just remove the corresponding rows.
103 We interpret the pruned edges $e \notin \mathcal{E}$ that have no corresponding row in the pruned \mathbf{I} as having zero
104 probability of being incident to any node.

105 **Loss function.** We apply the loss function only on the incidence probability matrix \mathbf{I} . For efficiency
106 purposes, we would like to train each incidence value separately as a binary classification and apply a
107 constant threshold (> 0.5) on \mathbf{I} at inference time to get a binary incidence matrix. In probabilistic
108 terms, this translates to a factorization of the joint distribution $p(\mathbf{I}|\mathbf{X})$ as, $\prod_{i,j} p(\mathbf{I}_{i,j}|\mathbf{X})$. In order to
109 still be able to model interrelations between different incidence probabilities, we impose a requirement
110 on the model f : the probability $\mathbf{I}_{i,j}$ produced by f depends on e_i and v_j . We could alternatively
111 factorize $p(\mathbf{I}|\mathbf{X})$ through the chain rule, but that is much less efficient as it introduces nm non-
112 parallelizable steps. This highlights the importance of the latent node and edge representations, which
113 enable us to model the dependencies in the output efficiently because predicting all $\mathbf{I}_{i,j}$ can happen
114 in parallel. Furthermore, this changes our assumption on the incidence probabilities from that of
115 independence to *conditional* independence on e_i and v_j , and we apply the binary cross-entropy loss
116 on each $\mathbf{I}_{i,j}$.

117 The binary classification over $\mathbf{I}_{i,j}$ highlights yet another reason for picking the incidence representa-
118 tion over the adjacency. Let us assume we are trying to learn a binary classifier that predicts for every
119 entry in the adjacency tensor whether it is 0 or 1. Removing all the 0 values (non-existing edges)
120 from the training set will clearly not work out in the adjacency case. In contrast, an existing edge in
121 the incidence matrix contains both 1s and 0s (except for the edge connecting all nodes), ensuring
122 that a binary incidence classifier sees both positive and negative examples. However, an adjacency
123 tensor has the advantage that the order of the entries is fully decided by the order of the nodes, which
124 are given by the input \mathbf{X} in our case. In the incidence matrix, the row order of the incidence matrix
125 depends on the edges, which are orderless.

126 When comparing the predicted incidence matrix with a ground-truth matrix, we need to account for
127 the orderless nature of the edges and the given order of the nodes. Hence, we require a loss function
128 that is invariant toward reordering over the rows of the incidence matrix, but equivariant to reordering
129 over the columns. We achieve this by matching every row in \mathbf{I} with a row in the pruned ground-truth
130 incidence matrix \mathbf{I}^* (containing the existing edges), such that the binary cross-entropy loss H over
131 all entries is minimal:

$$\mathcal{L}(\mathbf{I}, \mathbf{I}^*) = \min_{\pi \in \Pi} \sum_{i,j} H(\mathbf{I}_{\pi(i),j}, \mathbf{I}_{i,j}^*) \quad (1)$$

132 Finding a permutation π that minimizes the total loss is known as the linear assignment problem and
133 we solve it with an efficient variant of the Hungarian algorithm [13, 14]. We discuss the implications
134 on the computational complexity of this in [Appendix B](#).

135 Having established the training objective in Equation 1, we can now offer a more formal argument on
 136 the soundness of supervising existing edges only while *pruning the non-existing ones*, where \mathbf{J} can
 137 be understood as the full incidence matrix (proof in Appendix A):

138 **Proposition 1** (Supervising only existing edges). *Let $\mathbf{J} \in [\epsilon, 1]^{(2^n-1) \times n}$ be a matrix with at most
 139 m rows for which $\exists j: \mathbf{J}_{ij} > \epsilon$, with a small $\epsilon > 0$. Similarly, let $\mathbf{J}^* \in \{0, 1\}^{(2^n-1) \times n}$ be a matrix
 140 with at most m rows for which $\exists j: \mathbf{J}_{ij} = 1$. Let $\text{prune}(\cdot)$ denote the function that maps from a
 141 $(2^n - 1) \times n$ matrix to a $k \times n$ matrix, by removing $(2^n - 1) - k$ rows where all values are $\leq \epsilon$.
 142 Then, for a constant $c = (2^n - 1 - k)n \cdot H(\epsilon, 0)$ and all such \mathbf{J} and \mathbf{J}^* :*

$$\mathcal{L}(\mathbf{J}, \mathbf{J}^*) = \mathcal{L}(\text{prune}(\mathbf{J}), \text{prune}(\mathbf{J}^*)) + c \quad (2)$$

143 The matrix $\text{prune}(\mathbf{J})$ can be understood as the pruned incidence matrix that we defined earlier and
 144 $\text{prune}(\mathbf{J}^*)$ as the pruned ground-truth. In practice, the ϵ corresponds to a lower bound on the log in
 145 the entropy computation, like -100 in PyTorch [15]. Since the losses in Equation 2 are equivalent
 146 up to an additive constant, the gradients of the parameters are exactly equal in both the pruned and
 147 non-pruned cases. Thus, pruning the non-existing edges does not affect learning, while significantly
 148 reducing the asymptotic memory requirements.

149 **Summary.** In set-to-hypergraph tasks, the number of different edges that can be predicted grows
 150 exponentially with the input set size. We address this computational limitation by representing the
 151 edge connections with the incidence matrix and pruning all non-existing edges *before* explicitly
 152 deciding for each edge whether it exists or not. We show that pruning the edges is sound when the
 153 loss function respects the permutation symmetry in the edges.

154 3 Scaling by skipping the non-essential gradients

155 Next, we consider how to *learn* the pruned incidence matrix. Some tasks may require more compute
 156 than others, which can result in worse performance or intractable models if not properly addressed.
 157 A naive approach would increase the number of parameters, either by increasing the number of
 158 hidden dimensions or the depth of the neural network, which is clearly not scalable. Furthermore,
 159 the memory requirement of backprop would also grow with greater depth. Instead, we would like to
 160 increase the amount of sequential computation by reusing parameters. That is, we want the model
 161 f to be recurrent, $\mathcal{H}^{t+1} = f(\mathbf{X}, \mathcal{H}^t)$ with t denoting the t -th application of f starting from $t = 0$.
 162 Recurrent models are commonly applied to sequential data, where the input varies for each time step t
 163 [16], for example, the words in a sentence. In our case, we use the same input \mathbf{X} at every step. Using
 164 a recurrent model, we can increase the total number of iterations – to scale the number of sequential
 165 computation steps – without increasing the number of parameters. However, the recurrence does not
 166 address the growing memory requirements of backprop, since the activations of each iteration still
 167 need to be kept in memory.

168 **Iterative refinement.** In the rest of this section, we present an efficient training algorithm that
 169 can scale to any number of iterations at a constant memory cost. We build on the idea that if each
 170 iteration applies a small refinement, then it becomes unnecessary to backprop through every iteration.

171 We can define an iterative refinement as reducing the loss
 172 (by a little) after every iteration, $\mathcal{L}(\mathbf{I}^t, \mathbf{I}^*) < \mathcal{L}(\mathbf{I}^{t-1}, \mathbf{I}^*)$.
 173 Thus, the long-term dependencies between \mathcal{H}^t for t 's that
 174 are far apart can also be ignored, since f only needs to
 175 improve the current \mathcal{H}^t . We can encourage f to iteratively
 176 refine the prediction \mathcal{H}^t , by applying the loss \mathcal{L} after each
 177 iteration. This has the effect that f learns to move the \mathcal{H}^t in
 178 the direction of the negative gradient of \mathcal{L} , making it similar
 179 to a gradient descent update.

180 **Backprop with skips.** Similar to previous works that en-
 181 courage iterative refinement through (indirect) supervision
 182 on the intermediate steps [17], we expect the changes of each
 183 step to be small. Thus, it stands to reason that supervising
 184 *every* step is unnecessary and redundant. This leads us to
 185 a more efficient training algorithm that skips iterations in
 186 the backward-pass of backprop. Algorithm 1 describes the

Algorithm 1: Backprop with skips

Input: $\mathbf{X}, \mathbf{I}^*, S, B$
 $\mathcal{H} \leftarrow \text{initialize}(\mathbf{X})$
for s **in** S :
 with `no_grad()` :
 for t **in** `range(s)` :
 $\mathcal{H} \leftarrow f(\mathbf{X}, \mathcal{H})$
 $l \leftarrow 0$
 for t **in** `range(B)` :
 $\mathcal{H} \leftarrow f(\mathbf{X}, \mathcal{H})$
 $l \leftarrow l + \mathcal{L}(\mathcal{H}, \mathbf{I}^*)$
 `backward(l)`
 `gradient_step_and_reset()`

187 training procedure in pseudocode (in a syntax similar to PyTorch [15]). The argument S is a list of
 188 integers of length N that is the number of gradient updates per mini-batch. Each gradient update
 189 consists of two phases, first $s \in S$ iterations without gradient accumulation and then B iterations that
 190 add up the losses for backprop. Through these hyperparameters, we control the amount of resources
 191 used during training. Increasing hyperparameter B allows for models that do not strictly decrease the
 192 loss after every step and require supervision over multiple subsequent steps. Note that having the
 193 input \mathbf{X} at every refinement step is important so that the model does not forget the initial problem.

194 **Summary.** Motivated by the need for more compute to address complex problems, we propose a
 195 method that increases the amount of sequential compute of the neural network without increasing
 196 the memory requirement at training time. Our training algorithm requires the model f to perform
 197 iterative refining of the hypergraph, for which we present a method in the next section.

198 4 Scaling the set-to-hypergraph prediction model

199 In Section 2 and Section 3 we proposed two methods to overcome the memory scaling problems that
 200 appear for set-to-hypergraph tasks. To put these methods into practice, we need to specify a model f
 201 that fulfills the required properties. In what follows, we propose a specific implementation for each
 202 such property.

203 **Initialization.** As the starting point for the iterative refinement, we initialize the nodes \mathcal{V}^0 from the
 204 input set as $\mathbf{v}_i^0 = \mathbf{W} \mathbf{x}_i + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{d_{\mathcal{V}} \times d_{\mathbf{X}}}$, $\mathbf{b} \in \mathbb{R}^{d_{\mathcal{V}}}$ are learnable parameters. The affine
 205 map allows for hidden dimensions $d_{\mathcal{V}}$ that are different from the input feature dimensions $d_{\mathbf{X}}$. An
 206 informed initialization similar to the nodes is not available for the edges and the incidence matrix,
 207 since these are what we aim to predict. Instead, we choose an initialization scheme that respects
 208 the permutation symmetry of a set of edges while also ensuring that each edge starts out differently.
 209 The last point is necessary for permutation-equivariant operations to distinguish between different
 210 edges. The random initialization $\mathbf{e}_i^0 \sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$, with shared learnable parameters $\boldsymbol{\mu} \in \mathbb{R}^{d_{\mathcal{E}}}$
 211 and $\boldsymbol{\sigma} \in \mathbb{R}^{d_{\mathcal{E}}}$ fulfills both these properties, as it is highly unlikely for two samples to be identical.

212 **Conditional independence.** We want the incidence probabilities $\mathbf{I}_{i,j}$ to be conditionally independ-
 213 ent of each other given \mathbf{e}_i and \mathbf{v}_j . A straightforward way to model this is by concatenating both
 214 vectors (denoted with $[\cdot]$) and applying an MLP with a sigmoid activation on the scalar output:

$$\mathbf{I}_{i,j}^t = \text{MLP}([\mathbf{e}_i^{t-1}, \mathbf{v}_j^{t-1}]) \quad (3)$$

215 The superscripts point out that we produce a new incidence matrix for step t based on the edge and
 216 node vectors from the previous step. Note that we did not specify an initialization for the incidence
 217 matrix, since we directly replace it in the first step.

218 **Iterative refinement.** The training algorithm in Section 3 assumes that f performs iterative re-
 219 finement on \mathcal{H}^t , but leaves open the question on how to design such a model. Instead of iteratively
 220 refining the incidence matrix, i.e., the only term that appears in the loss (Equation 1), we focus on
 221 refining the edges and nodes.

222 A refinement step for some node $\mathbf{v}_i \in \mathcal{V}$ needs to account for the rest of the hypergraph, which
 223 also changes with each iteration. For this purpose, we apply the permutation-equivariant DeepSets
 224 [18] to produce node updates dependent on the full set of nodes from the previous iteration \mathcal{V}^{t-1} .
 225 The permutation-equivariance of DeepSets means that the output set retains the input order; thus
 226 it is sensible to refer to \mathbf{v}_i^t as the updated version of the *same* node \mathbf{v}_i^{t-1} from the previous step.
 227 Furthermore, we concatenate the aggregated neighboring edges weighted by the incidence probabili-
 228 ties $\rho_{\mathcal{E} \rightarrow \mathcal{V}}(j, t) = \sum_{i=1}^k \mathbf{I}_{i,j}^t \mathbf{e}_i^{t-1}$, to incorporate the relational structure between the nodes. This
 229 aggregation works akin to message passing in graph neural networks [19]. An indispensable input,
 230 required for adding skips in the backward pass during training, is the input features \mathbf{X} . Instead
 231 of directly concatenating the raw features \mathbf{x}_i , we use the initial nodes \mathbf{v}_i^0 . Finally, we express the
 232 refinement part for the nodes as:

$$\mathcal{V}^t = \text{DeepSets}(\{[\mathbf{v}_j^{t-1}, \rho_{\mathcal{E} \rightarrow \mathcal{V}}(j, t), \mathbf{v}_j^0] \mid j = 1 \dots n\}) \quad (4)$$

233 The updates to the edges \mathcal{E}^t mirror that of the nodes, except for the injection of the input set. Together
 234 with the aggregation function $\rho_{\mathcal{V} \rightarrow \mathcal{E}}(i, t) = \sum_{j=1}^n \mathbf{I}_{i,j}^t \mathbf{v}_j^{t-1}$, we can update the edges as:

$$\mathcal{E}^t = \text{DeepSets}(\{[\mathbf{e}_i^{t-1}, \rho_{\mathcal{V} \rightarrow \mathcal{E}}(i, t)] \mid i = 1 \dots k\}) \quad (5)$$

235 By sharing the parameters between different refinement steps, we naturally obtain a recurrent model.
 236 Previous works on recurrent models [20] saw improvements in training convergence by including
 237 layer normalization [21] between each iteration. Shortcut connections in ResNets [22] have been
 238 shown to encourage iterative refinement of the latent vector [17]. We add both shortcut connections
 239 and layer normalization to the updates in Equation 4 and Equation 5. Although we prune the negative
 240 edges, we still want to predict a variable number thereof. To achieve that we simply extend the
 241 incidence model in Equation 3 with an existence indicator:

$$\hat{I}_i^t = \sigma_i^t I_i^t \quad (6)$$

242 This can be seen as factorizing the probability into “ $p(e_i \text{ incident to } v_j) \cdot p(e_i \text{ exists})$ ” and replaces
 243 the aggregation weights in $\rho_{\mathcal{E} \rightarrow \mathcal{V}}$ and $\rho_{\mathcal{V} \rightarrow \mathcal{E}}$.

244 **Summary.** We propose a model that fulfills the requirements of our scalable set-to-hypergraph
 245 training framework from Section 2 and Section 3. By adding shortcut connections, we encourage it
 246 to perform iterative refinements on the hypergraph while being permutation equivariant with respect
 247 to both the nodes and the edges.

248 5 Experiments

249 In this section, we evaluate our approach on multiple set-to-hypergraph tasks, in order to compare to
 250 prior work and examine the main design choices. We refer to Appendix D for further details, results,
 251 and ablations. Code is included in the supplementary material.

252 5.1 Scaling Set-to-Hypergraph Prediction

253 First, we compare our model from Section 4 on three different set-to-hypergraph tasks against the
 254 state-of-the-art model. This allows us to see the difference between predicting the incidence matrix
 255 and predicting the adjacency tensors.

256 **Baselines.** Our main comparison is against Set2Graph [2], which is a strong and representative
 257 baseline for approaches that predict the adjacency structure, which we generally refer to as adjacency-
 258 based approaches. Serviansky et al. [2] modify the task in two of the benchmarks, to avoid storing
 259 an intractably large adjacency tensor. We explain in Appendix D how this affects the comparison.
 260 Additionally, we compare to Set Transformer [23] and Slot Attention [20], which we adapt to the
 261 set-to-hypergraph setting by treating the output as the pruned set of edges and producing the incidence
 262 matrix with the MLP from Equation 3. We refer to these two as incidence-based approaches that also
 263 include our model.

264 **Particle partitioning.** Particle colliders are an important tool for studying the fundamental particles
 265 of nature and their interactions. During a collision, several particles are emanated and measured by
 266 nearby detectors, while some particles decay beforehand. Identifying which measured particles share
 267 a common progenitor is an important subtask in the context of vertex reconstruction [24]. We can
 268 treat this as a set-to-hypergraph task: the set of measured particles is the input set and the common
 269 progenitors are the edges of the hypergraph. We use a simulated dataset of 0.9M data-sample with
 270 the default train/validation/test split [2, 24]. Each data-sample is generated from on one of three
 271 different distributions for which we report the results separately: *bottom jets*, *charm jets* and *light*
 272 *jets*. The ground-truth target is the incidence matrix that can also be interpreted as a partitioning
 273 of the input elements since every particle has exactly one progenitor (edge). In Table 1 we report
 274 the performances on each type of jet as the F1 score and Adjusted Rand Index (ARI). Our method
 275 outperforms all alternatives on bottom and charm jets while being competitive on light jets.

276 **Convex hull.** The convex hull of a finite set of d -dimensional points can be efficiently represented
 277 as the set of simplices that enclose all points. In the 3D case, each simplex consists of 3 points
 278 that together form a triangle. For the general d -dimensional case, the valid incidence matrices are
 279 limited to those with d incident vertices per edge. Finding the convex hull is an important and
 280 well-understood task in computational geometry, with optimal exact solutions [11, 25]. Nonetheless,
 281 predicting the convex hull for a given set of points poses a challenging problem for current machine
 282 learning methods, especially when the number of points increases [2, 5]. We generate an input set
 283 by drawing n 3-dimensional vectors from one of two distributions: Gaussian or spherical. For the
 284 Gaussian setting, points are sampled i.i.d. from a standard normal distribution. For the spherical

Table 1: Particle partitioning results. On three jet types performance measured in F1 score and adjusted rand index (ARI) for 11 different seeds. Our method outperforms the baselines on bottom and charm jets while being competitive on light jets.

Model	bottom jets		charm jets		light jets	
	F1	ARI	F1	ARI	F1	ARI
Set2Graph	0.646 \pm 0.003	0.491 \pm 0.006	0.747 \pm 0.001	0.457 \pm 0.004	0.972 \pm 0.001	0.931 \pm 0.003
Set Transformer	0.630 \pm 0.004	0.464 \pm 0.007	0.747 \pm 0.003	0.466 \pm 0.007	0.970 \pm 0.001	0.922 \pm 0.003
Slot Attention	0.600 \pm 0.012	0.411 \pm 0.021	0.728 \pm 0.008	0.429 \pm 0.016	0.963 \pm 0.002	0.895 \pm 0.009
Ours	0.679 \pm 0.002	0.548 \pm 0.003	0.763 \pm 0.001	0.499 \pm 0.002	0.972 \pm 0.001	0.926 \pm 0.002

Table 2: Convex hull results measured as F1 score. Our method outperforms all baselines considerably for all settings and set sizes (n).

Model	Spherical			Gaussian		
	$n=30$	$n=50$	$n \in [20 \dots 100]$	$n=30$	$n=50$	$n \in [20 \dots 100]$
Set2Graph	0.780	0.686	0.535	0.707	0.661	0.552
Set Transformer	0.773	0.752	0.703	0.741	0.727	0.686
Slot Attention	0.668	0.629	0.495	0.662	0.665	0.620
Ours	0.892	0.868	0.823	0.851	0.831	0.809

285 setting, we additionally normalize each point to lie on the unit sphere. Following Serviansky et al.
 286 [2], we use $n=30$, $n=50$ and $n \in [20 \dots 100]$, where in the latter case the input set size varies between
 287 20 and 100. Table 2 shows our results. Our method consistently outperforms all the baselines by a
 288 considerable margin. In contrast to Set2Graph, our model does not suffer from a drastic performance
 289 decline when increasing the input set size from 30 to 50. Furthermore, based on the results in the
 290 Gaussian setting, we also observe that all the incidence-based approaches handle varying input sizes
 291 much better than the adjacency-based approach.

292 **Delaunay triangulation.** A Delaunay triangulation of a finite set of 2D points is a set of triangles
 293 for which the circumcircles of all triangles have no point lying inside. When there exists more than
 294 one such set, Delaunay triangulation aims to maximize the minimum angle of all triangles. We
 295 consider the same learning task and setup as Serviansky et al. [2], who frame Delaunay triangulation
 296 as a mapping from a set of 2D points to the set of Delaunay edges, represented by the adjacency
 297 matrix. Since this is essentially a set-to-graph problem instead of a set-to-hypergraph one, we
 298 adapt our method for efficiency reasons, as we describe in Appendix D. We generate the input
 299 sets of size n , by sampling 2-dimensional vectors uniformly from the unit square, with $n=50$ or
 300 $n \in [20 \dots 80]$. In Table 3, we report the results for Set2Graph [2] and our adapted method. Since the
 301 other baselines were not competitive in convex hull finding, we do not repeat them here. Our method
 302 again outperforms Set2Graph on all metrics.

303 **Summary.** By predicting the positive edges only, we can significantly reduce the amount of required
 304 memory for set-to-hypergraph tasks. On three different benchmarks, we observe performance
 305 improvements when using this incidence-based approach, compared to the adjacency-based baseline.
 306 Interestingly, our method does *not* see a large discrepancy in performance between different input set
 307 sizes, both in convex hull finding and Delaunay triangulation. We attribute this to the recurrence of
 308 our iterative refinement scheme, which we look into next.

309 5.2 Ablations

310 **Effects of increasing (time) complexity.** The *intrinsic* complexity of finding a convex hull for
 311 a d -dimensional set of n points is in $\mathcal{O}(n \log(n) + n^{\lfloor \frac{d}{2} \rfloor})$ [11]. This scaling behavior offers an

Table 3: Delaunay triangulation results for different set sizes (n). Our method outperforms Set2Graph on all metrics.

Model	$n=50$				$n \in [20 \dots 80]$			
	Acc	Pre	Rec	F1	Acc	Pre	Rec	F1
Set2Graph	0.984	0.927	0.926	0.926	0.947	0.736	0.934	0.799
Ours	0.989	0.953	0.946	0.950	0.987	0.945	0.942	0.943

312 interesting opportunity to study the effects of increasing (time) complexity on model performance.
 313 The time complexity implies that *any* algorithm for convex hull finding scales super-linearly with
 314 the input set size. Since our learned model is not considered an algorithm that (exactly) solves the
 315 convex hull problem, the implications become less clear. In order to assess the relevancy of the
 316 problem’s complexity for our approach, we examine the relation between the number of refining steps
 317 and increases in the intrinsic resource requirement. The following experiments are all performed
 318 with standard backprop, in order to not introduce additional hyperparameters that may affect the
 319 conclusions.

320 First, we examine the performance of our approach with 3 iter-
 321 ations, trained on increasing set sizes $n \in [10 \dots 50]$. In Figure 2
 322 we observe a monotone drop in performance when training
 323 with the same number of iterations. The negative correlation
 324 between the set size and the performance confirms a relation-
 325 ship between the computational complexity and the difficulty
 326 of the learning task. Next, we examine the performance for
 327 varying numbers of iterations and set sizes. We refer to the
 328 setting, where the number of iterations is 3 and set size $n=10$,
 329 as the base case. All other set sizes and number of iterations
 330 are chosen such that the performance matches the base case as
 331 closely as possible. In Figure 2, we observe that the required
 332 number of iterations increases with the input set size, further
 333 highlighting that an increase in the number of iterations actually
 334 suffices in counteracting the performance decrease. Further-
 335 more, we observe that the number of refinement steps scales sub-linearly with the set size, different
 336 from what we would expect based on the complexity of the problem. We speculate this is due to the
 337 parallelization of our edge finding process, differing from incremental approaches that produce one
 338 edge at a time.

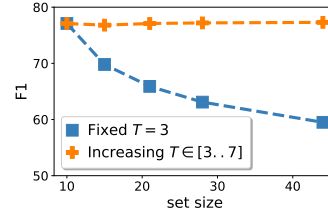


Figure 2: Increasing complexity. Increasing the iterations counteracts the performance decline from larger set sizes.

339 **Efficiency of backprop with skips.** To assess the
 340 efficiency of backprop with skips, we compare to
 341 truncated backpropagation through time (TBPTT)
 342 [26]. We consider two variants of our training algo-
 343 rithm: 1. Skipping iterations at fixed time steps and 2.
 344 Skipping randomly sampled time steps. In both the
 345 fixed and random skips versions, we skip half of the
 346 total iterations. We train all models on convex hull
 347 finding in 3-dimensions for 30 spherically distributed
 348 points. In addition, we include baselines trained with
 349 standard backprop that contingently inform us about
 350 performance degradation incurred by our method or
 351 TBPTT. Standard backprop increases the memory
 352 footprint linearly with the number of iterations T ,
 353 inevitably exceeding the available memory at some
 354 threshold. Hence, we deliberately choose a small set
 355 size in order to accommodate training with backprop
 356 for $T \in \{4, 16, 32\}$ number of iterations. We illustrate
 357 the differences between standard backprop, TBPTT and our backprop with skips in Figure 5 in the
 358 Appendix. The results in Figure 3 demonstrate that skipping half of the iterations in the backward-pass
 359 significantly decreases the training time without hurting predictive performance. When the memory
 360 budget is constricted to 4 iterations in the backward-pass, both TBPTT and backprop with skips
 361 outperform standard backprop considerably. We provide a detailed discussion of the computational
 362 complexity of our framework in Appendix B.

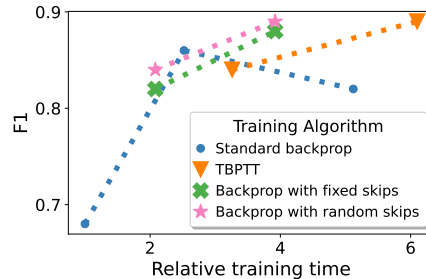


Figure 3: Training time of backprop with skips. Relative training time and performance for different $T \in \{4, 16, 32\}$. All runs require the same memory, except standard backprop $T \in \{16, 32\}$, which require more.

363 **Recurrent vs. stacked.** Recurrence plays a crucial role in enabling more computation without an
 364 increase in the number of parameters. By training the recurrent model using backprop with skips, we
 365 can further reduce the memory cost during training to a constant amount. Since our proposed training
 366 algorithm from Section 3 encourages iterative refinement akin to gradient descent, it is natural to
 367 believe that the weight-tying aspect of recurrence is a good inductive bias for modeling this. A reason
 368 for thinking so is that the “gradient” should be the same for the same I , no matter at which iteration

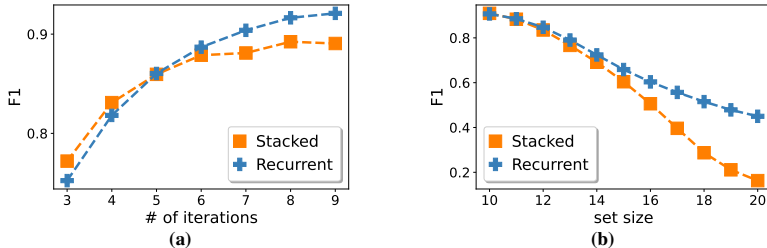


Figure 4: Recurrent vs. stacked. (a) Performance for different numbers of iterations. (b) Extrapolation performance on $n \in [11 \dots 20]$ for models trained with set size $n=10$. We stop training the recurrent model early, to match the validation performance of the stacked on $n=10$. The recurrent model derives greater benefits from adding iterations and generalizes better.

369 it is computed. Hence, we compare the recurrent model against a non-weight-tied (stacked) version
 370 that applies different parameters at each iteration. First, we compare the models trained for 3 to 9
 371 refinement steps. In Figure 4a, we show that both cases benefit from increasing the iterations. Adding
 372 more iterations beyond 6 only slightly improves the performance of the stacked model, while the
 373 recurrent version still benefits, leading to an absolute difference of 0.03 in F1 score for 9 iterations.
 374 Next, we train both versions with 15 iterations until they achieve a similar validation performance, by
 375 stopping training early on the recurrent model. The results in Figure 4b show that the recurrent variant
 376 performs better when tested at larger set sizes than trained, indicating an improved generalization
 377 ability.

378 6 Related Work

379 Set2Graph [2] is a family of maximally expressive permutation equivariant neural networks that
 380 map from an input set to (hyper)graphs. They show that their method outperforms many popular
 381 alternatives, including Siamese networks [27] and graph neural networks [28] applied to a k -NN
 382 induced graph. [2] extend the general idea, of applying a scalar-valued adjacency indicator function
 383 on all pairs of nodes [29], to the l -edge case (edges that connect l nodes). In Set2Graph, for each l the
 384 adjacency structure is modeled by an l -tensor, requiring memory in $\mathcal{O}(n^l)$. This becomes intractable
 385 already for small l and moderate set sizes. By pruning the negative edges, our approach scales in
 386 $\mathcal{O}(nk)$, making it applicable even when $l=n$. Recent works on set prediction map a learned initial
 387 set [23, 30] or a randomly initialized set [20, 31, 32] to the target space. Out of these, the closest one
 388 to our hypergraph refining approach is Slot Attention [20], which recurrently applies the Sinkhorn
 389 operator [33] in order to associate each element in the input set with a single slot (hyperedge). None
 390 of the prior works on set prediction consider the set-to-hypergraph task, but some can be naturally
 391 extended by mapping the input set to the set of positive edges, an approach similar to ours.

392 7 Conclusion and future work

393 By representing and supervising the set of positive edges only, we substantially improve the asymp-
 394 totic scaling and enable learning tasks with higher-order edges. On common benchmarks, we have
 395 demonstrated that our method outperforms previous works while offering a more favorable asymp-
 396 totic scaling behavior. In further evaluations, we have highlighted the importance of recurrence for
 397 addressing the intrinsic complexity of problems. We identify the Hungarian matching [13] as the
 398 main computational bottleneck during training. Replacing the Hungarian matched loss with a faster
 399 alternative, like a learned energy function [32], would greatly speed up training for tasks with a large
 400 maximum number of edges. Our empirical analysis is limited to datasets with low dimensional inputs.
 401 Learning on higher dimensional input data might require extensions to the model that can make larger
 402 changes to the latent hypergraph than is feasible with small iterative refinement steps. The idea here
 403 is similar to the observation from Jastrzebski et al. [17] for ResNets [22] that also encourage iterative
 404 refinement: earlier residual blocks apply large changes to the intermediate features while later layers
 405 perform (small) iterative refinements.

References

- 406
- 407 [1] Jonathan Shlomi, Peter Battaglia, and Jean-Roch Vlimant. Graph neural networks in particle
408 physics. *Machine Learning: Science and Technology*, 2020. 1, 14
- 409 [2] Hadar Serviansky, Nimrod Segol, Jonathan Shlomi, Kyle Cranmer, Eilam Gross, Haggai Maron,
410 and Yaron Lipman. Set2graph: Learning graphs from sets. In *Advances in Neural Information
411 Processing Systems*, 2020. 1, 2, 3, 6, 7, 9, 14, 15, 16
- 412 [3] Ivan Brugere, Brian Gallagher, and Tanya Y Berger-Wolf. Network structure inference, a survey:
413 Motivations, methods, and applications. *ACM Computing Surveys*, 51(2):1–39, 2018. 1
- 414 [4] Federico Battiston, Giulia Cencetti, Iacopo Iacopini, Vito Latora, Maxime Lucas, Alice Patania,
415 Jean-Gabriel Young, and Giovanni Petri. Networks beyond pairwise interactions: structure and
416 dynamics. *Physics Reports*, 2020. 1
- 417 [5] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural
418 Information Processing Systems*, 2015. 1, 6
- 419 [6] S. Feng, E. Heath, and B. Jefferson. Hypergraph models of biological networks to identify
420 genes critical to pathogenic viral response. *"BMC Bioinformatics"*, 22:287, 2021. 1
- 421 [7] O. Frank and D. Strauss. Markov graphs. *J. Am. Stat. Assoc.*, 81:832–842, 1986. 1
- 422 [8] Yibing Zhan, Zhi Chen, Jun Yu, BaoSheng Yu, Dacheng Tao, and Yong Luo. Hyper-relationship
423 learning network for scene graph generation. *arXiv preprint arXiv:2202.07271*, 2022. 1
- 424 [9] Alexander Trott, Caiming Xiong, and Richard Socher. Interpretable counting for visual question
425 answering. In *International Conference on Learning Representations*, 2018. 1
- 426 [10] Linyuan Lü and Tao Zhou. Link prediction in complex networks: A survey. *Physica A:
427 statistical mechanics and its applications*, 390(6):1150–1170, 2011. 1
- 428 [11] Bernard Chazelle. An optimal convex hull algorithm in any fixed dimension. *Discrete &
429 Computational Geometry*, 1993. 1, 6, 7
- 430 [12] Xavier Ouvrard, Jean-Marie Le Goff, and Stéphane Marchand-Maillet. Adjacency and tensor
431 representation in general hypergraphs part 1: e-adjacency tensor uniformisation using
432 homogeneous polynomials. *arXiv preprint arXiv:1712.08189*, 2017. 3
- 433 [13] Harold W Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics
434 Quarterly*, 2(1-2):83–97, 1955. 3, 9
- 435 [14] Roy Jonker and Anton Volgenant. A shortest augmenting path algorithm for dense and sparse
436 linear assignment problems. *Computing*, 1987. 3, 13
- 437 [15] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
438 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
439 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
440 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
441 high-performance deep learning library. In *Advances in Neural Information Processing Systems*,
442 pages 8024–8035, 2019. 4, 5
- 443 [16] Zachary C Lipton, John Berkowitz, and Charles Elkan. A critical review of recurrent neural
444 networks for sequence learning. *arXiv preprint arXiv:1506.00019*, 2015. 4
- 445 [17] Stanisław Jastrzebski, Devansh Arpit, Nicolas Ballas, Vikas Verma, Tong Che, and Yoshua
446 Bengio. Residual connections encourage iterative inference. In *International Conference on
447 Learning Representations*, 2018. 4, 6, 9
- 448 [18] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,
449 and Alexander J Smola. Deep sets. In *Advances in Neural Information Processing Systems*,
450 2017. 5
- 451 [19] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
452 message passing for quantum chemistry. In *International Conference on Machine Learning*,
453 2017. 5
- 454 [20] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg
455 Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with
456 slot attention. In *Advances in Neural Information Processing Systems*, 2020. 6, 9, 14

- 457 [21] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*
458 *arXiv:1607.06450*, 2016. 6
- 459 [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
460 recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016. 6, 9
- 461 [23] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh.
462 Set transformer: A framework for attention-based permutation-invariant neural networks. In
463 *International Conference on Machine Learning*, 2019. 6, 9
- 464 [24] Jonathan Shlomi, Sanmay Ganguly, Eilam Gross, Kyle Cranmer, Yaron Lipman, Hadar Servian-
465 sky, Haggai Maron, and Nimrod Segol. Secondary vertex finding in jets with neural networks.
466 *arXiv preprint arXiv:2008.02831*, 2020. 6, 14
- 467 [25] Franco P Preparata and Michael I Shamos. *Computational geometry: an introduction*. Springer
468 Science & Business Media, 2012. 6
- 469 [26] Ronald J Williams and Jing Peng. An efficient gradient-based algorithm for on-line training of
470 recurrent network trajectories. *Neural computation*, 1990. 8, 16
- 471 [27] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional
472 neural networks. In *Conference on Computer Vision and Pattern Recognition*, 2015. 9
- 473 [28] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,
474 Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural
475 networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019. 9
- 476 [29] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint*
477 *arXiv:1611.07308*, 2016. 9
- 478 [30] Yan Zhang, Jonathon Hare, and Adam Prügel-Bennett. Deep set prediction networks. In
479 *Advances in Neural Information Processing Systems*, 2019. 9
- 480 [31] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
481 Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on*
482 *Computer Vision*, 2020. 9
- 483 [32] David W Zhang, Gertjan J Burghouts, and Cees G M Snoek. Set prediction without impos-
484 ing structure as conditional density estimation. In *International Conference on Learning*
485 *Representations*, 2021. 9
- 486 [33] Ryan Prescott Adams and Richard S Zemel. Ranking via sinkhorn propagation. *arXiv preprint*
487 *arXiv:1106.1925*, 2011. 9
- 488 [34] David F Crouse. On implementing 2d rectangular assignment algorithms. *IEEE Transactions*
489 *on Aerospace and Electronic Systems*, 2016. 13
- 490 [35] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David
491 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.
492 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew
493 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.
494 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.
495 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul
496 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific
497 Computing in Python. *Nature Methods*, 2020. 13
- 498 [36] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
499 differential equations. *Advances in neural information processing systems*, 31, 2018. 13
- 500 [37] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
501 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
502 *preprint arXiv:2011.13456*, 2020. 13
- 503 [38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
504 *arXiv:1412.6980*, 2014. 14
- 505 [39] Stefano Rebay. Efficient unstructured mesh generation by means of delaunay triangulation and
506 bowyer-watson algorithm. *Journal of Computational Physics*, 106(1):125–138, 1993. 15

507 A Proof: Supervising positive edges only suffices

508 **Proposition 1** (Supervising only existing edges). *Let $\mathbf{J} \in [\epsilon, 1]^{(2^n-1) \times n}$ be a matrix with at most*
 509 *m rows for which $\exists j : \mathbf{J}_{ij} > \epsilon$, with a small $\epsilon > 0$. Similarly, let $\mathbf{J}^* \in \{0, 1\}^{(2^n-1) \times n}$ be a matrix*
 510 *with at most m rows for which $\exists j : \mathbf{J}_{ij} = 1$. Let $\text{prune}(\cdot)$ denote the function that maps from a*
 511 *$(2^n - 1) \times n$ matrix to a $k \times n$ matrix, by removing $(2^n - 1) - k$ rows where all values are $\leq \epsilon$.*
 512 *Then, for a constant $c = (2^n - 1 - k)n \cdot H(\epsilon, 0)$ and all such \mathbf{J} and \mathbf{J}^* :*

$$\mathcal{L}(\mathbf{J}, \mathbf{J}^*) = \mathcal{L}(\text{prune}(\mathbf{J}), \text{prune}(\mathbf{J}^*)) + c \quad (2)$$

513 *Proof.* We shorten the notation with $\mathbf{I} = \text{prune}(\mathbf{J})$ and $\mathbf{I}^* = \text{prune}(\mathbf{J}^*)$, making the relation to
 514 the incidence matrix \mathbf{I} defined in Section 2 explicit. Since \mathcal{L} is invariant to permutations over
 515 the rows of its input matrices, we can assume w.l.o.g. that the not-pruned rows are the first k
 516 rows, $\mathbf{J}_{:k} = \mathbf{I}$ and $\mathbf{J}_{:k}^* = \mathbf{I}^*$. For improved readability, let $\hat{H}(\mathbf{J}_{\pi(i)}, \mathbf{J}_i^*) = \sum_j H(\mathbf{J}_{\pi(i),j}, \mathbf{J}_{i,j}^*)$
 517 denote the element-wise binary cross-entropy, thus the loss in Equation 1 can be rewritten as
 518 $\mathcal{L}(\mathbf{J}, \mathbf{J}^*) = \min_{\pi \in \Pi} \sum_i \hat{H}(\mathbf{J}_{\pi(i)}, \mathbf{J}_i^*)$.

519 First, we show that there exists an optimal assignment between \mathbf{J}, \mathbf{J}^* that assigns the first
 520 k rows equally to an optimal assignment between \mathbf{I}, \mathbf{I}^* . More formally, for an optimal as-
 521 signment $\pi_I \in \arg \min_{\pi \in \Pi} \sum_i \hat{H}(\mathbf{I}_{\pi(i)}, \mathbf{I}_i^*)$ we show that there exists an optimal assignment
 522 $\pi_J \in \arg \min_{\pi \in \Pi} \sum_i \hat{H}(\mathbf{J}_{\pi(i)}, \mathbf{J}_i^*)$ such that $\forall 1 \leq i \leq k : \pi_J(i) = \pi_I(i)$. If $\pi_J(i) \leq k$ for all $1 \leq i \leq k$
 523 then the assignment for the first k rows is also optimal for \mathbf{I}, \mathbf{I}^* . So we only need to show that there
 524 exists a π_J such that $\pi_J(i) \leq k$ for all $1 \leq i \leq k$. Let π_J be an optimal assignment that maps an $i < k$ to
 525 $\pi_J > k$. Since π_J is a bijection, there also exists a $j < k$ that $\pi_J^{-1}(j) > k$ assigns to. The corresponding
 526 loss terms are lower bounded as follows:

$$\hat{H}(\mathbf{J}_i, \mathbf{J}_{\pi_J(i)}^*) + \hat{H}(\mathbf{J}_{\pi_J^{-1}(j)}, \mathbf{J}_j^*) \quad (7)$$

$$= \hat{H}(\mathbf{J}_i, \mathbf{0}) + \hat{H}(\epsilon, \mathbf{J}_j^*) \quad (8)$$

$$= - \sum_{l=1}^n \log(1 - \mathbf{J}_{i,l}) + \mathbf{J}_{j,l}^* \log(\epsilon) + (1 - \mathbf{J}_{j,l}^*) \log(1 - \epsilon) \quad (9)$$

$$\geq - \sum_{l=1}^n (1 - \mathbf{J}_{j,l}^*) \log(1 - \mathbf{J}_{i,l}) + \mathbf{J}_{j,l}^* \log(\epsilon) + (1 - \mathbf{J}_{j,l}^*) \log(1 - \epsilon) \quad (10)$$

$$\geq - \sum_{l=1}^n (1 - \mathbf{J}_{j,l}^*) \log(1 - \mathbf{J}_{i,l}) + \mathbf{J}_{j,l}^* \log(\mathbf{J}_{i,l}) + (1 - \mathbf{J}_{j,l}^*) \log(1 - \epsilon) \quad (11)$$

$$= \hat{H}(\mathbf{J}_i, \mathbf{J}_j^*) - \sum_{l=1}^n (1 - \mathbf{J}_{j,l}^*) \log(1 - \epsilon) \quad (12)$$

$$\geq \hat{H}(\mathbf{J}_i, \mathbf{J}_j^*) - \sum_{l=1}^n \log(1 - \epsilon) \quad (13)$$

$$= \hat{H}(\mathbf{J}_i, \mathbf{J}_j^*) + \hat{H}(\epsilon, \mathbf{0}) \quad (14)$$

527 Equality of Equation 8 holds since all rows with index $> k$ are ϵ -vectors in \mathbf{J} and zero-vectors
 528 in \mathbf{J}^* . The inequality in Equation 11 holds since all values in \mathbf{J} are lower bounded by ϵ . Thus,
 529 we have shown that either there exists no optimal assignment π_J that maps from a value $\leq k$ to
 530 a value $> k$ (which is the case when $\hat{H}(\mathbf{J}_i, \mathbf{J}_{\pi_J(i)}^*) + \hat{H}(\mathbf{J}_{\pi_J^{-1}(j)}, \mathbf{J}_j^*) > \hat{H}(\mathbf{J}_i, \mathbf{J}_j^*) + \hat{H}(\epsilon, \mathbf{0})$) or that
 531 there exists an equally good assignment that does not mix between the rows below and above k .
 532 Since the pruned rows are all identical, any assignment between these result in the same value
 533 $(2^n - 1 - k) \hat{H}(\epsilon, \mathbf{0}) = (2^n - 1 - k)n \cdot H(\epsilon, 0)$ that only depends on the number of pruned rows $2^n - 1 - k$
 534 and number of columns n . \square

535 B Computational complexity

536 The space complexity of the hypergraph representation presented in Section 2 is in $\mathcal{O}(nm)$, offering
 537 an efficient representation for hypergraphs when the maximal number of edges m is low, relative to

538 the number of all possible edges $m \ll 2^n$. Problems that involve edges connecting many vertices
 539 benefit from this choice of representation, as the memory requirement is independent of the maximal
 540 connectivity of an edge. This differs from the adjacency-based approaches that not only depend on
 541 the maximum number of nodes an edge connects, but scale exponentially with it. In practice, this
 542 improvement from $\mathcal{O}(2^n)$ to $\mathcal{O}(mn)$ is important even for moderate set sizes because the amount of
 543 required memory determines whether it is possible to use efficient hardware like GPUs. We showcase
 544 this in [Section D.4](#).

545 Backprop with skips, introduced [Section 3](#), further scales the memory requirement by a factor of B
 546 that is the number of iterations to backprop through in a single gradient update step. Notably, this
 547 scales constantly in the number of gradient update steps N and iterations skipped during backprop
 548 $\sum_i S_i$. Hence, we can increase the number of recurrent steps to adapt the model to the problem
 549 complexity (which is important, as we show in [Section 5.2](#)), at a constant memory footprint.

550 To compute the loss in [Equation 1](#), we apply a modified Jonker-Volgenant algorithm [[14](#), [34](#), [35](#)] that
 551 finds the minimum assignment between the rows of the predicted and the ground truth incidence
 552 matrices in $\mathcal{O}(m^3)$. In practice, this can be the main bottleneck of the proposed method when the
 553 number of edges becomes large. For problems with $m \ll n$, the runtime complexity is especially
 554 efficient since it is independent of the number of nodes.

555 **Comparison to Set2Graph.** When we compare Set2Graph and our method for hyperedges that
 556 connect two nodes (i.e., a normal graph) then Set2Graph has an efficiency advantage. For example,
 557 the authors report for Set2Graph a training time of 1 hour (on a Tesla V100 GPU) for Delaunay
 558 triangulations, while our method takes up to 9 hours (on a GTX 1080 Ti GPU). Nevertheless, the
 559 major bottleneck of the Set2Graph method is its memory and time complexity, which in practice
 560 can lead to intractable memory requirements even for small problems. For example, the convex
 561 hull problem in 10-dimensional space over 13 points requires more than 500GB ($\approx 4 * 13^10$) just
 562 for storing the adjacency tensor – that is without even considering the intermediate neural network
 563 activations. Even when we keep the space 3-dimensional as in the experiments reported by Set2Graph,
 564 it already struggles with memory requirements as the authors point out themselves. They circumvent
 565 this issue by considering an easier local version of the problem and restrict the adjacent nodes to
 566 the k-Nearest-Neighbors with $k=10$. In conclusion, when the hypergraph has relatively few edges,
 567 then our framework offers a much better scaling than Set2Graph. In practice, this does not merely
 568 translate to faster runtime but turns tasks that were previously intractable into feasible tasks.

569 C Backprop with skips

570 Our backprop with skips training consists of two aspects:

- 571 • Truncation of the backprop through time, and
- 572 • Skipping the gradient calculations for some intermediate $f(\mathbf{X}, \mathcal{H})$ as specified by the S .

573 To understand what the effect of truncation is, we first consider the case of standard backprop. When
 574 training a residual neural network with standard backprop it is possible to expand the loss function
 575 around \mathcal{H}^t from previous iterations $t < T$ as $\mathcal{L}(\mathcal{H}^T) = \mathcal{L}(\mathcal{H}^t) + \sum_{i=t}^{T-1} f(\mathbf{X}, \mathcal{H}^i) \frac{\partial \mathcal{L}(\mathcal{H}^i)}{\partial \mathcal{H}^i} +$
 576 $\mathcal{O}(f^2(\mathbf{X}, \mathcal{H}^i))$ (see equation 4 in [[13](#)]). Jastrzebski et al. [[13](#)] point out that the sum terms’ gradients
 577 point into the same half space as that of $\frac{\partial \mathcal{L}(\mathcal{H}^i)}{\partial \mathcal{H}^i}$, implying that $\sum_{i=t}^{T-1} \mathcal{L}(\mathcal{H}^i)$ is also a descent
 578 direction for $\mathcal{L}(\mathcal{H}_T)$. Thus, truncation can be interpreted as removing $\mathcal{L}(\mathcal{H}_i)$ terms from the loss
 579 function for earlier steps i . We remedy this, by explicitly adding (back) the $\mathcal{L}(\mathcal{H}_i)$ terms to the
 580 training objective. The second aspect of backprop with skips involves skipping gradients of some
 581 intermediate steps $f(\mathbf{X}, \mathcal{H}^i)$ entirely. Given a specific data point \mathbf{X} , we view f as approximating
 582 the gradient $\frac{\partial \mathcal{H}^t}{\partial t}$ where we denote the time step t as the argument of H instead of as the subscript.
 583 From this viewpoint, we are learning a vector field (analog to neural ODE [[36](#)] or diffusion models
 584 [[37](#)]) over the space of \mathcal{H} which we train by randomly sampling different values for the \mathcal{H}^t ’s.

585 D Experimental details

586 In this section, we provide further details on the experimental setup and additional results.

587 D.1 Particle partitioning

588 The problem considers the case where particles are collided at high energy, resulting in multiple
 589 particles shooting out from the collision. Each example in the dataset consists of the input set,
 590 which corresponds to the measured outgoing particles, and the ground truth partition of the input
 591 set. Each element in the partition is a subset of the input set and corresponds to some intermediate
 592 particle that was not measured, because it decayed into multiple particles before it could reach the
 593 sensors. The learning task consists of inferring which elements in the input set originated from the
 594 same intermediate particle. Note that the particle partitioning task bears resemblance to the classical
 595 clustering setting. It can be understood as a meta-learning clustering task, where both the number of
 596 clusters and the similarity function depend on the context that is given by \mathbf{X} . That is why clustering
 597 algorithms such as k -means cannot be directly applied to this task. For more information on how this
 598 task fits into the area of particle physics more broadly, we refer to Shlomi et al. [1].

599 **Dataset.** We use the publicly available dataset of 0.9M data-sample with the default
 600 train/validation/test split [2, 24]. The input sets consist of 2 to 14 particles, with each particle
 601 represented by 10 features. The target partitioning indicate the common progenitors and restrict the
 602 valid incidence matrices to those with a single incident edge per node.

603 **Setup.** While Set2Graph is only one instance of an adjacency-based approach, [2] show that it
 604 outperforms many popular alternatives: Siamese networks, graph neural networks and a non-learnable
 605 geometric-based baseline. All adjacency-based approaches incur a prohibitively large memory cost
 606 when predicting edges with high connectivity. In the case of particle partitioning, Set2Graph resorts
 607 to only predicting edges with at most 2 connecting nodes, followed by an additional heuristic to infer
 608 the partitions [2]. In contrast to that, all the incidence-based approaches do not require the additional
 609 post-processing step at the end.

610 We simplify the hyperparameter search by choosing the same number of hidden dimensions d for the
 611 latent vector representations of both the nodes $d_{\mathcal{V}}$ and the edges $d_{\mathcal{E}}$. In all runs dedicated to searching
 612 d , we set the number of total iterations $T=3$ and backpropagate through all iterations. We start with
 613 $d=32$ and double it, until an increase yields no substantial performance gains on the validation set,
 614 resulting in $d=128$. In our reported runs, we use $T=16$ total iterations, $B=4$ backprop iterations,
 615 $N=2$ gradient updates per mini-batch, and a maximum of 10 edges.

616 We apply the same $d=128$ to both the Slot Attention and Set Transformer baselines. Similar to the
 617 original version [20], we train Slot Attention with 3 iterations. Attempts with more than 3 iterations
 618 resulted in frequent divergences in the training losses. We attribute this behavior to the recurrent
 619 sinkhorn operation, that acts as a contraction map, forcing all slots to the same vector in the limit.

620 We train all models using the Adam optimizer [38] with a learning rate of 0.0003 for 400 epochs and
 621 retain the parameters corresponding to the lowest validation loss. All models additionally minimize a
 622 soft F1 score [2]. Since each particle can only be part of a single partition, we choose the one with
 623 the highest incidence probability at test time. Our model has 268162 trainable parameters, similar
 624 to 251906 for the Slot Attention baseline, but less than 517250 for Set Transformer and 461289 for
 625 Set2Graph [2]. The total training time is less than 12 hours on a single GTX 1080 Ti and 10 CPU
 626 cores.

627 The maximum number of edges is set to $m = 10$.

628 **Further results.** For completeness, we also report the results for the rand index (RI) in Table 4.

629 D.2 Convex hull finding

630 On convex hull finding in 3D, we compare our method to the same baselines as on the particle
 631 partitioning task.

632 **Setup.** Set2Graph learns to map the set of 3D points to the 3rd order adjacency tensor. Since storing
 633 this tensor in memory is not feasible, they instead concentrate on a local version of the problem,
 634 which only considers the k -nearest neighbors for each point [2]. We train our method with $T_{\text{total}}=48$,
 635 $T_{\text{BPTT}}=4$, $N_{\text{BPTT}}=6$ and set k equal to the highest number of triangles in the training data. At test
 636 time, a prediction admits an edge e_i if its existence indicator $\sigma_i > 0.5$. Each edge is incident to the
 637 three nodes with the highest incidence probability. We apply the same hyperparameters, architectures

Table 4: Additional particle partitioning results. On three jet types performance measured as rand index (RI). Our method outperforms the baselines on bottom and charm jets, while being competitive on light jets.

Model	bottom jets	charm jets	light jets
	RI	RI	RI
Set2Graph	0.736 \pm 0.004	0.727 \pm 0.003	0.970 \pm 0.001
Set Transformer	0.734 \pm 0.004	0.734 \pm 0.004	0.967 \pm 0.002
Slot Attention	0.703 \pm 0.013	0.714 \pm 0.009	0.958 \pm 0.003
Ours	0.781 \pm 0.002	0.751 \pm 0.001	0.969 \pm 0.001

638 and optimizer as in the particle partitioning experiment, except for: $T=48$, $B=4$, $N=6$. Since we
 639 do not change the model, the number of parameters remains at 268162 for our model. This notably
 640 differs to Set2Graph, which reports an increased parameter count of 1186689 [2]. We train our
 641 method until we observe no improvements on the F1 validation performance for 20 epochs, with a
 642 maximum of 1000 epochs. The set-to-set baselines are trained for 4000 epochs, and we retain the
 643 parameters resulting in the highest f1 score on the validation set. The training time is similar to our
 644 proposed method. The total training time is between 14 and 50 hours on a single GTX 1080 Ti and
 645 10 CPU cores.

646 We set the maximum number of edges m equal to the maximum number of triangles of any example
 647 in the training data. For the spherically distributed point sets, m is a constant that is $m = (n - 4)2 + 4$
 648 for $n \geq 4$. This can be seen from the fact that all points lie on the convex hull in this case. Note that
 649 the challenge lies not with finding which points lie on the convex hull, but in finding all the facets
 650 that constitute the convex hull. For the Gaussian distributed point sets, m varies between different
 651 samples. For $n = 30$ most examples have < 40 edges, for $n = 50$ most examples have < 50 edges,
 652 and for $n = 100$ most examples have < 60 edges.

653 D.3 Delaunay triangulation

654 The problem of Delaunay triangulation is, similar to convex hull finding a well-studied problem in
 655 computational geometry and has exact solutions in $\mathcal{O}(n \log(n))$ [39]. We consider the same learning
 656 task as Serviansky et al. [2], who frame Delaunay triangulation as mapping from a set of 2D points to
 657 the set of Delaunay edges, represented by the adjacency matrix. Note that this differs from finding
 658 the set of triangles, as an edge no longer remembers which triangles it is part of. Thus, this reduces to
 659 a set-to-graph task, instead of a set-to-hypergraph task.

660 **Model adaptation.** The goal in this task is to predict the adjacency matrix of an ordinary graph – a
 661 graph consisting of edges that connect two nodes – where the number of edges are greater than the
 662 number of nodes. One could recover the adjacency matrix based on the matrix product of $I^T I$, by
 663 clipping all values above 1 back to 1 and setting the diagonal to 0. This approach is inefficient, since
 664 in this case the incidence matrix is actually larger than the adjacency matrix. Instead of applying our
 665 method directly, we consider a simple adaptation of our approach to the graph setting. We replace the
 666 initial set of edges with the (smaller) set of nodes and apply the same node refinements on both sets.
 667 This change results in $\mathcal{E} = \mathcal{V}$ for the prediction and effectively reduces the incidence matrix to an
 668 adjacency matrix, since it is computed based on all pairwise combinations of \mathcal{E} and \mathcal{V} . We further
 669 replace the concatenation for the MLP modelling the incidence probability with a sum, to ensure that
 670 the predicted adjacency matrix is symmetric and represents an undirected graph. Two of the main
 671 design choices of our approach remain in this adaptation: Iterative refining of the complete graph
 672 with a recurrent neural network and BPTT with gradient skips. We train our model with $T=32$, $B=4$
 673 and $N=4$. At test-time, an edge between two nodes exists if the adjacency value is greater than 0.5.

674 **Setup.** We increase the latent dimensions to $d=256$, resulting in 595201 trainable parameters.
 675 This notably differs to Set2Graph, which increases the parameter count to 5918742 [2], an order of
 676 magnitude larger. The total training time is less than 9 hours on a single GTX 1080 Ti and 10 CPU
 677 cores.

678 D.4 Learning higher-order edges

679 The particle partitioning experiment exemplifies a case where a single edge can connect up to 14
 680 vertices. Set2Graph [2] demonstrates that in this specific case it is possible to approximate the
 681 hypergraph with a graph. They leverage the fact that any vertex is incident to exactly one edge and
 682 apply a post-processing step that constructs the edges from noisy cliques. Instead, we consider a task
 683 for which no straightforward graph based approximation exists. Specifically, we consider convex
 684 hull finding in 10-dimensional space for 13 standard normal distributed points. We train with $T=32$,
 685 $N=4$ and $B=4$. The test performance reaches an F1 score of 0.75, clearly demonstrating that the
 686 model managed to learn. This result demonstrates the improved scaling behavior can be leveraged
 687 for tasks that are computationally out of reach for adjacency-based approaches.

688 We demonstrated that the improved scaling behavior of our proposed method can be leveraged for
 689 tasks that are computationally out of reach for adjacency based approaches. The number of points and
 690 dimensions were chosen in conjunction, such that the corresponding adjacency tensor would require
 691 more storage than is feasible with current GPUs (available to us). For 13 points in 10 dimensions,
 692 explicitly storing the full adjacency tensor using 32-bit floating-point numbers would already require
 693 more than 500 GB. We intentionally kept the number of points and dimensions low, to highlight
 694 that the asymptotic scaling issue cannot be met by hardware improvements, since small numbers
 695 already pose a problem. Note that Set2Graph already struggles with convex hull finding in 3D, where
 696 the authors report that storing 3-rd order tensors in memory is not feasible. Instead, they consider a
 697 local version of the problem and take the k -Nearest-Neighbors out of the set of points that are part
 698 of the convex hull, with $k = 10$. While we limited our calculation of the storage requirement to
 699 the adjacency tensor itself, a typical implementation of a neural network also requires storing the
 700 intermediate activations, further exacerbating the problem for adjacency based approaches.

701 D.5 Backprop with skips

702 We compare backprop with skips to TBPTT [26] with $B=4$ every 4 iterations, which is the setting
 703 that is most similar to ours with regard to training time. In general, TBPTT allows for overlaps
 704 between subsequent BPTT applications, as we illustrate in Figure 5. We constrict both TBPTT and
 705 backprop with skips to a fixed memory budget, by limiting any backward pass to the most recent
 706 $B=4$ iterations, for $T \in \{16, 32\}$. The standard backprop results serve as a reference point to answer
 707 the question: “What if we apply backprop more frequently, resulting in a better approximation to the
 708 true gradients?”, without necessitating a grid search over all possible hyperparameter combinations
 709 for TBPTT. The results on standard backprop appear to indicate that performance worsens when
 710 increasing the number of iterations from 16 to 32. We observe that applying backprop on many
 711 iterations leads to increasing gradient norms in the course of training, complicating the training
 712 process. The memory limited versions did not exhibit a similar behavior, evident from the improved
 713 performance, when increasing the iterations from 16 to 32.

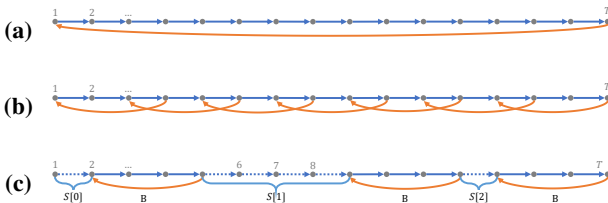


Figure 5: Adding gradient skips to backprop (a) Standard backprop (b) TBPTT, applying backprop on 4 iterations every 2nd iteration (c) Backprop with skips at iterations 1, 6, 7, 8, which effectively reduces the training time, while retaining the same number of refinement steps.