

# DNNs F - Compositional generalisation in group structured data

## 1 Introduction & Hypothesis

My mini-project focuses on compositional generalisation in group structured data. In this report, I analyze how transformers and graph neural networks (GNNs) generalise subgroups. I hypothesise that transformers may struggle because they may have difficulty handling examples when they aren't explicitly trained on a diverse set of rules. In contrast, GNNs might leverage their inherent structure to better capture subgroup properties, a hypothesis I rigorously evaluate through a carefully structured, logically detailed experiment on only symmetric groups as they are compositionally rich.

## 2 Data Construction & Representation Choices

In this exercise we focus mainly on the symmetric group  $S_n$  and the modular arithmetic group  $Z_p^k$  and we will focus on how we represent these in this section. (See [Appendix A.1](#) for full definitions)

### 2.1 Representation

Possible Cayley table representations include adjacency matrices, which encode operations directly as matrix entries (see [Appendix A.2.1](#)); positional encodings, which enhance the representation with positional information for better model interpretability ([A.2.2](#)); direct element pair representation ([A.2.3](#)); and graphical dependencies ([A.2.4](#)).

In terms of the choice of representation, this depends on the models that we use. For the transformer models, we use the direct element representation because it encodes each permutation element explicitly, preserving the inherent order and structure. This explicit encoding aligns well with the transformer's attention mechanism, enabling the model to learn relationships between elements more effectively.

Additionally, for GNNs we use the graphical dependencies. This approach leverages the natural graph structure of the permutation data, allowing the model to capture complex inter-element dependencies through message passing between nodes. The graphical representation is intrinsic to GNN architectures, enhancing their ability to model relational patterns in the data.

We encode each permutation as tokens rather than the actual elements. For example, consider  $[1, 2, 3]$ .

Two encoding schemes are available. In the first, the permutation is mapped to a unique integer  $n \in \mathbb{Z}$ . In the second, it is represented as  $[\pi(1), \pi(2), \pi(3)]$ , with  $\pi : \mathbb{Z} \rightarrow \mathbb{Z}$  reordering the elements.

Both methods encapsulate the entire permutation as tokens, each offering its own advantages. We still shuffle the numbers in the combinations because the transformer can learn the operations using simple addition for both the symmetric group and modular arithmetic groups [trivially] otherwise. <sup>1</sup>

In this report, we employ a single integer encoding for the group  $Z_p^k$  while retaining a list structure for the symmetric groups  $S_n$ . This decision is motivated by the following: for groups with repeated elements, using a list encoding would simply replicate the dataset and make learning trivial; whereas preserving the list structure for  $S_n$  (and its subgroups) maintains the sequential information necessary to capture complex relational patterns, paired with the empirical evidence that these methods train better.

### 2.2 Sampling Strategies

We compare structured subgroup-based sampling with random sampling. Subgroup-based sampling ensures controlled structural diversity, perhaps enhancing compositional generalization, whereas random sampling produces a less structured distribution that may limit interpretability or provide variety.

<sup>1</sup>By representing the permutation as tokens and shuffling the numbers, the transformer should learn the underlying arithmetic operations rather than rely on fixed positions. This forces the model to capture the compositional structure, which generalizes to both symmetric and modular arithmetic groups.

### 3 Generalisation of groups in transformers

This section aims to summarise the results with transformer experiments and hypothesise how well they generalise with subgroups. Here we experiment over the modular arithmetic group first before analysing the symmetric groups. From section E we know that transformer generalise well with more than 50% of training data therefore will do our experiments with smaller training splits.

#### 3.1 Methodology

We kept the model parameters (i.e., model size) constant for consistency while varying the hyperparameters through grid search and early stopping based on inference performance. A learning rate scheduler and validation set convergence were used to mitigate the influence of the initial hyperparameter settings (including optional warm-up LR). The primary hyperparameters searched were weight decay, dropout rate, batch size (and the number of epochs, learning rate). In general we found that dropout controlled the overfitting and we were looking for the highest value such that the training loss was still a decreasing function, this ended up around 0.5 on average, that batch size showed minimal changes as we were working in low data environments and then iterated over values for weight decay using known values from section E as a guide, ending up choosing in the order of  $1e-2$ .

#### 3.2 Experiment 1: $Z_6^2, Z_9^2$

Model	Accuracy (%)	Standard Deviation (%)
Random Sampling	80.44	24.21
Subgroup Sampling	67.96	35.67

Table 1: Accuracy and error margins for  $Z_6^2$ .<sup>2</sup>

In this experiment our method was to take all of the 1D and 2D subgroups.<sup>3</sup> We can immediately see that the model struggles to generalise but this is a very weak argument because we do not have statistical significance, seen with the high variability of the results and we are only testing one configuration and group.<sup>4</sup>

In an attempt to generalise our argument, we can look into a few other groups including  $Z_2^7$  and  $Z_3^3$  and in both of these we notice a consistent lack of generalisation behaviour. This would be described by the extreme where we have a lack of tokens observed leading to the models over-fitting on the training data and significantly struggle. See [Appendix B.1](#) for full experimental details. Through these experiments, we start to believe that transformer struggle to generalise while extrapolating from subgroups to groups compared to random sampling. We can now see what happens when we try the same experiments on symmetric groups instead.

#### 3.3 Experiment 3: Symmetric Groups

Group	Method	Accuracy $\pm$ 1 std (%)
$S_5$	Random Sampling	$100 \pm 0.1$
	Subgroup Sampling	$70.4 \pm 13.4$

Figure 1: Accuracy across methods

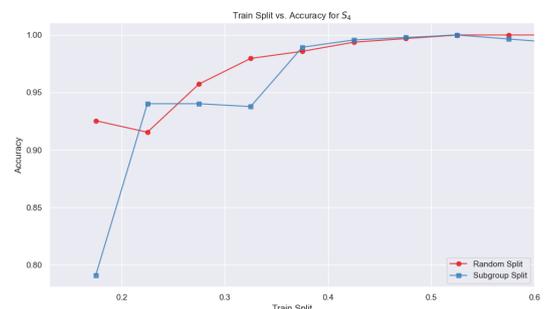


Figure 2:  $S_4$  accuracy vs training split

At first we work with  $S_5$  and we randomly chosen subgroups which are one and two-dimensional to meet a 30% data split. From this slice, we can notice that subgroups appear to hinder training performance within  $S_5$  like above. To do further analysis of this we move into  $S_4$  to try and get more detailed evaluation.

Initially, with  $S_4$  we did experiments with random selections of the subgroups in order to try and notice any patterns in the efficiency of subgroups. In the figure you can see how the training split seems to affect the accuracy across using

<sup>3</sup>The dimensionality of a subgroup defined by the number of basis permutations, these are generally the smaller subgroups

<sup>4</sup>Note: For this we use tiny transformers and use at least 7 repeats

random splits as opposed to random subgroups. From this figure we seem to notice that complete random sampling is slightly more effective than random subgroup selection, which reinforces our initial results from  $S_5$ . However, this graph is very noisy and we can hypothesise this is because some subgroups are much better than others. Therefore, now we can assign each individual subgroup a rank based on how often it appears in a training run that achieves a high accuracy. Now, using this total order on the subgroups, we can complete experiments while choosing each subgroup based on rank.

From the analysis we notice that training on just the first few ranked subgroup deems no significant improvement compared to random sampling. See [Appendix B.2](#) for full experimental details as we analyse the hypothesis test that the models are different and conclude that there is no significance showing difference using an independent t-test.<sup>5</sup>

### 3.4 Conclusions

From this, we can hypothesise that Transformers struggle to generalise when trained on subgroups. We initially observe that sampling smaller, random subgroups yields worse performance than training on the entire available sample set; moreover, hand-picking “best” subgroups does not ultimately outperform random selection in a statistically significant way.

This may be because Transformers are data-hungry and require highly representative, diverse training data to learn robust representations. Randomly sampled subgroups may fail to capture the full variability of the dataset, leading to poorer generalisation.<sup>6</sup> Conversely, even when we attempt to select the most “informative” subgroups, we do not see a significant improvement in low-data environments. To draw firm conclusions, one would need to explore a broader range of model configurations, since our findings are specific to this “tiny” Transformer.<sup>7</sup> Finally, training on only 10% of a massive dataset may not fully capture the benefits of scale—larger Transformers, trained on more extensive data, often tend to generalise far better in real-world applications.

## 4 Generalisation of symmetric groups in GNNs

Having explored the efficacy of Transformers in capturing complex dependencies, a logical next step is to consider Graph Neural Networks (GNNs), which are naturally suited to structured data.<sup>8</sup> We use GNNs that have message-passing mechanism to learn representations that capture both local and global structure, regardless of node labels. This makes them particularly effective for learning functions under symmetry groups such as the symmetric group  $S_n$  or its subgroups, where permuting node labels should not alter underlying properties.<sup>9</sup> By designing GNN architectures that are invariant or equivariant to these permutations, one can directly leverage group-theoretic insights—reducing redundancy, improving generalization, and ensuring that the model’s outputs remain consistent under re-labellings imposed by specific subgroups.

### 4.1 Methodology

For this particular task, I chose to use a Graphical Convolutional Network as my choice of GNN. This is a network that operates by letting each node update its features by aggregating information from its immediate neighbours. This approach is especially useful in graphs where edges can represent different relationships (e.g., in knowledge graphs or multi-relational data), and it supports more expressive modeling of how nodes influence one another. An example of GCNs’ ability to capture group-based relationships is seen in their application to citation networks, where research papers form natural communities that enhance classification performance<sup>10</sup>.

We train our RGCN-based model using mini-batches, computing gradients and updating weights each epoch via standard optimizers (Adam). We then evaluate performance on a validation set to monitor over-fitting and employ early stopping when the validation loss stops improving for a specified number of epochs [patience]. Hyperparameter tuning—covering factors such as learning rate, number of epochs, dropout rate, and hidden dimension size—is done via grid searches, using inference to control the change between runs, with the best configurations chosen based on validation accuracy and loss trends. This balanced approach helps the model converge efficiently and maintain good generalization.

See [Appendix C.2](#) for full setup and algorithmic details. Out of the (hyper)parameters here, we chose the number of convolutional layers [3] and hidden channels [128] at first by training on a network and finding the minimal values where we were able to successfully represent the data [see Experiment 0]. Next, we chose the LR scheduler and convergence parameters empirically by observing some full runs for the patience values and then chose ‘ReduceLROnPlateau’ as a

---

<sup>5</sup>Omitted from main section as trivial results

<sup>6</sup>Vaswani et al. (2017) introduced the original Transformer, underscoring the importance of attention in processing diverse training.

<sup>7</sup>Kaplan et al. (2020) observe that some models exhibit abrupt gains, potentially due to emergent capabilities at specific sizes.

<sup>8</sup>Foundational works in deep learning often treat GNNs within a unified framework that emphasizes the preservation of structure.

<sup>9</sup>Research in graph representation similarly emphasizes the importance of leveraging these symmetries to enhance generalization.

<sup>10</sup>Kipf, T. N., Welling, M. (2017). *Semi-Supervised Classification with Graph Convolutional Networks* (ICLR).

simple adaptive learning scheme that tries to avoid early stopping and encourage grokking. Finally, we did a grid search on the remaining tunable hyperparameters. <sup>11</sup>

## 4.2 Experiment 0

In our initial experiment, our goal was to identify model parameters that enable the learning of symmetric groups. To achieve this, we conducted a preliminary search—essentially an inference-based or binary search—over the number of channels and convolutional layers, aiming to determine the minimal architecture capable of learning  $S_4$  from randomly split data. Binary search is an efficient algorithm that repeatedly halves the search space at each step, allowing us to quickly converge on an optimal configuration.

As a result of this search, we trained on 90% of  $S_4$  and in particular with 128 hidden channels and 3 convolutional layers we achieved an accuracy of over 75% consistently on the test set. <sup>12</sup>

## 4.3 Experiment 1: $S_4$

The initial experiments on  $S_4$  seem to imply that subgroups tend to generalise better than random splits in low-resource environments.

This figure provides strong evidence that, across different training splits, subgroup-based sampling consistently achieves statistically significant performance compared to training on random splits of the Cayley table. See Appendix B.3.1 for full experimental setup and the statistical analysis demonstrating the significance of this difference.

The initial conclusion of the results from this suggests that we can say that the GNN will leverage the underlying algorithmic structure of the subgroups and use this to effectively learn the elements of the Cayley table. This is extremely useful as we are able to achieve performance of above 40% accuracy for a 6% training split, which is significantly better than the Transformer’s performance <sup>13</sup>. An additional note is that the error bars appear tighter for high training splits like at 30%. This shows us how sensitive deep learning models can be to training structure. The next steps for this experiment would be to analyse the broader applicability of these findings and how it works across domains. An attempt at  $S_5$  was tried but was unsuccessful mainly owing to limitations of GPU hours. See Appendix B.4.1 for details.

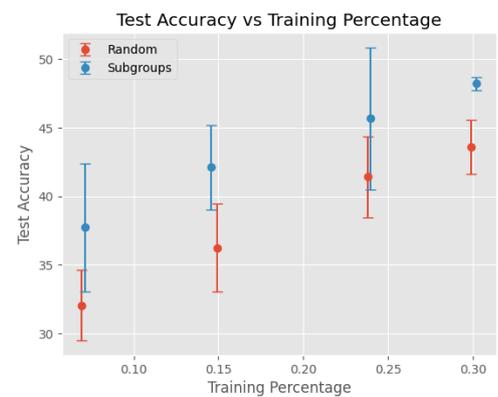


Figure 3: GNN results with  $S_4$

## 5 Summary & Further Analysis

In summary, our experiments suggest that when training on subgroups, the Graph Convolutional Network (GCN) outperforms transformers on symmetric groups. However, our analysis remains preliminary, as we have yet to evaluate the performance of the graph-based network on other group types, such as those found in modular arithmetic. The discrepancy in performance suggests Transformers rely heavily on statistical generalization—capturing frequent patterns explicitly—whereas GNNs inherently exploit compositional generalization, allowing them to generalize learned rules about subgroup interactions to new contexts. This architectural difference highlights the advantage of explicitly compositional inductive biases in structured data scenarios

Additionally, both architectures exhibit unique strengths that merit further investigation. For instance, the GCN’s superior performance on symmetric groups highlights its ability to capture structural properties effectively, whereas the Transformer’s design might be better suited for tasks involving more nuanced or sequential relationships. In total, our findings underscore the importance of tailoring model architectures to the specific characteristics of the subgroup being modeled.

As a next step, we should extend our analysis to include a broader range of group types, refine our evaluation metrics, and conduct a more comprehensive comparison between Transformer-based models and GNNs to better understand their respective advantages and limitations in subgroup learning.

<sup>11</sup>Choosing a dropout of 0.5, patience of around 1000 iterations, weight decay of 0.03 and an initial lr of 0.01.

<sup>12</sup>This is sufficient compared with 50% for 2 conv layers, similarly 55% for 64 channels. With larger models taking too long to run

<sup>13</sup>Achieved around 30% at optimal

# Appendix Content

**A: Data Construction & Groups background**

**B: Raw Results**

**C: GNN Details (including code)**

**D: Transformer Details (including code)**

**E: Group Generation Code**

## A Appendix A : Data Construction

### A.1 Definitions

We define the symmetric group  $S_n$  as the group consisting of all permutations of  $n$  distinct elements, with group operation given by permutation composition. Formally,

$$S_n = \{\sigma \mid \sigma \text{ is a bijection from } \{1, \dots, n\} \text{ to itself}\}.$$

We define the modular arithmetic addition group  $\mathbb{Z}_p^k$  as the direct product of  $k$  copies of the cyclic group of integers modulo  $p$ , i.e.,

$$\mathbb{Z}_p^k = \underbrace{\mathbb{Z}_p \times \mathbb{Z}_p \times \dots \times \mathbb{Z}_p}_{k \text{ times}},$$

where each component has group operation defined as addition modulo  $p$ .

#### A.1.1 Subgroup Definition and Calculation

We define subgroups  $H$  as subsets of a group  $G$  that satisfy the group axioms under the operation inherited from  $G$ . Formally, a subset  $H \subseteq G$  is a subgroup if for all  $a, b \in H$ , we have  $ab^{-1} \in H$ .

To identify and calculate subgroups, we:

- Fix a generating set and compute the generated elements,
- Verify closure under group operation and inverses,
- Ensure identity element inclusion.

## A.2 Summary of Representation Types

We summarize four distinct representation methods for symmetric groups ( $S_n$ ), using the symmetric group  $S_3$  as a running example. The permutations in  $S_3$  explicitly map positions to elements, for instance, the permutation  $(1, 3, 2)$  means element 1 maps to 1, element 2 maps to 3, and element 3 maps to 2.

The permutations of  $S_3$  in lexicographic order are:

$$\{(1, 2, 3), (1, 3, 2), (2, 1, 3), (2, 3, 1), (3, 1, 2), (3, 2, 1)\}$$

The Cayley table for the group operation (composition of permutations) on  $S_3$  is:

	(1,2,3)	(1,3,2)	(2,1,3)	(2,3,1)	(3,1,2)	(3,2,1)
(1,2,3)	(1,2,3)	(1,3,2)	(2,1,3)	(2,3,1)	(3,1,2)	(3,2,1)
(1,3,2)	(1,3,2)	(1,2,3)	(3,1,2)	(3,2,1)	(2,1,3)	(2,3,1)
(2,1,3)	(2,1,3)	(3,1,2)	(1,2,3)	(3,2,1)	(1,3,2)	(2,3,1)
(2,3,1)	(2,3,1)	(3,2,1)	(1,3,2)	(1,2,3)	(2,3,1)	(2,1,3)
(3,1,2)	(3,1,2)	(2,1,3)	(2,3,1)	(3,2,1)	(1,2,3)	(1,3,2)
(3,2,1)	(3,2,1)	(2,3,1)	(2,1,3)	(1,3,2)	(1,2,3)	(3,2,1)

Figure 4: Cayley Table for Permutation Composition

### A.2.1 Adjacency Matrix Representation

The adjacency matrix explicitly represents the composition of permutations as a table where each cell corresponds to the operation result.

The Cayley table shown above itself serves as an adjacency matrix for  $S_3$ , taken without the headers.

### A.2.2 Positional Encoding Representation

Positional encoding explicitly encodes position-element mappings within a permutation.

Example with permutation  $(3, 1, 2)$ :

$$[(1, 3), (2, 1), (3, 2)]$$

We would then input two of these lists into our models and look for the output in this form.

### A.2.3 Direct Element Pair Representation

Direct concatenation encodes two permutations explicitly as input pairs.

Example:

- Input permutations:  $(1, 2, 3)$  and  $(2, 3, 1)$
- Concatenated input:  $[1, 2, 3, 2, 3, 1]$
- Output permutation (composition):  $[3, 1, 2]$

### A.2.4 Graph-based Representation

Graph representation translates permutations into graph structures:

Example with permutations  $p = (1, 2, 3)$  and  $q = (2, 3, 1)$ :

- Nodes: elements  $\{1, 2, 3\}$ , encoded with one-hot vectors.
- Edges for permutation  $p$ :  $(1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3)$ .
- Edges for permutation  $q$ :  $(1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 1)$ .
- Edge types differentiate between input permutations.

This representation is particularly beneficial for graph neural network models due to its explicit relational structure.

[Backlink: Representation design choices](#)

## B Results Appendix

### B.1 Extra experiments with $Z$

After doing a through grid search through learning rate, weight decay, batch size, dropout, the following results were found. Note we kept the Transformer extremely small and did not change at all for these experiments as in each case we have tested that the transformer is capable of learning the operation [by testing on 90% training set] and making sure this learns effectively.

---

```
model = Transformer(embedding_dim=64, number_heads=4, num_layers=2, hidden_dim=64)
```

---

#### B.1.1 $Z_7^2$

Considering all the subgroups of this as our training set, we end up with a 14% training split. For random sampling at this percentage, we are able to achieve at least 15% accuracy on the test set expectedly due to the symmetry of the operator since this is an Abelian group. This can be supported by manual analysis of the sets. However, for the subgroups, we are able to learn effectively on the training set, achieving  $\sim 100\%$  even after using simple methods like increasing the dropout rates and weight decay as explained in the original Grokking paper, we are unable to achieve above 1% accuracy on the test set. A hypothesis as to why this occurs is that the model is struggling to extrapolate outside of this limited set it has been given. We also tried ...

#### B.1.2 $Z_3^3$

In this case we have a much richer model in terms of the training statistics. In this case for all the subgroups we end up using around 50% of the data points. As we know from section E, we would expect a random sample to get to 100% easily. [This has been verified]. However, after learning on only subgroups we only get an average performance of 80% across epochs, receiving [92, 65, 80] the first 3 times I ran this. This further supports evidence that subgroups hinder training.

### B.2 Ranked Experiments with $S_4$

#### B.2.1 How to form ranks

Top 3 ranked subgroups	Frequency
(1, 2, 4, 3), (1, 3, 2, 4), (1, 4, 3, 2), (1, 4, 2, 3), (1, 3, 4, 2), (1, 2, 3, 4)	11
(4, 1, 3, 2), (2, 4, 3, 1), (1, 2, 3, 4)	11
(1, 4, 3, 2), (3, 2, 1, 4), (3, 4, 1, 2), (1, 2, 3, 4)	9

Table 2: Occurrences of each subgroup.

These were generated by training the model on random selections of subgroups and incrementing the score of each subgroup based on its occurrence in the runs that converge on a solution above 90% accuracy. This was run for over 100 iterations and we get the table above as the head of that table and we use the first 2 subgroups from this for our further experiments.

We will be using the fact 6% is equivalent to training the model on  $S_4$  with only the top rank subgroup as the training set and 11% is equivalent to the top 2 subgroups.

#### B.2.2 Raw results

From the results in table 3, where we can see the experiments repeated twice out of the ten total. From these results, we can see a slight increase in performance of the sub-grouped model and can hypothesise this will be slightly better however, we can trivially see that this will not be statistically significant.

#### B.2.3 T test setup and results

For completion, we perform independent t-tests on these results to see if we can notice any significance in these results. The t-test yields a t-statistic of 0.338 and p-value of 0.7408. Since the p-value is much larger than any significance level we can say that there is no statistical difference between the accuracy scores of random sampling and highest rank subgroup sampling. This is assuming the accuracy score is noisy enough to fit a normal distribution and the variances

Method	Training Split (%)	Test Accuracy (%)
Subgroups	6	22.407
	6	22.870
	11	25.588
	11	25.488
Random	6	23.985
	6	25.461
	11	23.930
	11	24.708

Table 3: Initial Results across methods

of both sampling methods are roughly equivalent, the second of which may not hold as we scale but is assumed for this exercise.

The performance variation under subgroup-based sampling may reflect theoretical insights from group theory: certain subgroups generate richer, more generalizable rules, enabling better compositional generalization when selected appropriately. Our experiments suggest that subgroup richness (e.g., dimensionality and generator choice) directly impacts the compositional quality of learned representations.

### B.3 GNN Results

#### B.3.1 Experiment with $S_4$

These experiments show that random subgroup choices tend to perform better than random splits of the table when considering GCNs and keeping the percentage of training data consistent. Here we will outline how these results were acquired.

### B.4 Raw results

Method	Training Split (%)	Test Accuracy (%)
Subgroups	7	37.74 ± 4.65
	15	42.11 ± 3.09
	24	45.66 ± 5.17
	30	48.20 ± 0.48
Random	7	32.07 ± 2.56
	15	36.26 ± 3.18
	24	41.43 ± 2.95
	30	43.59 ± 1.95

Table 4: Results accuracy ± standard deviation for Different Training Splits

These results represent the results of training the GCN in low-resource environments. This was generated by choosing arbitrary subgroups and repeating at least 5 times to get the test accuracy with a standard deviation. More repeats were made for experiments with high standard deviation like the training split of 24%.

Then using these percentage splits we chose random splits. In this case we re-tuned some of the hyper-parameters like patience due to early convergence occurring occasionally as well as learning rate.

Additionally, we attempted this with many different hyper-parameters settings, keeping results from which the validation set converged as well as the training set as we noticed for many experiments the validation appeared to converge before the training had. Additionally, in some experiments we used the loss function paired with the validation accuracy in order to achieve a better notion of convergence and this achieved better performance especially for lower percentage splits.

We performed an independent t-test combining results across all training splits, obtaining a t-statistic of 1.91 with a p-value of 0.064. At a 10% significance level, we therefore reject the null hypothesis of equivalence, providing statistical evidence that subgroup-based sampling methods yield significantly different performance compared to random sampling. The independent t-test performed (Welch’s variant) assumes the independence and approximate normality of accuracy measurements but does not require equal variances between groups, which aligns with the observed data characteristics in our experiments.

Note: We combined all simulated accuracy samples into two larger groups—one for subgroup sampling and one for random sampling—and carried out an independent Welch's t-test (`scipy.stats.ttest_indwithequal_var = False`) to evaluate if the observed difference was statistically significant.

#### B.4.1 Experiment with $S_5$

With promising results for  $S_4$  let us see how this generalises to  $S_5$ . We are not able to do as extensive an evaluation as  $S_4$  due to limitations of GPU hours, but we can take a training split as a data point and evaluate this. <sup>14</sup>

Since  $S_5$  has a more subgroup-rich structure, we should be able to take lower training percentages than  $S_4$  and still expect good performance. In this case, I have chosen 14% training sets. I initially ran this and got 32% for random sampling and 25% for subgroup sampling, however, upon further analysis for these results, training had not converged even though we had finished the `max_epochs` of 7000.

---

<sup>14</sup>We are not able to do as extensive hyper-parameter tuning or as many repeats

## C Graph Neural Network

Here I am outlining all my design decisions and the reasoning behind them here along with pseudocode derived from my notebook.

### C.1 Input format

#### C.1.1 Pseudocode

Input dataset is the symmetricdataset( $n$ ) defined

---

**Algorithm 1** create\_graph\_dataset(tensor\_dataset,  $n$ )

---

```
1: Input: tensor_dataset, integer  $n$ 
2: Output: graph_data_list
3: graph_data_list  $\leftarrow$  empty list
4: for each  $(x, y)$  in tensor_dataset do
5:    $p, q \leftarrow x$  ; p_edge_src, p_edge_tgt  $\leftarrow$  empty lists
6:   for  $i \leftarrow 0$  to  $n - 1$  do
7:     Append  $i$  to p_edge_src
8:     Append  $(p[i] - 1)$  to p_edge_tgt
9:   end for
10:  q_edge_src, q_edge_tgt  $\leftarrow$  empty list
11:  for  $i \leftarrow 0$  to  $n - 1$  do
12:    Append  $i$  to q_edge_src
13:    Append  $(q[i] - 1)$  to q_edge_tgt
14:  end for
15:  edge_index  $\leftarrow$  Tensor([p_edge_src || q_edge_src, p_edge_tgt || q_edge_tgt])
16:  edge_type  $\leftarrow$  Tensor(0, 0, ..., 1, 1)
17:  x_features  $\leftarrow$  Identity matrix of size  $n \times n$ 
18:  y_target  $\leftarrow y - 1$ 
19:  data_obj  $\leftarrow$  torch.Data( $x = x\_features$ ,  $edge\_index = edge\_index$ ,  $edge\_type = edge\_type$ ,  $y = y\_target$ )
20:  Append data_obj to graph_data_list
21: end for
22: return graph_data_list
```

---

#### C.1.2 Explanation

For this network, a graph is required as input, which consists of node features and connectivity information. For clarity, we assume 1-based indexing and illustrate the process with an example.

Suppose we are composing two permutations,  $p = [3, 1, 4, 2]$  and  $q = [2, 4, 1, 3]$ . The operation of interest is the composition of  $p$  with  $q$  in  $S_4$ .

For each permutation, directed edges are created such that the source node corresponds to the position in the permutation and the target node corresponds to the value at that position. For instance, in permutation  $p$ :

$p[1] = 3$  implies an edge from node 1 to node 3,  $p[2] = 1$

Complete edges for  $p$  is: (1, 3), (2, 1), (3, 4), (4, 2). For  $q$ , they are: (1, 2), (2, 4), (3, 1), (4, 3).

After constructing the edges for each permutation they are merged into a single structure : Edge Index Tensor. A 2-row tensor that specifies the connectivity of the graph. The first row contains the source nodes for each edge and the second row contains the corresponding target nodes.

We then annotate each edge with an edge type to indicate its origin from the respective permutation graph:

Next, node features are created. Each node is assigned a feature vector that describes it. One-hot encoding is used, meaning node 1 is represented as [1, 0, 0, 0], node 2 as [0, 1, 0, 0], and so on. This encoding aids the network in identifying the nodes, while subsequent layers can learn more complex representations based on connectivity.

Finally, the target graph is prepared, which would be [1, 2, 3, 4] from  $p$  and  $q$

### C.2 Choice of type of GNN

There are many different types of GNNs that we could have chosen for this task. In particular for this we chose to use a Graph Convolution Network. This is a type of Graph Neural Network that works by letting each node in a graph update its own feature representation by gathering information from its immediate neighbours.

The main operations are:

1. Aggregation: Collect feature information from neighbouring nodes in the graph
2. Transformation: Aggregated features are combined using a set of learned weights. Along with a non-linear activation function

This is then repeated over several layers allowing each node to add information from a wider part of the graph.

This type of GNN has a lower complexity compared to other models and is effective for learning local structures which is why it was chosen for this task.

### C.2.1 Pseudocode

---

#### Algorithm 2 RGCNModel: Relational Graph Convolutional Network

---

```

1: Input:  $n$ , hidden_channels, dropout, in_channels =  $n$ , out_channels =  $n$ , num_relations = 2
2: Output: Predicted node embeddings
3: Define class RGCNModel(nn.Module):
4:     Initialize:
5:         conv1  $\leftarrow$  RGCNConv(in_channels, hidden_channels, num_relations)
6:         bn1  $\leftarrow$  BatchNorm1d(hidden_channels)
7:         conv2  $\leftarrow$  RGCNConv(hidden_channels, hidden_channels, num_relations)
8:         bn2  $\leftarrow$  BatchNorm1d(hidden_channels)
9:         conv3  $\leftarrow$  RGCNConv(hidden_channels, hidden_channels, num_relations)
10:        fc_dropout  $\leftarrow$  Dropout(p=dropout)
11:        fc  $\leftarrow$  Linear(hidden_channels, out_channels)
12:        self.dropout  $\leftarrow$  dropout
13:    Define forward function:
14:        Extract ( $x$ ,  $edge\_index$ ,  $edge\_type$ ) from input data
15:         $x \leftarrow$  bn1(conv1( $x$ ,  $edge\_index$ ,  $edge\_type$ ))
16:         $x \leftarrow$  ReLU( $x$ )
17:         $x \leftarrow$  Dropout( $x$ )
18:         $x \leftarrow$  bn2(conv2( $x$ ,  $edge\_index$ ,  $edge\_type$ ))
19:         $x \leftarrow$  conv3( $x$ ,  $edge\_index$ ,  $edge\_type$ )
20:         $x \leftarrow$  ReLU( $x$ )
21:         $x \leftarrow$  Dropout( $x$ )
22:         $x \leftarrow$  fc_dropout( $x$ )
23:         $x \leftarrow$  fc( $x$ )
24:    return  $x$ 

```

---

### Main Hyperparameters & Parameters

1. Number of relations is trivially 2 since we have 2 input graphs
2. Number of convolutional layers: 3 (RGCNConv)
3. Hidden channels: *hidden\_channels*
4. Learning rate, max\_epochs, weight decay, dropout
5. Learning rate scheduler and patience:
  - Scheduler: ReduceLROnPlateau
  - Patience: 100 epochs for the scheduler, 700 for early stopping

---

**Algorithm 3** Train and Evaluate RGCN Model

---

```
1: (Hyper)Parameters  $n, b\_size \leftarrow 32, num\_relations \leftarrow 2, in\_chan, out\_chan \leftarrow n, hid\_chan \leftarrow 128, dropout$ 
2:  $num\_epochs \leftarrow 5000, patience \leftarrow 700, wd, lr$ 
3:  $train\_dataset, val\_dataset, test\_dataset \leftarrow create\_graph\_dataset(dataset, n)$  ▷ Create datasets
4: Initialize  $train\_loader, val\_loader, test\_loader$ 
5:  $model \leftarrow RGCNModel(n, in\_chan, hid\_chan, out\_chan, num\_relations, dropout)$  ▷ Initialize the model
6:  $criterion \leftarrow CrossEntropyLoss()$ 
7:  $optimizer \leftarrow AdamW(model.parameters(), lr, wd)$  ▷ Define loss and optimizer
8:  $scheduler \leftarrow ReduceLROnPlateau(optimizer, mode='min', factor, patience=100)$ 
9:
10: Initialize:  $best\_val\_loss, best\_val\_acc, patience\_counter, best\_model\_state$ 
11: for  $epoch \leftarrow 0$  to  $num\_epochs$  do ▷ Training loop
12:    $model.train()$ 
13:    $train\_loss \leftarrow 0.0$ 
14:   for each data in  $train\_loader$  do
15:      $optimizer.zero\_grad()$ 
16:      $out \leftarrow model(data)$ 
17:      $y \leftarrow data.y.view(-1)$ 
18:      $loss \leftarrow criterion(out, y)$ 
19:      $loss.backward()$ 
20:      $optimizer.step()$ 
21:      $train\_loss \leftarrow train\_loss + loss.item()$ 
22:   end for
23:    $avg\_train\_loss \leftarrow train\_loss / len(train\_loader)$ 
24:    $val\_loss, val\_acc, \leftarrow evaluate(model, val\_loader, criterion)$  ▷ Validation step
25:    $scheduler.step(val\_loss)$ 
26:
27:   if  $val\_loss < best\_val\_loss$  then ▷ Early stopping and checkpointing
28:      $best\_val\_loss \leftarrow val\_loss$  ; Reset  $patience\_counter$ 
29:   else
30:      $patience\_counter ++$ 
31:   end if
32:   if  $val\_acc > best\_val\_acc$  then  $best\_val\_acc \leftarrow val\_acc$ ; Save  $best\_model\_state$ 
33:   end if
34:   if  $patience\_counter \geq patience$  then break
35:   end if
36: end for
37:
38:  $model.load\_state\_dict(best\_model\_state)$  ▷ Load the best model state
39:
40:  $test\_loss, test\_acc, test\_preds, test\_labels \leftarrow evaluate(model, test\_loader, criterion)$  ▷ Evaluation
```

---

## D Transformer Implementation

### D.1 Transformer Code

---

**Algorithm 4** Encoder Only Transformer Model

---

```
1: Input: input_len, output_dim, dropout, vocab_size, embed_dim, num_heads, hidden_dim, num_layers
2: Output: Transformed token representations
3: Define class Transformer(nn.Module):
4:   Initialize:
5:     embed  $\leftarrow$  Embedding(vocab_size, embed_dim)
6:     encoder_layer  $\leftarrow$  TransformerEncoderLayer(**kwargs)
7:     encoder  $\leftarrow$  TransformerEncoder(encoder_layer, num_layers)
8:     fc  $\leftarrow$  Linear(embed_dim, vocab_size)
9:     positional_encoding  $\leftarrow$  Learnable Parameter of shape (input_len, embed_dim)
10:  Define forward function:
11:    Input: x of shape (batch_size, 2n)
12:    x  $\leftarrow$  embed(x)
13:    x  $\leftarrow$  x + positional_encoding (optional)
14:    x  $\leftarrow$  permute(x, (1, 0, 2)) ▷ Switch dimensions around for Encoder
15:    x  $\leftarrow$  encoder(x)
16:    x  $\leftarrow$  permute(x, (1, 0, 2))
17:    n  $\leftarrow$  x.shape[1]//2
18:    x  $\leftarrow$  view(x, (batch_size, n, 2, embed_dim))
19:    x  $\leftarrow$  mean(x, dim=2)
20:    x  $\leftarrow$  fc(x)
21:    return x
```

---

### D.2 Invocation / Training Loop

---

**Algorithm 5** Training a Transformer Model on Dataset

---

```
1: input_len  $\leftarrow$  2k/2n, output_dim  $\leftarrow$  k/n, vocab_size  $\leftarrow$  p
2: dataset, sorted_subgroups, map_to_int, map_to_tokens  $\leftarrow$  repeated_cyclic_modulo_dataset_subgroups_flat(p, k)
3: total_samples  $\leftarrow$  len(dataset)
4: train_dataset, val_dataset, test_dataset  $\leftarrow$  split(dataset, [train_size, val_size, test_size])
5: Initialize train_loader, test_loader, val_loader
6: model  $\leftarrow$  Transformer(**kwargs)
7: criterion  $\leftarrow$  CrossEntropyLoss()
8: optimizer  $\leftarrow$  AdamW(model.parameters(), lr, wd) ▷ Initialize Model, Loss, and Optimizer
9:
10: for epoch  $\leftarrow$  0 to num_epochs do ▷ Training Loop
11:   model.train()
12:   train_loss  $\leftarrow$  0
13:   for batch_inputs, batch_targets in train_loader do
14:     optimizer.zero_grad()
15:     logits  $\leftarrow$  model(batch_inputs)
16:     logits  $\leftarrow$  reshape(logits, -1, vocab_size)
17:     batch_targets  $\leftarrow$  reshape(batch_targets, -1)
18:     loss  $\leftarrow$  criterion(logits, batch_targets)
19:     loss.backward()
20:     optimizer.step()
21:     train_loss  $\leftarrow$  train_loss + loss.item()
22:     clip_gradients(model)
23:   end for
24:   avg_train_loss  $\leftarrow$  train_loss / len(train_loader) ▷ Validation Step
25:   avg_val_loss, val_acc  $\leftarrow$  evaluate(model, val_loader, criterion)
26: end for
```

---

## E Dataset pseudocode

### E.1 Generating $Z_p^n$

#### E.1.1 Generating group

---

**Algorithm 6** Repeated Cyclic Modulo Dataset for  $Z_p^k$

---

```
1: function repeated_cyclic_modulo_dataset_subgroups_flat( $p, k$ )    ▷ Generate dataset for  $Z_p^k$  with cyclic addition
2:
3:   alphabet  $\leftarrow$  torch.arange( $p$ )                                ▷ Create the alphabet for  $Z_p$ 
4:   elements  $\leftarrow$  cartesian product of  $k$  copies of alphabet    ▷ Generate all group elements in  $Z_p^k$ 
5:   group_elements  $\leftarrow$  cartesian product of  $k$  copies of alphabet
6:   if  $k = 1$  then
7:     group_elements  $\leftarrow$  group_elements.unsqueeze(1)
8:   end if
9:    $n \leftarrow$  number of elements in group_elements                ▷ Generate input pairs
10:  indices  $\leftarrow$  torch.arange( $n$ )
11:  prod_indices  $\leftarrow$  cartesian product of indices
12:   $X_1 \leftarrow$  group_elements[prod_indices[:, 0]]
13:   $X_2 \leftarrow$  group_elements[prod_indices[:, 1]]
14:   $X \leftarrow$  concatenate( $X_1, X_2$  along dim=1)
15:   $Y \leftarrow (X_1 + X_2) \bmod p$ 
16:  tokens  $\leftarrow$  random permutation of  $p$                         ▷ Create a mapping from integers to tokens
17:  map_to_tokens  $\leftarrow$   $i$ : tokens[ $i$ ] for  $i$  in range( $p$ )
18:  map_to_int  $\leftarrow$  tokens[ $i$ ]:  $i$  for  $i$  in range( $p$ )
19:
20:   $X \leftarrow$  tensor of mapped token values from original  $X$       ▷ Convert dataset to tokenized format
21:   $Y \leftarrow$  tensor of mapped token values from original  $Y$ 
22:  dataset  $\leftarrow$  TensorDataset( $X, Y$ )
23:  return dataset, sorted_subgroups, map_to_tokens, map_to_int
24: end function
```

---

## E.1.2 Generating the subgroup

---

### Algorithm 7 Subgroup Generation

---

```

1: function closure_of_generators(generators,  $m$ ,  $k$ )    ▷ Compute the subgroup generated by given generators  $\mathbb{Z}_m^k$ 
2:   identity  $\leftarrow$  tuple of zeros
3:   subgroup  $\leftarrow$  set(generators)  $\cup$  {identity}
4:   changed  $\leftarrow$  True
5:   while changed do
6:     changed  $\leftarrow$  False
7:     new_elements  $\leftarrow$  empty set
8:     for  $a$  in subgroup do
9:       for  $b$  in subgroup do
10:         $s \leftarrow$  tuple( $a[i] + b[i] \bmod m$  for  $i$  in  $0, \dots, k - 1$ )
11:        if  $s \notin$  subgroup then
12:          new_elements  $\leftarrow$  new_elements  $\cup$  { $s$ }
13:        end if
14:      end for
15:    end for
16:    if new_elements  $\neq \emptyset$  then
17:      subgroup  $\leftarrow$  subgroup  $\cup$  new_elements
18:      changed  $\leftarrow$  True
19:    end if
20:  end while
21:  return sorted(subgroup)
22: end function
23: function generate_all_subgroups_nonprime( $m$ ,  $k$ )    ▷ Generate all subgroups of  $\mathbb{Z}_m^k$  for non-prime  $m$ 
24:   all_elements  $\leftarrow$  all  $k$ -tuples with entries in  $\{0, \dots, m - 1\}$ 
25:   identity  $\leftarrow$  tuple of  $k$  zeros
26:   group_order  $\leftarrow m^k$ 
27:   seen  $\leftarrow$  empty set
28:   trivial  $\leftarrow$  frozenset({identity})    ▷ Yield the trivial subgroup
29:   seen  $\leftarrow$  seen  $\cup$  {trivial}
30:   yield trivial
31:   nonzero  $\leftarrow$  all_elements  $\setminus$  {identity}
32:   for  $d \leftarrow 1$  to  $k$  do
33:     for each generators in combinations(nonzero,  $d$ ) do
34:       subgroup  $\leftarrow$  closure_of_generators(generators,  $m$ ,  $k$ )
35:       if |subgroup| = group_order then
36:         continue    ▷ Skip if subgroup is the full group
37:       end if
38:       if subgroup  $\notin$  seen then
39:         seen  $\leftarrow$  seen  $\cup$  {subgroup}
40:         yield subgroup    ▷ Yield so that we are generating on the fly
41:       end if
42:     end for
43:   end for
44: end function

```

---

## E.2 Generating $S_n$

### E.2.1 Generating the group

---

**Algorithm 8** Generate Dataset for Symmetric Group  $S_n$

---

```
1: function symmetric_dataset( $n$ )                                ▷ Generate dataset for  $S_n$ , representing the Cayley table
2:   elements  $\leftarrow \{1, 2, \dots, n\}$                         ▷ Generate all permutations of  $\{1, 2, \dots, n\}$ 
3:   permutations  $\leftarrow$  all permutations of elements
4:    $X\_alphabet \leftarrow []$ ,  $Y\_alphabet \leftarrow []$ 
5:   for perm1 in permutations do
6:     for perm2 in permutations do
7:       perm_out  $\leftarrow$  apply perm2 to perm1 to elements
8:       Append [perm1, perm2] to  $X\_alphabet$ 
9:       Append [perm_out] to  $Y\_alphabet$ 
10:    end for
11:  end for
12:   $X \leftarrow$  torch.tensor( $X\_alphabet$ )
13:   $Y \leftarrow$  torch.tensor( $Y\_alphabet$ )
14:
15:  perm  $\leftarrow$  random permutation of  $\{1, \dots, n\}$ 
16:  map_to_tokens  $\leftarrow \{i + 1: \text{perm}[i] \text{ for } i \text{ in range}(n)\}$ 
17:  map_to_int  $\leftarrow \{\text{perm}[i]: i + 1 \text{ for } i \text{ in range}(n)\}$ 
18:   $X \leftarrow$  apply map_to_tokens to  $X$                                 ▷ Map dataset to tokenized format
19:   $Y \leftarrow$  apply map_to_tokens to  $Y$ 
20:  dataset  $\leftarrow$  TensorDataset( $X, Y$ )
21:  return dataset, map_to_tokens, map_to_int
22: end function
```

---

## E.2.2 Generating subgroups

---

**Algorithm 9** Subgroup Closure and Enumeration

---

```
1: function closure_from_generators(generators, comp_table, cache) ▷ Compute subgroup closure
2:   key ← frozenset(generators)
3:   if key in cache then return cache[key]
4:   end if
5:   closure ← set(generators), changed ← True
6:   while changed do
7:     new_elements ← {comp_table[a][b] for a, b in closure}
8:     new_closure ← closure ∪ new_elements
9:     if new_closure = closure then
10:      changed ← False
11:    else
12:      closure ← new_closure
13:    end if
14:  end while
15:  cache[key] ← frozenset(closure)
16:  return cache[key]
17: end function
18: function enumerate_subgroups(generators, start, group_elements, comp_table, cache, seen, result) ▷ Recursive
subgroup enumeration
19:   H ← closure_from_generators(generators ∪ {0}, comp_table, cache)
20:   if H ∈ seen then return
21:   end if
22:   seen, result ← seen ∪ {H}, result ∪ {H}
23:   for i ← start to len(group_elements) do
24:     g ← group_elements[i]
25:     if g ∉ H then
26:       enumerate_subgroups(generators ∪ {g}, i + 1, group_elements, comp_table, cache, seen, result)
27:     end if
28:   end for
29: end function
30: function all_subgroups_ordered(perms, comp_table)
31:   enumerate_subgroups((), 0, range(1, len(perms)), comp_table, cache, seen, result)
32:   return result
33: end function
```

---

## E.3 Split Function

---

**Algorithm 10** Dataset Splitting Based on Subgroups or Random Sampling

---

```
1: perms, mapping ← symmetric_group(n)
2: comp_table ← precompute_composition_table(perms, mapping)
3: subgroups ← all_subgroups_ordered(perms, comp_table)
4: formatted_subgroups ← [
5:   convert_subgroup_indices_to_permutations(sg, perms)
6:   for sg in sorted(subgroups) ]
7: subgroups, potentials ← split_symmetric_subgroup_(formatted_subgroups)
8: train_dataset, test ← split_cayley_dataset_by_subgroups(dataset, subgroups, [0,1])
9: if len(test) = 0 or len(train_dataset) = 0 then
10:  continue
11: else
12:  flag ← False
13:  break
14: end if
15: val_dataset, test_dataset ← random_split(test, [len(test)//2, len(test) - len(test)//2])
```

---