

Group-Invariant Machine Learning on the Kreuzer-Skarke Dataset

Christian Ewert^a, Sumner Magruder^b, Vera Maiboroda^{c,d}, Yueyang Shen^e, Pragya Singh^f, Daniel Platt^{g,*}

^a*German Center for Neurodegenerative Diseases (DZNE), Bonn, Germany*

^b*Yale University, New Haven, CT, USA*

^c*University Paris-Saclay, Gif-sur-Yvette, France*

^d*Alternative Energies and Atomic Energy Commission (CEA) Paris-Saclay, Irfu/DPhP, Gif-sur-Yvette, France*

^e*Statistics Online Computational Resource, Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI, USA*

^f*Department of Computing, Imperial College London, London, UK*

^g*Department of Mathematics, Imperial College London, London, UK*

Abstract

We use supervised machine learning to predict Hodge numbers for Calabi-Yau threefolds encoded by reflexive polyhedra. The Hodge number is invariant to the order of the vertices and the swapping of axes. Incorporating these properties, i.e. the invariance of column and row permutations for a matrix containing the polyhedron's vertices, promises better performance for Hodge number prediction. In fact, we find that machine learning models that incorporate these symmetries outperform models that do not. We compare approaches incorporating different degrees of invariance and improve the state-of-the-art on the dataset.

Keywords: Machine Learning, Invariant Machine Learning, Calabi–Yau, String Compactifications, Hodge Number, Kreuzer-Skarke

1. Introduction

Calabi-Yau threefolds are widely studied objects in pure mathematics and string theory. Heterotic string theory is ten-dimensional and Calabi-Yau threefolds are candidates to construct string compactifications down to four dimensions, with the hope of modeling real-world phenomena. However, most such Calabi-Yau manifolds do not satisfy certain consistency conditions dictated by physics (He et al., 2014, Section 3) and are therefore not suitable as models for string theory. Through the definition of “favourable manifold”, one datum that enters the checks of the consistency conditions are the Hodge numbers of the Calabi-Yau manifold (He et al., 2014, p.7). Many Calabi-Yau threefolds can be obtained as codimension one submanifolds of toric 4-folds, which are encoded by a reflexive polytope obtained as the image of their moment map. There are 473,800,776 reflexive polytopes in dimension 4, which give rise to an even larger number of Calabi-Yau threefolds and a complete list of these polytopes was published by Kreuzer and Skarke. The Hodge numbers of all Calabi-Yau threefolds corresponding to reflexive polytopes have been computed explicitly, but other string theoretic consistency conditions have only been checked for a few Calabi-Yau manifolds. Checking these consistency conditions is a computationally hard problem and it is desirable to have a fast method to identify the most promising candidates before running full computations. We study the problem of predicting Hodge numbers as a model case to demonstrate that machine learning is suitable as a fast method for this type of check.

Deep learning architectures are often built to make use of the data structure of the inputs to generate an output. For example, multi-layer perceptrons (MLPs) are used for vector-valued inputs and a trained network relies on the fixed order of the vector components for accurate inference. Convolutional neural networks (CNNs) have been designed to implement translation invariance on image- or tensor-valued inputs. While reflexive polytopes are also encoded as matrices (a stack of their vertices), the Hodge number is invariant to column and row permutations corresponding to a reordering of vertices and a swapping of axes. Incorporating invariance of the output with respect to some action on the inputs promises better performance. In a project at the *London Geometry and Machine Learning (LOGML) Summer School 2022*, we tested various models with different degrees of invariance. In this paper, we compare the performance of the different approaches. Python code for the experiments in this article is available at www.github.com/danielplatt/kreuzer-skarke-ML.

2. Related Work

Several supervised machine learning models have been trained to predict Hodge numbers of Calabi-Yau manifolds. This line of research was started by He (2017), where a shallow, fully-connected neural network was used to learn Hodge numbers of complete intersection Calabi-Yau threefolds. Following this, more sophisticated neural networks (Bull et al., 2018, 2019; Erbin and Finotello, 2021) were applied to the problem, improving the original prediction accuracy. Like the reflexive polytopes dataset, this problem is invariant under the action of a large symmetry group. Aslan et al. (2023) applied group-invariant machine learning models to the same dataset further

*D.P. is the corresponding author (E-Mail: daniel.platt.berlin@gmail.com)

60 improving performance. Berglund et al. (2021) trained a neural
61 network to predict the Hodge numbers of Calabi-Yau threefolds
62 arising from polytopes of the dataset by Kreuzer and Skarke.

63 3. Dataset and Methods

64 3.1. Dataset

65 In Kreuzer and Skarke (2000), the 473,800,776 reflexive
66 four-dimensional polyhedra were classified, giving rise to (at
67 least) this many Calabi-Yau threefolds. For each reflexive poly-
68 tope, the dataset lists some information about the dual polytope
69 alongside the Euler characteristic and Hodge numbers of the
70 corresponding Calabi-Yau threefold. In this paper, we consider
71 the subset of around 78,000 reflexive polyhedra in dimension
72 four with 26 vertices and their corresponding first Hodge num-
73 bers. The polyhedra are encoded as matrices $M \in \mathbb{R}^{4 \times 26}$ (26
74 vertices stacked in column vectors) and the considered subset
75 includes 35 different Hodge numbers $h^{1,1}(M)$, i.e. 6-39 and
76 42. We use the subset as an example but we expect similar re-
77 sults for other numbers of vertices. We randomly partition the
78 dataset into 80% training data and 20% test data once and use
79 this split for every training and evaluation. The Hodge number
80 $h^{1,1}$ is invariant to column and row permutations of M . On one
81 hand, applying a column permutation merely changes the order
82 of vertices in M but leaves the polyhedron unchanged. On the
83 other hand, applying a row permutation yields an isomorphic
84 polyhedron. Thus, the function $M \mapsto h^{1,1}(M)$ is invariant under
85 the action of the group $S_4 \times S_{26}$, where S_n denotes the symmet-
86 ric group on a set of n elements. For mathematical background
87 on the dataset, we refer the interested reader to Section A.2.

88 3.2. Preprocessing

89 In addition to the dataset above which we refer to as *origi-*
90 *nal dataset*, we obtain group-invariant features from the dataset
91 by applying the preprocessing step from Aslan et al. (2023,
92 Section 2.2), denoted by π and referred to as *preprocessed*
93 *dataset*.¹ The rationale behind this is that precomposing any
94 model $\phi : \mathbb{R}^{4 \times 26} \rightarrow \mathbb{R}$ with an approximately group-invariant
95 map $\pi : \mathbb{R}^{4 \times 26} \rightarrow \mathbb{R}^{4 \times 26}$ yields an approximately group-
96 invariant composition $\phi \circ \pi$. In practice, we implemented this
97 as follows: instead of training our machine learning models ϕ
98 on matrices $M \in \mathbb{R}^{4 \times 26}$, we instead used the training data $\pi(M)$.
99 We also build a permuted version of the original dataset where
100 columns and rows of the matrices are randomly permuted (de-
101 noted as M_p) and refer to this as the *permuted dataset*. For the
102 last dataset, we apply the aforementioned preprocessing π to
103 obtain group-invariant features on the permuted dataset ($\pi(M_p)$)
104 and refer to this as *preprocessed permuted dataset*.

105 3.3. Algorithms

106 *Random Guess (baseline)*. For each polytope, we randomly as-
107 sign one of the 35 Hodge numbers appearing in the dataset.

108 *Majority Class Only (baseline)*. For each polytope, we choose
109 the most frequent Hodge number in the dataset, i.e. the Hodge
110 number 16.

111 *MLP with reduced input (baseline)*. Berglund et al. (2021)
112 trained a neural network to predict $h^{1,1}$ from features invariant
113 to column and row permutations, e.g. the number of vertices
114 in the polytope and its dual. While they trained their method
115 on matrices $M \in \mathbb{R}^{4 \times k}$ with arbitrary k , we use their results on
116 the entire dataset as the first baseline in our subproblem where
117 $k = 26$ and refer to it as *MLP with reduced input*.

118 *XGBoost*. We used the gradient-boosted decision tree imple-
119 mentation of Chen and Guestrin (2016). All inputs were flat-
120 tened, we trained with a maximum depth of 6 and used 50%
121 of the $4 * 26 = 104$ features to build each tree. The *L1* regu-
122 larisation factor was 1 and the model was trained for a maxi-
123 mum of 60000 epochs and with early stopping using patience
124 100. The model was trained on the regression objective with
125 Tweedie loss.

126 *CNN*. This neural network consists of a convolutional neural
127 network (CNN) and subsequent multi-layer perceptron (MLP).
128 The CNN consists of 2D convolution layers with kernel size 2
129 and 32 and 64 channels respectively, ReLU activations, and 2×2
130 max-pooling layers. The MLP consists of two hidden linear
131 layers with 128 and 50 neurons respectively and ReLU activa-
132 tions. The network was trained for 1000 epochs using the Adam
133 optimiser with early stopping using patience 100. The model is
134 trained on the classification objective with cross-entropy loss.

135 *Vision Transformer*. We use the transformer neural network
136 from Dosovitskiy et al. (2021) which has approximately 25 mil-
137 lion parameters over 19 layers and uses 4 attention heads. For
138 the attention heads, we use a patch size of 4×4 . The $\mathbb{R}^{4 \times 26}$ input
139 matrix is augmented to $\mathbb{R}^{52 \times 52}$ via stacking of random row and
140 column permutations. The model was trained for 60 epochs
141 using the Adam optimiser on the classification objective with
142 cross-entropy loss.

143 *Invariant MLP*. This model (Hartford et al., 2018) con-
144 tains three equivariant layers with 256 channels and
145 cross-connections between layers followed by sum-pooling
146 of each channel and five fully-connected layers with
147 (256, 256, 256, 256, 35) neurons. We trained the model with the
148 Adam optimiser and batch size 32 for 300 epochs with early
149 stopping using patience 12. This model is invariant to column
150 and row permutations of the input and uses a mean squared er-
151 ror loss in a regression setting.

152 *PointNet++*. The PointNet++ model (Qi et al., 2017b) is a neu-
153 ral network architecture designed for the processing of point
154 clouds extending PointNet (Qi et al., 2017a) with local instead
155 of global information aggregation via graphs. We used the im-
156 plementation from Fey. The model consists of four linear layers
157 with a message-passing mechanism for local feature aggrega-
158 tion. Each of these layers is followed by a ReLU activation and
159 the block of these layers is followed by a global aggregation of

¹We use the implementation from github.com/danielplatt/Fundamental-Domain-Projections.

160 features which are used by a single subsequent linear layer to
 161 make predictions. We trained the model using the Adam
 162 optimiser and batch size 64 for 300 epochs with early stopping
 163 using patience 15. This model is invariant under column per-
 164 mutations, but it is *not* invariant under row permutations, as the
 165 coordinates of the vertices are the inputs to the first layer of the
 166 network. It was trained with a cross-entropy loss function.

167 *MLP with inv. features.* This model consists of a column- and a
 168 row-encoder and a single decoder. Each encoder contains a sin-
 169 gle MLP that is applied independently to each column or row of
 170 the input. Per column or row, the respective encoder obtains a
 171 feature of length 256. Max-pooling operations are applied sep-
 172 arately across column and row features similar to the pooling in
 173 PointNet (Qi et al., 2017a). The column-encoder obtains a fea-
 174 ture invariant to column permutations while the row-encoder
 175 obtains a feature invariant to row permutations. Both features
 176 are added and processed by the decoder into the Hodge num-
 177 ber prediction. The MLPs in the encoder contain two hidden
 178 and one output layer with 256 neurons while the MLP in the
 179 decoder contains three hidden layers with 256 neurons and the
 180 output layer. The decoder MLP uses batch normalization before
 181 the activation. Except for the very last layer, all layers have a
 182 ReLU activation. The model was trained in a classification set-
 183 ting with a cross-entropy loss function, Adam optimizer, and
 184 batch size 32 for 300 epochs with early stopping using patience
 185 12.

Model	Accuracy Original Dataset	Accuracy Permuted Dataset
Random Guess	2.86%	2.86%
Majority Class Only	11.28%	11.28%
MLP with reduced input*	46.89%	46.89%
XGBoost	55.02%	29.70%
XGBoost+ π	70.63%	59.84%
CNN	$56.71 \pm 0.38\%$	$30.78 \pm 0.49\%$
CNN+ π	$71.98 \pm 0.61\%$	$61.28 \pm 0.35\%$
Vision Transformer	$62.06 \pm 0.44\%$	$43.70 \pm 0.60\%$
Vision Transformer+ π	$69.00 \pm 0.48\%$	$61.02 \pm 0.63\%$
Invariant MLP	$79.30 \pm 0.90\%$	$78.45 \pm 0.92\%$
Invariant MLP+ π	$78.15 \pm 1.73\%$	$77.96 \pm 1.50\%$
PointNet++	$95.04 \pm 0.39\%$	$86.56 \pm 0.52\%$
PointNet++ π	$95.15 \pm 0.97\%$	$90.75 \pm 0.44\%$
MLP with inv. features	$96.86 \pm 0.31\%$	$92.04 \pm 0.54\%$
MLP with inv. features+ π	$96.66 \pm 0.30\%$	$95.37 \pm 0.37\%$

* This model was trained and evaluated on a superset containing polytopes with a wide variety of vertices, not just 26 vertices. Furthermore, the inputs contain information about the polytopes' dual (see method description for details).

Table 1: Test accuracies for the task of predicting the first Hodge number of a Calabi-Yau threefold given by a reflexive polytope. Given are the mean performance and sample standard deviation over five runs. “+ π ” denotes that the machine learning model was trained on preprocessed data as outlined in Aslan et al. (2023, Section 2.2). We evaluate the methods’ performance on the original dataset and a version of the dataset where the columns and rows of each input matrix are randomly permuted.

4. Results and Discussion

In this section, we compare the methods described in Section 3.3 on the test sets of the datasets described in Section 3.2: the original and permuted datasets as well as group-invariant features based on these datasets, i.e. the preprocessed (denoted by + π) and preprocessed permuted datasets. For every method, we train one instance on the original dataset and one on the preprocessed dataset and report the accuracies on the respective corresponding test-sets. In addition and without retraining, we compare these instances on the permuted versions of the datasets. We list the accuracies in Table 1 and visualise these results in Figure 1. Note that the two models *MLP with reduced input* and *Invariant MLP* are already invariant under $S_4 \times S_{26}$, so composing with π makes no difference.

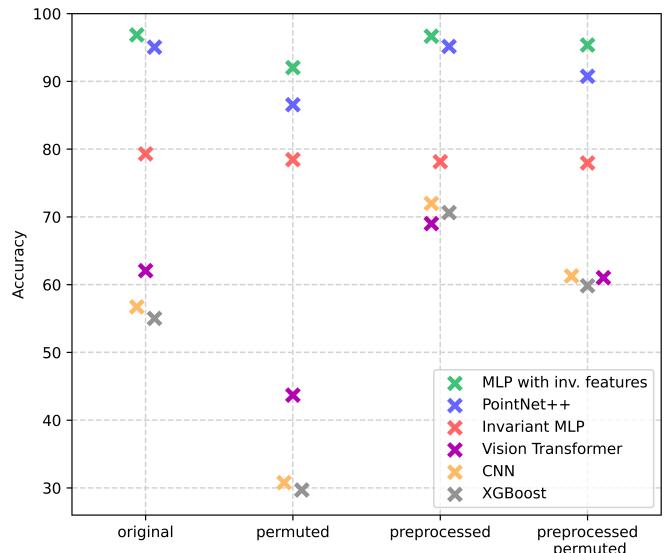


Figure 1: Performance summary of models across the four different datasets: the *original* dataset, a random permutation of the original dataset (*permuted*), the group-invariant preprocessing applied to the original (*preprocessed*) and applied to the randomly permuted dataset (*preprocessed permuted*).

We make the following observations:

1. Improved accuracy: The baseline achieved 46.89%; our best model achieves 96.86%.
2. Classification seems to perform better than regression: While it might seem more natural to formulate the task of predicting Hodge numbers as a regression task, some teams found that a formulation as a classification task, i.e. a multi-label output, performed better. Our best methods also use a classification setting.
3. Accuracy declines significantly (29.6-46%) if inputs are permuted and models are not permutation-invariant: Models such as *XGBoost*, *CNN*, and *Vision Transformer* are not by design permutation-invariant and have also not been trained with suitable data augmentation. As expected, their performance declines when permuted versions of the inputs are presented (see Table 1).

- 217 4. Building in group invariance improves performance: We
 218 notice that the models *MLP with inv. features*, *Invariant*
 219 *MLP*, and *PointNet++* perform better than other models
 220 that have no group-invariance built in. More specifically,
 221 the *Invariant MLP* is invariant under row and column per-
 222 mutations, *PointNet++* is invariant under column permu-
 223 tations, and the *MLP with inv. features* contains parts to
 224 extract features invariant under row and column permu-
 225 tations but not both.
- 226 5. Preprocessing step π improves performance: We note that
 227 making the models more group-invariant by applying the
 228 preprocessing step π increases performance for all mod-
 229 els that are not already fully group-invariant on the task in
 230 which input matrices were randomly permuted. For the
 231 models *XGBoost*, *CNN*, and *Vision Transformer* perfor-
 232 mance improves more than two times the sample standard
 233 deviation, for the other models the performance change (a
 234 decrease in two cases) is less than two times the sample
 235 standard deviation. For the models *MLP with inv. features*
 236 and *PointNet++*, performance is further improved by ap-
 237 plying the preprocessing step π for the scenario with ran-
 238 dom permutations.
- 239 6. Partially invariant models outperform fully invariant mod-
 240 els: We observe that the two partially invariant models
 241 *MLP with inv. features* and *PointNet++* perform bet-
 242 ter than the fully group-invariant model *Invariant MLP*.
 243 We assume that partially invariant models are a compro-
 244 mise between *performance* and *generalizability* (to input
 245 changes that they are invariant to). On one hand, design-
 246 ing a neural network with built-in invariance involves en-
 247 forcing the retrieval of invariant features which can intro-
 248 duce bottlenecks and over-constrain its architecture which,
 249 in turn, negatively affects its performance. On the other
 250 hand, neglecting the invariances, architectures can have
 251 sufficient capacity without bottlenecks but lack the desired
 252 invariances and generalizability. Therefore, partially in-
 253 variant models perform well because they are not over-
 254 constrained and additionally benefit from increased gen-
 255 eralizability caused by their partial invariance.

256 5. Conclusion

257 Our experiments show that the Hodge numbers of Calabi-
 258 Yau manifolds defined by reflexive polytopes can be predicted
 259 with good accuracy using machine learning.

260 Furthermore, we observe that building group invariance into
 261 the machine learning models improves performance: we see
 262 that models that are by definition invariant or partially invariant
 263 perform better than models that are not. We also see that mod-
 264 els that are not group-invariant benefit from applying a group-
 265 invariant preprocessing step in most cases.

266 Interestingly, our strongest models are those that are partially
 267 invariant but not invariant under the full group. We assume that
 268 this is the case because the partially invariant models are a com-
 269 promise between *performance* and *generalizability* (to the input
 270 changes that they are invariant to). More specifically, partially

invariant models perform well as their architecture is not over-
 271 constrained to implement full invariance and benefit from in-
 272 creased generalizability through their partial invariance.

273 Acknowledgments

274 This project was started during the *London Geometry and*
Machine Learning 2022 summer school, and we thank the orga-
 275 nisers for creating such a positive and productive environ-
 276 ment. We thank Benjamin Aslan and David Sheard for help-
 277 ful conversations. DP gratefully acknowledges funding from
 278 the Simons Collaboration in Special Holonomy. YS acknowl-
 279 edges funding from Statistical Online Computational Resources
 280 (SOCR). PS thanks Imperial College London and Goldman
 281 Sachs, where she worked previously, for their support. The
 282 work was carried out independently and neither Goldman Sachs
 283 nor Imperial was involved as an official partner.

284 Declaration of Competing Interests

285 The authors have no competing interests to declare.

286 References

- 287 Aslan, B., Platt, D., Sheard, D., 2023. Group invariant machine learning by
 288 fundamental domain projections, in: Sanborn, S., Shewmake, C., Azeglio,
 289 S., Di Bernardo, A., Miolane, N. (Eds.), Proceedings of the 1st NeurIPS
 290 Workshop on Symmetry and Geometry in Neural Representations, pp. 181–
 291 218. URL: <https://proceedings.nlr.press/v197/aslan23a.html>.
- Batyrev, V.V., 1994. Dual polyhedra and mirror symmetry for calabi-yau hy-
 292 persurfaces in toric varieties. *Journal of Algebraic Geometry* 3, 493–535.
- Berglund, P., Campbell, B., Jejjala, V., 2021. Machine Learning Kreuzer-
 293 Skarke Calabi-Yau Threefolds. doi:10.48550/arXiv.2112.09117.
- Bull, K., He, Y.H., Jejjala, V., Mishra, C., 2018. Machine learning CICY three-
 294 folds. *Physics Letters B* 785, 65–72. doi:10.1016/j.physletb.2018.
 295 08.008.
- Bull, K., He, Y.H., Jejjala, V., Mishra, C., 2019. Getting CICY high. *Physics*
 296 Letters B 795, 700–706. doi:10.1016/j.physletb.2019.06.067.
- Chen, T., Guestrin, C., 2016. Xgboost: A scalable tree boosting system, in: Pro-
 297 ceedings of the 22nd ACM SIGKDD International Conference on Knowl-
 298 edge Discovery and Data Mining, Association for Computing Machinery.
 299 pp. 785–794. doi:10.1145/2939672.2939785.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner,
 300 T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit,
 301 J., Houlsby, N., 2021. An Image is Worth 16x16 Words: Transformers for
 302 Image Recognition at Scale, in: International Conference on Learning Rep-
 303 resentations. URL: <https://openreview.net/pdf?id=YicbFdNTTy>.
- Erbin, H., Finotello, R., 2021. Inception neural network for complete inter-
 304 section Calabi-Yau 3-folds. *Machine Learning: Science and Technology* 2,
 305 02LT03. doi:10.1088/2632-2153/abda61.
- Fey, M., . Point Cloud Classification. URL: https://colab.research.google.com/drive/1D45E5bUK3gQ40YpZo65ozs7hg51-eo_U.
- Fulton, W., 1993. Introduction to Toric Varieties. volume 131 of *Annals of Mathematics Studies*. Princeton University Press. doi:10.1515/9781400882526.
- Hartford, J., Graham, D., Leyton-Brown, K., Ravanbakhsh, S., 2018. Deep
 306 Models of Interactions Across Sets, in: Dy, J., Krause, A. (Eds.), Proceed-
 307 ings of the 35th International Conference on Machine Learning, pp. 1909–
 308 1918. URL: <https://proceedings.nlr.press/v80/hartford18a/hartford18a.pdf>.
- He, Y.H., 2017. Machine-learning the string landscape. *Physics Letters B* 774,
 309 564–568. doi:10.1016/j.physletb.2017.10.024.
- He, Y.H., Lee, S.J., Lukas, A., Sun, C., 2014. Heterotic model building: 16
 310 special manifolds. *Journal of High Energy Physics* 2014. doi:10.1007/jhep06(2014)077.

330 Kreuzer, M., Skarke, H., . CY/4d: reflexive polyhedra in 4 dimensions. URL:
 331 <http://hep.itp.tuwien.ac.at/~kreuzer/CY/>.
 332 Kreuzer, M., Skarke, H., 2000. Complete classification of reflexive polyhedra
 333 in four dimensions. Advances in Theoretical and Mathematical Physics 4,
 334 1209–1230. doi:10.4310/atmp.2000.v4.n6.a2.
 335 Moon, K.R., van Dijk, D., Wang, Z., Gigante, S., Burkhardt, D.B., Chen, W.S.,
 336 Yim, K., van den Elzen, A., Hirn, M.J., Coifman, R.R., Ivanova, N.B., Wolf,
 337 G., Krishnaswamy, S., 2019. Visualizing structure and transitions in high-
 338 dimensional biological data. Nature Biotechnology 37, 1482–1492. doi:10.
 339 1038/s41587-019-0336-3.
 340 Qi, C.R., Su, H., Mo, K., Guibas, L.J., 2017a. PointNet: Deep Learning on
 341 Point Sets for 3D Classification and Segmentation, in: IEEE Conference on
 342 Computer Vision and Pattern Recognition 2017, IEEE Computer Society.
 343 pp. 77–85. doi:10.1109/cvpr.2017.16.
 344 Qi, C.R., Yi, L., Su, H., Guibas, L.J., 2017b. PointNet++: Deep
 345 Hierarchical Feature Learning on Point Sets in a Metric Space,
 346 in: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus,
 347 R., Vishwanathan, S., Garnett, R. (Eds.), Advances in Neural In-
 348 formation Processing Systems, Curran Associates, Inc. URL:
 349 https://proceedings.neurips.cc/paper_files/paper/2017/file/d8bf84be3800d12f74d8b05e9b89836f-Paper.pdf.
 350 Yau, S.T., 2008. A survey of Calabi-Yau manifolds. Surveys in differential
 351 geometry. Vol. XIII. Geometry, analysis, and algebraic geometry: forty
 352 years of the Journal of Differential Geometry, Int. Press. doi:10.4310/
 353 sdg.2008.v13.n1.a9.
 354

388 construction to four-dimensional reflexive polyhedra, then the
 389 resulting Calabi-Yau manifold is of complex dimension 3.
 389

355 A. Appendix

356 A.1. Visualisation

357 To better understand the reflexive polytopes dataset we em-
 358 ploy dimensionality reduction techniques to visualize the low
 359 dimensional embedding of input matrices. In Figure A, we pro-
 360 vide plots of input data under the *Potential of Heat-diffusion*
 361 *for Affinity-based Trajectory Embedding* (PHATE) (Moon et al.,
 362 2019). Data points are coloured according to their first Hodge
 363 number, though this Hodge number is *not* used in the com-
 364 putation of the dimensionality reduction. PHATE was cho-
 365 sen specifically due to its ability to preserve both the high-
 366 dimensional global and local distances in low dimensions.

367 In the embeddings, Hodge numbers generally increase along
 368 the first axis. For the original dataset, PHATE creates a thin
 369 region of similar input matrices with low Hodge numbers. This
 370 structure disappears when PHATE embeddings of input data
 371 with random row and column permutations applied are com-
 372 puted. For both the original and the permuted dataset, applying
 373 the group-invariant preprocessing step from Aslan et al. (2023,
 374 Section 2.2) only has a minor effect on the plot of the dimen-
 375 sionality reduction. It seems to lead to areas of matrices with
 376 low Hodge numbers that are slightly better separated from the
 377 rest of the input data, though we are unable to quantify this.
 378 For example, for the original data, the thin region of matrices
 379 corresponding to low Hodge numbers is elongated.

380 A.2. Mathematical background

381 Calabi-Yau threefolds are complex manifolds of complex di-
 382 mension 3. Many ways to construct them are known (see survey
 383 by Yau (2008, Section 2.5)). One construction is the following:
 384 a reflexive polyhedron defines a toric variety, see e.g. the text-
 385 book by Fulton (1993, Section 1.5). A hypersurface in this toric
 386 variety, suitably desingularised, is then a Calabi-Yau manifold,
 387 which was first observed by Batyrev (1994). If one applies this

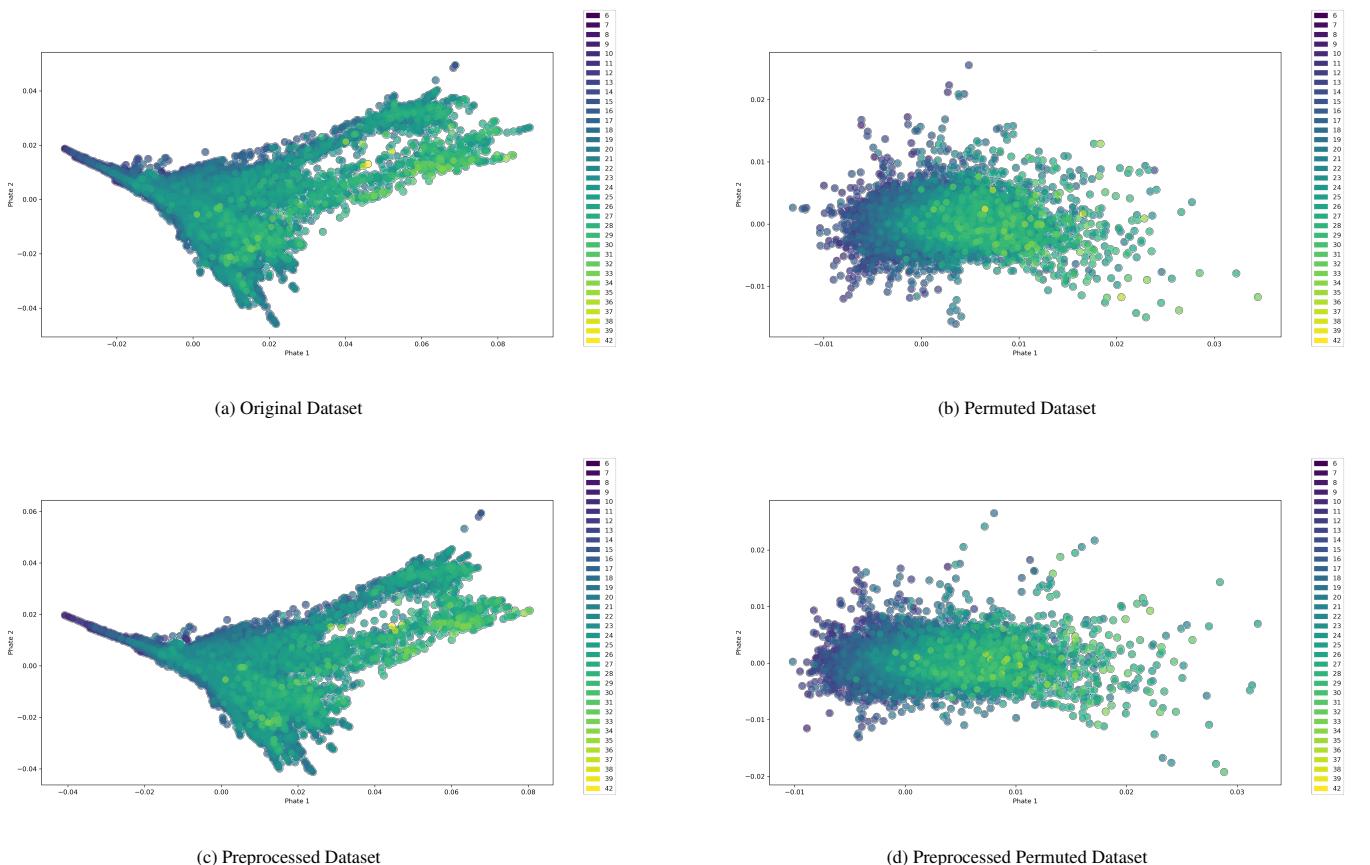


Figure A: The PHATE embeddings of inputs. Each point represents one $\mathbb{R}^{4 \times 26}$ matrix, which in turn encodes a Calabi-Yau manifold. Points are coloured according to the first Hodge number of the corresponding Calabi-Yau manifold. The Hodge numbers are not used during the computation of the embeddings. (a) Embedding of the original dataset, (b) embedding of the original dataset after random row and column permutations were applied to the input matrices, (c) embedding after a group-invariant preprocessing step was applied to the original data, (d) embedding after a group-invariant preprocessing step was applied to the randomly permuted data.