

Attention to the strengths of physical interactions: Transformer and graph-based event classification for particle physics experiments

Luc Bultjes¹, Sascha Caron^{1,2*}, Polina Moskvitina^{1,2†}, Clara Nellist^{2,3‡}, Roberto Ruiz de Austri^{4°}, Rob Verheyen^{5§} and Zhongyi Zhang^{1,2,6||}

¹ High Energy Physics, Radboud University Nijmegen, Heyendaalseweg 135, 6525 AJ Nijmegen, the Netherlands

² Nikhef, Science Park 105, 1098 XG Amsterdam, the Netherlands

³ Institute of Physics, University of Amsterdam, 1090 GL Amsterdam, The Netherlands

⁴ Instituto de Física Corpuscular, IFIC-UV/CSIC, Valencia, Spain

⁵ Department of Physics and Astronomy, University College London, London, WC1E 6BT, UK

⁶ The Bethe Center for Theoretical Physics, Bonn University, 53115 Bonn, Germany

* scaron@nikhef.nl, † p.moskvitina@nikhef.nl, ‡ c.nellist@nikhef.nl, ° r.ruiz@ific.uv.es, § r.verheyen@ucl.ac.uk, || zhongyi@th.physik.uni-bonn.de

Abstract

A major task in particle physics is the measurement of rare signal processes. These measurements are highly dependent on the classification accuracy of these events in relation to the huge background of other Standard Model processes. Reducing the background by a few tens of percent with the same signal efficiency can already increase the sensitivity considerably. This work demonstrates the importance of incorporating physical information into deep learning-based event selection. The paper includes this information into different methods for classifying events, in particular Boosted Decision Trees, Transformer Architectures (Particle Transformer) and Graph Neural Networks (Particle Net). In addition to the physical information previously proposed for jet tagging, we add particle measures for energy-dependent particle-particle interaction strengths as predicted by the leading order interactions of the Standard Model (SM). We find that the integration of physical information into the attention matrix (transformers) or edges (graphs) notably improves background rejection by 10% to 40% over baseline models (a graph network), with about 10% of this improvement directly attributable to what we call the SM interaction matrix. In a simplified statistical analysis, we find that such architectures can improve the significance of signals by a significant factor compared to a graph network (our base model).

Copyright attribution to authors.

This work is a submission to SciPost Physics.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

3 **1 Introduction**

2

4 **2 Data description**

3

5	2.1	Data generation	3
6	2.2	Event and object selection	4
7	2.3	Data format	4
8	3	Machine learning models	5
9	3.1	Boosted decision tree	5
10	3.2	Fully-connected network	6
11	3.3	Convolutional network	6
12	3.4	Particle Net	7
13	3.5	Particle Transformer	7
14	3.6	Particle Transformer as Set Transformer	8
15	3.7	Particle Transformer with Focal Loss	8
16	4	Informing the ML models about physics	9
17	4.1	Pairwise kinematic features based on 4-vectors	9
18	4.2	Pairwise kinematic features and the Standard Model Interaction Matrix	9
19	5	Results	11
20	5.1	Summary of Model details	11
21	5.2	A search for 4 top production	11
22	5.3	The top-top-Higgs searches	17
23	6	Conclusions	19
24	A	Additional Plots and Tables	21
25	A.1	AUC	21
26	A.2	ttH	23
27		References	23

28
29

30 1 Introduction

31 With the unprecedented amount of data provided by the upcoming runs of the Large Hadron
32 Collider (LHC), one can start to measure rare processes with very small cross-sections. Ex-
33 amples are the recent observations of four top quarks originating from a single proton-proton
34 collision event [1, 2]. At the heart of this endeavor is the difficult task of detecting and mea-
35 suring rare signaling processes amidst the overwhelming background noise generated by the
36 multitude of Standard Model processes. Accurate classification of these events is crucial, as
37 even a small reduction in background noise on the order of a few tens of percent while main-
38 taining the same signal detection efficiency can lead to a profound increase in sensitivity.

39 At the same time, the recent deep learning revolution has found a large variety of applica-
40 tions in high energy physics (see e.g. [3] for a review). One of these is the development of a
41 large variety of architectures for the purpose of classification of particle physics data, includ-
42 ing improved BDTs [4], convolutional networks [5], graph neural networks [6] and attention-
43 based architectures [7, 8]. These methods are mainly used in the context of the classification
44 of jet data. On the other hand, their application to event-level data has not yet been explored
45 to the same degree, and BDTs are still the most commonly used method.

46 In this work, we perform a comparison of different event classification methods and high-
 47 light a crucial aspect: the inclusion of physical information and inductive biases in machine
 48 learning architectures. In addition to the already proposed data used for jet tagging, we intro-
 49 duce a new approach by incorporating particle measures that capture the subtleties of particle-
 50 particle interaction strengths as predicted by the Feynman rules of the Standard Model ¹ With
 51 this in mind, we present a publicly available dataset to test classification methods with four
 52 top and top-top-Higgs events. We perform wide hyperparameter scans over all models and
 53 compute several performance metrics to evaluate their performance.

54 In Section 2, we describe the data generation and the data format. Section 3 briefly de-
 55 scribes the Machine Learning (ML) models used in the comparison and their optimization. In
 56 Section 4, we explore how to inform the ML models about physics. We discuss our results in
 57 Section 5, followed by the conclusions in Section 6.

58 2 Data description

59 In this section, we describe the data generation and the data format used for this work.

60 2.1 Data generation

61 We simulated proton-proton collisions at a center-of-mass energy of 13 TeV. The relevant pro-
 62 duction processes for this work consist of the backgrounds $t\bar{t} + X$, where $X = Z, W^+$ and
 63 W^+W^- and signal processes including the four top production process and $t\bar{t}H$ production.
 64 The hard scattering process generation was performed at leading order, where up to two addi-
 65 tional jets are added to the final state in the case of the background processes, and up to one
 66 for the signal. The cross-sections for all the processes and the corresponding total number of
 67 events generated are depicted in Tab. 1.

68 The hard scattering was generated with MG5_aMC@NLO version 2.7 [9] with the NNPDF31_1.0
 69 parton distribution functions set [10], using the 5 flavor scheme. Parton showering was per-
 70 formed with Pythia version 8.239 [11], while the MLM merging scheme [12] was used to
 71 merge high-multiplicity hard scattering events with the parton shower. A fast detector simula-
 72 tion was performed with Delphes version 3.4.2 [13] using the ATLAS detector map. Finally,
 73 an $H_T = \sum_{jets} E_T > 400$ GeV restriction was imposed at the parton level during event gener-
 74 ation. The purpose of this restriction is generation efficiency, since our signal region imposes
 75 $H_T > 500$ GeV on the reconstructed objects (see Section 2.2 for details.).

Physics process	σ (pb)	N_{tot}	ϵ
$pp \rightarrow t\bar{t}t\bar{t}$ (+1 j)	0.01	32463742	0.007
$pp \rightarrow t\bar{t}h$ (+2 j)	0.022	29783343	0.001
$pp \rightarrow t\bar{t}W^\pm$ (+2 j)	0.045	8954246	0.005
$pp \rightarrow t\bar{t}W^+W^-$ (+2 j)	0.0096	20160377	0.003
$pp \rightarrow t\bar{t}Z$ (+2 j)	0.034	10605846	0.011

Table 1: Signal and background processes with the corresponding LO (Leading Order) cross-section σ in pb (second column), the total number of generated events N_{tot} (third column) and the efficiency ϵ of the cuts applied (fourth column).

¹Pairwise particle-particle interaction strengths and 4-vector correlations are called “pairwise kinematic features” in the following text, see chapter 4 for details.

76 2.2 Event and object selection

77 Collision events consist of objects such as jets, b-jets, leptons, and photons, each with their
78 corresponding kinematic variables (see Section 2.3). Following the strategy in Ref. [14], an
79 event is saved if at least one of the following conditions is met:

- 80 • At least one jet or a b -jet with transverse momentum $p_T > 60$ GeV and pseudorapidity
81 $|\eta| < 2.8$,
- 82 • At least one electron with $p_T > 25$ GeV and $|\eta| < 2.47$, except for $1.37 < |\eta| < 1.52$,
- 83 • At least one muon with $p_T > 25$ GeV and $|\eta| < 2.7$,
- 84 • At least one photon with $p_T > 25$ GeV and $|\eta| < 2.37$.

85 The subsequent object selection then follows in Ref. [15], meaning that individual objects
86 are kept only if they pass the following requirements:

- 87 • Electron candidates satisfying $p_T > 28$ GeV and $|\eta| < 2.47$ are selected. In the region
88 $1.37 < |\eta| < 1.52$, known as the LAr crack region, electrons are rejected in order to
89 reduce the contribution from non-prompt and fake electrons due to detector design in
90 the liquid Argon calorimeter.
- 91 • Muon candidates are required to pass the Medium quality working point, with $p_T > 28$
92 GeV, and $|\eta| < 2.5$.
- 93 • Jet candidates satisfying $p_T > 25$ GeV and $|\eta| < 2.5$ are selected.

94 Finally, to remove as much background as possible with respect to the signal events, we
95 define a signal region [15] that requires at least six jets, at least two of which are b-tagged,
96 $H_T > 500$ GeV, and two leptons of the same sign, or at least three leptons for each event. The
97 resulting efficiencies are shown in Table 1.

98 2.3 Data format

99 The generated Monte Carlo data were saved as ROOT files and then processed into CSV files
100 using the event selection presented in Section 2.2. Each line in the CSV files is of variable
101 length and contains three event specifiers followed by the kinematic features for each object
102 in the event. The specific line format follows the event format used in the Dark Machines
103 challenges [14, 16], given by

$$\text{event ID; process ID; weight; } \cancel{E}_T; \phi_{\cancel{E}_T}; \text{obj}_1, E_1, p_{T_1}, \eta_1, \phi_1; \text{obj}_2, E_2, p_{T_2}, \eta_2, \phi_2; \dots$$

104 such that each object is represented by a string that starts with an identifier obj_n ², followed
105 by its kinematic properties in the form of a four-vector containing the full energy E and the
106 transverse momentum p_T in units of MeV, as well as the pseudo-rapidity η and the azimuthal
107 angle ϕ . The other relevant quantities are \cancel{E}_T and $\phi_{\cancel{E}_T}$, which represent the magnitude of
108 E_T^{miss} and the azimuthal angle $\phi_{E_T^{\text{miss}}}$ of the missing transverse energy. The other three com-
109 ponents represent the identity of an event, the corresponding physical process, and the event
110 weight, which is given by the cross-section of the process divided by the total number of events
111 generated.

112 Since the length of the events is variable, the data is zero-padded to the largest number of
113 objects found in the events within in the entire dataset. The dataset includes 302 072 events,

²j: jet, b: b-jet, e-: electron, e+: positron, μ^- : muon, μ^+ : antimuon, g: photon.

114 half of which correspond to the four tops signal and half of which are background processes.
 115 All background processes have an equal number of events. Finally, the dataset was split in 80%
 116 used for training, 10% used for validation and 10% used for testing. The data are available in
 117 CSV format in Ref. [17].

118 3 Machine learning models

119 In this section we provide a brief summary of the models and methods used in this work.

120 3.1 Boosted decision tree

121 We use Light Gradient Boosting Machine (LightGBM³) for this study to test the performance
 122 of Boosted Decision Tree (BDT).

123 BDTs combine a series of weak classifiers (decision trees) into a stronger classifier through
 124 gradient boosting. The boosting strategy is defined with respect to a series of previous deci-
 125 sion trees f_1, f_2, \dots, f_{t-1} which remain fixed, while the t -th tree f_t is calculated. This process is
 126 made highly efficient in LightGBM by converting the input data to histograms, and using gra-
 127 dient based sampling to focus on the data that are not well modelled. This procedure reduces
 128 memory usage and is optimized for both CPU and GPU performance. LightGBM uses first- and
 129 second-order derivatives to minimize the loss for the next iteration for gradient boosting

$$\begin{aligned} \text{Loss}^{(t)} &= \sum_{i=1}^n l(y_i, (\hat{y}_i^{(t-1)} + f_t(x_i))) + \sum_{i=1}^t \omega(f_i) \\ &\approx \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \omega(f_t) + \text{constant} \\ g_i &= \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)}), \quad h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)}). \end{aligned} \quad (1)$$

130 where l is a reconstruction loss functions, e.g. Mean Square Error, Binary Cross Entropy, etc.,
 131 f_t is the t -th tree, and $\hat{y}^{(t-1)}$ is the class label predicted by f_1, f_2, \dots, f_{t-1} . The term $\omega(f_i)$
 132 represents tree complexity terms that involve properties such as depth, number of leaves, etc.
 133 LightGBM uses a depth-first algorithm to add branches to the tree f_t with limitation on largest
 134 depth while minimizing Equation (1).

135 A common property of collision data in particle physics is a wide variability in object num-
 136 bers and types. A structured formatting of the data would thus lead to a large degree of
 137 sparsity, which we found to significantly degrade the BDT performance. We thus choose to
 138 pre-process the data by limiting the maximum number of (jets, b-jets, e^- , e^+ , μ^- , μ^+) to the
 139 **(4, 4, 1, 1, 1, 1)** hardest objects respectively.

140 One useful feature of this BDT is the fact that training is much faster than is the case for
 141 all the other architectures described below. Therefore, fine-tuning of the hyperparameters and
 142 adjustment of the data format is easy and efficient. Another attractive advantage of BDTs is
 143 their capacity to indicate feature importance, which details which input value is most impor-
 144 tant for its performance. In Section 4.1 we discuss the inclusion of high-level features beyond
 145 the raw four-vectors, where this feature is especially useful.

³<http://github.com/microsoft/LightGBM> with binary cross entropy as loss function, auc as early stop metric, 5000 estimators, 500 leaves, 0.01 learning rate, gbdt boost type, and max depth equaling to 15.

146 3.2 Fully-connected network

147 Fully-connected neural networks (FCNs) [18] are deep neural networks in their most basic
 148 form. They consist of several layers of neurons, each of which is connected to every neuron in
 149 the following layer. The connections represent a linear transformation with trainable parame-
 150 ters, which are followed by a non-linear activation function. After the final layer with a single
 151 node, a sigmoid activation is applied to produce a classification score.

152 The hidden layers all use ReLU activations, and the first five layers were followed by a
 153 Dropout layer with probability **0.5**. The network was trained with the Adam optimizer [19]
 154 with default parameters. An exponential learning rate schedule with $\gamma = \mathbf{0.95}$ was used, as
 155 well as early stopping by monitoring the performance on the validation set. The batch size
 156 and learning rates, as well as the network parameters were optimized using Optuna [20].

157 As is the case for BDTs, we found that the FCN performance generally deteriorates when
 158 applied to large, sparse input data. The data is thus pre-processed following the prescription
 159 given in Section 3.1.

160 3.3 Convolutional network

161 Convolutional Neural Networks (CNNs) [21] are mainly applied to analyze data where adja-
 162 cent items have a causal relationship, i.e. image data. CNNs apply convolutional operations
 163 through trainable filter matrices that slide over the data to produce output that is translation-
 164 ally equivariant. The convolutional layers are usually followed by pooling operations to reduce
 165 the dimension of the data in the inner network layers. Here, we utilize a one-dimensional vari-
 166 ant of such an architecture (1D CNN) called the DeepAK8 algorithm, originally used for jet
 167 tagging [22].

168 We incorporate 11 particle features as given in Table 2. Each feature is then represented by
 169 an array of size $N_{\max} = \mathbf{18}$, the maximum number of objects in an event. The event-wide fea-
 170 tures E_T^{miss} and $\phi_{E_T^{\text{miss}}}$ are added to the \mathbf{p}_T and ϕ feature vectors respectively. Following [22],
 171 the network consists of a set of 1D convolutional blocks that pass over each of the feature
 172 vectors separately. The output of these blocks is concatenated and passed to a FCN with ReLU
 173 activations. The blocks are composed of two sub-blocks which consist of a set of convolutional
 174 layers with a ReLU activation function followed by a max pooling layer and a Dropout layer
 175 with dropout probability **0.2**. The model was trained with the Adam optimizer with default
 176 parameters, with a learning rate scheduler and early stopping which both monitor the valida-
 177 tion AUC (Area Under the Curve) to prevent overfitting. The number of convolutional layers,
 178 the number of filters and the kernel size, as well as the FCN parameters were optimized with
 179 Optuna.

$$\overline{\overline{\text{Variables per particle}}}$$

$$E, \mathbf{p}_T, \eta, \phi, \text{jet}_{\text{tag}}, \text{b-jet}_{\text{tag}}, e_{\text{tag}}^-, e_{\text{tag}}^+, \mu_{\text{tag}}^-, \mu_{\text{tag}}^+, \gamma_{\text{tag}}$$

Table 2: Particle input variables for the 1D CNN, Particle Net and Particle Transformer. FCNs and BDTs use the same variables, but limit the information of 4 momenta. For leptons only the lepton with the highest \mathbf{p}_T per lepton type (including charge) is considered, while for jets and \mathbf{b} -jets, only the four jets with highest \mathbf{p}_T are considered. As a result, only 12 objects are used for FCNs and BDTs.

180 3.4 Particle Net

181 Particle Net (PN) [23] is a graph-based architecture based on Dynamic Graph Convolutional
 182 Neural Networks [24]. It treats events as particle cloud inspired by a point cloud [24] in
 183 Computer Vision challenges. Every final-state particle, encoded by the variables shown in
 184 Table 2, is represented by an individual node in the graph, carrying $(E, \mathbf{p}_T, \eta, \phi)$ as node
 185 values. Edges are constructed by connecting these particles with their k -nearest neighbours
 186 (kNN), where distances are defined as $\Delta R_{ij} = \sqrt{(\Delta\eta)_{ij}^2 + (\Delta\phi)_{ij}^2}$. The graph representing
 187 the event thus has N (number of final state particles) nodes and kN edges.

188 Messages are passed to every node i by all k neighbouring nodes j in the graph by applying
 189 the operation

$$x'_i = \frac{1}{k} \sum_{j=1}^k \text{FCN}(x_i, x_i - x_{i_j}) \quad (2)$$

190 to every node. Here, j runs over the k nearest neighbors of i and the weights of the
 191 FCN are the same for every node and edges combination. We performed experiments with an
 192 attention-weighted procedure rather than the simple averaging over edges of Equation (2),
 193 but found no difference in performance.

194 The node features are then updated to x'_i . Multiple layers of the above procedure are
 195 applied consecutively, and the node features after every step are concatenated, averaged over
 196 the nodes, and then processed by another FCN which also receives E_T^{miss} and $\phi_{E_T^{\text{miss}}}$ to obtain
 197 a classification.

198 We performed a wide hyperparameter scan over the ParticleNet architecture and found no
 199 significant difference in performance as long as sufficient capacity is available. We thus choose
 200 to use the hyperparameter settings recommended in Ref. [23] and the training procedure
 201 of [25]. Our implementation is based on Ref. [26].

202 3.5 Particle Transformer

203 Particle Transformer (ParT) [8] is a transformer-based architecture originally developed for jet
 204 tagging. It is inspired by the success of similar architectures in fields such as natural language
 205 processing [27], embedding individual particles rather than words. At its core lies the repeated
 206 application of the self-attention mechanism

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{T}) = \text{SoftMax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d})\mathbf{V}, \quad (3)$$

207 where \mathbf{Q} , \mathbf{K} and \mathbf{V} are trainable d -dimensional linear projections of the particle embedding
 208 based on the variables of Table 2.

209 The application of the attention mechanism serves to correlate every particle with all oth-
 210 ers. Furthermore, it is applicable to vary numbers of particles and is explicitly permutation
 211 invariant. Classification is obtained by appending a classification token to the list of parti-
 212 cle embeddings before the last few layers of the transformer. This token is a trainable set of
 213 weights that is identical for every event. The attention mechanism then correlates the classi-
 214 fication token with the event, after which it is processed by an FCN, which also receives E_T^{miss}
 215 and $\phi_{E_T^{\text{miss}}}$, to produce a classification label. Our implementation is based on Ref. [28].

216 As was the case for Particle Net, a hyperparameter scan over the ParT architecture does
 217 not lead to significant differences in performance. We thus choose to use the hyperparameter
 218 settings and training procedure recommended in Ref. [8].

219 3.6 Particle Transformer as Set Transformer

220 We incorporate the Set Transformer architecture [29] into the ParT model. In a Set Trans-
 221 former, the matrix \mathbf{Q} in Eq. (3) is no longer a projection of the input states, but rather a set of
 222 so-called inducing points. These are parameters that are jointly optimized with the rest of the
 223 parameters of the transformer. The model should thus learn to use \mathbf{Q} to effectively summarize
 224 the information contained in \mathbf{V} for any possible state.

225 This modification leads to a self-attention mechanism that is permutation *invariant* [30],
 226 meaning that permuted inputs produce exactly the same output. The usual self-attention
 227 mechanism is permutation *equivariant*, meaning that the outputs permute along with the in-
 228 puts. Since collision data presents as an unordered set of particles, the performance of the
 229 model may benefit from the former, as it imposes a stricter constraint.

230 In our experiments, we explored various configurations for the number of inducing points,
 231 specifically testing sets of {18, 20, 30, 40, 50, 100, 200} points. While the performance of the
 232 transformer model without pairwise features increased slightly with increasing number of in-
 233 ducing points, we did not observe any improvement with increasing number of inducing points
 234 in the model with pairwise features. Once again, the best Set Transformer model proved to be
 235 the one containing all pairwise kinematic interactions as explained in the following chapter
 236 (labelled ‘SetT_{int. SM}’).

237 3.7 Particle Transformer with Focal Loss

238 For the particle transformer we perform experiments using the focal loss [31] in place of the
 239 usual cross-entropy loss. It is given by

$$\text{Focal Loss} = -\alpha_t(1 - p_t)^\gamma \log(p_t). \quad (4)$$

240 where:

- 241 • α_t is a balancing factor that weights the importance of the different classes t . The alpha
 242 parameter essentially adjusts the importance given to each class and can handle class
 243 imbalances.
- 244 • p_t is the model’s estimated probability for the class label t ,
- 245 • γ is the focusing parameter, which adjusts the rate at which easy-to-classify examples
 246 are down-weighted. High γ values would decrease the contribution of events which are
 247 very much signal-like, i.e. p_t were is large.

248 While the focal loss was originally developed to handle class imbalance, it can still enhance
 249 model performance in other cases. The scaling factor $(1 - p_t)^\gamma$ increases the weight of difficult
 250 training samples, where the model does not yet assign large probability, while attenuating the
 251 loss for well-classified samples.

252 We performed a comprehensive hyperparameter scan over the focal loss parameters using
 253 the extended ParT model (with SM running coupling constants including the pairwise kine-
 254 matic features with the (third) SM interaction matrix as explained in the following chapter).
 255 Scans were performed over $\alpha \in \{0.25, 0.5, 0.75, 1\}$ and $\gamma \in \{0, 1, 2, 3, 4, 5, 6\}$. The best re-
 256 sults were achieved for $\alpha = 0.75$ and $\gamma = 3$ (the model is labelled ‘ParT_{int. SM (FL)}’). Even
 257 at these optimal values, the overall model performance was not better than that of models
 258 trained with the usual cross-entropy loss. Thus, results presented below pertain to models
 259 trained with cross-entropy loss, unless otherwise specified. However, this conclusion is highly
 260 dependent on the mixture of background processes. In particular, the focal loss leads to better
 261 performance in separating some backgrounds.

262 4 Informing the ML models about physics

263 4.1 Pairwise kinematic features based on 4-vectors

264 Previous work has highlighted that the inclusion of information beyond raw four-vector data,
 265 such as correlations of four-vectors (called here pairwise features), can improve deep learning
 266 classifier performance [8, 25] in jet physics. These pairwise 4-vector correlations are invariant
 267 masses or distances between two objects, which are typically known from jet physics.

268 Similarly, it is common practice to include high-level features in the training of BDTs to
 269 improve event classification, see e.g. [32]. The work of Ref. [25] suggests that this increase in
 270 performance is due to the resulting implicit embedding of Lorentz symmetry in the network
 271 architecture through features that adhere to (sub)symmetries. Lorentz's symmetry has previ-
 272 ously been shown to function as a strong inductive bias for neural network design [33–36].

273 For the BDT, we perform experiments with the inclusion of a variety of high-level features,
 274 which are treated on the same footing as the low-level ones. Similarly, we follow [8, 25]
 275 and include pairwise features in Particle Net and Particle Transformer through a trainable
 276 embedding U_{ij} for particles i and j . They are then included in Particle Net by replacing
 277 Equation (2) with

$$x'_i = \frac{1}{k} \sum_{j=1}^k \text{FCN}(x_i, x_i - x_j + U_{ij}) \quad (5)$$

278 and in Particle Transformer by replacing Equation (3) with

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{T}) = \text{SoftMax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d} + \mathbf{U})\mathbf{V}. \quad (6)$$

279 In all three above cases, we evaluated the performance of a wide variety of kinematic pairwise
 280 features, including m_{ij} , $\Delta\mathbf{R}_{ij}$, the jet-based features used in Ref. [8] and three-body invariant
 281 masses. Using the feature importance indicator of the BDT, and empirically for Particle Net
 282 and Particle Transformer, we find that for all architectures the performance is saturated by the
 283 inclusion of only m_{ij} and $\Delta\mathbf{R}_{ij}$. Furthermore, the BDT indicates that we find that the pairwise
 284 invariant masses lead to the biggest gain in performance. This result is in line with the findings
 285 of [25]. In the next section, we experiment with adding further information through dynamics,
 286 while maintaining the above kinematic information.

287 4.2 Pairwise kinematic features and the Standard Model Interaction Matrix

288 The Standard Model (SM) of particle physics provides the most comprehensive framework for
 289 understanding the electromagnetic, weak and strong nuclear interactions between elementary
 290 particles. We explore incorporating the dynamics of particle interactions described by the SM
 291 through the inclusion of a separate interaction matrix in the embedding U_{ij} for the PN and
 292 ParT models. The interaction matrix consists of entries indicating the significance of pairwise
 293 particle interactions.

294 To systematically investigate the effect of adding dynamic information to the models, we
 295 explore the use of three types of interaction matrices with increasing amounts of physical
 296 information. In the first matrix (abbreviation SMids and called SM matrix[1] in Table 3),
 297 an entry '1' indicates an interaction possible at leading order in the SM, while a '0' indicates
 298 interactions that only appear at higher orders. The following pairwise interactions are assigned
 299 a '1' in the matrix: jet–jet, jet–b-jet, jet– γ , b-jet–b-jet, b-jet– γ , $e^- - e^+$, $e^- - \gamma$, $e^+ - \gamma$, $\mu^- - \mu^+$, μ^-
 300 $-\gamma$, $\mu^+ - \gamma$. The omission of other particle interactions with a "0" does not mean that they are
 301 physically impossible, but is a practical limitation for the model. The simplified representation
 302 does not take into account the full complexity of the SM, but should provide a computationally
 303 tractable method for learning high-level interaction features.

304 In the second iteration of the interaction matrix (abbreviation SMconst and called SM ma-
 305 trix[2] in Table 3) we use the coupling constants of the SM as fixed parameters: $g_z = 0.758$ for
 306 the weak force for leptons, $g_s = 1.22$ for the strong force in jet interactions, and $g_e = 0.31$ for
 307 the electromagnetic force in photon interactions. The interactions between jets and photons as
 308 well as between b-jets and photons are determined by the electromagnetic coupling constant
 309 g_e , since photons have no colour charge. Consequently, the interactions are characterized as
 310 follows:

- 311 • For the jet- γ interactions, the modified coupling constant is $g_e \times 0.5$ to reflect the assump-
 312 tion that jets originate mainly from quarks for the signals investigated in this work. The
 313 factor **0.5** in the jet- γ coupling constant comes from the average charge of the quarks,
 314 which is calculated as $(\frac{1}{3} + \frac{2}{3})/2$, assuming an equal distribution of the quark charges
 315 of $\frac{1}{3}$ and $\frac{2}{3}$.
- 316 • For the b-jet- γ interactions, we take into account the electric charge of the b-quarks by
 317 using $g_e \times \frac{1}{3}$.

318 For the third interaction matrix (abbreviation SM and called SM matrix[3] in Table 3) we
 319 take the energy dependence of the coupling constants into account:

- 320 • For QED, the running of the fine-structure constant α is described by the Renormalization
 321 Group Equation (RGE). At one-loop level for a given pair of particle types (i, j) , it can
 322 be approximated as:

$$\alpha(Q^2) = \frac{\alpha(\mu_0^2)}{1 - \frac{n\alpha(\mu_0^2)}{3\pi} \cdot \ln\left(\frac{Q^2}{\mu_0^2}\right)}, \quad (7)$$

$$g_e = \sqrt{4\pi\alpha}$$

323 The factor n approximates the contribution of the different particles in the loop. We
 324 used $n = 3$ and considered only leptons. Other choices did not have much influence.

- 325 • For QCD, the running of α_s is more complex due to the non-Abelian nature of the theory.
 326 The one-loop RGE for a given pair of particle types (i, j) , α_s is:

$$\alpha_s(Q^2) = \frac{\alpha_s(\mu_0^2)}{1 + \frac{\alpha_s(\mu_0^2)(33-2n_f)}{12\pi} \ln\left(\frac{Q^2}{\mu_0^2}\right)}, \quad (8)$$

$$g_s = \sqrt{4\pi\alpha_s}$$

327 Here $\mu_0 = 91.1876$ GeV, $\alpha(\mu_0) = \frac{1}{127.5}$, $\alpha_s(\mu_0) = 0.118$, $n_f = 6$ is the number of quark
 328 flavors that are active at the energy scale Q^2 . To calculate these constants at a specific
 329 scale, such as the average transverse momentum \bar{p}_t of a particle pair in an event, we set
 330 $Q^2 = \bar{p}_t^2 = \left(\frac{p_t^i + p_t^j}{2}\right)^2$ as the energy scale in the RGEs to calculate $\alpha(Q^2)$ and $\alpha_s(Q^2)$.

331 g_e gives the effective coupling strength for electromagnetic interactions, while g_s gives
 332 the effective coupling strength for strong interactions at a given energy scale Q^2 . We
 333 used g_z as a constant value from the previous version of the matrix.

334 The interaction matrix provides a structured approach to encode SM-particle interactions
 335 for training machine learning models, especially models such as ParT. By simplifying the wide
 336 range of possible interactions into a prioritized scheme, the matrix allows learning to focus on
 337 the most important interactions. The interaction matrices here are structured in such a way
 338 that large negative numbers (-10k) are used if no particle exist (i.e. masked). This masking is
 339 achieved by the softmax activation function, which exponentiates the values in the attention
 340 matrix and thus pushes the irrelevant values towards zero.

341 5 Results

342 5.1 Summary of Model details

343 Table 3 summarizes the details of the ML models and the sessions that we found after the
 344 hyperparameter studies discussed above. In the following sections, we will discuss the perfor-
 345 mance of these models applied to 4 top signals with different backgrounds and to signals with
 346 top-top-Higgs events.

NN structure	Pairwise kinematic features	Loss function
BDT		Cross-entropy
BDT _{int.}	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij}$	
FCN		
CNN		
PN		
PN _{int.}	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij}$	
PN _{int.} SMids	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[1]$	
PN _{int.} SM const	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[2]$	
PN _{int.} SM	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[3]$	
ParT		
ParT _{int.}	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij}$	Focal [$\alpha = 0.75, \gamma = 3$]
ParT _{int.} SM (FL)	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[3]$	
ParT _{int.} SMids	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[1]$	Cross-entropy
ParT _{int.} SM const	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[2]$	
ParT _{int.} SM	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[3]$	
SetT _{int.} SM	$\mathbf{m}_{ij}, \Delta\mathbf{R}_{ij} + \text{SM matrix}[3]$	

Table 3: Summary of Machine Learning (ML) model details, including neural network (NN) structures and their respective loss functions. This table also highlights the inclusion of pairwise kinematic features in certain models. The particle input variables for these models are detailed in Table 2.

347 5.2 A search for 4 top production

348 In order to investigate the relationship between the amount of training data and the model's
 349 performance, we plotted learning curves on the Fig. 1 that shows the area under the ROC
 350 curve (AUC) scores as a function of training size. The x-axis represents the size of the training
 351 set, while the y-axis denotes the AUC score achieved by the model on a test set.

352 As illustrated in the figure, there is a clear trend of improving AUC scores with an increase in
 353 the training set size, affirming the hypothesis that larger datasets enhance model performance.
 354 Notably, this improvement is more pronounced in the initial stages of increasing the data

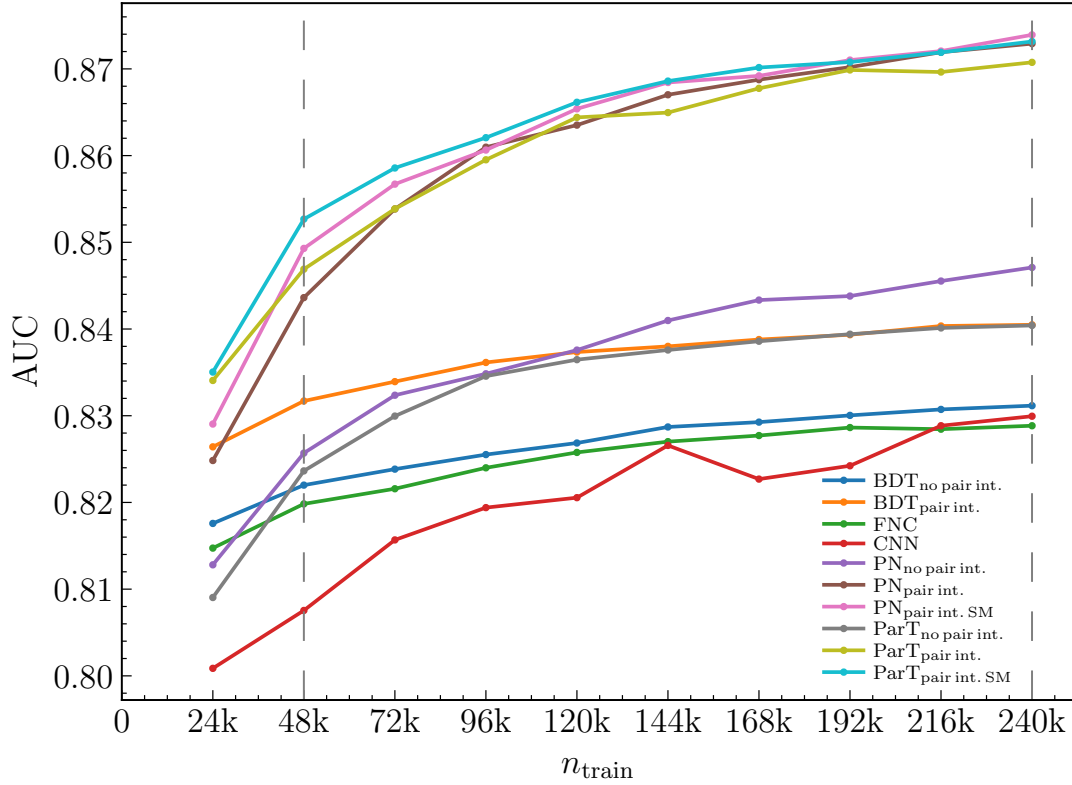


Figure 1: Learning curves of various machine learning models. This plot illustrates the relationship between the size of the training dataset and the AUC (Area Under the Curve) for each model. The vertical dashed lines represent the referenced training sizes of 48k and 240k event in the dataset.

355 volume. Beyond a certain point, however, the rate of improvement in AUC scores begins
 356 to plateau. This observation suggests that while additional training data is beneficial, the
 357 marginal gains in model accuracy diminish after reaching a certain dataset size.

358 Furthermore, the learning curves also provide insights into the data efficiency and learning
 359 capacity of the different models. Models like PN and ParT demonstrate a steeper ascent in the
 360 AUC scores with fewer data, indicating better data efficiency, while others show a more gradual
 361 improvement, reflecting their need for larger datasets to achieve comparable performance.

362 Our analysis and its key findings, based on the 48k training dataset, are summarized in
 363 Table 4. Full details covering the entire 240k training dataset can be found in Table 8 in the
 364 Appendix. This table shows the AUC values together with the background efficiencies (ϵ_B) at
 365 signal efficiencies (ϵ_S)⁴ of 30% and 70% for each evaluated method. In particular, the PN and
 366 ParT architectures with the inclusion of pairwise features and SM running coupling constants
 367 (labelled ‘int. SM’) consistently achieved the best performance in all metrics and at different
 368 signal efficiencies.

369 The background efficiencies at a signal efficiency of 30% vary among the models. Gener-
 370 ally, a lower background efficiency at this signal efficiency level indicates a model’s strength in
 371 maintaining signal detection while effectively rejecting a significant portion of the background.
 372 At a higher signal efficiency of 70%, the background efficiencies increase for all models, which
 373 is expected as increasing signal efficiency typically comes at the cost of allowing more back-

⁴ $\epsilon_S \equiv \frac{TP}{TP+FN}$ and $\epsilon_B \equiv \frac{FP}{TN+FP}$.

		BDT	BDT _{int.}	FCN	CNN
$t\bar{t} + h$	AUC	0.825(0)	0.831(0)	0.821(2)	0.778(6)
	$\epsilon_B(\epsilon_S = 0.7)$	0.206(0)	0.192(0)	0.203(1)	0.272(11)
	$\epsilon_B(\epsilon_S = 0.3)$	0.026(1)	0.026(0)	0.026(1)	0.037(1)
$t\bar{t} + W$	AUC	0.891(0)	0.895(0)	0.887(0)	0.867(5)
	$\epsilon_B(\epsilon_S = 0.7)$	0.099(0)	0.092(0)	0.103(1)	0.125(8)
	$\epsilon_B(\epsilon_S = 0.3)$	0.011(0)	0.011(0)	0.010(0)	0.011(1)
$t\bar{t} + WW$	AUC	0.740(0)	0.746(0)	0.737(1)	0.745(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.347(0)	0.339(0)	0.342(5)	0.335(3)
	$\epsilon_B(\epsilon_S = 0.3)$	0.050(0)	0.051(0)	0.054(0)	0.051(0)
$t\bar{t} + Z$	AUC	0.833(0)	0.856(0)	0.836(0)	0.839(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.191(0)	0.163(0)	0.192(0)	0.190(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.026(0)	0.019(0)	0.023(0)	0.021(1)
		PN	PN _{int.}	PN _{int. SM}	ParT _{int. SM (FL)}
$t\bar{t} + h$	AUC	0.824(0)	0.842(1)	0.846(1)	0.844(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.199(0)	0.176(3)	0.171(2)	0.176(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.025(0)	0.019(1)	0.020(1)	0.020(1)
$t\bar{t} + W$	AUC	0.887(0)	0.895(2)	0.900(1)	0.902(4)
	$\epsilon_B(\epsilon_S = 0.7)$	0.102(1)	0.097(1)	0.091(1)	0.091(5)
	$\epsilon_B(\epsilon_S = 0.3)$	0.011(0)	0.011(0)	0.010(0)	0.011(0)
$t\bar{t} + WW$	AUC	0.742(0)	0.760(1)	0.765(0)	0.768(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.335(2)	0.311(1)	0.297(2)	0.294(7)
	$\epsilon_B(\epsilon_S = 0.3)$	0.051(0)	0.044(1)	0.044(1)	0.044(1)
$t\bar{t} + Z$	AUC	0.851(0)	0.879(1)	0.887(1)	0.892(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.168(4)	0.136(1)	0.126(2)	0.119(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.020(0)	0.016(1)	0.016(0)	0.016(0)
		ParT	ParT _{int.}	ParT _{int. SM}	SetT _{int. SM}
$t\bar{t} + h$	AUC	0.824(0)	0.837(2)	0.846(1)	0.845(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.197(3)	0.179(6)	0.174(1)	0.176(3)
	$\epsilon_B(\epsilon_S = 0.3)$	0.023(0)	0.020(0)	0.020(0)	0.020(0)
$t\bar{t} + W$	AUC	0.896(1)	0.899(1)	0.905(2)	0.898(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.097(2)	0.090(1)	0.089(3)	0.094(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.010(0)	0.010(0)	0.009(0)	0.011(0)
$t\bar{t} + WW$	AUC	0.737(0)	0.767(1)	0.769(0)	0.763(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.354(3)	0.295(5)	0.288(2)	0.301(5)
	$\epsilon_B(\epsilon_S = 0.3)$	0.050(1)	0.040(0)	0.042(0)	0.047(1)
$t\bar{t} + Z$	AUC	0.839(1)	0.885(0)	0.891(1)	0.886(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.182(2)	0.130(1)	0.119(3)	0.129(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.021(1)	0.016(0)	0.015(0)	0.014(0)

Table 4: The areas under the ROC curve and the background efficiencies, at signal efficiencies of 70% and 30% respectively, correspond to the 48k training dataset. Quoted uncertainties are extracted from three independent runs for each network architecture. Numbers in bold indicate the best performance. In cases where the performances of multiple architectures are the best within the uncertainty, the results are both indicated.

374 ground events. Certain models, particularly PN_{int. SM} and ParT_{int. SM} with pairwise kinematic
375 features and the SM interaction matrix, manage to maintain relatively lower background effi-

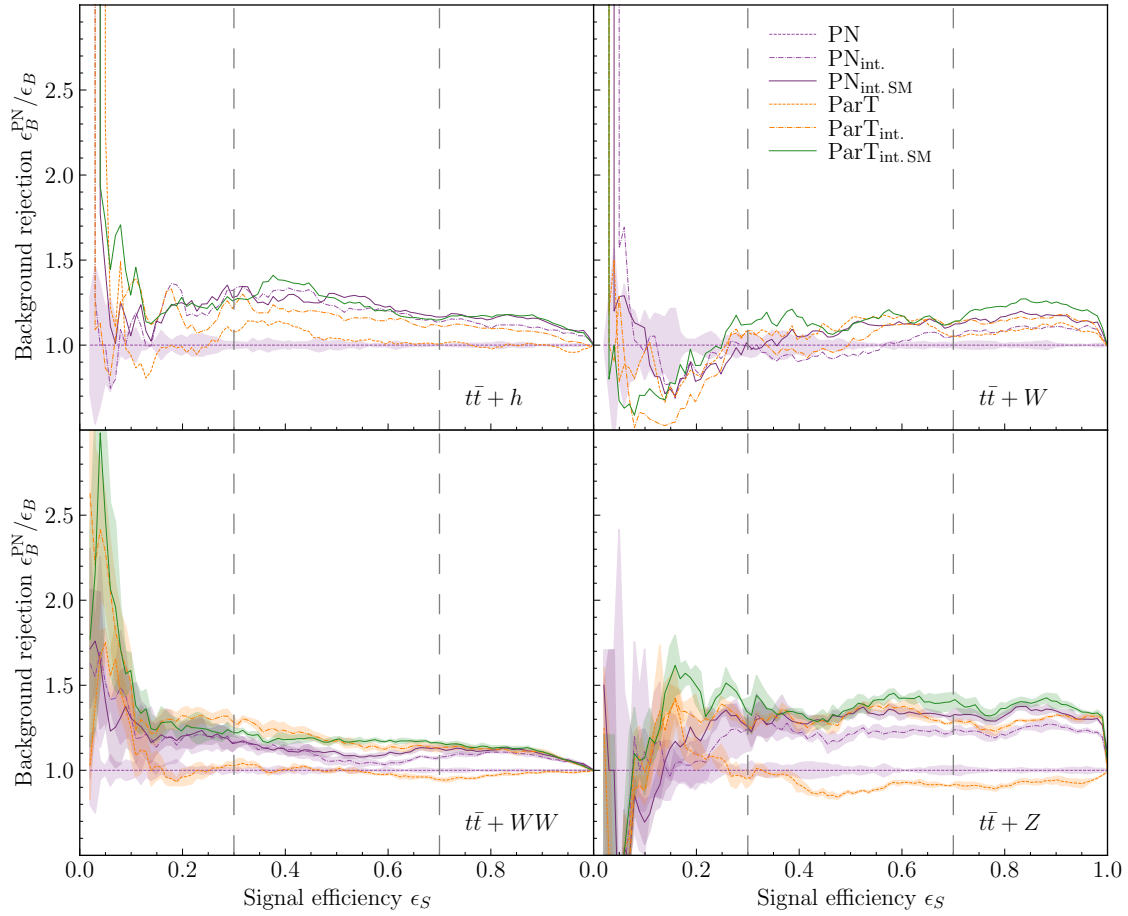


Figure 2: Signal efficiency versus background rejection plot for the four background processes corresponding to the 48k training dataset.

376 efficiencies, underscoring their efficiency in handling a more challenging balance between signal
 377 and background.

378 Figure 2 presents an alternative metric that more effectively illustrates the significance
 379 of these differences. In this comparison, we evaluate the background rejection as a function
 380 of signal efficiencies to a PN baseline model using the 48k training dataset. The PN and ParT
 381 architectures including physical information, particularly $\text{PN}_{\text{int. SM}}$ and $\text{ParT}_{\text{int. SM}}$ with pairwise
 382 kinematic features and the SM interaction matrix, demonstrate an improvement in background
 383 reduction compared to the PN baseline model of 10-40% for signal efficiencies between 30
 384 and 90%. Best performance is found for models that include the SM interaction matrix.

385 Fig. 3 shows the signal and background distributions as a function of the classifier score,
 386 normalized to the total cross-section. This figure, with its solid lines and error bands, contains
 387 the mean and standard deviation observed over three independent runs for each architecture
 388 across the entire dataset. A critical observation here is the tendency of the best performing
 389 architectures to concentrate large portions of the background at lower classifier values, espe-
 390 cially for background processes with higher cross-sections such as $t\bar{t} + W$ and $t\bar{t} + Z$. This
 391 property is of crucial importance for the discrimination of backgrounds in signal fits in LHC
 392 experiments.

393 Table 5 compares the performance of various models at two distinct signal efficiency levels,
 394 $\epsilon_S = 0.3$ and $\epsilon_S = 0.7$. Significance, denoted as σ , is defined as the signal count s divided by
 395 the square root of the background count b . This calculation is done under the assumption of

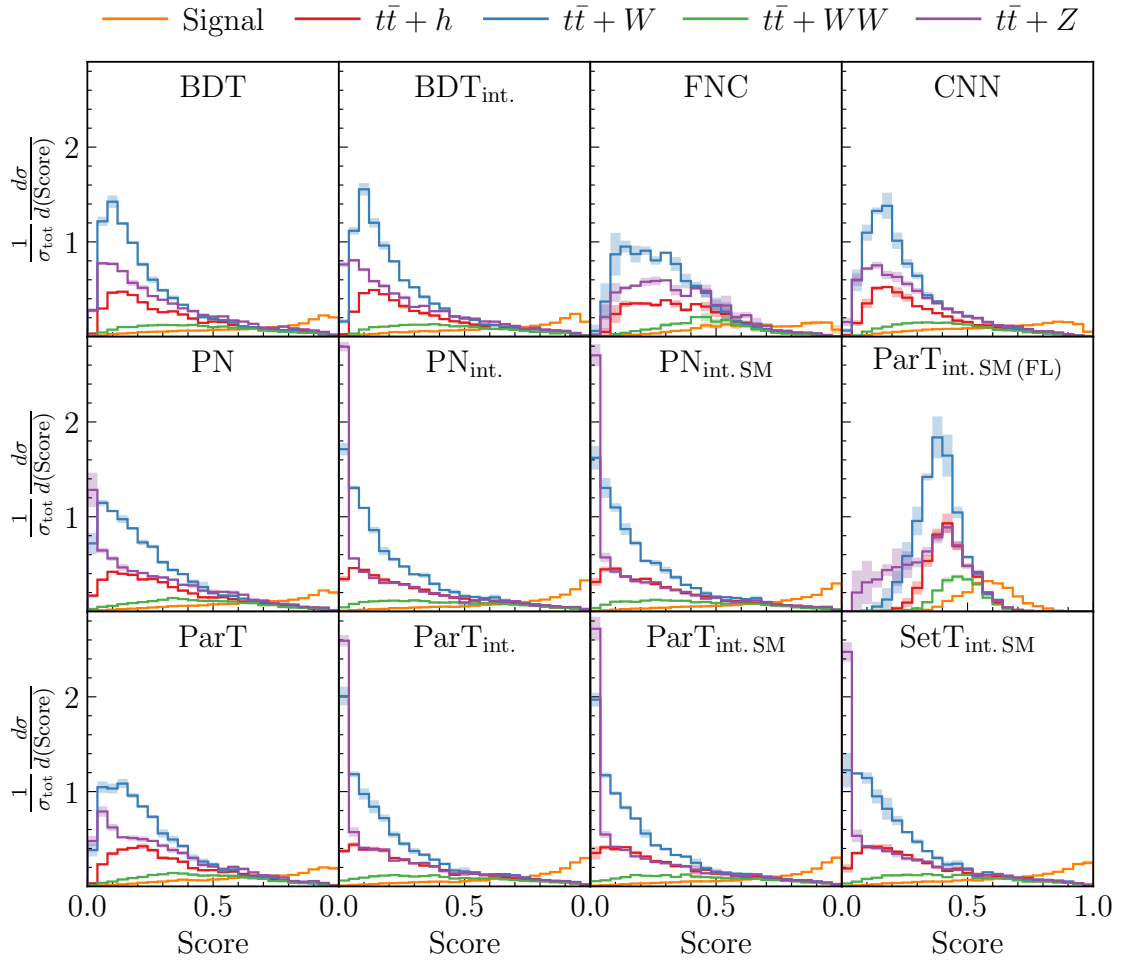


Figure 3: The distribution of the signal and the backgrounds as a function of the classifier score, normalized to the total cross-section. The solid lines and error bands correspond to the mean and standard deviation over three independent runs for each architecture across the entire dataset.

396 a luminosity of 100 fb^{-1} and incorporates the LO (Leading Order) cross-section from Table 1.
 397 We also consider the impact of systematic errors on the significance, represented by $\sigma_{\delta_{\text{sys}}=0.2}$,
 398 which is calculated as s divided by the square root of the effective background count b_{sys} ,
 399 where $b_{\text{sys}} = b + (b \cdot \delta_{\text{sys}})^2$ and δ_{sys} is set to **0.2**.

400 At $\epsilon_S = \mathbf{0.3}$, the model captures **30%** of the true signal events. The significance without
 401 systematic errors at this level suggests effective discrimination between signal and background.
 402 However, introducing a systematic error of **20%** noticeably reduces the significance, under-
 403 scoring the influence of factors like instrumental or theoretical uncertainties. At $\epsilon_S = \mathbf{0.7}$,
 404 the model identifies **70%** of the true signal events. While this higher efficiency captures more
 405 signal events, the corresponding raw significance drops. The impact of systematic errors is
 406 more pronounced at this efficiency, as evidenced by a further decrease in $\sigma_{\delta_{\text{sys}}=0.2}$.

		σ	$\sigma_{\delta_{\text{sys}}=0.2}$
BDT	$\epsilon_S = \mathbf{0.3}$	20.77	6.79
	$\epsilon_S = \mathbf{0.7}$	16.82	2.01
BDT _{int.}	$\epsilon_S = \mathbf{0.3}$	21.93	7.53
	$\epsilon_S = \mathbf{0.7}$	17.51	2.17
FCN	$\epsilon_S = \mathbf{0.3}$	20.31	6.51
	$\epsilon_S = \mathbf{0.7}$	16.67	1.97
CNN	$\epsilon_S = \mathbf{0.3}$	20.88	6.86
	$\epsilon_S = \mathbf{0.7}$	16.73	1.98
PN	$\epsilon_S = \mathbf{0.3}$	23.09	8.29
	$\epsilon_S = \mathbf{0.7}$	17.68	2.21
PN _{int.}	$\epsilon_S = \mathbf{0.3}$	25.30	9.83
	$\epsilon_S = \mathbf{0.7}$	20.51	2.97
PN _{int. SM}	$\epsilon_S = \mathbf{0.3}$	25.65	10.09
	$\epsilon_S = \mathbf{0.7}$	20.50	2.97
ParT	$\epsilon_S = \mathbf{0.3}$	22.37	7.82
	$\epsilon_S = \mathbf{0.7}$	17.72	2.23
ParT _{int.}	$\epsilon_S = \mathbf{0.3}$	24.54	9.29
	$\epsilon_S = \mathbf{0.7}$	20.21	2.89
ParT _{int. SM}	$\epsilon_S = \mathbf{0.3}$	25.36	9.88
	$\epsilon_S = \mathbf{0.7}$	20.53	2.98
ParT _{int. SM (FL)}	$\epsilon_S = \mathbf{0.3}$	26.19	10.48
	$\epsilon_S = \mathbf{0.7}$	20.28	2.91
SetT _{int. SM}	$\epsilon_S = \mathbf{0.3}$	25.58	10.03
	$\epsilon_S = \mathbf{0.7}$	20.18	2.88

Table 5: Significance table calculated for the entire dataset.

407 Comparative analysis reveals that the different versions of the particle transformer with
 408 SM interaction matrix (PartT_{int. SM} with and without focal loss and as Set Transformer) achieve
 409 the highest significance without systematic errors at $\epsilon_S = \mathbf{0.3}$. In addition, ParT_{int. SM} attains
 410 the highest significance, accounting for systematic errors at $\epsilon_S = \mathbf{0.7}$. These findings highlight
 411 the crucial role of model selection based on specific analytical requirements and the significant
 412 impact of systematic errors, especially at higher signal efficiency levels.

413 Compared to the baseline graph network (PN), it is interesting to estimate how much the
 414 sample statistic (or integrated luminosity) would have to be increased in order to achieve a
 415 similar increase in significance, neglecting systematic errors. An increase of significance e.g.
 416 from **2.21 σ** (baseline PN model at 70% signal efficiency) to **2.98 σ** (ParT_{int. SM}) corresponds

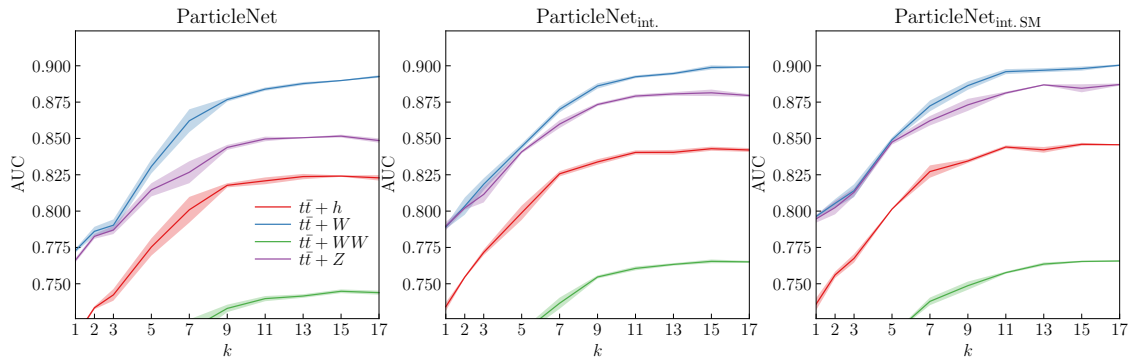


Figure 4: The performance of Particle Net on the four background processes as a function of k , the number of nearest neighbours. The solid lines and error bands correspond to the mean and standard deviation over three independent runs for every value of k correspond to the 48k training dataset.

417 to an increase in integrated luminosity of approximately 82%. An increase of significance e.g.
 418 from 8.29σ (baseline PN model at 30% signal efficiency) to 9.88σ (ParT_{int. SM}) corresponds
 419 to an increase in integrated luminosity of approximately 42% and an increase from 8.29σ to
 420 10.48σ (ParT_{int. SM (FL)}) corresponds to an increase in integrated luminosity of approximately
 421 60%.

422 Finally, in Fig. 4 shows the performance corresponding to the 48k training dataset using the
 423 AUC metric for PN as a function of k , the number of nearest neighbors. As one might expect,
 424 performance improves with k , eventually saturating when k approaches the limit where every
 425 particle is connected to all other particles. Note that this limit would require more significant
 426 computational overhead in the context of jet physics, as the number of objects can grow much
 427 larger and the complexity scales like $\mathcal{O}(kn)$. Here, the number of objects is limited and setting
 428 $k = n$ is unproblematic. Note that this setup converts PN into an architecture that is quite
 429 similar to ParT, as is reflected in the results.

430 5.3 The top-top-Higgs searches

431 To evaluate the impact of including ongoing coupling constants on the efficiency of neural
 432 networks, a second analysis was performed focusing on the search for top-top-Higgs signals.
 433 The main results of this study are presented in Table 6. The table shows the AUC for both 4 top
 434 and top-top-Higgs signal detection. The first row of the table shows the AUC values obtained
 435 for the 4 top signal, followed by the second row explaining the results for the top-top-Higgs
 436 signal. It is important to emphasize that the entire data set was used when analysing the 4
 437 top signal, while the data set for the top-top-Higgs signal, although identical in composition,
 438 intentionally excluded the data of the 4 top signal. This study shows three different ML models:
 439 the standard ParT architecture, the extended ParT with integrated pairwise features (labelled
 440 ‘int.’) and the extended ParT with SM running coupling constants (labelled ‘int. SM’). The
 441 model containing both the pairwise features and the SM interaction matrix performs best,
 442 which again confirms our earlier results. Here too, the background can be significantly reduced
 443 by about 30% compared to a PN baseline model.

444 Table 7 shows the significance estimate for the simplified top-top Higgs analysis. Here
 445 we have not included all backgrounds (compared to a full ATLAS or CMS analysis), and the
 446 values should only be used to see how important background reduction can be. For example,
 447 an increase in significance from 3.80σ (base PN model at 70% signal efficiency) to 5.02σ
 448 (PN_{int. SM}) corresponds to an increase in integrated luminosity of about 75%, which confirms

449 our previous results.

450 Comprehensive results covering other versions of the SM interaction matrices are presented
 451 in Table 9 in the Appendix.

		PN	PN _{int.}	PN _{int. SM}
$t\bar{t}t\bar{t}$	AUC	0.8471(1)	0.8729(0)	0.8739(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.1758(3)	0.1387(1)	0.1369(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0207(0)	0.0182(0)	0.0176(0)
$t\bar{t} + h$	AUC	0.8146(2)	0.8505(0)	0.8523(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.2292(1)	0.1787(0)	0.1733(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0471(1)	0.0345(0)	0.0340(0)
		ParT	ParT _{int.}	ParT _{int. SM}
$t\bar{t}t\bar{t}$	AUC	0.8404(0)	0.8708(0)	0.8732(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.1842(3)	0.1394(0)	0.1366(0)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0230(0)	0.0172(0)	0.0169(0)
$t\bar{t} + h$	AUC	0.8058(1)	0.8507(0)	0.8532(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.2399(2)	0.1794(1)	0.1748(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0502(0)	0.0357(0)	0.0351(0)

Table 6: Results for the 4 top and top-top-Higgs signals: the areas under the ROC curve and the background efficiencies, at signal efficiencies of 70% and 30% respectively, correspond to the entire training dataset. Quoted uncertainties are extracted from three independent runs for each network architecture. Numbers in bold indicate the best performance.

		σ	$\sigma_{\delta_{sys} = 0.2}$
PN	$\epsilon_S = 0.3$	32.18	7.63
	$\epsilon_S = 0.7$	34.30	3.80
PN _{int.}	$\epsilon_S = 0.3$	37.53	10.27
	$\epsilon_S = 0.7$	38.75	4.84
PN _{int. SM}	$\epsilon_S = 0.3$	37.86	10.44
	$\epsilon_S = 0.7$	39.50	5.02
ParT	$\epsilon_S = 0.3$	30.24	6.76
	$\epsilon_S = 0.7$	33.48	3.62
ParT _{int.}	$\epsilon_S = 0.3$	36.82	9.90
	$\epsilon_S = 0.7$	38.39	4.75
ParT _{int. SM}	$\epsilon_S = 0.3$	37.20	10.10
	$\epsilon_S = 0.7$	39.05	4.91

Table 7: Significance table calculated for the top-top-Higgs signal.

452 6 Conclusions

453 In this work, we present a novel approach to event classification in particle physics by inte-
454 grating physical information, in particular energy-dependent SM interactions, into advanced
455 machine learning models. Our study focuses on improving transformer models with an at-
456 tention matrix and graph networks such as PN with edge features, both of which reflect the
457 dynamical nature of SM interactions.

458 The results show that PN and ParT exhibit superior performance when these pairwise fea-
459 tures and interaction matrices are integrated. This integration improves background suppres-
460 sion by **10 – 40%** over the baseline models (PN without other physical information), with
461 approximately **10%** of this improvement directly attributable to the SM interaction matrix. In
462 a simplified statistical analysis, we find that these ML models increase significance by up to
463 **30%** compared to the baseline model. To achieve a similar improvement in significance by
464 increasing the luminosity L , one needs to increase L by about **70%**, assuming that significance
465 improvement scales with \sqrt{L} when signal and background events are proportional to L . We
466 conclude that embedding SM interactions as physical information in network structures is an
467 important avenue in this field that could lead to more accurate and efficient event classification
468 in particle physics.

469 Acknowledgements

470 The author(s) gratefully acknowledges the computer resources at Artemisa, funded by the Eu-
471 ropean Union ERDF and Comunitat Valenciana as well as the technical support provided by the
472 Instituto de Física Corpuscular, IFIC (CSIC-UV). R. RdA is supported by PID2020-113644GB-
473 I00 from the Spanish Ministerio de Ciencia e Innovación and by the PROMETEO/2022/69
474 from the Spanish GVA. RV is supported by the European Research Council (ERC) under the
475 European Union's Horizon 2020 research and innovation programme (grant agreement No.
476 788223, PanScales).

477 Appendix

478 A Additional Plots and Tables

479 This appendix complements the main text by providing an additional plot and two compre-
 480 hensive tables. It summarizes the results for the entire 240k dataset, providing a complete
 481 perspective on the data's scope and the analysis outcomes.

482 A.1 AUC

483 Fig. 5 displays the ROC curves for all architectures against various backgrounds, providing a
 484 visual perspective on their comparative performances throughout the entire dataset.

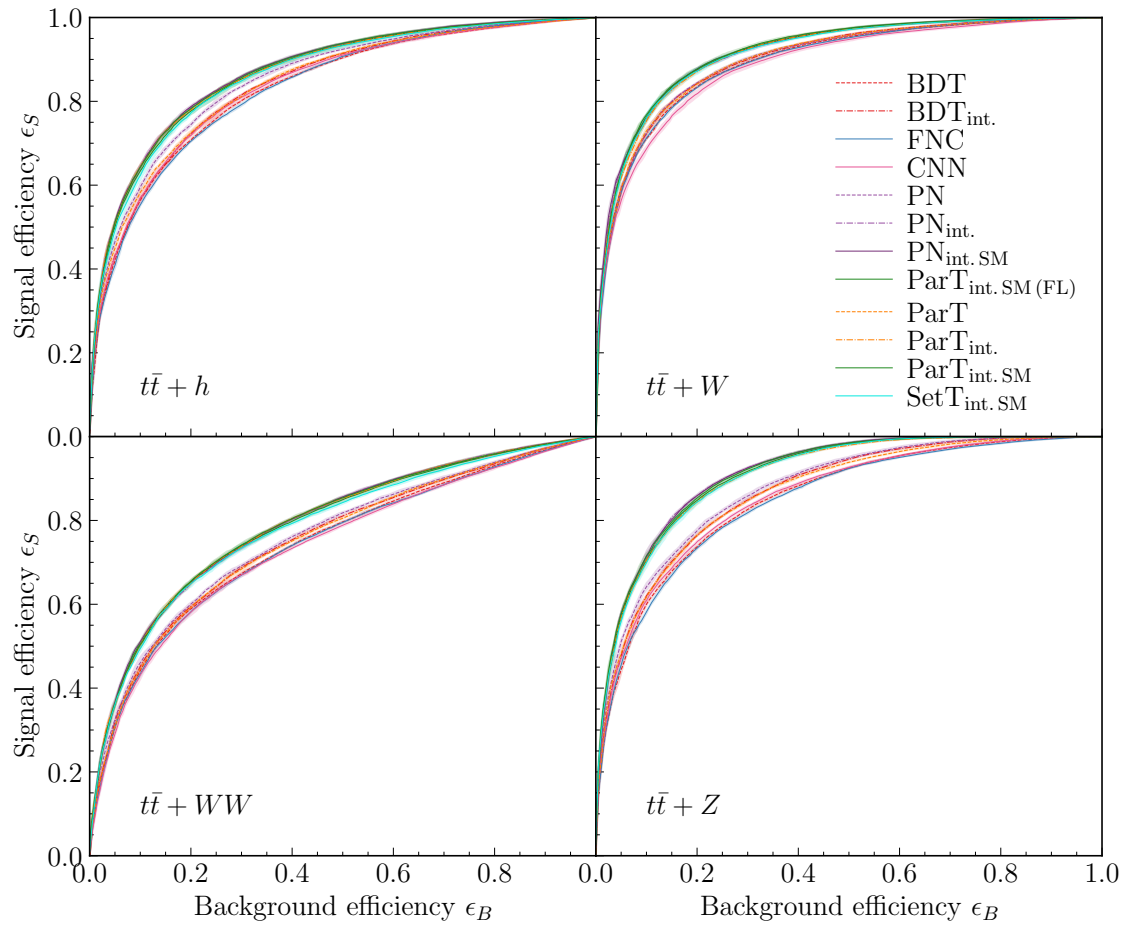


Figure 5: Receiver Operating Characteristic (ROC) curves for all architectures for the 4 top signal across the four background processes. The solid lines and error bands represent the mean and standard deviation of three independent runs for each architecture over the entire training dataset.

		BDT	BDT _{int.}	FCN	CNN
$t\bar{t} + h$	AUC	0.833(0)	0.840(0)	0.832(0)	0.838(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.193(0)	0.183(0)	0.195(0)	0.182(3)
	$\epsilon_B(\epsilon_S = 0.3)$	0.022(0)	0.022(0)	0.023(1)	0.021(2)
$t\bar{t} + W$	AUC	0.896(0)	0.900(0)	0.895(0)	0.888(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.093(0)	0.087(0)	0.093(1)	0.107(3)
	$\epsilon_B(\epsilon_S = 0.3)$	0.009(0)	0.009(1)	0.011(0)	0.009(0)
$t\bar{t} + WW$	AUC	0.745(0)	0.754(0)	0.742(0)	0.739(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.339(0)	0.317(2)	0.341(0)	0.344(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.048(0)	0.045(0)	0.048(0)	0.052(0)
$t\bar{t} + Z$	AUC	0.852(1)	0.869(0)	0.848(0)	0.857(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.167(1)	0.149(2)	0.170(2)	0.161(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.020(0)	0.018(0)	0.020(0)	0.018(0)
		PN	PN _{int.}	PN _{int. SM}	ParT _{int. SM (FL)}
$t\bar{t} + h$	AUC	0.854(1)	0.871(0)	0.872(0)	0.867(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.161(2)	0.129(1)	0.129(3)	0.138(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.016(0)	0.017(0)	0.017(0)	0.016(0)
$t\bar{t} + W$	AUC	0.901(0)	0.917(1)	0.919(0)	0.918(1)
	$\epsilon_B(\epsilon_S = 0.7)$	0.089(1)	0.072(1)	0.071(2)	0.071(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.008(0)	0.007(0)	0.007(0)	0.007(0)
$t\bar{t} + WW$	AUC	0.759(2)	0.791(0)	0.793(0)	0.791(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.312(4)	0.256(2)	0.252(2)	0.249(7)
	$\epsilon_B(\epsilon_S = 0.3)$	0.043(1)	0.036(1)	0.035(2)	0.035(1)
$t\bar{t} + Z$	AUC	0.876(2)	0.913(0)	0.913(0)	0.909(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.139(5)	0.095(1)	0.094(2)	0.097(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.015(0)	0.013(0)	0.012(0)	0.010(0)
		ParT	ParT _{int.}	ParT _{int. SM}	SetT _{int. SM}
$t\bar{t} + h$	AUC	0.843(1)	0.869(0)	0.871(0)	0.864(3)
	$\epsilon_B(\epsilon_S = 0.7)$	0.179(3)	0.131(2)	0.132(0)	0.141(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.019(0)	0.015(0)	0.015(0)	0.016(1)
$t\bar{t} + W$	AUC	0.901(0)	0.915(0)	0.918(1)	0.915(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.087(3)	0.078(1)	0.072(1)	0.074(2)
	$\epsilon_B(\epsilon_S = 0.3)$	0.008(0)	0.009(0)	0.008(0)	0.009(0)
$t\bar{t} + WW$	AUC	0.753(1)	0.792(1)	0.792(1)	0.786(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.318(5)	0.250(2)	0.248(2)	0.257(5)
	$\epsilon_B(\epsilon_S = 0.3)$	0.047(1)	0.032(0)	0.034(0)	0.036(1)
$t\bar{t} + Z$	AUC	0.866(0)	0.907(1)	0.912(0)	0.907(2)
	$\epsilon_B(\epsilon_S = 0.7)$	0.150(2)	0.098(2)	0.093(3)	0.100(4)
	$\epsilon_B(\epsilon_S = 0.3)$	0.017(0)	0.012(1)	0.011(0)	0.011(0)

Table 8: The areas under the ROC curve and the background efficiencies, at signal efficiencies of 70% and 30% respectively, correspond to the entire training dataset (240k events). Quoted uncertainties are extracted from three independent runs for each network architecture. Numbers in bold indicate the best performance. In cases where the performances of multiple architectures are the best within the uncertainty, the results are both indicated.

485 **A.2 ttH**

486 Comprehensive results covering other versions of the SM interaction matrices are presented in
 487 Table 9, which details the Area Under the Curve (AUC) for both the 4 top and top-top-Higgs
 488 signals. Specifically, the first row illustrates the AUC results for the 4 top signal, while the
 489 subsequent row delineates the outcomes related to the top-top-Higgs signal. It is important
 490 to note that, although the full dataset was employed for the 4 top signal analysis, the dataset
 491 used for the top-top-Higgs signal analysis was the same, yet it explicitly excluded data from
 492 the 4 top.

493 Five scenarios are provided for the PN and the ParT: the standard ParT/PN architecture,
 494 ParT/PN with the inclusion of pairwise features (int.), ParT/PN with the first iteration of the
 495 interaction matrix (SMids), ParT/PN with the second iteration of the interaction matrix (where
 496 the coupling constants on the SM are fixed parameters) and ParT/PN with the inclusion of SM
 497 running coupling constants (int. SM, which is the third iteration).

		PN	PN _{int.}	PN _{int. SMids}	PN _{int. SM const}	PN _{int. SM}
$t\bar{t}t\bar{t}$	AUC	0.8471(1)	0.8729(0)	0.8725(0)	0.8727(0)	0.8739(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.1758(3)	0.1387(1)	0.1377(0)	0.1384(0)	0.1369(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0207(0)	0.0182(0)	0.0178(0)	0.0178(0)	0.0176(0)
$t\bar{t} + h$	AUC	0.8146(2)	0.8505(0)	0.8489(1)	0.8505(0)	0.8523(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.2292(1)	0.1787(0)	0.1785(1)	0.1764(3)	0.1733(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0471(1)	0.0345(0)	0.0343(1)	0.0350(0)	0.0340(0)
		ParT	ParT _{int.}	ParT _{int. SMids}	ParT _{int. SM const}	ParT _{int. SM}
$t\bar{t}t\bar{t}$	AUC	0.8404(0)	0.8708(0)	0.8715(0)	0.8717(0)	0.8732(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.1842(3)	0.1394(0)	0.1389(2)	0.1372(1)	0.1366(0)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0230(0)	0.0172(0)	0.0180(0)	0.0167(0)	0.0169(0)
$t\bar{t} + h$	AUC	0.8058(1)	0.8507(0)	0.8473(0)	0.8497(0)	0.8532(0)
	$\epsilon_B(\epsilon_S = 0.7)$	0.2399(2)	0.1794(1)	0.1836(3)	0.1801(1)	0.1748(1)
	$\epsilon_B(\epsilon_S = 0.3)$	0.0502(0)	0.0357(0)	0.0355(1)	0.0367(0)	0.0351(0)

Table 9: Results for the 4 top and top-top-Higgs signals: the areas under the ROC curve and the background efficiencies, at signal efficiencies of 70% and 30% respectively, correspond to the entire training dataset. Quoted uncertainties are extracted from three independent runs for each network architecture. Numbers in bold indicate the best performance.

498 **References**

- 499 [1] G. Aad *et al.*, *Observation of four-top-quark production in the multilepton final state with*
 500 *the ATLAS detector*, *Eur. Phys. J. C* **83**(6), 496 (2023), doi:[10.1140/epjc/s10052-023-](https://doi.org/10.1140/epjc/s10052-023-11573-0)
 501 [11573-0](https://doi.org/10.1140/epjc/s10052-023-11573-0), [2303.15061](https://arxiv.org/abs/2303.15061).
- 502 [2] A. Hayrapetyan *et al.*, *Observation of four top quark production in proton-proton collisions*
 503 *at $s=13\text{TeV}$* , *Phys. Lett. B* **847**, 138290 (2023), doi:[10.1016/j.physletb.2023.138290](https://doi.org/10.1016/j.physletb.2023.138290),
 504 [2305.13439](https://arxiv.org/abs/2305.13439).
- 505 [3] G. Karagiorgi, G. Kasieczka, S. Kravitz, B. Nachman and D. Shih, *Machine Learning in*
 506 *the Search for New Fundamental Physics*, doi:[10.48550/arXiv.2112.03769](https://doi.org/10.48550/arXiv.2112.03769) (2021), [2112.](https://arxiv.org/abs/2112.03769)
 507 [03769](https://arxiv.org/abs/2112.03769).

- 508 [4] Y. Coadou, *Boosted Decision Trees*, p. 9–58, WORLD SCIENTIFIC,
509 doi:[10.1142/9789811234033_0002](https://doi.org/10.1142/9789811234033_0002) (2022).
- 510 [5] L. de Oliveira, M. Kagan, L. Mackey, B. Nachman and A. Schwartzman, *Jet-images — deep*
511 *learning edition*, JHEP **07**, 069 (2016), doi:[10.1007/JHEP07\(2016\)069](https://doi.org/10.1007/JHEP07(2016)069), [1511.05190](https://arxiv.org/abs/1511.05190).
- 512 [6] J. Shlomi, P. Battaglia and J.-R. Vlimant, *Graph neural networks in particle physics*,
513 *Machine Learning: Science and Technology* **2**(2), 021001 (2020), doi:[10.1088/2632-](https://doi.org/10.1088/2632-2153/abbf9a)
514 [2153/abbf9a](https://doi.org/10.1088/2632-2153/abbf9a).
- 515 [7] V. Mikuni and F. Canelli, *Point cloud transformers applied to collider physics*, *Mach. Learn.*
516 *Sci. Tech.* **2**(3), 035027 (2021), doi:[10.1088/2632-2153/ac07f6](https://doi.org/10.1088/2632-2153/ac07f6), [2102.05073](https://arxiv.org/abs/2102.05073).
- 517 [8] H. Qu, C. Li and S. Qian, *Particle Transformer for Jet Tagging*,
518 doi:[10.48550/arXiv.2202.03772](https://doi.org/10.48550/arXiv.2202.03772) (2022), [2202.03772](https://arxiv.org/abs/2202.03772).
- 519 [9] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, O. Mattelaer, H. S. Shao,
520 T. Stelzer, P. Torrielli and M. Zaro, *The automated computation of tree-level and next-to-*
521 *leading order differential cross sections, and their matching to parton shower simulations*,
522 JHEP **07**, 079 (2014), doi:[10.1007/JHEP07\(2014\)079](https://doi.org/10.1007/JHEP07(2014)079), [1405.0301](https://arxiv.org/abs/1405.0301).
- 523 [10] R. D. Ball *et al.*, *Parton distributions from high-precision collider data*, *Eur. Phys. J. C*
524 **77**(10), 663 (2017), doi:[10.1140/epjc/s10052-017-5199-5](https://doi.org/10.1140/epjc/s10052-017-5199-5), [1706.00428](https://arxiv.org/abs/1706.00428).
- 525 [11] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna, S. Pres-
526 tel, C. O. Rasmussen and P. Z. Skands, *An introduction to PYTHIA 8.2*, *Comput. Phys.*
527 *Commun.* **191**, 159 (2015), doi:[10.1016/j.cpc.2015.01.024](https://doi.org/10.1016/j.cpc.2015.01.024), [1410.3012](https://arxiv.org/abs/1410.3012).
- 528 [12] M. L. Mangano, M. Moretti, F. Piccinini, R. Pittau and A. D. Polosa, *ALPGEN, a gen-*
529 *erator for hard multiparton processes in hadronic collisions*, JHEP **07**, 001 (2003),
530 doi:[10.1088/1126-6708/2003/07/001](https://doi.org/10.1088/1126-6708/2003/07/001), [hep-ph/0206293](https://arxiv.org/abs/hep-ph/0206293).
- 531 [13] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens and M. Sel-
532 vaggi, *DELPHES 3, A modular framework for fast simulation of a generic collider experi-*
533 *ment*, JHEP **02**, 057 (2014), doi:[10.1007/JHEP02\(2014\)057](https://doi.org/10.1007/JHEP02(2014)057), [1307.6346](https://arxiv.org/abs/1307.6346).
- 534 [14] T. Aarrestad *et al.*, *The Dark Machines Anomaly Score Challenge: Benchmark Data and*
535 *Model Independent Event Classification for the Large Hadron Collider*, *SciPost Phys.* **12**(1),
536 043 (2022), doi:[10.21468/SciPostPhys.12.1.043](https://doi.org/10.21468/SciPostPhys.12.1.043), [2105.14027](https://arxiv.org/abs/2105.14027).
- 537 [15] G. Aad *et al.*, *Evidence for $t\bar{t}t\bar{t}$ production in the multilepton final state in proton–proton*
538 *collisions at $\sqrt{s} = 13$ TeV with the ATLAS detector*, *Eur. Phys. J. C* **80**(11), 1085 (2020),
539 doi:[10.1140/epjc/s10052-020-08509-3](https://doi.org/10.1140/epjc/s10052-020-08509-3), [2007.14858](https://arxiv.org/abs/2007.14858).
- 540 [16] G. Brooijmans *et al.*, *Les Houches 2019 Physics at TeV Colliders: New Physics Working*
541 *Group Report*, In *11th Les Houches Workshop on Physics at TeV Colliders: PhysTeV Les*
542 *Houches*, doi:[10.48550/arXiv.2002.12220](https://doi.org/10.48550/arXiv.2002.12220) (2020), [2002.12220](https://arxiv.org/abs/2002.12220).
- 543 [17] DarkMachines Community, *The 4tops dataset*, Zenodo, doi:[10.5281/zenodo.7277951](https://doi.org/10.5281/zenodo.7277951)
544 (2022).
- 545 [18] S. Haykin, *Neural Networks: A Comprehensive Foundation*, Prentice Hall PTR, USA, 1st
546 edn., ISBN 0023527617 (1994).
- 547 [19] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*,
548 doi:[10.48550/arXiv.1412.6980](https://doi.org/10.48550/arXiv.1412.6980) (2017), [1412.6980](https://arxiv.org/abs/1412.6980).

- 549 [20] T. Akiba, S. Sano, T. Yanase, T. Ohta and M. Koyama, *Optuna: A next-generation hyperpa-*
550 *rameter optimization framework*, doi:[10.48550/arXiv.1907.10902](https://doi.org/10.48550/arXiv.1907.10902) (2019), [1907.10902](https://arxiv.org/abs/1907.10902).
- 551 [21] G. Agarwal, L. Hay, I. Iashvili, B. Mannix, C. McLean, M. Morris, S. Rappoccio and
552 U. Schubert, *Explainable AI for ML jet taggers using expert variables and layerwise relevance*
553 *propagation*, JHEP **05**, 208 (2021), doi:[10.1007/JHEP05\(2021\)208](https://doi.org/10.1007/JHEP05(2021)208), [2011.13466](https://arxiv.org/abs/2011.13466).
- 554 [22] A. M. Sirunyan *et al.*, *Identification of heavy, energetic, hadronically decaying particles*
555 *using machine-learning techniques*, JINST **15**(06), P06005 (2020), doi:[10.1088/1748-](https://doi.org/10.1088/1748-0221/15/06/P06005)
556 [0221/15/06/P06005](https://doi.org/10.1088/1748-0221/15/06/P06005), [2004.08262](https://arxiv.org/abs/2004.08262).
- 557 [23] H. Qu and L. Gouskos, *ParticleNet: Jet Tagging via Particle Clouds*, Phys. Rev. D **101**(5),
558 056019 (2020), doi:[10.1103/PhysRevD.101.056019](https://doi.org/10.1103/PhysRevD.101.056019), [1902.08570](https://arxiv.org/abs/1902.08570).
- 559 [24] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein and J. M. Solomon, *Dynamic graph*
560 *cnn for learning on point clouds*, doi:[10.48550/arXiv.1801.07829](https://doi.org/10.48550/arXiv.1801.07829) (2019), [1801.07829](https://arxiv.org/abs/1801.07829).
- 561 [25] C. Li, H. Qu, S. Qian, Q. Meng, S. Gong, J. Zhang, T.-Y. Liu and Q. Li,
562 *Does Lorentz-symmetric design boost network performance in jet physics?*,
563 doi:[10.48550/arXiv.2208.07814](https://doi.org/10.48550/arXiv.2208.07814) (2022), [2208.07814](https://arxiv.org/abs/2208.07814).
- 564 [26] Hqucms/weaver, *Streamlined Neural Network Training*, GitHub repository, Available
565 online: <https://github.com/hqucms/weaver>.
- 566 [27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser and
567 I. Polosukhin, *Attention is all you need*, doi:[10.48550/arXiv.1706.03762](https://doi.org/10.48550/arXiv.1706.03762) (2023), [1706.](https://arxiv.org/abs/1706.03762)
568 [03762](https://arxiv.org/abs/1706.03762).
- 569 [28] Jet-universe/particle_transformer, *Official implementation of 'Particle Transformer for Jet*
570 *Tagging'*, GitHub repository, Available online: [https://github.com/jet-universe/particle_](https://github.com/jet-universe/particle_transformer)
571 [transformer](https://github.com/jet-universe/particle_transformer).
- 572 [29] J. Lee, Y. Lee, J. Kim, A. R. Kosior, S. Choi and Y. W. Teh, *Set trans-*
573 *former: A framework for attention-based permutation-invariant neural networks*,
574 doi:[10.48550/arXiv.1810.00825](https://doi.org/10.48550/arXiv.1810.00825) (2019), [1810.00825](https://arxiv.org/abs/1810.00825).
- 575 [30] Y. Tang and D. Ha, *The sensory neuron as a transformer: Permutation-invariant neural net-*
576 *works for reinforcement learning*, doi:[10.48550/arXiv.2109.02869](https://doi.org/10.48550/arXiv.2109.02869) (2021), [2109.02869](https://arxiv.org/abs/2109.02869).
- 577 [31] T.-Y. Lin, P. Goyal, R. Girshick, K. He and P. Dollár, *Focal loss for dense object detection*,
578 doi:[10.48550/arXiv.1708.02002](https://doi.org/10.48550/arXiv.1708.02002) (2018), [1708.02002](https://arxiv.org/abs/1708.02002).
- 579 [32] M. Drees, M. Shi and Z. Zhang, *Machine learning optimized search for the z' from*
580 *$u(1)_{L_\mu-L_\tau}$ at the lhc*, doi:[10.48550/arXiv.2109.07674](https://doi.org/10.48550/arXiv.2109.07674) (2022), [2109.07674](https://arxiv.org/abs/2109.07674).
- 581 [33] A. Butter, G. Kasieczka, T. Plehn and M. Russell, *Deep-learned Top Tagging with a Lorentz*
582 *Layer*, SciPost Phys. **5**(3), 028 (2018), doi:[10.21468/SciPostPhys.5.3.028](https://doi.org/10.21468/SciPostPhys.5.3.028), [1707.08966](https://arxiv.org/abs/1707.08966).
- 583 [34] M. Erdmann, E. Geiser, Y. Rath and M. Rieger, *Lorentz Boost Networks: Autonomous*
584 *Physics-Inspired Feature Engineering*, JINST **14**(06), P06006 (2019), doi:[10.1088/1748-](https://doi.org/10.1088/1748-0221/14/06/P06006)
585 [0221/14/06/P06006](https://doi.org/10.1088/1748-0221/14/06/P06006), [1812.09722](https://arxiv.org/abs/1812.09722).
- 586 [35] A. Bogatskiy, B. Anderson, J. T. Offermann, M. Roussi, D. W. Miller and R. Kondor, *Lorentz*
587 *Group Equivariant Neural Network for Particle Physics*, doi:[10.48550/arXiv.2006.04780](https://doi.org/10.48550/arXiv.2006.04780)
588 (2020), [2006.04780](https://arxiv.org/abs/2006.04780).

- 589 [36] S. Gong, Q. Meng, J. Zhang, H. Qu, C. Li, S. Qian, W. Du, Z.-M. Ma and T.-Y. Liu, *An*
590 *efficient Lorentz equivariant graph neural network for jet tagging*, JHEP **07**, 030 (2022),
591 doi:[10.1007/JHEP07\(2022\)030](https://doi.org/10.1007/JHEP07(2022)030), [2201.08187](https://arxiv.org/abs/2201.08187).