

Jet energy calibration with deep learning as a Kubeflow pipeline

Daniel Holmberg¹, Dejan Golubovic² and Henning Kirschenmann³

¹Department of Computer Science, University of Helsinki, 00560 Helsinki, Finland.

²CERN, 1211 Geneva 23, Switzerland.

³Helsinki Institute of Physics, 00560 Helsinki, Finland.

Contributing authors: daniel.holmberg@helsinki.fi; dejan.golubovic@cern.ch; henning.kirschenmann@cern.ch;

Abstract

Precise measurements of the energy of jets emerging from particle collisions at the LHC are essential for a vast majority of physics searches at the CMS experiment. In this study, we leverage well-established deep learning models for point clouds and CMS open data to improve the energy calibration of particle jets. To enable production-ready machine learning based jet energy calibration an end-to-end pipeline is built on the Kubeflow cloud platform. The pipeline allowed us to scale up our hyperparameter tuning experiments on cloud resources, and serve optimal models as REST endpoints. We present the results of the parameter tuning process and analyze the performance of the served models in terms of inference time and overhead, providing insights for future work in this direction. The study also demonstrates improvements in both flavor dependence and resolution of the energy response when compared to the standard jet energy corrections baseline.

Keywords: LHC, CMS, Open Data, Jet Energy, Kubeflow, MLOps, GNN

1 Introduction

The adoption of machine learning methods has had a profound impact on the field of high energy physics, greatly increasing the discovery potential in the data measured by the particle detectors at the Large Hadron Collider (LHC) [1]. Deep learning especially has proven very useful [2] with graph neural networks being one of the most expressive and versatile architectures to choose for many tasks [3] ranging from reconstructing particle tracks [4] to classifying complete events [5].

In this paper, we study the application of deep learning for calibrating the energy of particle jets

at the Compact Muon Solenoid (CMS) experiment [6]. Jets in this context originate from high energy proton-proton collisions producing color charged partons that undergo hadronization forming collimated sprays of color neutral particles. Calibrating the energy of jets is an involved process, split into several factorized steps, some based only on simulations and some on comparisons with data [7]. A precise calibration of the jet energy scale is crucial for a wide variety of physics analyses, most prominently e.g. measurements of the top quark mass [8] and inclusive jet cross-section measurements.

Deep learning has been successfully applied by the CMS collaboration in the past to increase the

energy resolution of bottom jets with a feedforward neural network [9]. These efforts are here extended to all jet flavors in a QCD-jet data sample publicly accessible on the CERN OpenData portal [10]. Additionally, by adopting recent advancements in representation learning, specifically ones made for jet classification [11, 12], more information about jet constituents can be included in the training process, which has proven beneficial for jet calibration [13].

Furthermore, since operationalizing machine learning workflows is a challenge in itself for many organizations [14], we introduce a cloud native pipeline for running jet energy calibration experiments. It runs on the Kubeflow platform [15] that comes with readily available components for hyperparameter tuning and model serving among others. Kubeflow has been used by researchers in various domains such as bioinformatics to achieve rapid scaling with containers [16], or as a means to create automated machine learning workflows for a service-aware 5G network model adapting to drift in the input data [17]. Adopting a cloud native workflow enables the workload to be smoothly deployed also on public cloud resources, as was done for fast simulation of electromagnetic showers using generative deep learning at CERN [18].

The rest of this paper is structured as follows. In Section 2, we explain the contents of the CMS open dataset used for this study. In Section 3, we go through the data distribution, feature sets and models used to calibrate jet energy. Section 4 introduces the Kubeflow pipeline used for training and serving our models on internal cloud resources. Section 5 shows the results that our models yield, and lastly in Section 6 we draw some final conclusions.

2 Dataset definition and conventional jet energy calibration

In this study, we utilize a dataset prepared in the context of the CMS OpenData effort [10]. The dataset consists of particle jets extracted from simulated proton-proton (pp) collision events at $\sqrt{s} = 13$ TeV. These events are generated at leading-order perturbative QCD with

PYTHIA 8 [19], and include CMS detector simulation and event reconstruction.

“Particle-level jets” or “generator jets” are clustered using the anti- k_T algorithm [20] with radius parameter of 0.4 from stable (decay length $c\tau > 1$ cm) final-state particles resulting from the hadronization of partons originating from the pp collisions. As these particles propagate through the detector, they leave signals in detector components such as the tracker and the electromagnetic and hadronic calorimeter. The modeling of the interaction with the material and detector response relies on the CMS Full simulation, based on GEANT4 [21].

The types of quarks or gluons that initiate the formation of the jet determines the flavor of a jet. Flavor labeling is done using a technique called “ghost association” [22], where heavy flavor hadrons and light quark and gluon partons are added as “ghost particles” to the clustering process. If heavy flavor hadron ghosts are found in the jet, it is labeled as a heavy flavor (b or c) jet. If no ghost hadron is found, the jet is checked for light flavor (uds) or gluon (g) partons, identifying it as a corresponding jet [23].

The Particle Flow (PF) approach at CMS [24] is a method that attempts to reconstruct each particle in the event individually, prior to the jet clustering, based on information from all relevant sub-detectors, resulting in a list of “PF candidates”. Various methods for per-particle pileup (additional pp collisions occurring in the same bunch crossing as the event of interest) mitigation can be applied [25]. The charged hadron subtraction (CHS) is the default for narrow radius jets in Run 2 CMS data, meaning that charged particles associated with pileup vertices are removed prior to jet clustering, thereby reducing the impact of pileup on jets. The remaining reconstructed PF candidates are then used as input to jet clustering algorithms to form a “reconstructed jet” (anti- k_T , $R=0.4$) that is supposed to be as close to the particle-level jet in terms of kinematic quantities as possible.

The aim of jet energy corrections is to correct the energy of the reconstructed jets—on average—back to the jet energy of particle-level jets. The measured jet energy is affected by various effects, such as energy loss in the detector material, non-linear response of the detector, and pileup. CMS factorizes the jet energy corrections into levels

to individually correct for various effects [7]: The L1 correction corrects for offset energy induced by pileup. It is determined from simulated samples with/without pileup, parametrizing the offset energy as a function of median event energy density, jet area, p_T , and η . The L2L3 correction corrects for remaining detector response dependence and is parametrized as a function of jet p_T and η . The L2L3Residual correction corrects for any remaining differences between experimental data and simulation, but is not applicable in this MC-only study.

These conventional jet energy corrections do not take into account the substructure of the jets. The PF approach already leads to a reduced dependence of, e.g. the jet response on the flavor of a jet in comparison to calorimeter-only reconstruction. However, taking the substructure of jets into account promises potential for reduced flavor response differences and associated uncertainties as well as an improved jet-energy resolution.

3 Jet energy regression

Standard jet energy corrections can be further improved upon by using supervised machine learning. An important step in this direction has been taken by the CMS Collaboration for b jets specifically [9]. For these jets a significant part of the jet energy is carried by semileptonically decaying b-hadrons, decaying to charged leptons and neutrinos. Neutrinos escape detection by the CMS detector since they only interact via the weak force leading to jet energy being underestimated. However, using a neural network trained on a sample of simulated top quarks event decaying into b jets and W bosons significant improvements in the energy resolution for b jets were achieved.

These efforts can be generalized to other jet flavors too. Using a QCD sample such as the one described in Section 2 enables the training of a regression model for multiple jet flavors. This approach can potentially address any flavor discrepancies in the energy response. Notably, the response of light quark jets and gluon jets can vary significantly due to the higher color charge of gluon jets, which results in more and softer particles with a lower calorimeter response. Although the PF algorithm can reduce the response difference substantially by replacing the non-linear calorimeter measurement of

charged hadron energy with the corresponding track momentum, 15% of jet energy is still carried by neutral hadrons subject to the calorimeter response non-linearity [24].

3.1 Data distribution and input features

The jets in the CMS open dataset span a large spectrum both in terms of p_T and pseudorapidity η as seen in Figure 1. However, the very low p_T region experiences high pileup, resulting in lower-quality jets. Additionally, the forward region of the detector lacks tracking information leading to a worse event reconstruction quality and less reliable measurements to train on. To address these issues, the dataset is filtered by $p_T^{\text{gen}} > 20$ GeV and $|\eta| < 2.5$. In total 1.42M jets were used, 60% of which were allocated to the training set. The model was validated at the end of every training epoch on a separate set with 20% of all jets. Once the training had finished, model performance was evaluated on a test set with the remaining 20% of all jets.

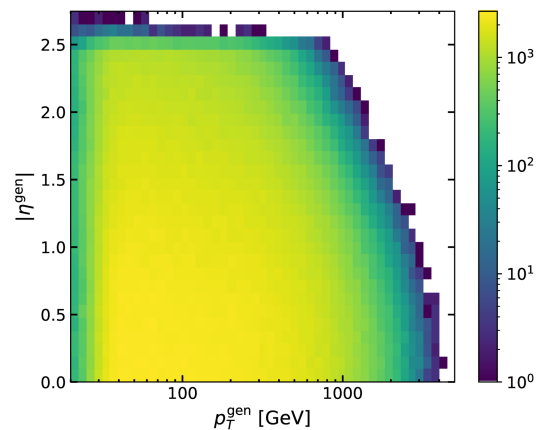


Fig. 1: A heatmap illustrating the distribution of jets in the dataset with respect to the generated transverse momentum (p_T^{gen}) and the absolute value of generated pseudorapidity ($|\eta^{\text{gen}}|$)

Table 1 presents the features used for training our regression model selected from a comprehensive list of variables available in the CMS open dataset [10]. Eight of the features describe jets as a whole. The reconstructed p_T is log-transformed to

reduce the width of the otherwise long-tailed distribution. In addition to p_T , jet coordinates η and ϕ as well as jet mass and catchment area are also included in the set of training features. The last three jet-level variables $p_T D$, σ_2 and multiplicity can help to discriminate between quark and gluon jets [26].

Every particle within a jet has six associated features. The two most statistically significant for the regression task at hand are: 1) the log-transformed p_{T_i} of a PF candidate indexed with i , and 2) the same variable, but relative to the p_T of the whole jet, also log-transformed. Additionally, each PF candidate has four location-based features: the η_i , ϕ_i and θ_i detector coordinates, along with the distance R_i from the particle to the center of the jet.

3.2 Regression target and loss function

The aim of the regression is to train a model to map a set of features describing a reconstructed jet towards the transverse momentum of the corresponding generator-level jet. However, by itself the particle-level p_T^{gen} follows a decreasing exponential distribution that covers many orders of magnitude in the energy spectrum as seen in Figure 1. To counteract this, p_T^{gen} can be divided by the similarly distributed p_T^{reco} that is part of the training set. This gives a target distribution on the order of one with a reduced variance compared to the original target distribution [9]. To further correlate the target with the input features the logarithm is taken yielding the final regression target as $y = \log(p_T^{\text{gen}}/p_T^{\text{reco}})$. The distribution of the target is narrow and centralized around zero as seen in Figure 2. In order to get the corrected transverse momentum p_T^{corr} , the exponential function is applied to the prediction \hat{y} , and the resulting correction factor is multiplied by p_T^{reco} . The per-jet energy response can then be defined as $R = p_T^{\text{corr}}/p_T^{\text{gen}}$.

The mean absolute error (MAE) is selected as the loss function to minimize for this problem, as it assigns less importance to outliers compared to the more commonly used mean squared error (MSE) loss. In line with the previous study on b jet energy regression [9], a loss function with reduced sensitivity to the tails of the target distribution is

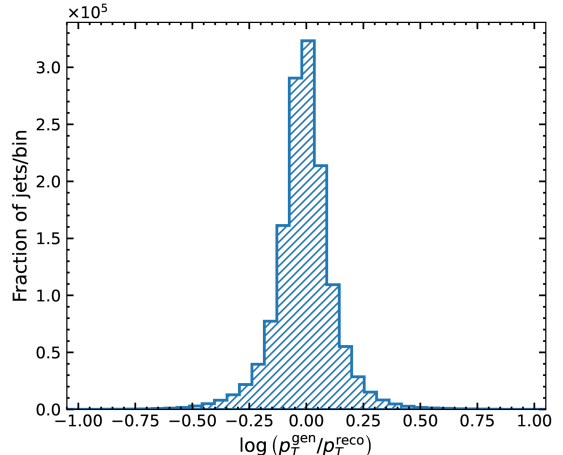


Fig. 2: Distribution of the regression target

preferred. To prevent potential spikes in the training loss, jets with a target value smaller than -1 or larger than 1 are excluded, as they are considered too poorly reconstructed to be taken into account. The loss function used is thus:

$$L = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| I_{|y_i| < 1}. \quad (1)$$

A further motivation for the use of the MAE loss function is that the statistic it learns is the median of the target distribution, whereas for example the minimum of the MSE loss lies on the function that maps input features to the expected value of the target. Predicting the median in this case can be seen as beneficial since it is a robust measure of central tendency. By learning the median of the target distribution, the model is more likely to make accurate predictions even in the presence of outliers.

3.3 Deep learning models

During recent years many new approaches of applying deep learning in jet physics have emerged, especially for the purpose of jet tagging where the aim is to classify jets based on the particle initiating them. Some proposed approaches to do this is for example to treat the jets as images [27], sequences [28] or trees [29]. While

Table 1: Overview of the input features used for training the regression model, categorized into jet-level features and PF candidate features. Jet features describe the overall properties of the jet, while PF features characterize the individual particles within the jet

Category	Variable	Description
Jet features	$\log p_T$	Logarithm of a jet's p_T
	η	Pseudorapidity of a jet
	ϕ	Azimuthal angle of a jet
	m	Mass of a jet
	A	Catchment area of a jet
	$p_T D$	Fragmentation distribution for a jet's p_T
	σ_2	Minor ellipse axis of a jet
	multiplicity	Jet constituent multiplicity
PF features	$\log p_{T_i}$	Logarithm of a particle's p_T
	$\log \frac{p_{T_i}}{p_T}$	Logarithm of the fractional p_T of a particle
	η_i	Pseudorapidity of a particle
	ϕ_i	Azimuthal angle of a particle
	θ_i	Polar angle of a particle
	R_i	Distance from a particle to the center of the jet

these methods perform well and surpass traditional multivariate methods, the way they represent jet constituents is not ideal. A particle jet can contain up to $\mathcal{O}(100)$ particles whereas an image of a jet will contain $\mathcal{O}(1000)$ pixels, and as noted in [27] the images are indeed very sparse with 5-10% of pixels being active. When instead considering a sequence or tree of particles as the representation a notable issue that arises is that the particles must be ordered in some fashion to be used by the deep learning model, e.g. a recurrent neural network or recursive neural network. However, the constituents of a particle jet have no intrinsic ordering.

More natural ways of representing particles have been found by adopting point cloud based formalisms from the wider machine learning community as particle clouds [12]. This kind of representation treats a collection of particles as a graph structure $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where each individual particle serves as a node $\mathcal{V} = \{1, \dots, n\}$ with potential edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ connecting them. In this work, we compare two separate models that operate on data represented in this way.

The simplest case of representation learning on particle clouds is when the set of edges is empty $\mathcal{E} = \emptyset$. This will lead to the particles taking on the form of an unordered set, and was originally

proposed for the Particle Flow Network (PFN) model [11] adapting from the Deep Sets [30] framework. The key idea here is that input feature vectors $\mathbf{x}_i \in \mathbb{R}^F$ are mapped with an equivariant function into a latent feature vector $\mathbf{h}_i = \psi(\mathbf{x}_i)$. In practice that would be a multilayer perceptron (MLP) with shared weights for all elements of the set.

To make predictions for the particle jet as a whole the latent feature vectors must be aggregated using a permutation invariant pooling, such as summing, averaging or taking the max value. Following the PFN and Deep Sets papers, summation $\sum_{i \in \mathcal{V}}$ is chosen as the global pooling operation. Figure 3 (b) in conjunction with Figure 3 (c) show the complete network architecture where the output of the Deep Sets block connects to global pooling in the network head. The global particle representation is concatenated with jet features before being passed into a final MLP mapping towards the regression target. Note that the rectified linear unit (ReLU) [31] is used as activation function, and dropout [32] is applied in the head of the network.

Spatial information can be used to further increase the expressivity of a point cloud based model. The ParticleNet [12] architecture uses edge convolution (EdgeConv), first introduced as a

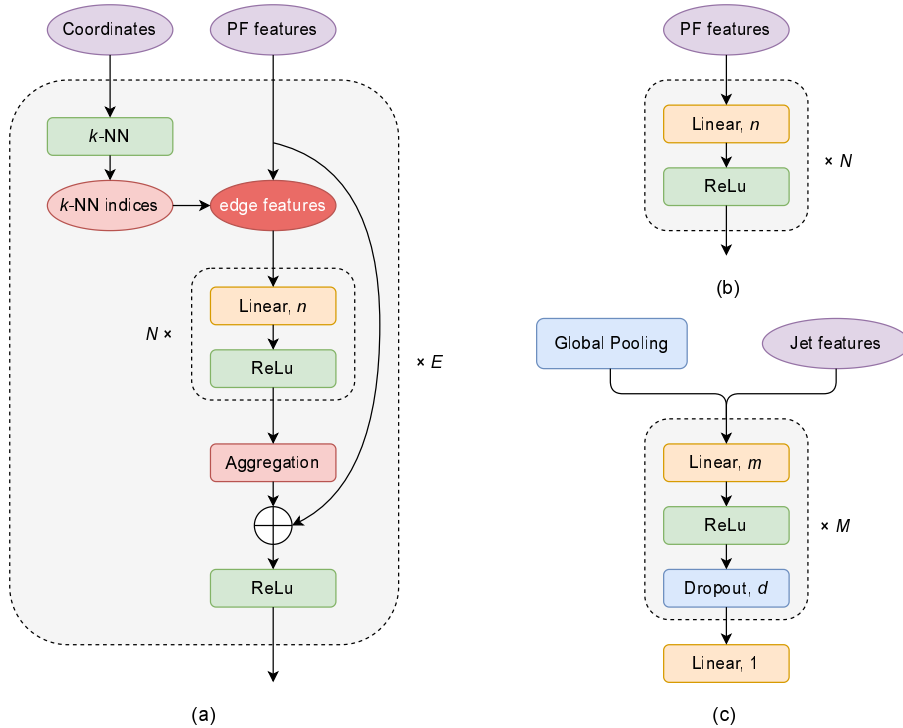


Fig. 3: Illustration of model architectures: (a) EdgeConv block for ParticleNet, (b) Deep Sets block for PFN, and (c) Network head shared by both models

building block of dynamic graph convolutional neural networks (DGCNN) [33], to incorporate information on the local neighborhood of each particle. Detector coordinates in the (η, ϕ) -plane are used to calculate the Euclidean distance matrix from pairwise distances between particles. The k -NN algorithm is then applied to construct edges connecting each particle to its k nearest neighboring particles.

Messages between every point \mathbf{x}_i and its neighbors \mathbf{x}_j are learned using an asymmetric edge function $\psi(\mathbf{x}_i, \mathbf{x}_j - \mathbf{x}_i)$ implemented as an MLP with shared weights. Permutation invariant aggregation in the form of averaging $\frac{1}{k} \sum_{i \in \mathcal{N}_i^k}$ over the learned edge features for the k nearest neighbors is used to update every node in the particle cloud. A shortcut connection [34] from the original node features is added to the output of the aggregation before being passed through the ReLU activation function. This concludes the EdgeConv block shown in Figure 3 (a).

If multiple EdgeConv blocks are stacked after one another the input graphs are dynamically updated by calculating the pairwise distance

matrix from the latent feature space learned by the previous block. Global average pooling $\frac{1}{n} \sum_{i \in \mathcal{V}}$ is applied on the output of the last EdgeConv block as it is passed into the network head in Figure 3 (c). An identical procedure to that in the PFN model is applied, where the jet features are concatenated with the pooled particle features, and finally passed through one last MLP mapping towards the regression target.

The models are implemented in PyTorch [35] as part of the weaver deep learning framework for high energy physics [36]. Weaver supports handling common particle physics data formats such as ROOT [37] or Awkward Array [38], as well as distributed training, and model inference. To scale up the training on cluster resources where GPUs are distributed over separate nodes the code must support collective communications to synchronize gradients over machines. Different backends can be chosen in PyTorch for this purpose such as Message Passing Interface (MPI) [39] for CPU parallelization, or NVIDIA Collective Communication Library (NCCL) [40] for multi-GPU communication.

4 KubeFlow pipeline

The analysis was carried out on the KubeFlow-based machine learning platform at CERN [15]. KubeFlow is an open-source machine learning toolkit that supports the entire machine learning lifecycle by providing features such as notebooks, pipelines, hyperparameter optimization, distributed training, model serving, and model monitoring. KubeFlow is built on top of Kubernetes [41], a container orchestrator, leveraging the scalability, ease of use and integration of cutting-edge infrastructure technologies. Deployed as a set of Kubernetes resources, KubeFlow application code runs as a collection of micro-services that communicate with each other and process user workloads.

Machine learning workflows typically take the form of a directed acyclic graph that begins with data processing, proceeds through model training, and ends with an inference phase of the trained model. KubeFlow facilitates developing such workflows by offering several features, including a web interface for managing, tracking, and running pipelines, an engine for scheduling pipeline steps, a software development kit (SDK) for defining and running pipelines using Python, and higher-level abstraction tools like KALE [42], which convert notebooks to pipelines. Regardless of how a pipeline is defined, it is always converted to a Kubernetes YAML [43] definition file before being submitted for execution. A pipeline runs as a sequence of pods (containers), with each pipeline step waiting for its dependencies to complete successfully before proceeding. The pipeline developed for this study is shown as part of the KubeFlow user interface (UI) in Figure 4.

Pipelines offer benefits in resource utilization. By allowing users to define hardware requirements for each step, pipelines ensure that GPUs are only utilized when needed, for example during training and inference steps. Other steps that use CPU-only resources make GPUs available for other users in the cluster. Additional features include support for pipeline scheduling which enables automatic execution of repeated or periodic workflows, running pipelines with different input parameters without any code changes, and grouping pipeline runs into experiments making it easier to track and compare similar runs.

4.1 AutoML experiment

The KubeFlow Katib component [44] offers a streamlined process for automated machine learning (AutoML) supporting hyperparameter tuning, early stopping and neural architecture search. Here we use hyper-parameter optimization, that includes three main steps: 1. implementing a script that trains a model and takes hyperparameters as command line arguments, 2. building a docker image with all dependencies to run the script, and 3. specifying a YAML file with the definition of the hyper-parameters. The YAML file defines the search algorithm, early stopping options, maximum number of parallel jobs, hardware resources, and other options. The YAML file can be generated manually, using an SDK, or using the high-level KALE tool.

Katib schedules search jobs in the form of multiple trials, each trial corresponding to a unique combination of hyperparameters as seen in Figure 5. A trial can be a Kubernetes job running a script in a single pod, a pipeline where multiple pods run sequentially or a distributed training job where multiple pods run synchronously utilizing multiple GPUs to train the model. Being framework-agnostic, KubeFlow supports running search jobs with any machine learning framework.

In this study, Katib trials were executed using the PyTorchJob Kubernetes custom resource [45]. The KubeFlow training operator component provides distribution on the level of containerization; training a single model using multiple cluster GPUs that can be located in different machines. In addition to implementing the training to use distribution strategies, it is necessary to specify a YAML definition of a distributed job. The YAML definition includes many attributes, including the number of worker replicas, memory, and command line arguments to the machine learning code. The PyTorchJob also includes S3 credentials to access the training data stored in a bucket on CERN object storage.

To optimize the performance of ParticleNet and PFN, Katib was configured to use the Random Search algorithm [46] to find the optimal set of parameters for minimizing the test loss. The total number of trials was set to 60 for both models with 10 trials running in parallel using one GPU each to achieve fast scheduling of pods on the cluster. Each trial ran for 50 epochs, with the batch

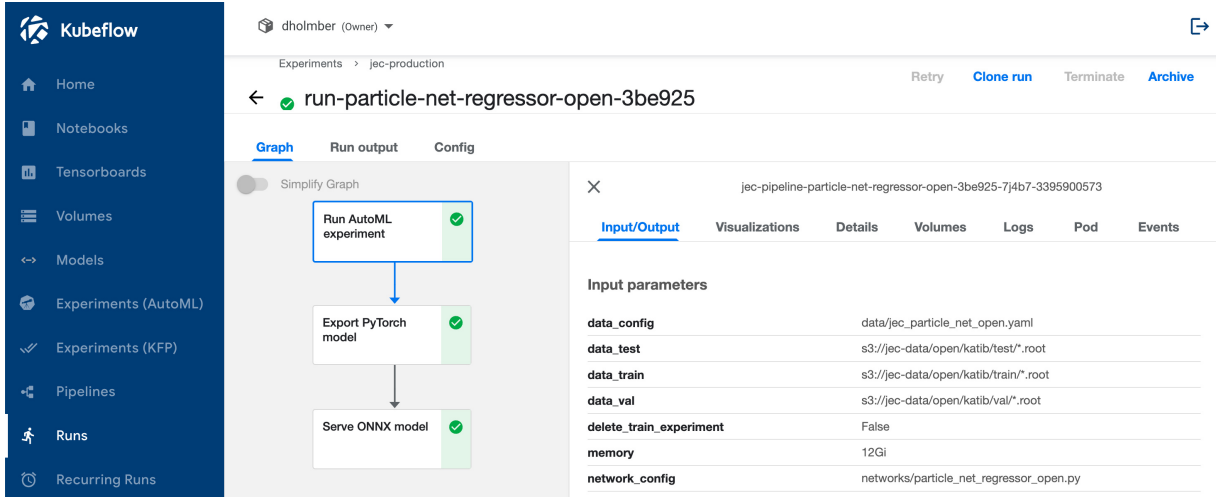


Fig. 4: Kubeflow UI for a jet energy regression pipeline run. The pipeline consists of three steps: 1. hyperparameter tuning using Kubeflow’s AutoML component Katib, 2. exporting the optimal PyTorch model to the ONNX format, and 3. serve the exported model over HTTP with KServe

size set to 500. The initial learning rate is set as part of the hyperparameter optimization process. After 70% of epochs, 35 in this case, a scheduler starts decreasing the learning rate exponentially on a per-epoch basis down to 1% of the initial value at the end of training. The model with the lowest validation loss is chosen as the final model for each trial, and is then run on the evaluation set to get the test loss that we are trying to minimize.

With reference to Figure 3 (b) and 3 (c), the search space for the hyperparameter tuning was defined as follows for PFN:

- linear layers: $N \in \{1, 2, 3, 4, 5\}$
- linear layer units: $n \in \{50, 100, 200, 400\}$
- linear layers: $M \in \{1, 2, 3, 4, 5\}$
- linear layer units: $m \in \{50, 100, 200, 400\}$.

The search space for ParticleNet, with respect to Figure 3 (a) and 3 (c), was defined as:

- EdgeConv blocks: $E \in \{1, 2, 3\}$
- nearest neighbors: $k \in \{4, 8, 16\}$
- linear layers: $N \in \{1, 2, 3\}$
- linear layer units: $n \in \{50, 100, 200\}$
- linear layers: $M \in \{1, 2, 3\}$
- linear layer units: $m \in \{50, 100, 200\}$.

Lastly, the mutual hyperparameters considered were:

- dropout rate: $d \in [0; 0.5]$

- initial learning rate: $lr \in [10^{-5}; 10^{-2}]$
- optimizer: $optim \in \{\text{AdaGrad, Adam, AdamW, Ranger, RMSProp}\}$.

The completion of trials can be viewed directly from the Katib UI, while Kubeflow’s TensorBoard component allows tracking the progress of individual training runs. To set up a TensorBoard server, it is necessary to specify a model output path, either on a persistent volume claim (PVC) within the cluster or an S3 object storage endpoint. As long as the model training is writing to the specified location, the model performance can be monitored in real-time using TensorBoard servers.

4.2 Exporting the optimal model

Once the optimal hyperparameters for both PFN and ParticleNet models have been determined through the hyperparameter tuning process, the optimal PyTorch model is exported to the ML framework agnostic Open Neural Network Exchange (ONNX) format [47]. This enables seamless integration with other ML tools, such as NVIDIA Triton Inference Server [48] for model serving, and eases the deployment process.

The second step of the Kubeflow pipeline involves running a PyTorchJob to carry out this conversion. It retrieves the optimal PFN and ParticleNet models from the S3 bucket, converts

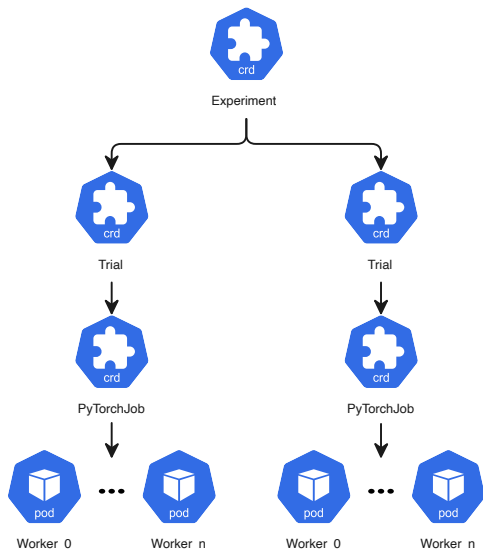


Fig. 5: The structure of an Experiment Kubernetes custom resource. An AutoML experiment consists of multiple trials, with each trial representing a unique combination of hyperparameters. Each trial monitors a PyTorchJob, which can submit one or more workers to train the model. Once completed, the Experiment resource retains the outcomes of all trials and the optimal trial result based on predefined metrics

them to ONNX format, and stores the resulting ONNX models back into the same S3 bucket. The PyTorchJob can be defined using a YAML file that specifies the necessary configurations, such as the PyTorch model input path and the ONNX model output path, network configuration, and hardware resources to run the export job.

When exporting a model to ONNX, a configuration file is created alongside the model file to facilitate serving the model using Triton. We made a schema in the Protocol Buffers (protobuf) [49] format with model input/output dimensions, data types (32-bit float) and graph optimization level for ONNX Runtime. This was compiled into a Python file that can be used to automatically generate model configuration in protobuf text format with the correct input and output dimensions when exporting a model in PyTorch.

ONNX Runtime defines a static computational graph for the model as opposed to the dynamic one used by PyTorch during training, which allows for various graph optimizations that

can improve inference performance, such as graph-level transformations, node eliminations, node fusions, and layout optimizations. An extended optimization level is available that enables complex node fusions. However, these optimizations were found to cause issues when serving ParticleNet, and as a result they were only applied to PFN. A more basic graph optimization level with semantics-preserving graph rewrites that removes redundant nodes and redundant computation was chosen for ParticleNet.

To allow Triton to accept dynamically varying batch sizes, we configured the maximum batch size to be 100k in the model configuration file. Batch requests that large are not necessarily recommended due to the spiky network load they would produce and the excessive amount of memory that must be allocated to the inference server.

The export job produces an output directory structure, as shown in Figure 6, that follows Triton’s specifications. The base model repository is in our case an S3 bucket path and a unique id for every pipeline run. The top-level repository can contain many subdirectories (or pseudo-folders since object storage has a flat address space), each representing distinct models. The optimal ONNX model is placed in a numeric sub-folder signifying model version during exportation, and the automatically generated model configuration file is placed alongside that folder.

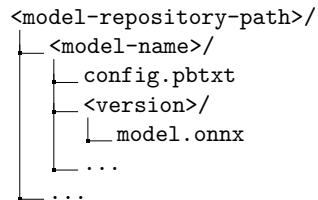


Fig. 6: Model repository layout for Triton Inference Server with ONNX backend

4.3 Model serving

After exporting the optimal PFN and ParticleNet PyTorch models to the ONNX format and storing them in an S3 bucket, the models are served using custom InferenceService resources. An InferenceService is the interface used for deploying models on KubeFlow’s inference platform KServe [50].

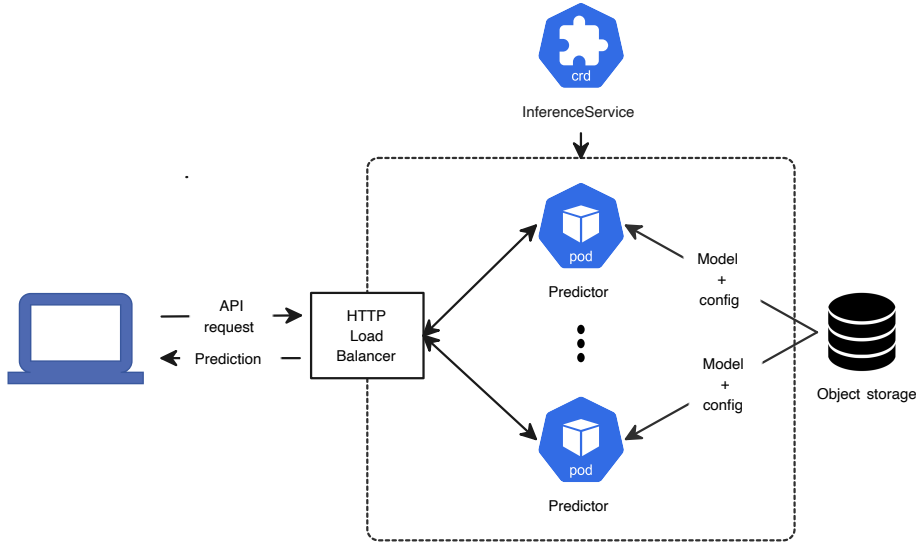


Fig. 7: A diagram depicting model serving using KServe, highlighting load balancing of user inference requests and the scalability of predictor pods

The InferenceService can be specified using a YAML file with various configuration options, such as hardware allocation for the server (CPU, GPU and memory), runtime version of the predictor, and the path to the model repository. In order to make authenticated requests to S3 storage, a Kubernetes ServiceAccount with the required access rights was deployed on the cluster and attached to every InferenceService.

KServe relies on Knative [51] for scaling serverless workloads and supports scale-to-zero, optimizing cost efficiency. Istio [52], another key technology in KServe, acts as a service mesh that uses Kubernetes sidecars (containers deployed alongside a main container in a pod) for network traffic management, providing features such as progressive “canary” model rollouts, traffic routing, ingress management, logging, load balancing, and security.

Triton is utilized as the predictor for the InferenceService. It is an open-source inference server capable of serving multiple models concurrently, supporting various machine learning frameworks. We have specified the ONNX Runtime as platform in the model configuration file telling Triton explicitly which backend to use.

The complete inference workflow is illustrated in Figure 7. The InferenceService creates Triton

pods that retrieves the ONNX model and configuration from S3 object storage. The pods act as REST endpoints [53] that can be queried over HTTP. When user sends inference requests, a load balancer is responsible for receiving them and distribute them to available inference server pods. The pods can scale up or down dynamically based on the volume of incoming requests. Servers pass the input data through the deep learning model and return model output to the load balancer that sends it back to the user.

5 Results

5.1 Hyperparameter optimization

The results of the hyperparameter tuning using Random Search can be analyzed to some extent with Pearson’s correlation coefficient. Table 2 shows the correlation of all continuous and ordinal hyperparameters with the inverted test loss for the models with the lowest validation loss during each training run. Because the choice of hyperparameters is stochastic for every trial it is difficult to isolate the impact of a single hyperparameter on the model’s performance, and thus particularly high correlation scores are not expected here. Note also that the correlation is limited to the search space laid out in Section 4.1.

Table 2: A two-fold presentation of the results from hyperparameter tuning. The upper part of the table shows the Pearson correlation between hyperparameters and the inverted test loss for PFN and ParticleNet. The bottom part lists the top three best sets of hyperparameters found for both models with the corresponding test loss values for those trials. Gain signifies the relative improvement in loss for each trial compared to the standard corrections baseline ($L_{\text{Baseline}} = 9.427\text{e-}2$), and is computed as $\Delta L/L_{\text{Model}}$

	Model	E	k	N	n	M	m	d	lr	$optim$	L	Gain [%]
$\rho_{hp,1/L}$	PFN			0.27	-0.21	-0.08	0.33	-0.54	-0.39			
	ParticleNet	0.22	-0.04	0.17	-0.21	0.02	0.40	-0.39	-0.34			
top trials	PFN			5	50	3	400	9.45e-3	1.08e-3	Ranger	8.768e-2	7.52
				5	50	4	100	4.13e-2	7.60e-4	Adam	8.770e-2	7.49
				3	200	4	400	8.58e-2	7.33e-4	RMSProp	8.770e-2	7.49
	ParticleNet	3	16	3	50	3	200	1.16e-2	8.72e-3	Ranger	8.746e-2	7.79
		3	16	2	100	2	100	1.14e-1	1.84e-3	Ranger	8.752e-2	7.71
		3	8	2	50	3	200	8.70e-2	2.50e-3	AdamW	8.755e-2	7.68

The table suggests that for PFN, having more linear layers with fewer units in the Deep Sets block is weakly associated with a lower loss. For ParticleNet, having more EdgeConv blocks with additional linear layers and fewer units offers an advantage. The number of nearest neighbors k in ParticleNet’s particle graph appears uncorrelated with a lower loss. The remaining correlation values exhibit similar behavior for both models. The loss is relatively unaffected by the number of linear layers in the network head, but more units in these layers tend to yield a lower loss. Increased dropout negatively impacts the regression task with the most certainty among all correlation results. Finally, a lower learning rate tends to produce better results.

The best hyperparameters found for both models, as listed in the lower section of Table 2, do tend to align with the Pearson correlation scores. For the PFN model, the top three trials all used a configuration with more linear layers (3–5) in the Deep Sets block and fewer units (50–200), which aligns with the correlation scores observed for these parameters. Similarly for ParticleNet, the top three trials incorporate more EdgeConv blocks (3) with a higher number of linear layers (2–3) and fewer units (50–100), reflecting the corresponding correlations. The initial learning rate and dropout across the top trials are, with exception for the optimal ParticleNet trial’s learning rate, notably low as suggested by the Pearson correlation.

To further highlight the impact of poorly adjusted dropout and learning rate we can compare the average of those parameters for trials that

fall in the upper and lower quartiles ranked by test loss. For PFN, the average initial learning rate for trials in the lower quartile of losses was $3.0\text{e-}3$, while for trials in the upper quartile it was higher, at $5.6\text{e-}3$. Similarly, the average dropout for trials in the lower quartile was 0.11, compared to a significantly higher 0.31 in the upper quartile. A similar pattern was observed with the ParticleNet model, with the average initial learning rate for trials in the lower quartile being $3.3\text{e-}3$, compared to $6.3\text{e-}3$ in the upper quartile, and the average dropout for trials in the lower quartile being 0.16, compared to 0.31 in the upper quartile.

A high learning rate allows the model to learn quickly, but it may also cause the model to overshoot the optimal solution and not converge well. The aim of dropout on the other hand is for the model to learn more robust, generalizable representations of the data. However, if the dropout rate is too high, as is the case for many trials in the upper loss quartile, the model may struggle to learn from the data at all, leading to underfitting.

The optimizer algorithm as a nominal variable falls outside the scope of the correlation analysis. However, Ranger proved to be the most successful. It combines LookAhead [54] with $k = 6$ and $\alpha = 0.5$, and an inner RAdam optimizer [55] with $\beta_1 = 0.95$, $\beta_2 = 0.999$ and $\epsilon = 10^{-5}$. RAdam can help to stabilize the learning rate and adapt it based on the variance of the gradient, making it a robust option when the learning rate is ill-adjusted. Furthermore, LookAhead has empirically been shown to improve convergence

by considering multiple directions in the parameter space. It also mitigates the impact of poorly chosen hyperparameters on training by smoothing out noisy gradient updates.

The gain measure in Table 2 represents the percentage improvement in test loss achieved by the PFN and ParticleNet models over the standard jet energy corrections loss value. This metric serves as an indicator of the models’ relative performance, providing a quantifiable measure of the benefits realized through hyperparameter optimization. The gain observed for the optimal PFN configuration is 7.52% whereas the optimal ParticleNet model achieves a 7.79% improvement over the baseline.

In this study, Random Search proved straightforward to set up and it showcased favorable practical properties. Random state is the only input parameter, the algorithm allows for trials to be discontinued or restarted without jeopardizing the experiment, and compared to grid search, it is more efficient for a given computational budget [46]. However, Katib offers several other AutoML algorithms such as Bayesian optimization [56] or Hyperband [57] that could be considered for future work since they have great potential to more effectively find an optimal set of hyperparameters.

It should be noted that more advanced algorithms often rely on additional input parameters that may alter the outcome which adds a level of complexity in the setup. Furthermore, while Random Search is embarrassingly parallel, Bayesian Optimization uses Gaussian process regression to iteratively model the search space and is therefore inherently sequential. Trials can still be queued in parallel, but the choice of parameter configuration for back-to-back trials is less informed than when the algorithm is run sequentially. Hyperband, on the other hand, is an extension of Random Search, and offers more efficient resource allocation to trials that matter by invoking early stopping on poorly performing configurations. However, adjusting resource allocation when candidate configurations have different convergence rates is an open challenge [57], a circumstance occurring in our experiment with varying learning rates and models with differing numbers of layers and hidden units.

5.2 Model complexity and inference performance

Table 3 compares the complexity of the optimal PFN and ParticleNet models in terms of loss, number of parameters, and Multiply-Accumulate operations (MACs). MACs represent the number of multiplications and additions performed during a single forward pass, indicating computational complexity. While ParticleNet achieves a lower test loss than PFN due to the inclusion of particle locality information, it has significantly higher computational complexity. A smaller number of nearest neighbors in the particle graph and fewer channels in the linear layers can be considered for reducing the complexity while still maintaining good performance [12].

Table 3: Comparison of model complexity for the optimal configuration of PFN and ParticleNet

Model	Loss	# Params	MACs
PFN	8.768e-2	355.45k	1.43M
ParticleNet	8.746e-2	123.47k	47.59M

The optimal PFN and ParticleNet models are served as REST endpoints using the Triton Inference Server running on top of KServe. We used the Python Triton client to request predictions for different batch sizes to evaluate how these models compare, and how request batch size affects roundtrip time, inference time, and overhead. The roundtrip time encompasses the total duration for a request to be processed, including both HTTP request time and inference time. Meanwhile, the overhead, calculated as the difference between roundtrip time and inference time, represents additional delay induced by factors such as data serialization / deserialization and network latency. The results of these tests for both models when served either on a CPU or a Tesla V100 GPU are displayed in Figure 8.

In the context of model performance, PFN achieve lower roundtrip times than ParticleNet due to its much lower computational complexity. However, the difference is less pronounced for smaller batch sizes. The larger overhead at small batch sizes affects both models similarly,

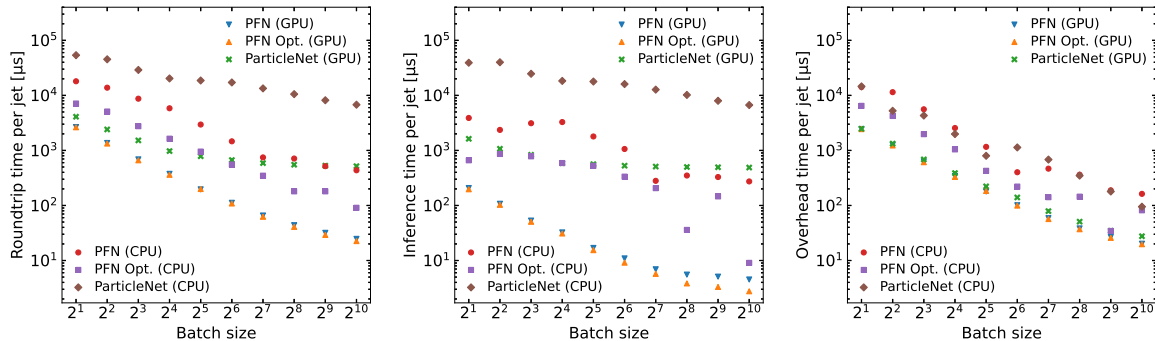


Fig. 8: Comparison of prediction request roundtrip time (left), device inference time (middle) and overhead (right) for PFN and ParticleNet ONNX models served with Triton. The represented values are based on the average processing time for 1k repetitions with randomly selected jets, covering a range of batch sizes from 2 to 1024

thereby reducing the relative performance difference between the models.

For very fast models such as PFN running on GPU the inference time is minimal due to the low model complexity and the high parallel processing capabilities of the GPU. This leads to overhead being the dominating factor in the roundtrip time. Conversely, for slower models like ParticleNet, the inference time especially when processing large batch sizes is substantially longer than the overhead time. As a result, the overhead comprises a fraction of the roundtrip time in that scenario.

We can note that PFN with extended ONNX graph optimization (PFN Opt.) shows performance enhancements compared to PFN with basic graph optimizations. Extended graph optimizations go beyond the simple, semantics-preserving transformations of basic optimizations, and applies complex node fusions after graph partitioning, tailoring the computation to the specific execution provider (CPU or GPU). This results in more efficient computations that are better suited to the architecture of the hardware on which the model runs.

The faster inference time on GPU compared to CPU for both PFN and ParticleNet models can be primarily attributed to the difference in the underlying architecture of these hardware platforms. GPUs are specifically designed for high-throughput, parallel processing and are capable of executing thousands of threads simultaneously. This characteristic is particularly beneficial for

inference tasks which involve large-scale matrix operations that can be parallelized effectively. In contrast, CPUs have fewer cores and are optimized for sequential tasks. Furthermore, the performance gap can be widened when the batch size is large, as larger batch sizes enable better utilization of the GPU’s parallel processing capabilities, leading to a faster per-jet inference time.

GPU inference and overhead times are also more consistent compared to CPU times. When allocating a processor in a Kubernetes cluster the instance is assigned a virtual processor (vCPU) shared across different processes, and due to their general-purpose nature, CPUs are typically tasked with managing a wider variety of processes, including system operations and other applications running in the background. A scheduler managed by the kernel has to coordinate time slots on the physical CPUs, which can lead to more variability in the availability of resources for the inference tasks. In contrast, GPU allocations on KubeFlow will currently result in a dedicated GPU for the task, resulting in more consistent performance.

When faced with the choice between CPU and GPU for deep learning inference it really depends on the hardware resources available and application requirements. Hardware accelerators such as GPUs provide superior throughput but often come with increased costs and power consumption. For applications with moderate inference workloads and less stringent response time requirements,

CPU-based inference may be more cost-effective. In contrast, high-throughput, low-latency applications can benefit significantly from GPU investment. It is also noteworthy that while a dedicated GPU will yield good results when benchmarking, virtualization of these resources as vGPUs could lead to better resource utilization [58], which is especially important as demand for hardware accelerators continues to increase.

5.3 Jet energy response flavor dependence

In this section, we present the analysis of the flavor dependence of the calibrated energy response produced by the optimal PFN and ParticleNet models. Figure 9 displays the median response for each jet flavor obtained using the two deep learning models and the standard corrections baseline. Both models exhibit a reduction in the differences between jet flavors compared to the baseline. A notable improvement in flavor dependence of the energy calibration is observed between light quark jets and gluon jets, which has been a known shortcoming in standard jet energy corrections.

The uncertainty for each data point, represented by the error bars in Figure 9, is calculated using the statistical bootstrapping method. By randomly sampling all response values 30 times, new sets of response values are generated similar in magnitude to the original dataset. The median is computed for each sample, followed by the standard deviation of the 30 median values, which provides an uncertainty for each point. There are fewer s, c and b jets in the QCD sample compared to the amount of u, d and g jets contributing to them having a higher bootstrapped uncertainty. The sample also has fewer jets in the endcap region compared to barrel region resulting in a higher uncertainty for the median response in the right plot in Figure 9 compared to the left one.

The sum of absolute errors (SAE) is used to evaluate the improvement in flavor dependence. It can be computed directly from the points in Figure 9 by summing the absolute difference between the median response for each flavor and the mean of the same points. Mathematically, it

can be expressed as:

$$\text{SAE} = \sum_{\text{flavor}} |R_{50\%, \text{flavor}} - \frac{1}{n} \sum_{\text{flavor}} (R_{50\%, \text{flavor}})| \quad (2)$$

where $\text{flavor} = \{u, d, s, c, b, g\}$ represents the complete set of jet flavors in the QCD sample. The relative improvement in flavor dependence for a model compared to the standard JEC is denoted by α and is defined as:

$$\alpha = 1 - \text{SAE}_{\text{Model}} / \text{SAE}_{\text{Baseline}}. \quad (3)$$

The values that α can take on ranges from any negative value to one, where a negative value indicates that the model performs worse than the standard correction, zero corresponds to no improvement, and $\alpha = 1$ would mean that the median response produced by the deep learning model is identical for all jet flavors.

As seen in Table 4, the improvement in flavor dependence varies across different p_T intervals for both models. In the very low p_T region ($30 \text{ GeV} < p_T^{\text{gen}} < 100 \text{ GeV}$), the improvements in flavor dependence are modest. As the p_T intervals become larger, both models demonstrate more significant improvements. Which of the two models perform better in a certain p_T interval or detector region varies. However, the differences in performance between them become less pronounced as p_T increases. It is also worth noting that the improvement in flavor dependence is not uniform across the barrel and endcap regions. For both models, the improvement is generally larger in the endcap region, especially in the intermediate and high p_T intervals.

Flavor response differences and even more so their differences between different generators are an input to both ATLAS and CMS flavor-related uncertainties. If the flavor responses become more alike, because the underlying jet properties are taken into account, as demonstrated in Table 4 and Figure 9, then one can also expect the uncertainties based on generator differences to be decreased.

5.4 Jet energy resolution

The performance of the regression models can also be assessed by examining the relative jet energy

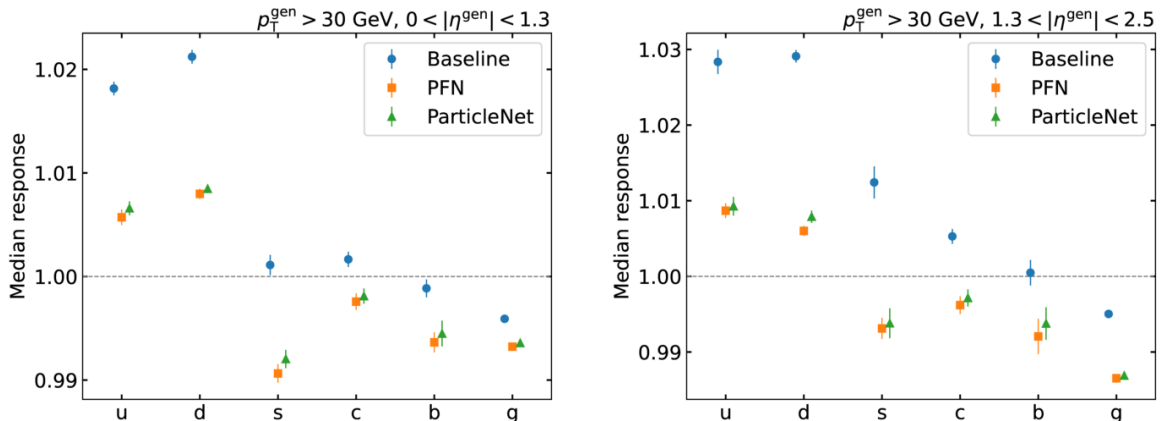


Fig. 9: Median energy response separated by jet flavor in the barrel region (left) and the endcap region (right). The improvement over the baseline is about 40% in the barrel region and 50% in the endcap region of the detector for both models

Table 4: Summary of jet energy regression results. Improvements in flavor dependence and energy resolution produced by PFN and ParticleNet compared to baseline JEC are presented in multiple p_T intervals and detector regions

Interval [GeV]	Model	α_{barrel} [%]	β_{barrel} [%]	α_{endcap} [%]	β_{endcap} [%]
$30 < p_T^{\text{gen}} < 100$	PFN	8.04	0.23	12.01	0.97
	ParticleNet	6.71	0.64	17.97	1.25
$100 < p_T^{\text{gen}} < 300$	PFN	23.25	1.43	49.59	6.93
	ParticleNet	24.52	1.61	45.07	7.07
$300 < p_T^{\text{gen}} < 1000$	PFN	57.52	4.53	68.05	12.90
	ParticleNet	56.11	4.81	70.02	11.93
$p_T^{\text{gen}} > 1000$	PFN	68.62	7.95	37.91	4.97
	ParticleNet	68.34	7.75	37.56	9.37

resolution. We define it here as the interquartile range (IQR) divided by the median for the response:

$$\bar{s} = \frac{R_{75\%} - R_{25\%}}{R_{50\%}}. \quad (4)$$

The IQR serves as a measure of response resolution. Both the median and IQR are robust statistics, meaning that they are less affected by outliers compared to the mean and standard deviation respectively. The uncertainty of the relative resolution is measured using the same bootstrapping technique as in the previous section and

results are shown in Figure 10. As the high p_T and endcap region are less populated as indicated in Figure 1, the uncertainties are sizeable for endcap high p_T jets.

The improvement in relative resolution for a model with respect to standard corrections is denoted here as β . We define it as one minus the ratio of relative jet energy resolution between the models and the baseline:

$$\beta = 1 - \bar{s}_{\text{Model}} / \bar{s}_{\text{Baseline}}. \quad (5)$$

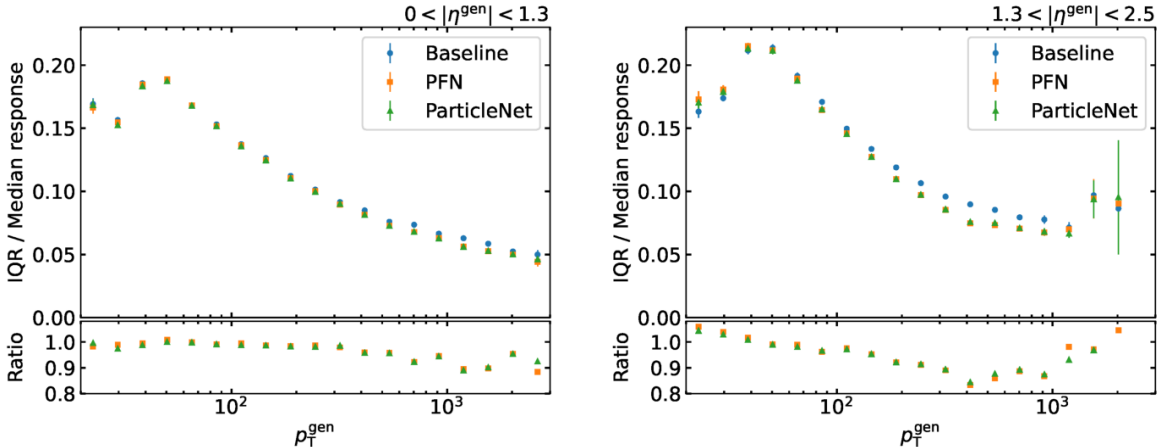


Fig. 10: Relative jet energy resolution in the barrel region (left) and the endcap region (right) binned logarithmically. The bottom panels show the ratio between the relative resolution produced by the deep learning models and the relative resolution after standard jet energy corrections

When examining the results presented in Table 4, we observe that both PFN and ParticleNet models achieve improvements in the energy resolution compared to the baseline. The improvements vary across different p_T intervals, with the largest improvements observed in the intermediate to high p_T intervals for both models. This behavior can be attributed to several factors. High p_T jets generally have more complex substructures and are more likely to undergo hard parton splittings, resulting in a higher multiplicity for the jet. This increased complexity leaves more room for improvement for machine learning-based approaches. The effect of pileup also diminishes at higher p_T resulting in less noisy data to train on. Owing to the limited amount of training data in the endcap region for jets with $p_T^{\text{gen}} > 1000$ GeV and reaching the kinematic limit of phase space in that regime, the improvements achieved through deep learning are comparatively smaller there.

6 Conclusion

In this paper, we presented a deep learning based workflow for calibrating the energy of particle jets in the CMS detector. By utilizing advancements in learning on particle clouds in the form of the PFN and ParticleNet models, we managed to improve upon standard jet energy corrections derived solely from kinematic quantities.

The results, categorized into jet energy resolution and flavor dependence, suggest that the performance of both networks is generally comparable, with larger improvements at higher p_T . The most notable difference between the two models is that the inclusion of locality information in ParticleNet results in a slightly better energy resolution at the expense of higher model complexity and inference time.

We have also demonstrated the potential of the Kubeflow platform for operationalizing ML workflows in high energy physics. As the field is witnessing a growing integration of ML techniques, the capabilities offered by Kubeflow, supporting the continual development of scalable ML solutions, are becoming increasingly more relevant. The pipeline we developed for this work enabled us to efficiently scale up our AutoML experiments on cloud resources and serve the optimal models as easily queryable REST endpoints. Having each step in the pipeline defined using Kubernetes custom resources allows for fine-grained access to hardware resources on the cloud and well-versioned, reusable machine learning workflows.

Acknowledgments. We wish to thank the CMS collaboration and the CERN OpenData group for publishing high-quality simulated data under an open-access policy. We acknowledge the support of the CERN IT department for providing the computational resources for this work.

Funding. Corresponding author D.H. is supported by the Academy of Finland under the ICT 2023: Frontier AI Technologies program (Grant No. 345635).

Data availability statement. The simulated dataset used for this study is hosted on the CERN OpenData portal. The instructions and code to replicate the studies in this paper are available at: <https://zenodo.org/record/7799179>

References

- [1] Radovic, A., Williams, M., Rousseau, D., Kagan, M., Bonacorsi, D., Himmel, A., Aurisano, A., Terao, K., Wongjirad, T.: Machine learning at the energy and intensity frontiers of particle physics. *Nature* **560**(7716), 41–48 (2018). <https://doi.org/10.1038/s41586-018-0361-2>
- [2] Guest, D., Cranmer, K., Whiteson, D.: Deep learning and its application to LHC physics. *Annual Review of Nuclear and Particle Science* **68**(1), 161–181 (2018). <https://doi.org/10.1146/annurev-nucl-101917-021019>
- [3] Shlomi, J., Battaglia, P., Vlimant, J.-R.: Graph neural networks in particle physics. *Machine Learning: Science and Technology* **2**(2), 021001 (2020). <https://doi.org/10.1088/2632-2153/abbf9a>
- [4] Ju, X., Farrell, S., Calafiura, P., Murnane, D., Gray, L., Klijnsma, T., Pedro, K., Cerati, G., Kowalkowski, J., Perdue, G., *et al.*: Graph neural networks for particle reconstruction in high energy physics detectors. In: *Advances in Neural Information Processing Systems*, vol. 32 (2019). <https://doi.org/10.48550/arXiv.2003.11603>
- [5] Choma, N., Monti, F., Gerhardt, L., Palczewski, T., Ronaghi, Z., Prabhat, P., Bhimji, W., Bronstein, M.M., Klein, S.R., Bruna, J.: Graph Neural Networks for IceCube Signal Classification. In: *IEEE International Conference on Machine Learning and Applications*, vol. 17, pp. 386–391 (2018). <https://doi.org/10.1109/ICMLA.2018.00064>
- [6] The CMS Collaboration: The CMS experiment at the CERN LHC. *Journal of Instrumentation* **3**(08), 08004 (2008). <https://doi.org/10.1088/1748-0221/3/08/S08004>
- [7] The CMS Collaboration: Jet energy scale and resolution in the CMS experiment in pp collisions at 8 TeV. *Journal of Instrumentation* **12**(02), 02014 (2017). <https://doi.org/10.1088/1748-0221/12/02/P02014>
- [8] The CMS Collaboration: Measurement of the top quark mass using a profile likelihood approach with the lepton+jets final states in proton-proton collisions at $\sqrt{s}=13$ TeV. Technical report, CERN, Geneva (2023). <https://cds.cern.ch/record/2848244>
- [9] The CMS Collaboration: A deep neural network for simultaneous estimation of b jet energy and resolution. *Computing and Software for Big Science* **4**(1), 10 (2020). <https://doi.org/10.1007/s41781-020-00041-z>
- [10] Kallonen, K.: Sample with jet properties for jet-flavor and other jet-related ML studies. *JetNTuple.QCD.RunIL13TeV.MC*. CERN Open Data Portal (2019). <https://doi.org/10.7483/OPENDATA.CMS.RY2V.T797>
- [11] Komiske, P.T., Metodiev, E.M., Thaler, J.: Energy flow networks: deep sets for particle jets. *Journal of High Energy Physics* **2019**(1), 121 (2019). [https://doi.org/10.1007/JHEP01\(2019\)121](https://doi.org/10.1007/JHEP01(2019)121)
- [12] Qu, H., Gouskos, L.: Jet tagging via particle clouds. *Phys. Rev. D* **101**, 056019 (2020). <https://doi.org/10.1103/PhysRevD.101.056019>
- [13] Gambhir, R., Nachman, B., Thaler, J.: Learning uncertainties the frequentist way: Calibration and correlation in high energy physics. *Physical review letters* **129**(8), 082001 (2022). <https://doi.org/10.1103/PhysRevLett.129.082001>
- [14] Mäkinen, S., Skogström, H., Laaksonen, E., Mikkonen, T.: Who Needs MLOps: What Data Scientists Seek to Accomplish and How Can MLOps Help? In: *IEEE/ACM Workshop*

- on AI Engineering - Software Engineering for AI, vol. 1, pp. 109–112 (2021). <https://doi.org/10.1109/WAIN52551.2021.00024>
- [15] Golubovic, D., Rocha, R.: Training and Serving ML workloads with Kubeflow at CERN. In: 25th International Conference on Computing in High-Energy and Nuclear Physics, vol. 251, p. 02067 (2021). <https://doi.org/10.1051/epjconf/202125102067>
- [16] Yuan, D.Y., Wildish, T.: Bioinformatics application with Kubeflow for batch processing in clouds. In: High Performance Computing, pp. 355–367 (2020). https://doi.org/10.1007/978-3-030-59851-8_24
- [17] Tsourdinis, T., Chatzistefanidis, I., Makris, N., Korakis, T.: AI-driven Service-aware Real-time Slicing for beyond 5G Networks. In: IEEE Conference on Computer Communications Workshops, vol. 41, pp. 1–6 (2022). <https://doi.org/10.1109/INFOCOMWKSHPS54753.2022.9798391>
- [18] Carminati, F., Khattak, G., Loncar, V., Nguyen, T.Q., Pierini, M., Rocha, R.B.D., Samaras-Tsakiris, K., Vallecorsa, S., Vlimant, J.-R.: Generative adversarial networks for fast simulation. *Journal of Physics: Conference Series* **1525**(1), 012064 (2020). <https://doi.org/10.1088/1742-6596/1525/1/012064>
- [19] Sjöstrand, T., Ask, S., Christiansen, J.R., Corke, R., Desai, N., Ilten, P., Mrenna, S., Prestel, S., Rasmussen, C.O., Skands, P.Z.: An introduction to PYTHIA 8.2. *Computer physics communications* **191**, 159–177 (2015). <https://doi.org/10.1016/j.cpc.2015.01.024>
- [20] Cacciari, M., Salam, G.P., Soyez, G.: The anti-kt jet clustering algorithm. *Journal of High Energy Physics* **2008**(04), 063 (2008). <https://doi.org/10.1088/1126-6708/2008/04/063>
- [21] Agostinelli, S., Allison, J., Amako, K.a., Apostolakis, J., Araujo, H., Arce, P., Asai, M., Axen, D., Banerjee, S., Barrand, G., *et al.*: GEANT4—a simulation toolkit. *Nuclear instruments and methods in physics research section A: Accelerators, Spectrometers, Detectors and Associated Equipment* **506**(3), 250–303 (2003). [https://doi.org/10.1016/S0168-9002\(03\)01368-8](https://doi.org/10.1016/S0168-9002(03)01368-8)
- [22] The CMS Collaboration: Identification of heavy-flavour jets with the CMS detector in pp collisions at 13 TeV. *Journal of Instrumentation* **13**(05), 05011 (2018). <https://doi.org/10.1088/1748-0221/13/05/P05011>
- [23] The CMS Collaboration: Jet algorithms performance in 13 TeV data. Technical report, CERN, Geneva (2017). <http://cds.cern.ch/record/2256875>
- [24] The CMS collaboration: Particle-flow reconstruction and global event description with the CMS detector. *JINST* **12**(10), 10003 (2017). <https://doi.org/10.1088/1748-0221/12/10/P10003>
- [25] The CMS Collaboration: Pileup mitigation at CMS in 13 TeV data. *JINST* **15**(09), 09018 (2020). <https://doi.org/10.1088/1748-0221/15/09/P09018>
- [26] The CMS Collaboration: Performance of quark/gluon discrimination in 8 TeV pp data. Technical report, CERN, Geneva (2013). <https://cds.cern.ch/record/1599732>
- [27] de Oliveira, L., Kagan, M., Mackey, L., Nachman, B., Schwartzman, A.: Jet-images—deep learning edition. *Journal of High Energy Physics* **2016**(7), 1–32 (2016). [https://doi.org/10.1007/JHEP07\(2016\)069](https://doi.org/10.1007/JHEP07(2016)069)
- [28] Guest, D., Collado, J., Baldi, P., Hsu, S.-C., Urban, G., Whiteson, D.: Jet flavor classification in high-energy physics with deep neural networks. *Physical Review D* **94**(11), 112002 (2016). <https://doi.org/10.1103/PhysRevD.94.112002>
- [29] Louppe, G., Cho, K., Becot, C., Cranmer, K.: QCD-aware recursive neural networks for jet physics. *Journal of High Energy Physics* **2019**(1), 1–23 (2019). [https://doi.org/10.1007/JHEP01\(2019\)057](https://doi.org/10.1007/JHEP01(2019)057)

- [30] Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R.R., Smola, A.J.: Deep Sets. In: *Advances in Neural Information Processing Systems*, vol. 30 (2017). <https://doi.org/10.48550/arXiv.1703.06114>
- [31] Glorot, X., Bordes, A., Bengio, Y.: Deep sparse rectifier neural networks. In: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research*, vol. 15, pp. 315–323 (2011). <https://proceedings.mlr.press/v15/glorot11a>
- [32] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**(56), 1929–1958 (2014). <https://jmlr.org/papers/v15/srivastava14a.html>
- [33] Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph CNN for learning on point clouds. *ACM Transactions On Graphics (TOG)* **38**(5), 1–12 (2019). <https://doi.org/10.1145/3326362>
- [34] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 29, pp. 770–778 (2016). <https://doi.org/10.1109/CVPR.2016.90>
- [35] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., *et al.*: Pytorch: An imperative style, high-performance deep learning library. In: *Advances in Neural Information Processing Systems*, vol. 32 (2019). <https://doi.org/10.48550/arXiv.1912.01703>
- [36] Qu, H.: Weaver: A machine learning R&D framework for high energy physics applications. [Online]. Available: <https://github.com/hqucms/weaver>, Accessed on: March 10, 2023.
- [37] Brun, R., Rademakers, F.: ROOT — An object oriented data analysis framework. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* **389**(1), 81–86 (1997). [https://doi.org/10.1016/S0168-9002\(97\)00048-X](https://doi.org/10.1016/S0168-9002(97)00048-X)
- [38] Pivarski, J., Elmer, P., Lange, D.: Awkward Arrays in Python, C++, and Numba. In: *24th International Conference on Computing in High Energy and Nuclear Physics*, vol. 245, p. 05023 (2020). <https://doi.org/10.1051/epjconf/202024505023>
- [39] Gabriel, E., Fagg, G.E., Bosilca, G., Angskun, T., Dongarra, J.J., Squyres, J.M., Sahay, V., Kambadur, P., Barrett, B., Lumsdaine, A., *et al.*: Open MPI: Goals, concept, and design of a next generation MPI implementation. In: *Recent Advances in Parallel Virtual Machine and Message Passing Interface: 11th European PVM/MPI Users’ Group Meeting*, pp. 97–104 (2004). https://doi.org/10.1007/978-3-540-30218-6_19
- [40] NVIDIA Corporation: NCCL: Optimized Primitives for Collective Multi-GPU Communication. [Online]. Available: <https://github.com/nvidia/ncl>, Accessed on: April 13, 2023.
- [41] Burns, B., Grant, B., Oppenheimer, D., Brewer, E., Wilkes, J.: Borg, Omega, and Kubernetes. *ACM Queue* **14**, 70–93 (2016). <https://doi.org/10.1145/2890784>
- [42] Pavlou, C.S., Kessler, F.B., Katsakioris, I., Kostis, L., Stefano, F., Alexiou, T., Valerio, M.: KALE: KubeFlow Automated pipeLines Engine. [Online]. Available: <https://github.com/kubeflow-kale/kale>, Accessed on: April 13, 2023.
- [43] The YAML Project: YAML Ain’t Markup Language. [Online]. Available: <https://yaml.org>, Accessed on: April 13, 2023.
- [44] George, J., Gao, C., Liu, R., Liu, H.G., Tang, Y., Pydipaty, R., Saha, A.K.: A scalable and cloud-native hyperparameter tuning system. *arXiv eprint* (2020). <https://doi.org/10.48550/arXiv.2006.02085>
- [45] The KubeFlow Project: Kubernetes

- Custom Resource and Operator for PyTorch Jobs. [Online]. Available: <https://github.com/kubeflow/pytorch-operator>, Accessed on: April 13, 2023.
- [46] Bergstra, J., Bengio, Y.: Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research* **13**(10), 281–305 (2012). <http://jmlr.org/papers/v13/bergstra12a.html>
- [47] ONNX Runtime developers: ONNX Runtime: A cross-platform, high performance ML inferencing and training accelerator. [Online]. Available: <https://onnxruntime.ai>, Accessed on: April 13, 2023.
- [48] NVIDIA Corporation: Triton Inference Server: An Optimized Cloud and Edge Inferencing Solution. [Online]. Available: <https://github.com/triton-inference-server>, Accessed on: April 13, 2023.
- [49] Google Inc.: Protocol Buffers: a language-neutral, platform-neutral extensible mechanism for serializing structured data. [Online]. Available: <https://protobuf.dev>, Accessed on: May 30, 2023.
- [50] The KServe Project: KServe: Standardized Serverless ML Inference Platform on Kubernetes. [Online]. Available: <https://github.com/kserve/kserve>, Accessed on: April 13, 2023.
- [51] The Knative Project: Knative: Kubernetes-based platform to build, deploy, and manage modern serverless workloads. [Online]. Available: <https://knative.dev>, Accessed on: April 13, 2023.
- [52] The Istio Project: Istio: Connect, secure, control, and observe services. [Online]. Available: <https://istio.io>, Accessed on: April 13, 2023.
- [53] Fielding, R.T., Taylor, R.N.: Principled design of the modern web architecture. *ACM Trans. Internet Technol.* **2**(2), 115–150 (2002). <https://doi.org/10.1145/514183.514185>
- [54] Zhang, M., Lucas, J., Ba, J., Hinton, G.E.: Lookahead optimizer: k steps forward, 1 step back. In: *Advances in Neural Information Processing Systems*, vol. 32 (2019). <https://doi.org/10.48550/arXiv.1907.08610>
- [55] Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., Han, J.: On the variance of the adaptive learning rate and beyond. In: *International Conference on Learning Representations*, vol. 8 (2020). <https://doi.org/10.48550/arXiv.1908.03265>
- [56] Brochu, E., Cora, V.M., de Freitas, N.: A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. Technical Report UBC TR-2009-023, University of British Columbia, Department of Computer Science (2009). <https://doi.org/10.48550/arXiv.1012.2599>
- [57] Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., Talwalkar, A.: Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research* **18**(185), 1–52 (2018). <http://jmlr.org/papers/v18/16-558.html>
- [58] Golubovic, D., Gaponic, D., Guerra, D., Rocha, R.: Efficient Access to Shared GPU Resources Part 1: Mechanisms, Motivations and Use Cases for GPU Concurrency on Kubernetes (2023). <https://kubernetes.web.cern.ch/blog/2023/01/09/efficient-access-to-shared-gpu-resources-part-1>