# Accelerating GAN training using highly parallel hardware on public cloud

*Renato* Cardoso[1,*], *Dejan* Golubovic[1], *Ignacio* Peluaga Lozada[1], *Ricardo* Rocha[1], *João* Fernandes[1], and *Sofia* Vallecorsa[1]

[1]CERN, 1 Esplanade des Particules, Geneva, Switzerland

**Abstract.** With the increasing number of Machine and Deep Learning applications in High Energy Physics, easy access to dedicated infrastructure represents a requirement for fast and efficient R&D. This work explores different types of cloud services to train a Generative Adversarial Network (GAN) in a parallel environment, using Tensorflow data parallel strategy. More specifically, we parallelize the training process on multiple GPUs and Google Tensor Processing Units (TPU) and we compare two algorithms: the TensorFlow built-in logic and a custom loop, optimised to have higher control of the elements assigned to each GPU worker or TPU core. The quality of the generated data is compared to Monte Carlo simulation. Linear speed-up of the training process is obtained, while retaining most of the performance in terms of physics results. Additionally, we benchmark the aforementioned approaches, at scale, over multiple GPU nodes, deploying the training process on different public cloud providers, seeking for overall efficiency and cost-effectiveness. The combination of data science, cloud deployment options and associated economics allows to burst out heterogeneously, exploring the full potential of cloud-based services.

## 1 Introduction

In recent years, several studies have demonstrated the benefits of using Deep Learning (DL) to solve different tasks related to data processing in High Energy Physics (HEP): for example, generative models are being tested as fast alternatives to Monte Carlo based simulation and anomaly detection algorithms are being explored to design searches for rare new-physics processes.

Training of such models has been made tractable with the improvement of optimization methods and the advent of dedicated infrastructure. High Performance Computing (HPC) and storage technologies are often required by these types of projects, together with the provision of high processing capacity on multi-architectures, ranging from large multi-core systems to hardware accelerators like GPUs or TPUs. These requirements directly translate into the necessity to extend the HEP computing model that, so far, has relied on on-premise resources and grid computing. A hybrid model has the potential to seamlessly complement and extend local infrastructure, using widely supported techniques, allowing flexible and efficient resource bursting to public cloud for CERN research teams and science beyond. In

---

*e-mail: Renato.Cardoso@cern.ch

this context, our work explores the deployment of HEP DL models on native cloud-based frameworks, with well-established industry support and large communities, whilst, at the same time, integrate it with the ongoing work to transparently and generically burst CERN on-premises infrastructure to public clouds. For this reason, in addition to the optimisation of the specific physics use case, we provide relevant feedback on the most efficient models to develop an end-to-end integrated strategy that includes cost prediction and optimization when consuming public cloud resources at scale, for scientific research. The development and optimisation of deep generative models as fast simulation solutions is one of the most important research directions as far as DL for HEP is concerned. HEP experiments rely heavily on simulations in order to model complex processes and describe detectors response: the classical Monte Carlo approach can reproduce theoretical expectations with a high level of precision, but it is extremely demanding in terms of computing resources [1]. For this reason, different strategies trading some physics accuracy for speed are being developed as fast simulation techniques and deep generative models are today among the most promising solutions [2–5]. This work focuses on 3DGAN [6]: a three-dimensional convolutional Generative Adversarial Network (GAN) that is being optimised for the simulation of electromagnetic calorimeters. This work summarises the specific development required by the 3DGAN architecture and the adversarial training process for efficient data parallel training on Google TPUs and multiple GPUs setups. It also details the effect of distributed and mixed-precision training (on TPUs) on 3DGAN physics performance: this aspect is essential to ensure that the performance of complex simulations tasks (such as the GAN generation process) is preserved at an acceptable level. In addition, we describe the 3DGAN training deployment on public cloud providers (Google Cloud Platform (GCP) and Microsoft Azure) comparing different approaches, based on Kubernetes, Kubeflow and the Azure Machine Learning service. Our deployments are among the first examples of establishing an innovative hybrid platform successfully deploying scientific deep learning workloads on public clouds at large scale [7].

## 2  Generative Adversarial Networks in High Energy Physics

Generative models represent a fundamental part of deep learning and GANs [8], in particular, can generate sharp and realistic images with high resolution [9, 10]. In HEP, most GAN applications address the problem of calorimeters simulation. The LaGAN [11] and CaloGAN [12] models first introduced GAN for simulating simplified detectors; since then, multiple prototypes have been developed [4, 5, 13]. While distributed deep neural networks training is a well-established field, with multiple approaches thoroughly investigated and state-of-the-art solutions widely spread, the specific GAN use case is still a relatively new research subject: [14] represents a notable example at large scale, while additional work can be found in [15, 16]. Low precision computation (int8, float16, bfloat16) is another active research field, aimed at reducing computing resources, mostly targeted to speed-up inference: models are trained in float32 precision, on multi-node clusters, quantized to lower precision and then loaded on light-weight devices such as edge-devices. Lower precision formats can be used at training time as well, in combination with full precision float32 in order to maintain the level of accuracy (mixed precision training) [17, 18].

Details about the 3DGAN model can be found in [6]: it represents the first proof of concept on using 3D convolutional GANs to simulate high granularity calorimeters. The 3DGAN model uses an ACGAN-like [19] approach to introduce physics constraints and improve image fidelity: in particular, two quantities, representing the particle energy, $E_P$, and its trajectory, measured by the angle $\theta$, are used as input to the generation step. In general terms, simulating samples using generative models is much faster than using a Monte Carlo approach. In the case of 3DGAN, a several orders of magnitude speed-up is observed with

respect to the HEP Monte Carlo state-of-the-art [20]: a 67000x speed-up has been recently achieved using int8 quantized inference on Intel Xeon Cascade Lake [21]. The training process is however very time consuming: an entire week is needed in order to train the model to convergence using a single NVIDIA V100 GPU and therefore a distributed training approach is essential [22].

## 3 The adversarial training process acceleration

The Tensorflow 2.3 *distribute.Strategy* API runs the training process across a single node with multiple GPUs, multiple nodes, as well as multiple TPU cores, while ensuring portability with minimal code changes; the downside is that only graph mode is available on TPUs, whereas time penalties are introduced in eager mode on GPUs. This work relies on the *mirrored strategy*, that implements synchronous parallel training over multiple GPUs on a single node; the *multi-worker mirrored strategy*, which extends it across multiple nodes and the *TPU strategy*, that implements synchronous training on TPUs, using its own, optimised, *All_reduce* and collective operations [23].

The 3DGAN adversarial training, shown in Algorithm 1, trains, alternately, the discriminator and generator networks an equal number of times. It includes several steps to initialise the generation process, train the discriminator network on real and fake data and, finally, the generator itself.

---
**Algorithm 1** The 3DGAN training process

---
**for** $N_{epochs}$ **do**
    **for** $N_{batches}$ **do**
        get 3D image batch
        get labels: $E_p$, $E_{CAL}$, $\theta$
        generator input: sample latent noise batch and concatenate(noise, $E_p$, $\theta$)
        generate fake 3D image batch
        calculate fake $E_{CAL}$ batch
        train discriminator on real batch
        train discriminator on fake batch
        **for** 2 **do**
            get labels: $E_p$, $E_{CAL}$, $\theta$
            generator input: sample latent noise batch and concatenate(noise, $E_p$, $\theta$)
            train the generator
        **end for**
    **end for**
**end for**

---

Algorithm 1 can be implemented using standard training APIs, such as *keras.train_on_batch* or custom training loops. In the first case, the distribution logic is built-in and only the batch size is manually updated according to the number of parallel replicas. Taking a look at the pseudo-code in Algorithm 1, it is evident that the generator input initialization steps are not distributed by *keras.train_on_batch*: they are sequentially run and end up being a bottleneck when the number of replicas increases. This effect is shown in Figure 1: the discriminator training time (per batch) stays the same while increasing the number of GPUs, with a small increase due to reduce and broadcast operations. On the other hand, the generator initialisation time increases linearly with the number of GPUs, since it is sequentially run.

In order to reduce this bottleneck, we rewrite Algorithm 1 as a custom training loop, using Tensorflow *tf.function* to explicitly redefine the forward step in training mode. In addition, the *tf.function* includes all previously sequential steps. At this point, the remaining bottleneck is represented by the data pre-processing and distribution steps. In order to address this issue, we rely on some important advantages related to using a customised training loop: loading
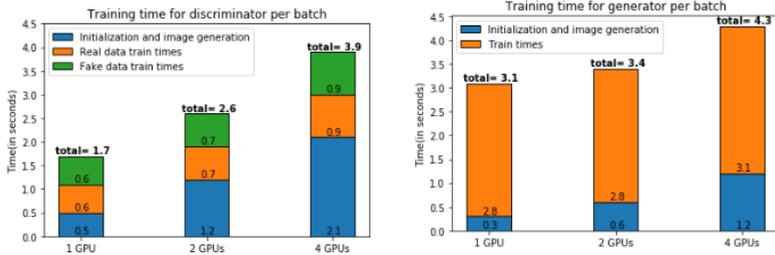
Figure 1: (left) Discriminator training times divided between processing, real data training and fake data training. (right) Generator training times. Times are quoted for one data batch.

data directly to the accelerators (GPUs and TPUs) and running data preparation (batching and shuffling) on the CPU host while the GPUs/TPUs are training. More specifically, we convert the original HDF5 data to the native tensor format, TF Records, and introduce iterators instead of manually instantiating batches on the CPU host. This final optimisation step yields results we discuss in the next section for TPUs and in section 5 for a larger number of GPUs on Google Cloud Platform (GCP) and Microsoft Azure clouds.

## 4 Training on Google Tensor Processing Units

Application-specific integrated circuits (ASICs), such as Google Tensor Processing Units (TPUs) dedicated to ML/DL workloads, are being developed with the objective of accelerating matrix multiplication. In order to have efficient communication across TPUs, Google implements a 2D Torus topology by using Inter-Core Interconnect (ICI) links that enable a direct connection between chips and simplify the rack-level deployment. In addition, the v2 and v3 TPUs use two cores per chip, further decreasing the latency with respect to single-core chips. In particular, each TPU core contains a 128x128 matrix multiplier unit (MXU) which uses systolic arrays.

Google TPUs use bfloat16 for MXU calculations, in order to accelerate data processing while, at the same time, maintaining a high level of accuracy. To fully exploit bfloat16, the Tensorflow XLA compiler maps TPU dependent instructions while also targeting CPU and GPUs instructions. The XLA compiler analyses the code and optimises the data transfer among different TPU elements, maps concurrent instructions in different parts of the TPU hardware and increases the ICI communication capabilities between TPU cores. All these optimisation steps occur during the first *tf.function* iteration, which in consequence, is slower. For this reason we quote average results excluding the first batch.

### 4.1 Results

Our initial test is run on 8 cores TPUs, using a batch size $BS = 128$ in order to match the design MXU dimension. Results for v2 and v3 TPUs are compared in figure 2: using the double core v3 TPU we obtain 2x speed-up on the epoch training time. Figure 2 also shows, on the center panel, the batch size impact on training time on 8 cores v3 TPU: a batch size $BS = 64$ is under utilising the MXU capacity and, therefore, it is processed in the same amount of time as a 128 batch; on the other hand, $BS = 256$, which is double the optimal size, is processed using twice the time. Figure 2, right panel, shows weak scaling results, obtained using a $BS = 128$ batch size on v3 TPUs with a number of cores ranging from 8 to
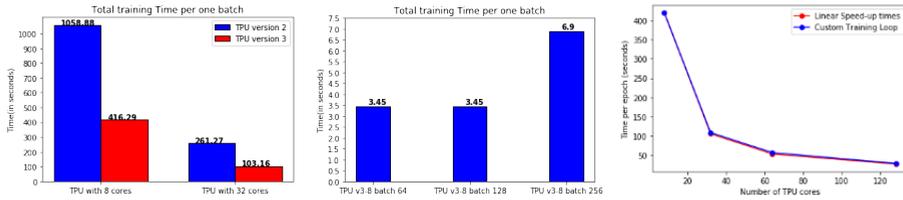
Figure 2: (left) Time per epoch (in seconds) on 8 cores v2 and v3 TPU, $BS = 128$. (center) Time per epoch (in seconds) on 8 cores v3 TPU for batch sizes $BS = 64, 128, 256$. (right) Measured training time per epoch (in seconds) compared to the theoretical linear speed-up.

128: the time per epoch scales linearly. Running on a 128 cores TPU, reduces the training time per epoch to about 30 seconds, allowing 3DGAN to train to convergence in a little more than one hour. As mentioned above, preserving the fidelity of the synthetic images is critical, when accelerating training. In order to verify this constraint, we validate the quality of the generated images against Monte Carlo simulation: the left panel on figure 3 shows the calorimeter energy response along the longitudinal plane. It can be noticed, that while the GAN model seems to slightly distort the response toward the left side of the distribution, the agreement remains overall very good.
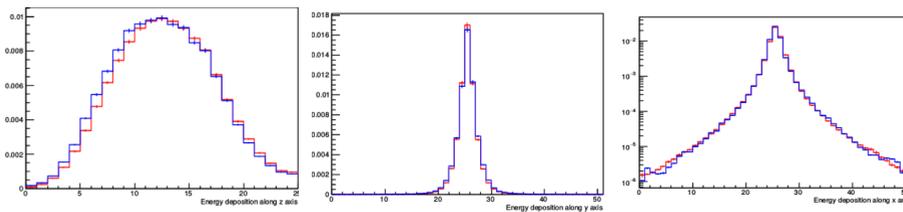


Figure 3: Calorimeter energy response. The GAN prediction is in blue, Monte Carlo in red. (left) Along the detector longitudinal plane. (center) Along the detector transverse plane in linear scale and (right) logarithmic scale.

The energy response along the transverse plane, shown in the center panel of figure 3, is also optimally reproduced, both in the central calorimeter cells (the distribution bulk high-lighted in the linear scale plot) and at the volume edges (shown in the logarithmic scale plot).

To further analyse the results obtained, a profiling tool was employed using tensorboard. The program is mostly compute bound, in fact the TPUs are only idle 0.7% of the time, with 38% of the time being used for the forward-pass (represented as the functional models), 61% for back-propagation (the gradient tapes) and 0.2% for *AllReduce*.

On the other hand, analysis on the GPUs shows an increase on the idle time to around 2.9%, for a higher number of GPUs, this is due to an increase in the time it takes to do the *AllReduce* operation, which occurs after every gradient tape operation shown in the profile.

## 5  Automated deployment using public cloud services

Training deep generative models for simulation is a workload that fits well the public cloud for multiple reasons. First of all, training is done periodically and can greatly benefit from

the use of on-demand resources when a longer-term investment on-premises cannot be justified. Secondly the availability of a large number of GPUs, often a scarce resource on data centers on-premises, as well as access to the most recent architectures. Finally, the access to specialized accelerators such as Google TPUs and Graphcore IPUs, that are vendor specific, allowing for a generic unlocked deployment strategy, whilst at the same time, benefiting from a side-by-side comparison.

The deployment can take different approaches, including the low level management of virtual machines which keeps the burden of a significant fraction of the infrastructure management on the end user; on the other hand, one can also use open platforms such as Kubernetes, which offer APIs abstracting the infrastructure and delegating most of the operations to the cloud service itself, being available in all major public clouds; another possibility is to use MLaaS provided by the vendor, that normally consists in a PaaS optimised and dedicated to machine learning, such as the Azure Machine Learning service. Below, we describe and present results for the last two approaches.

## 5.1 Google Cloud Platform with Kubeflow

Our first setup relies on Kubeflow [24]. Our existing on-premises deployment and configuration was reused and directed to the Google Kubernetes Engine (GKE), Google Cloud's managed Kubernetes service offering, benefiting from the resources under CERN IT's Cloud Broker Pilot project (CloudBank EU)[25]. To monitor resource usage during training we rely on Prometheus and custom Grafana dashboards showing GPU power usage, memory usage, temperature and overall utilization.

The Kubernetes cluster includes multiple node groups, each of them with a profile where the number of GPUs per node varies from 1 to 8. This allows the evaluation of the penalty introduced from spreading the training across a larger number of nodes, as well as the penalty of increasing the density in a single node. Instances are of the n1-standard family, from n1-standard-4 for single GPU nodes to n1-standard-32 for 8 GPU nodes. All tests rely on preemptible instances, which are low SLA resources living for a maximum of 24 hours that can be deleted at any point if capacity is reclaimed - but are over 3 times cheaper compared to reserved instances for V100 GPUs. No reclaim was experienced during the full set of tests performed.

### 5.1.1 Experiments and Results

In the initial set of runs, the goal is to determine an optimal batch size for training the neural network. Higher batch sizes mean less steps to complete the entire data set, which results in faster completion of the training. Though smaller batch sizes can yield to better generalization performance, in this case we tolerate a small loss of generalization performance (as discussed in section 5.2). From the hardware perspective, higher batch sizes offer better utilization of GPUs by exploiting advantages of the hardware. For the purpose of the paper, as models are being trained on enterprise machine learning platforms, it is important to provide maximum utilization of resources in order to minimize the cost for training jobs. This experiment is conducted by fixing the number of nodes (4), GPUs per node (8) and the number of epochs (5). The only hyper-parameter varying is the batch size. Batch sizes are selected to be multiples of 2 due to the GPU Single Instruction Multiple Data (SIMD) organization. Minimum batch size is 16, since it is sufficiently small to provide good generalization performance. The maximum batch size (96) is selected to fit the GPU memory. As expected, the training times are reduced when using larger batches. The increase in performance is
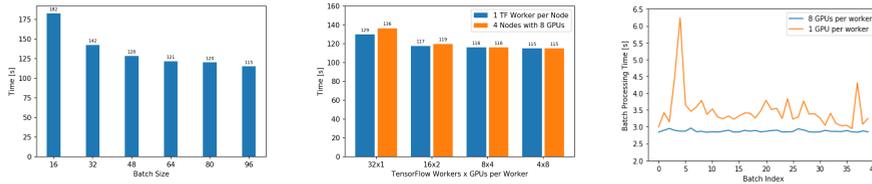
Figure 4: (left) Impact of batch size on time to perform one epoch of training. (center) Impact of various combinations of hardware layout and TF worker configuration on time to perform one epoch of training. Orange represents the first set of runs, where hardware configuration is fixed with 4 nodes with 8 GPUs each, while number of workers and GPUs per worker change. Blue represents the second set of runs where number of workers and GPUs per worker correspond to number of nodes and GPUs per node. (right) Batch processing times for the second epoch of training using hardware layout with 4 nodes with 8 GPUs per node. Orange represents the batch processing times when number of workers is 32 and each worker has 1 GPU, Blue represents the batch processing times when number of workers is 4 and each worker has 8 GPUs.

prominent between batch size of 16 and 32. With batch sizes of 64 and 80 the difference is almost negligible.

As described in the infrastructure chapter, the cluster being used can be configured with a different number of nodes and GPUs per node. In this experiment, we select the optimal layout with respect to the TensorFlow distributed training setup, including the number of GPUs per node, TensorFlow workers and number of GPUs in each worker. The objective is to understand the impact of communication overhead in the TensorFlow distributed setup. The experiment has two sets of runs using a total of 32 GPUs. In the first set the number of nodes (4) and the number of GPUs per node (8) are fixed. Different combinations of workers are selected: 32 workers with 1 GPU each, 16 workers with 2 GPUs each, 8 workers with 4 GPUs each and 4 workers with 8 GPUs each. In the second set the number of nodes and GPUs per node varies: 32 nodes with 1 GPU each, 16 nodes with 2 GPUs each, 8 nodes with 4 GPUs each and 4 nodes with 8 GPUs each. The workers configuration follows the nodes configuration and is the same as in the previous run: 32 workers with 1 GPUs each, 16 workers with 2 GPUs each, 8 workers with 4 GPUs each and 4 workers with 8 GPUs each. Essentially, the same set of worker configurations is run on a different hardware layout. The results in figure 4, on the right, show that the impact of different hardware configurations is not huge across the runs. The last three runs from both sets show times that differ by less than 3 seconds. The first run in both sets (32 workers with 1 GPU each) shows a significant increase in time, meaning the optimal configuration should not strive for less GPUs per node and less GPUs per worker. The best results are achieved when the number of workers matches the number of nodes and the number of GPUs per worker matches the number of GPUs per node. In subsequent experiments this setup will be used.

It was also observed that a higher number of workers introduces instability in the processing time per one batch, leading to instability in the processing time per one epoch and a less predictable time for the entire run completion. With a run with 32 workers, the variation in batch processing time is significant compared to the run with 4 workers. This goes in favor of avoiding higher number of workers and focusing on using nodes with as many GPUs as possible.
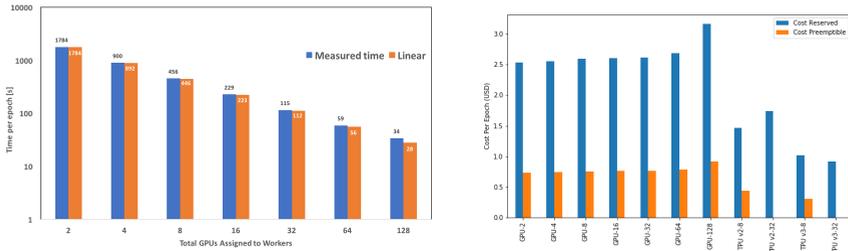
Figure 5: (left) The impact of increasing number of total GPUs assigned to workers on time to process one epoch of training. (right) Cost of running one epoch on the Google Cloud Platform. The cost per epoch remains similar for all GPU combinations, though the training time improves almost linearly.

We next determine how the performance changes with the increased number of GPUs available to workers. This set of runs strives for maximum utilization of the maximum number of GPUs available. Building on the previous tests a batch size of 96 and one worker per node are selected, varying only the number of nodes and GPUs per node. The performance improvement is close to linear with a slight drop when moving to 128 GPUs, which are positive results regarding TensorFlow's distributed training capabilities to scale workloads horizontally (see figure 5). This is particularly true when combined with Kubeflow and the ease of configuration and automation of the training process in large infrastructures.

### 5.1.2  Cost Analysis Overview

An important point when using public cloud resources is to understand upfront the costs involved in order to reach optimal service settings and configurations and, consequently, maximize cost effectiveness to reach a sustainable position. The infrastructure costs during the tests on Google Cloud Platform are driven by the cost of the GPUs, with virtual machines accounting for less than 5% of the total. Below, we present the numbers for using 2 to 128 GPUs and for a subset of TPU configurations. The calculation is done taking the training time of one epoch against the hourly GPU costs in region *europe-west4*  and include both reserved and preemptible options. It should be noted that when using reserved instances it is possible to apply for sustained or committed-use discounts which reduce the difference to preemptible instances costs. A few points are worth noting (figure 5, right panel): the cost per epoch remains similar when increasing the number of GPUs, while the training time is reduced up to 52 times for 128 GPUs (almost linear up to 64 GPUs); using preemptible GPUs can allow for significant cost savings and should be used whenever possible; finally, the best results are achieved using preemptible TPU v3-8, which are 2.4 times cheaper than their GPU equivalent, considering the training time. This is an interesting result that positions TPUs as a cost-effective option combined with their potential to reduce the overall training time. Preemptible TPUs are only available up to 8 cores, with the reserved instance costs being less attractive. As an example, the overall training cost with 128 GPUs is the same as TPU v3-32, but it takes half the time. Future work will include a comparison with a larger number of TPU cores.

## 5.2  Microsoft Azure with Azure's Machine Learning Service

In the case of Microsoft Azure cloud, the 3DGAN deployment makes use of the Azure Machine Learning service [26] offered through CERN openlab's Azure cloud research enroll-
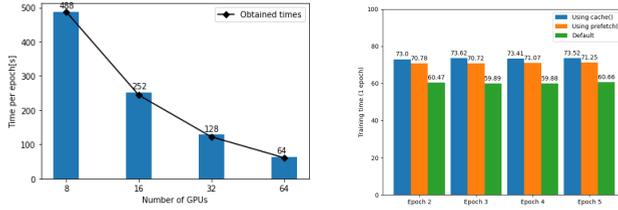
Figure 6: (left) Microsoft Azure training time per epoch for a $BS = 64$ batch size. (right) Azure ML service automatic optimisation performance compared to manual cache and prefetch implementation. All times are measured in seconds.
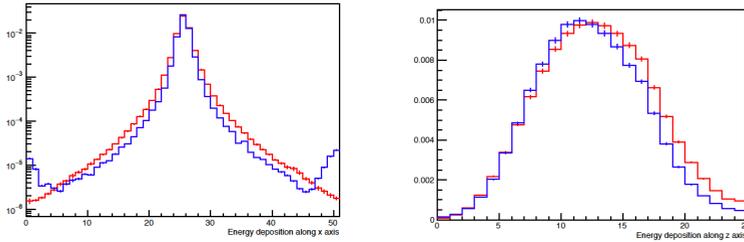


Figure 7: Calorimeter energy response as predicted by 64 GPUs GAN (blue) and Monte Carlo (red). (left) Distribution along the detector transverse plane in logarithmic scale.(right) Distribution along the calorimeter depth in linear scale.

ment. The service allows the user to be abstracted from the underlying infrastructure detail, since hardware is a managed component of the service stack. This means, essentially, that the job needs to be configured while the provisioning of the compute cluster is entirely operationalised by the Azure Machine Learning service, as part of the created deployment workflow. The interaction with the Azure ML service can be done using different ways such as Web interface, CLI and SDK. In this case, the SDK approach is preferred because of its wider support. It provides a space, called Workspace, where the user can manage the different run submissions. A detailed description of the general architecture can be found in [26]. For our tests, we use GPU-powered nodes with 24 vCPU cores, 448 GiB memory and 4 V100 GPU. The 3DGAN training time scales very close to linear up to a total of 16 nodes (64 GPUs). Azure ML automatically optimises the data set management, in aspects such as caching, pre-fetching and parallel data loading according to the hardware infrastructure setup. In figure 6, right panel, we verify that manually modifying caching, pre-fetching and auto-tuning features in Tensorflow does not improve on the train time performance. In terms of image fidelity, the work in [15] has shown that data parallel training can affect the quality of the images especially at the border of the detector sensitive region. Indeed, in figure 7, on the left, the calorimeter energy response in the detector transverse plane shows a similar performance degradation at the edges of the energy distribution. The amount of energy deposited in this region is orders of magnitudes smaller than in the central pixels; it is therefore much harder for the GAN to correctly predict it. It should be noted, however, that this behaviour does not critically affect physics applications, since such small energy deposits are usually neglected. On the other hand, the slight shift of the longitudinal energy distribution, on the right in figure 7, is significant and additional optimisation is needed in order to reduce this effect.

The experiments on The CERN openlab Azure enrollment were performed using an Azure credits grant, currently evolving to provide capabilities for realistic analysis and overall resource costing optimisation. These capabilities were not yet ready at the time this paper was produced and for this reason we defer any cost analysis to future studies.

## 6 Conclusions and Future Plans

With Deep Learning models in the HEP domain becoming more complex, computational requirements increase, triggering the need to consider HPCaaS and MLaaS offers in the public cloud. This work presents results of the first deployment of a three-dimensional convolutional GAN for detector simulation on TPUs, demonstrating an efficient parallelization of the adversarial training process: the 3DGAN training time is brought down from about a week to around one hour. This result enables large architecture hyper-parameter scans to run in just a few days (instead of weeks) and, therefore, it greatly extends the range of detector geometries than can be simulated by the 3DGAN model. Validation against Monte Carlo shows a slight over-estimation of the energy distribution along the edges of the detector sensitive volume for a large number of nodes (above 64 GPUs): simple strategies, typically employed to ease the effect of increasing global batch sizes are not sufficient to improve convergence at scale. We are currently investigating several different directions, such as the role played by the batch normalisation layer, which is one of the components greatly affecting the 3DGAN model performance [6]. In addition to the physics result, this work demonstrates the deployment of scientific DL workloads using public cloud services, complementing the on-premises infrastructure offers in use by research organisations. This is achieved by exploring innovative, open hybrid provisioning and orchestration models, profiting from the technological potential of each cloud provider, unlocked from any particular commercial vendor or provisioning model. Although public cloud services are offered in the general market as commodity services, they must be validated for research use cases technically but also in terms of cost optimisation requirements. These elements combined, as part of a global strategy when scaling out on-demand heterogeneously, can profit from the full potential of cloud-based services that are evolving from basic elastic provisioning of virtual resources, to a transparent and adaptive smart continuum. These possibilities must be assessed with close alignment across performance, usability and economics, considering complexities such as type of resources, data movement patterns, optimisation of architectures adapted to the scientific workloads, provisioning lock-in risks in and also data processing aspects to achieve a long-term sustainable position. Future work will continue to explore models of resources combination from specific vendors and explore seamless orchestration of resources federated across multiple cloud providers, combining technical and cost-effectiveness requirements. In addition, data governance aspects will also be considered, ensure the presence of a set of technical measures to guarantee adequacy with current legislation (e.g. GDPR and Free Flow of data).

## References

[1] J. Albrecht et al., Computing and Software for Big Science **3**, 7 (2019)
[2] M. Erdmann et al., arXiv preprint arXiv:1807.01954 (2018)
[3] V. Chekalina et al., arXiv preprint arXiv:1812.01319 (2018)
[4] Erdmann et al., Comput Softw Big Sci 3 (2019)
[5] A. Ghosh (ATLAS Collaboration), Tech. rep., CERN, Geneva (2019)
[6] G. Khattak et al., in *18th IEEE ICMLA Conference* (2019)
[7] R. Rocha, L. Heinrich, Kubecon - CloudNativeCon Europe 2019 (2019)

[8] I.J. Goodfellow et al., arXiv preprint arXiv:1406.2661 (2014)

[9] M. Arjovsky et al. (2017), `1701.07875`

[10] T. Karras et al., arXiv preprint arXiv:1710.10196 (2017)

[11] L. de Oliveira et al., arXiv preprint arXiv:1701.05927 (2017)

[12] M. Paganini et al., arXiv preprint arXiv:1705.02355 (2017)

[13] G.R. Khattak et al., in *25th IEEE ICIP conference* (2018)

[14] M. Mustafa et al., Computational Astrophysics and Cosmology **6**, 1 (2019)

[15] S. Vallecorsa et al., in *High Performance Computing* (2018), Vol. 11203

[16] A. Author1, *Title3: removed because of ISC review policy* (Address2, 2019)

[17] S.R. Nandakumar et al., Frontiers in Neuroscience **14** (2020)

[18] J. Osorio et al. (2020)

[19] A. Odena et al., arXiv preprints arXiv:1610.09585 (2016)

[20] *Geant* (2020), accessed Dec 20, 2020, `http://geant.cern.ch/`

[21] F. Rehm et al., in *Proceedings of the 10th ICPRAM conference* (2021)

[22] J.R. Vlimant et al., in *EPJ Web of Conferences* (2019), Vol. 214

[23] M. Abadi et al., accessed Dec 20, 2020, `https://www.tensorflow.org/`

[24] *Kubeflow* (2020), [Online; accessed Dec 21 2020], `https://www.kubeflow.org/`

[25] *Cloud Broker Project (CloudBank EU)* (2020), accessed Dec 20, 2020, `https://indico.cern.ch/event/919839/contributions/4147766/attachments/2162052/3648772/Cloud_Broker_Pilot.pdf`

[26] *Microsoft Azure Machine Learning service* (2020), accessed Dec 20, 2020, `https://azure.microsoft.com/en-us/services/machine-learning/`