

Introduction

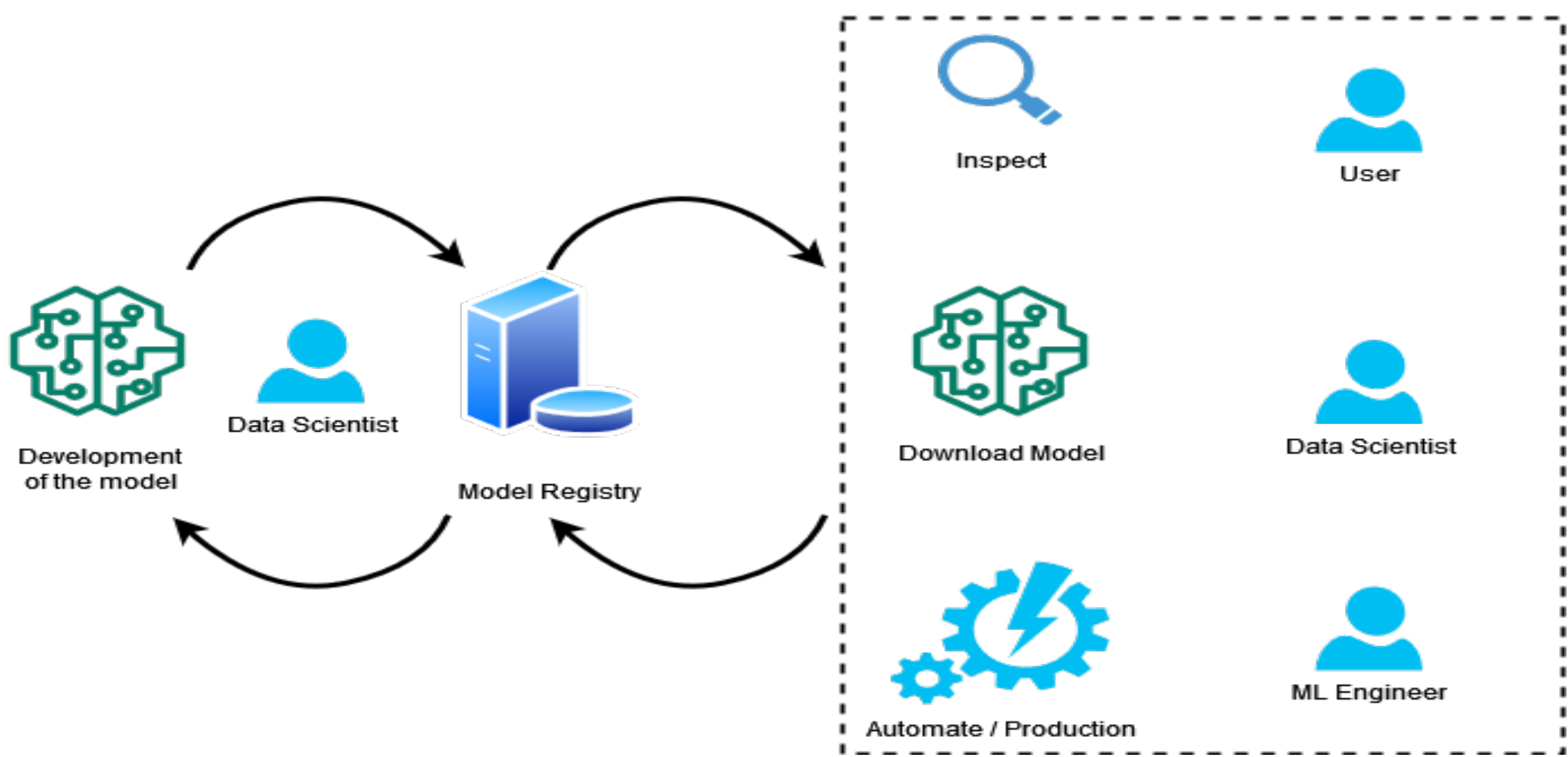
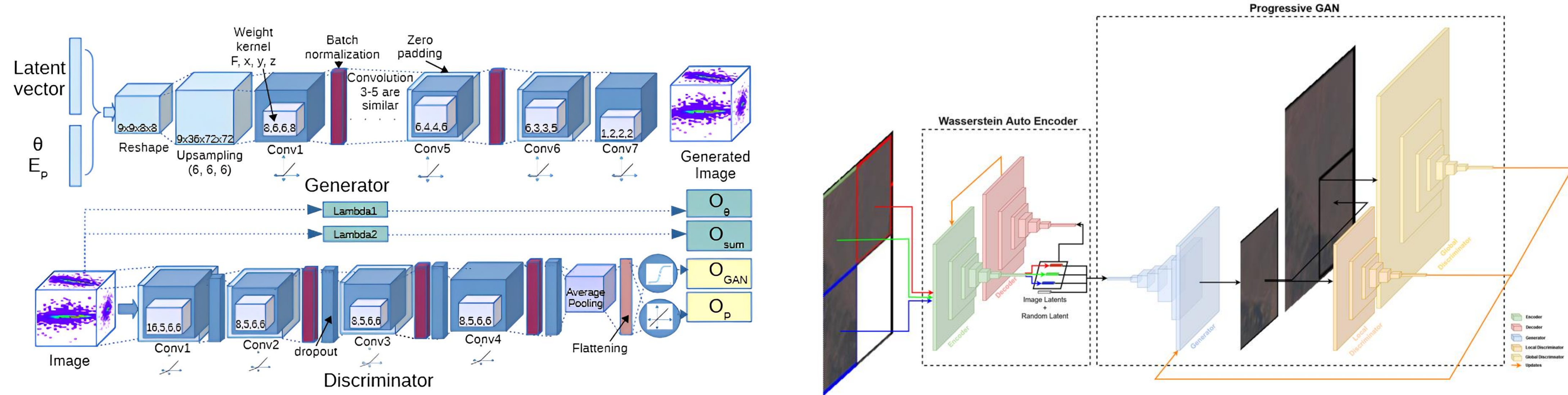
The use of machine learning at CERN is increasing and moving toward more complex algorithms. **Development and optimization become more challenging**, while **model generalization and re-usability** turn critical. To address these challenges, it is crucial to have efficient and flexible ways to **track changes and monitor performance metrics**.

The Oracle Accelerated Data Science SDK (ADS) is a Python library that is part of the OCI Data Science and provides intuitive access to a **Model Catalog** to **store, evaluate, monitor and train machine learning models**.

Models

We study two use cases, both developed in Tensorflow 2 and train using data parallelism:

- A 3D convolutional GAN (**3DGAN**) to generate calorimeter images [1], composed of **a few million parameters**, that trains in **4 days to train** on 1 GPU.
- A **hybrid Variational AutoEncoder GAN models (VAE-ProGAN)** that generates satellite images [2], contains **~80 million parameters**, and trains in **~1 month** on 1 GPU (using close to 32 GB of GPU memory).



What is a Model Catalog and How is it used

A Model Catalog is a centralized repository that enables the publication of production-ready models through a central user interface to facilitates model search, review, and documentation access. Furthermore, it serves as a collaborative platform for managing the lifecycle of all models, accelerating model deployment, streamlining lifecycle management, and improve governance. Thus, utilization of a Model Catalog should result in resource saving, faster turnaround and efficient knowledge transfer.

We used the Model Catalog available on Oracle ADS to enable energy cost analysis for the 2 use cases.

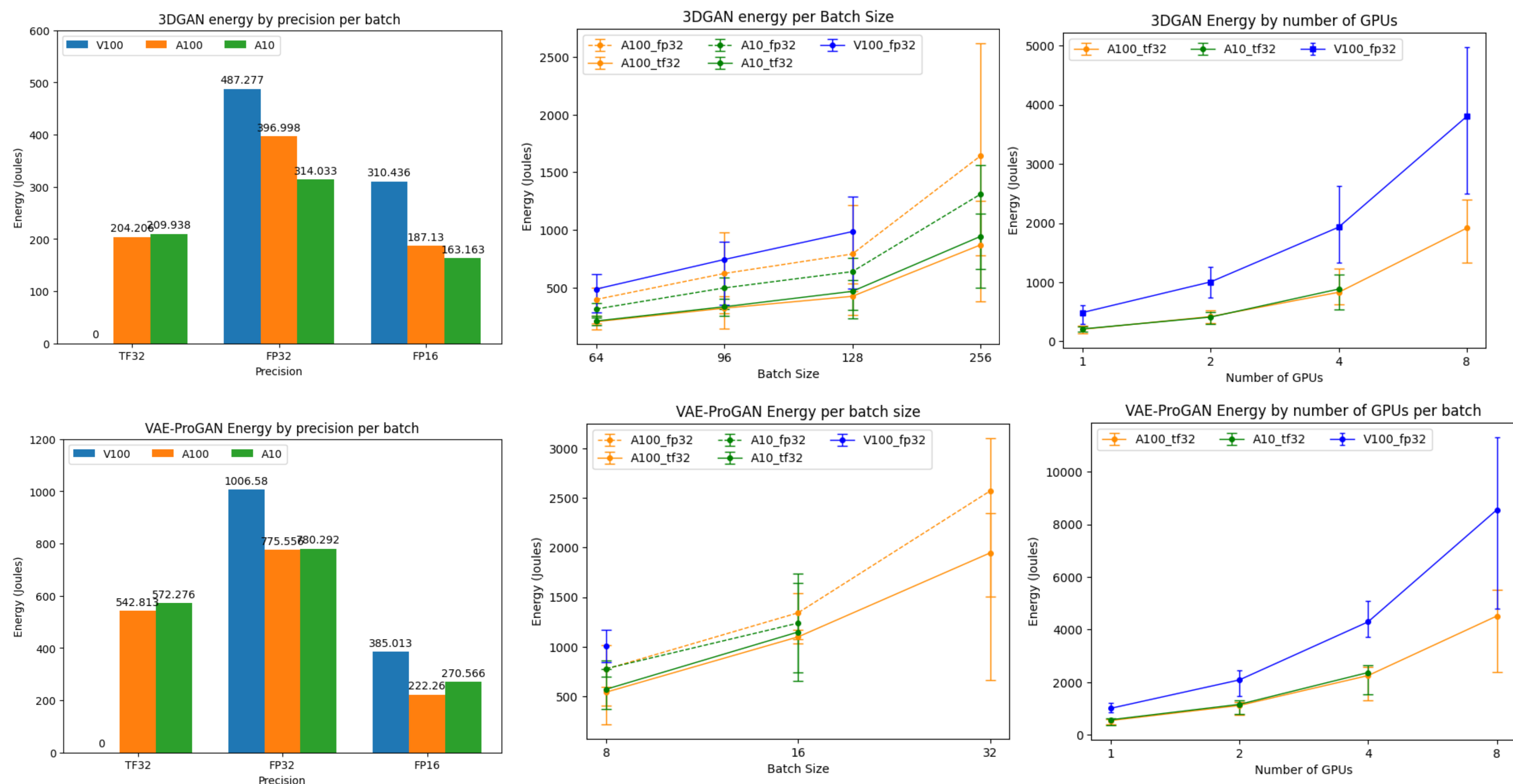
Energy cost analysis

We have previously studied [3] training times and associated costs has crucial metric for models' re-usability. With environmental concerns [4] and growing computational power utilization [5, 6], it is essential to also investigate **energy consumption** during training [7, 8].

Following [7], we analyzed **different GPU types, precision, batch size, and number of GPUs**: in particular we calculate the energy cost per batch averaging over 10 batches.

Results show similar values for A100 and A10 with **A100 having a slight advantage when using TF32 while being substantially faster**. Mixed FP16 shows lower consumption than TF32, but it could induce a loss in accuracy. We used the power obtained from the GPUs in Watts and the time for batch training to calculate **the total energy required for full training**, using **A100 with TF32: 46.3MJ** for 3DGAN and **778.9 MJ** for VAE-ProGAN.

In summary, A100 with TF32 delivers superior energy efficiency while maintaining loss accuracy and fast training times. It also allows for larger batch sizes and larger models to be trained. Further testing is required to verify the results for Physics Accuracy.



Energy measurements with respect to **Precision** (left), **Batch Size** (Middle), **Number of GPUs** (Right), for V100, A100 and A10. All measurements were calculated for a single batch size taking the average of multiple calculations.

GPU RAM limits on V100 (16GiB) and A10 (24 GiB) prevented proportional batch size increases compared to A100 (40 GiB).

Conclusion and Future Plans

The increase of computational power and its environmental effects have become a pressing issue for ML development, that could be eased by further advancement of model reusability and generalization. Tools that enable this strategies in complex, collaborative environments exist. We have used the Oracle ADS model catalog, to evaluate the energy cost of training our models, identifying the most energy efficient way of deployment the training process for the two GAN models, which exhibit very different computation requirements.

References

- [1] Khattak, G. R., Vallecorsa, S., Carminati, F., & Khan, G. M. (2021). Fast Simulation of a High Granularity Calorimeter by Generative Adversarial Networks. ArXiv [Physics.Ins-Det]. Retrieved from <http://arxiv.org/abs/2109.07388>
- [2] Cardoso, R., Vallecorsa, S., & Nemni, E. (2022). Conditional Progressive Generative Adversarial Network for satellite image generation. ArXiv [Cs.CV]. Retrieved from <http://arxiv.org/abs/2211.15303>
- [3] Cardoso, R., Golubovic, D., Lozada, I. P., Rocha, R., Fernandes, J., & Vallecorsa, S. (2021). Accelerating GAN training using highly parallel hardware on public cloud. ArXiv [Cs.LG]. Retrieved from <http://arxiv.org/abs/2111.04628>
- [4] Anthony, L. F. W., Kanding, B., & Selvan, R. (2020). Carbontracker: Tracking and Predicting the Carbon Footprint of Training Deep Learning Models. ArXiv [Cs.CY]. Retrieved from <http://arxiv.org/abs/2007.03051>
- [5] Amodei, D. and Hernandez, D. (2018) Ai and Compute, AI and compute. Available at: <https://openai.com/research/ai-and-compute> (Accessed: May 4, 2023).
- [6] Hernandez, D., & Brown, T. B. (2020). Measuring the Algorithmic Efficiency of Neural Networks. CoRR, abs/2005.04305. Retrieved from <https://arxiv.org/abs/2005.04305>
- [7] Desislavov R, Martínez-Plumed F, Hernández-Orallo J. Trends in AI inference energy consumption: Beyond the performance-vs-parameter laws of deep learning. Sustainable Computing: Informatics and Systems. 2023;38:100857. doi:10.1016/j.suscom.2023.100857
- [8] Garcia-Martin E, Rodrigues CF, Riley G, Grahm H. Estimation of energy consumption in machine learning. Journal of Parallel and Distributed Computing. 2019;134:75-88. doi:10.1016/j.jpdc.2019.07.007