# CaloChallenge 2022 — Final Evaluation and Lessons Learned
## — ML4Jets 2024 Paris —

Claudius Krause

Institute of High Energy Physics (HEPHY), Austrian Academy of Sciences (OeAW)
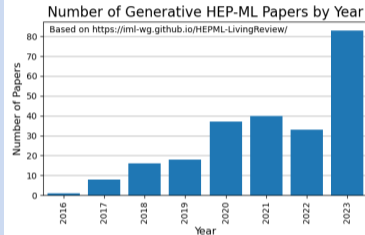
November 4, 2024

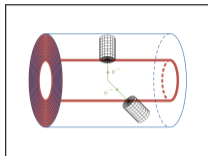$$\Rightarrow \text{arXiv:2410.21611} \Leftarrow$$

# It all started in 2021 …

- … the LHC-Olympics had just concluded.

- Generative AI was kicking off in HEP in 2020.

- Applications to Detector Simulation, as major bottleneck, were gaining popularity.

- However, $\mathcal{O}(10)$ architectures used $\mathcal{O}(8)$ datasets.

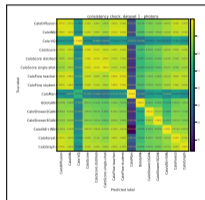$\Rightarrow$ We created the CaloChallenge to benchmark and trigger new developments.



Number of Generative HEP-ML Papers by Year
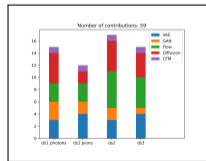Based on https://iml-wg.github.io/HEPML-LivingReview/
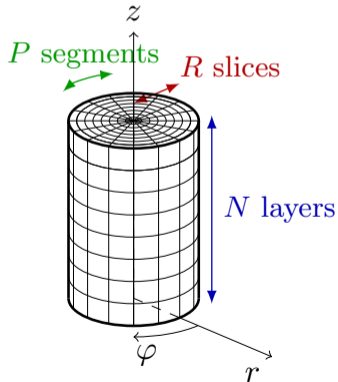
I: Datasets



II: Evaluation Metrics



III: Results

# CaloChallenge Showers are voxelized in cylindrical coordinates.

- There 4 datasets in increasing complexity / dimensionality.
- Particles enter perpendicular to front surface:

# CaloChallenge Showers are voxelized in cylindrical coordinates.

- Showers are usually sparse.
- Energy depositions span several orders of magnitude.

Photon shower at E = 1.0 GeV



Layer 0    Layer 1    Layer 2    Layer 3    Layer 12

$10^{-2}$    $10^{-1}$    $10^{0}$    $10^{1}$

# CaloChallenge Showers are voxelized in cylindrical coordinates.

- Showers are usually sparse.
- Energy depositions span several orders of magnitude.

Photon shower at E = 1048.6 GeV



Layer 0  Layer 1  Layer 2  Layer 3  Layer 12

$10^{-2}$  $10^{-1}$  $10^{0}$  $10^{1}$  $10^{2}$  $10^{3}$  $10^{4}$

# The Fast Calorimeter Simulation Challenge 2022

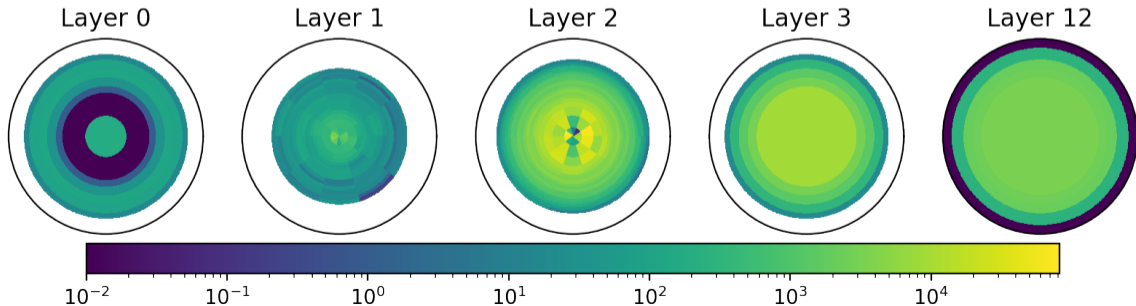The main task: | Develop a model that samples from $p(\text{shower}|E_{\text{incident}})$

https://calochallenge.github.io/homepage/

Michele Faucci Giannelli, Gregor Kasieczka, **CK**, Ben Nachman, Dalila Salamani, David Shih, and Anna Zaborowska

- Dataset 1:  AtlFast3 trainig data  ($\gamma$: 368, $\pi$: 533 voxels)
  [2109.02551, Comput.Softw.Big Sci.]  $E_{\text{inc}} \in [256\ \text{MeV}, 4.2\ \text{TeV}]$

- Dataset 2:  Par04 simulated detector  ($e^-$: 6480 voxels)  $E_{\text{inc}} \in [1\ \text{GeV}, 1\ \text{TeV}]$

- Dataset 3:  Par04 simulated detector  ($e^-$: 40500 voxels)  $E_{\text{inc}} \in [1\ \text{GeV}, 1\ \text{TeV}]$

I: Datasets



II: Evaluation Metrics



III: Results

# How to evaluate generative models?

In text / image / video generation: "by eye".
⇒ Our brains are incredible good at this task, but it doesn't scale.



imagined with Meta AI.

In high-energy physics: need to find something better!
⇒ We want to correctly cover $p(x)$ of the entire phase space.

1. Can look at histograms of derived features / observables.
⇒ To quantify, we use the *separation power* of high-level feature histograms:

$$S(h_1, h_2) = \frac{1}{2} \sum_{i=1}^{n_{\text{bins}}} \frac{(h_{1,i} - h_{2,i})^2}{h_{1,i} + h_{2,i}}$$

But: this is just a 1-dim projection!

# A Classifier provides the "ultimate metric".

According to the Neyman-Pearson Lemma we have:

- The likelihood ratio is the most powerful test statistic to distinguish two samples.
- A powerful classifier trained to distinguish the samples should therefore learn (something monotonically related to) $w = \frac{p_{\text{data}}}{p_{\text{model}}}$.
- If this classifier is confused, we conclude $\Rightarrow$ $p_{\text{data}}(x) = p_{\text{model}}(x)$
- $\Rightarrow$ This captures the full phase space incl. correlations. CK/D. Shih [2106.05285, PRD]

❷ Now, the AUC provides a single number to compare different models.

But: are AUCs of different models really comparable?

# A Classifier tells us much more about the model.

Failure modes of the model can now be seen in the $w = \frac{p_{\text{data}}}{p_{\text{model}}}$ histogram:

Data manifold over-
populated by model:
$\Rightarrow$ missmodeled feature



Data manifold not
populated by model:
$\Rightarrow$ missed feature

R. Das et al. [2305.16774, SciPost]

Cluster plots show where events lie in phase space:

figures by B. Schmidthaler / M. Rosendorf





small weights:

large weights:

# Other quality metrics we looked at.

**④** *KPD/FPD.*  <span style="color:gray">Kansal et al. [2211.10295, Phys.Rev.D]</span>

- Fréchet physics distance (FPD): Fréchet distance between Gaussian fits to obs.
- Kernel physics distance (KPD): kernel-based MMD between observables.

**⑤** *Pearson Correlation* between layer energies.  <span style="color:gray">Ahmad et al. [2406.12898]</span>

**⑥** *Precision / Recall / Density / Coverage.*  <span style="color:gray">Naeem et al. [2002.09797]</span>

- How many "real" samples are close to "fake" manifold.
- How many "fake" samples are close to "real" manifold.

# Other important metrics to look at.

$\Rightarrow$ The *generation time*.
- on CPU/GPU architectures
- for batch sizes 1 / 100 / 10000

$\Rightarrow$ The *number of trainable parameters*.
- as proxy for model size
- in training / generation

# Other important metrics to look at.

$\Rightarrow$ The *generation time*.
- on CPU/GPU architectures
- for batch sizes 1 / 100 / 10000

$\Rightarrow$ The *number of trainable parameters*.
- as proxy for model size
- in training / generation

- start singularity container
- load model weights + biases
- generate samples
- save them to `.hdf5`

# CaloChallenge 2022 — Final Evaluation and Lessons Learned

I: Datasets



II: Evaluation Metrics



III: Results

# The ~~preliminary~~ final! results of the CaloChallenge

I will only be able to share some highlights of the results of the CaloChallenge.

The final write-up, **arXiv:2410.21611**, has a lot more content!



- We received 59 submissions for all datasets.

- They were generated by 23 different models.

- All types of generative AI architectures were used.

# Comparing different quality metrics: high-level features



Correlation of high-level binary AUC to sum of separation powers, dataset 2

Legend:
- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDREAM

Scores correlate strongly.

# Comparing different quality metrics: classifier input



Correlation of low-level binary AUC to high-level binary AUC, dataset 1 - pions

Legend:
- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- DNN CaloSim
- CaloShowerGAN
- CaloVAE+INN
- CaloForest
- CaloGraph

↙ better

Scores correlate strongly, but 2 lines form.
Interestingly: along the type of architecture!

# Comparing different quality metrics: classifier architecture



Correlation of low-level binary AUC to CNN ResNet binary AUC, dataset 2

Legend:
- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDiT
- CaloDREAM

CNN ResNet is much better classifier, but correlation is still strong.

# Comparing different quality metrics: binary vs. multiclass



Correlation of low-level binary AUC to multiclass log-posterior, dataset 2

Legend:
- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDiT
- CaloDREAM

Very clear correlation, both seem to capture similar things!

# Comparing different timing metrics: CPU vs. GPU

Correlation of CPU to GPU generation times at batch size 100, dataset 3

GPU geneneration time [ms] vs. CPU geneneration time [ms]

↙ better

Legend:
- CaloDiffusion
- L2LFlows MAF
- conv. L2LFlows
- MDMA
- CaloClouds
- Calo-VQ
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- GEANT4 transformer
- CaloPointFlow
- CaloVAE+INN
- Calo-VQ(norm)
- CaloDREAM

GPU much faster, but times correlate.

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - photons

Legend:
- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- CaloShower2GAN
- CaloShower3GAN
- CaloVAE+INN
- CaloGraph

Axis labels: multiclass log-posterior (y-axis); GPU generation time, batch size 100, in ms (x-axis); ↖ better

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - photons

- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- CaloShower2GAN
- CaloShower3GAN
- CaloVAE+INN
- CaloGraph

Diffusion models are good, but slow.

Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - photons

- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- CaloShower2GAN
- CaloShower3GAN
- CaloVAE+INN
- CaloGraph

VAEs and GANs are fast, but not as good

ÖAW
AUSTRIAN
ACADEMY OF
SCIENCES

HEPHY
INSTITUTE OF
HIGH ENERGY PHYSICS

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - photons

Legend:
- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- CaloShower2GAN
- CaloShower3GAN
- CaloVAE+INN
- CaloGraph

Normalizing Flows sit in the sweet spot!

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - pions

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - pions

Legend:
- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- DNN CaloSim
- CaloShowerGAN
- CaloVAE+INN
- CaloGraph

y-axis: multiclass log-posterior

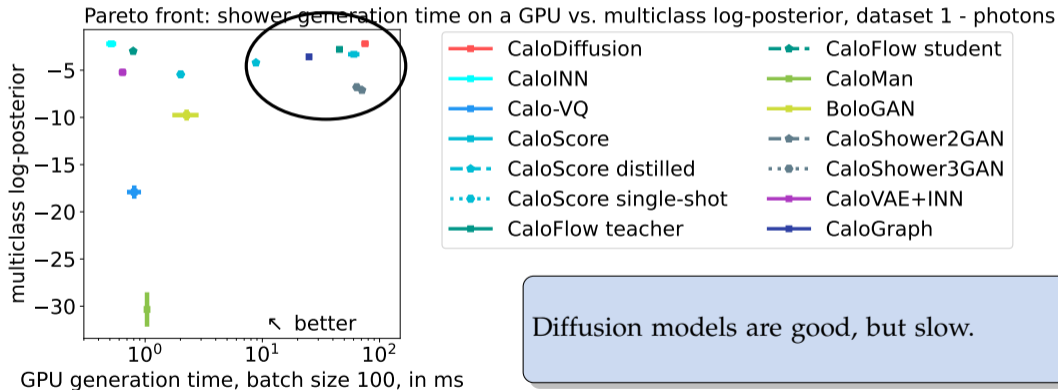x-axis: GPU generation time, batch size 100, in ms

↖ better

Diffusion models are again good, but slow.

# Pareto Fronts: Quality vs. Generation Time

Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 1 - pions



Legend:
- CaloDiffusion
- CaloINN
- Calo-VQ
- CaloFlow teacher
- CaloFlow student
- CaloMan
- BoloGAN
- DNN CaloSim
- CaloShowerGAN
- CaloVAE+INN
- CaloGraph

Normalizing Flows are still strong, but a VAE wins.
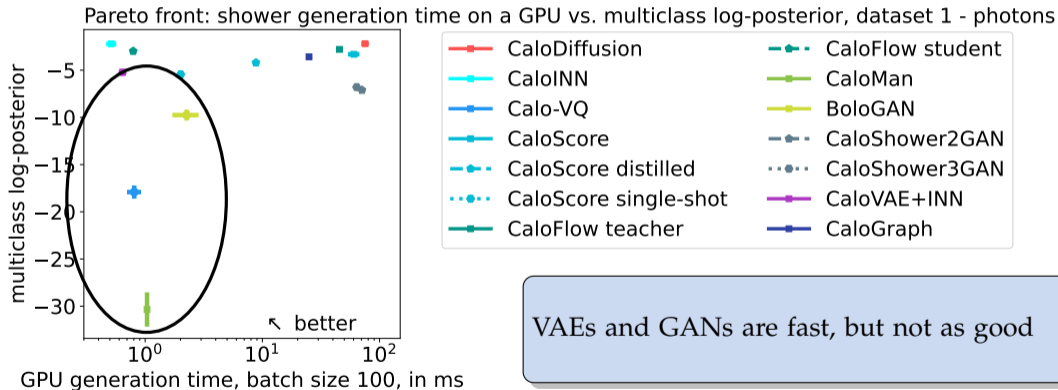
# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 2
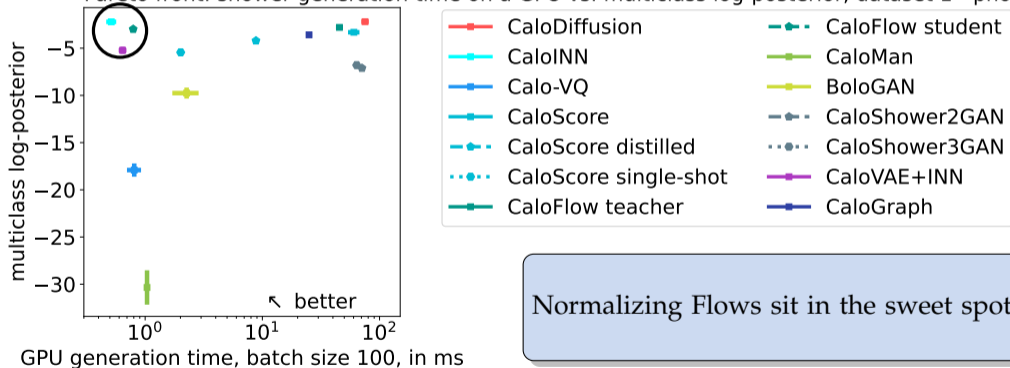
# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 2

Legend:
- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDiT
- CaloDREAM

Again, a similar cluster for Diffusion models up here.

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 2

Legend:
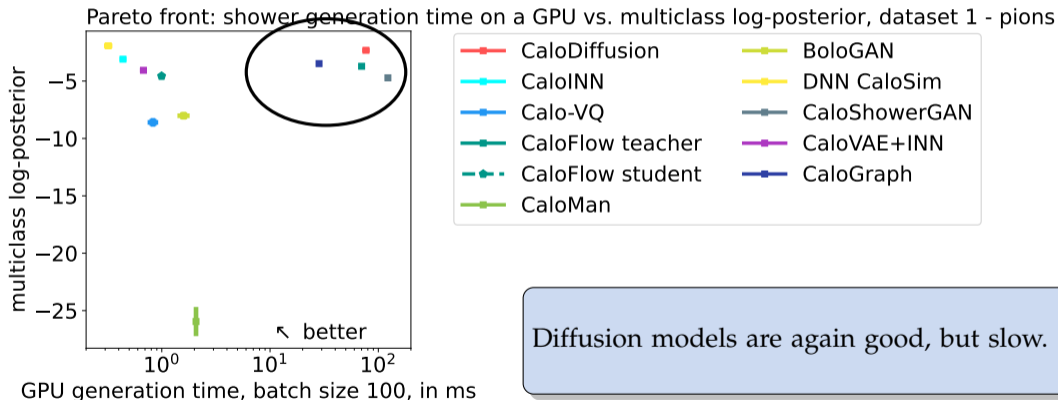- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDiT
- CaloDREAM

And a group of VAEs and GANs here

# Pareto Fronts: Quality vs. Generation Time

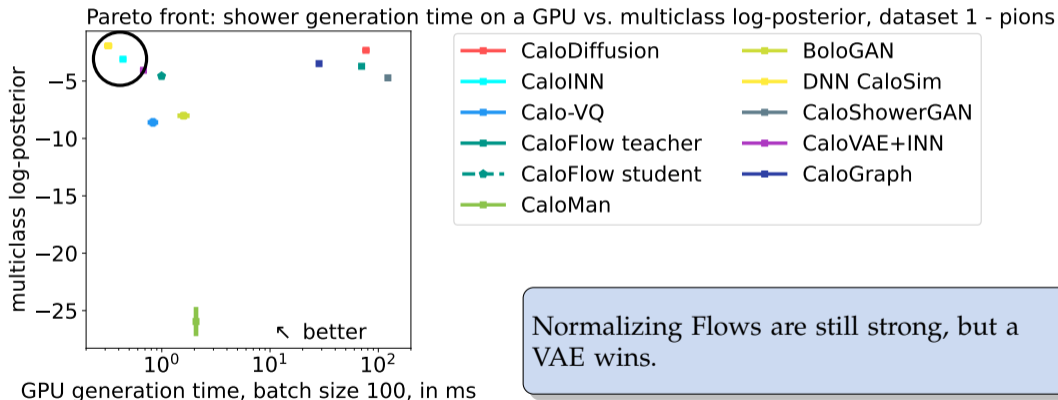Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 2



Legend:
- CaloDiffusion
- conv. L2LFlows
- CaloINN
- MDMA
- Calo-VQ
- CaloScore
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- SuperCalo
- DeepTree
- CaloPointFlow
- CaloVAE+INN
- CaloLatent
- CaloDiT
- CaloDREAM
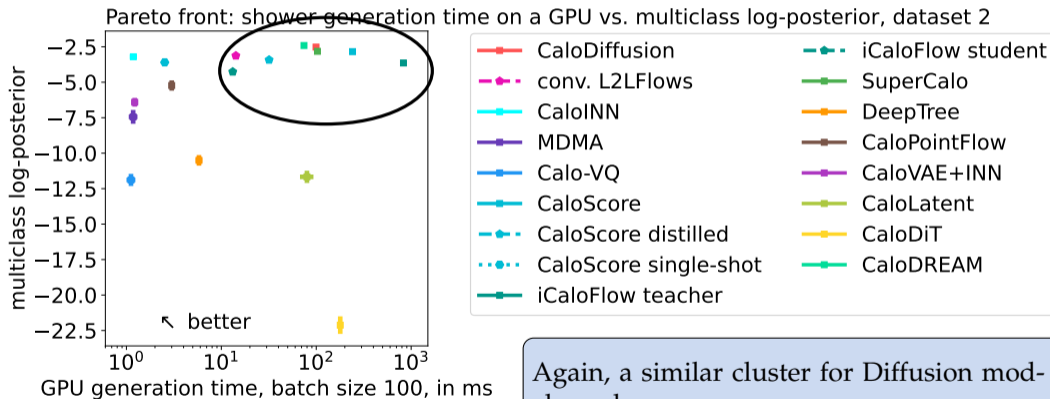
And a Normalizing Flows in the corner

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 3

Legend:
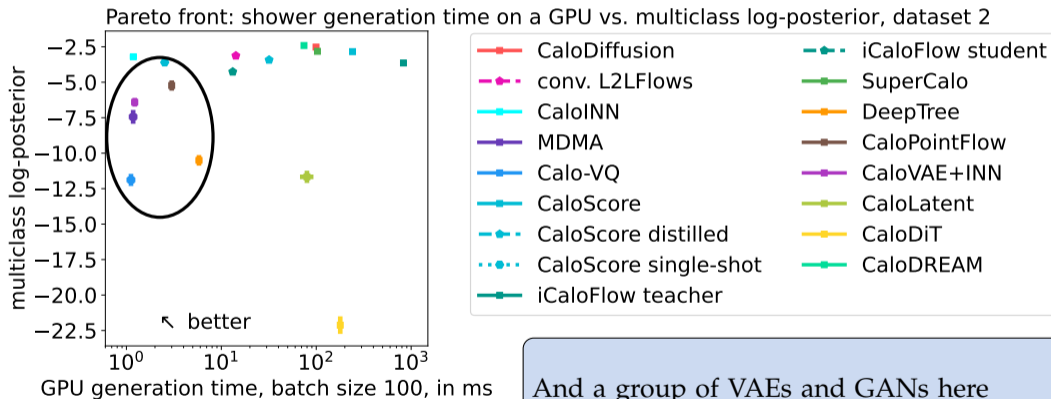- CaloDiffusion
- L2LFlows MAF
- conv. L2LFlows
- MDMA
- CaloClouds
- Calo-VQ
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- GEANT4 transformer
- CaloPointFlow
- CaloVAE+INN
- Calo-VQ(norm)
- CaloDREAM

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 3

Legend:
- CaloDiffusion
- L2LFlows MAF
- conv. L2LFlows
- MDMA
- CaloClouds
- Calo-VQ
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- GEANT4 transformer
- CaloPointFlow
- CaloVAE+INN
- Calo-VQ(norm)
- CaloDREAM

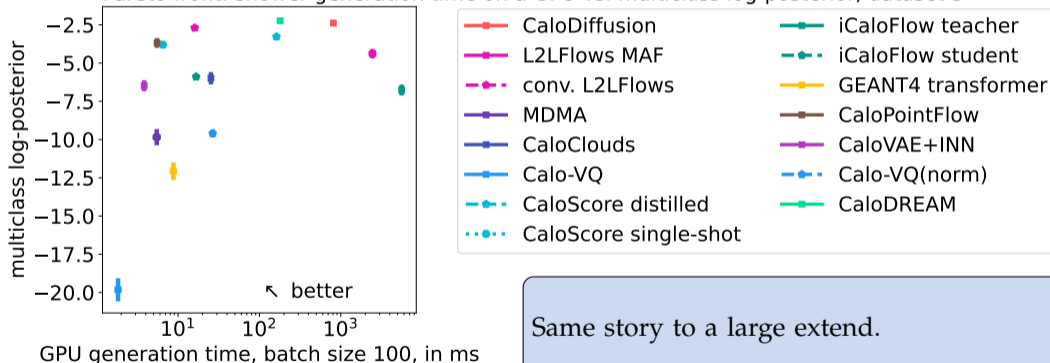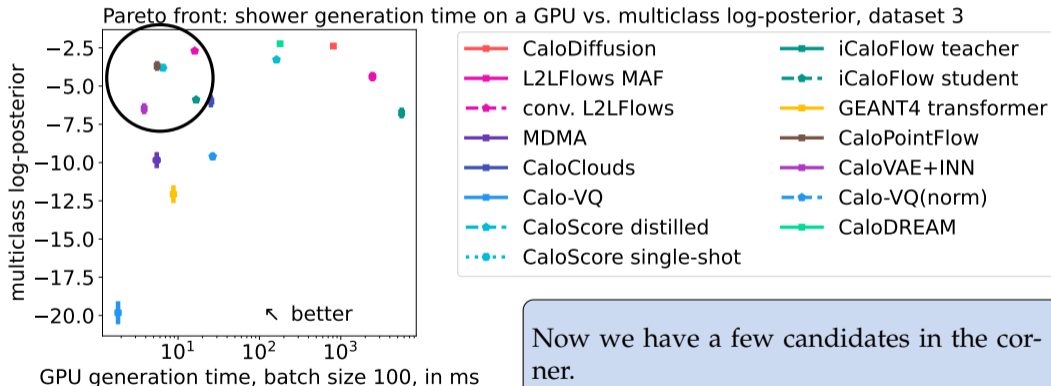Same story to a large extend.

# Pareto Fronts: Quality vs. Generation Time



Pareto front: shower generation time on a GPU vs. multiclass log-posterior, dataset 3

Legend:
- CaloDiffusion
- L2LFlows MAF
- conv. L2LFlows
- MDMA
- CaloClouds
- Calo-VQ
- CaloScore distilled
- CaloScore single-shot
- iCaloFlow teacher
- iCaloFlow student
- GEANT4 transformer
- CaloPointFlow
- CaloVAE+INN
- Calo-VQ(norm)
- CaloDREAM

Now we have a few candidates in the corner.

# CaloChallenge 2022 — Final Evaluation and Lessons Learned

The CaloChallenge was well-received in the community:

- 20+ papers
- Even more talks at ML4Jets / ML 4 Physical Sciences@NeurIPS / CHEP / …
- Many discussions and feedback on evaluation metrics etc.
- All repositories are public!

# CaloChallenge 2022 — Final Evaluation and Lessons Learned

The CaloChallenge was well-received in the community:
- 20+ papers
- Even more talks at ML4Jets / ML 4 Physical Sciences@NeurIPS / CHEP / …
- Many discussions and feedback on evaluation metrics etc.
- All repositories are public!

Final evaluation:
- Quality: Diffusion and CFM better than NF better than GAN/VAE.
- Speed: GAN/VAE faster than NF faster than Diffusion and CFM.

# CaloChallenge 2022 — Final Evaluation and Lessons Learned

The CaloChallenge was well-received in the community:
- 20+ papers
- Even more talks at ML4Jets / ML 4 Physical Sciences@NeurIPS / CHEP / …
- Many discussions and feedback on evaluation metrics etc.
- All repositories are public!

Final evaluation:
- Quality: Diffusion and CFM better than NF better than GAN/VAE.
- Speed: GAN/VAE faster than NF faster than Diffusion and CFM.

Lessons Learned:
- Various correlations between quality metrics for all datasets.
- Next step: embedding models in full fast simulation to see how trade-offs play out.

# CaloChallenge 2022 — Final Evaluation and Lessons Learned

The CaloChallenge was well-received in the community:
- 20+ papers
- Even more talks at ML4Jets / ML 4 Physical Sciences@NeurIPS / CHEP / …
- Many discussions and feedback on evaluation metrics etc.
- All repositories are public!

Thank you!

Final evaluation:
- Quality: Diffusion and CFM better than NF better than GAN/VAE.
- Speed: GAN/VAE faster than NF faster than Diffusion and CFM.

Lessons Learned:
- Various correlations between quality metrics for all datasets.
- Next step: embedding models in full fast simulation to see how trade-offs play out.