

ML4Jets 2024 Trip Report

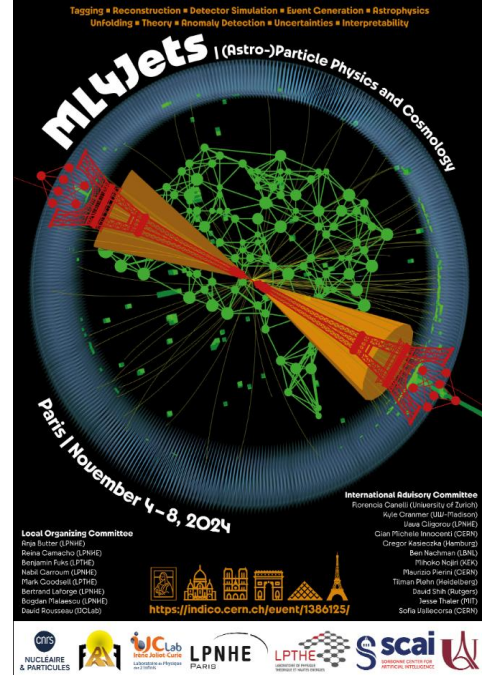
Peter McKeown, Piyush Raikwar

02.12.2024



About the Workshop

- 5 day workshop in Paris (at LPNHE, on the Sorbonne Uni Campus)
- 140 people in person + 150 online
- Tracks:
 - Anomaly detection
 - Astro & Cosmo
 - Detector simulation
 - Event generation
 - Foundation models
 - Reconstruction
 - Tagging
 - Theorie
 - Uncertainties & Interpretation
 - Unfolding & Inference

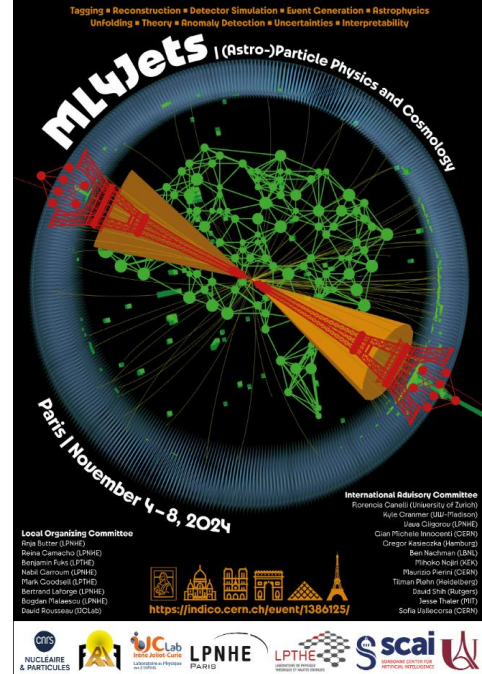


About the Workshop

- 5 day workshop in Paris (at LPNHE, on the Sorbonne Uni Campus)
- 140 people in person + 150 online
- Tracks:

- Anomaly detection
- Astro & Cosmo
- Detector simulation
- Event generation
- Foundation models
- Reconstruction
- Tagging
- Theorie
- Uncertainties & Interpretation
- Unfolding & Inference

Disclaimer: we will show a *heavily* curated and biased selection

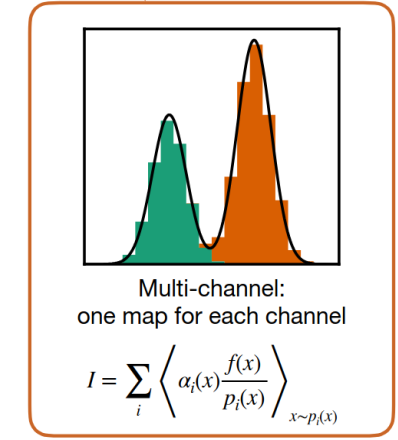
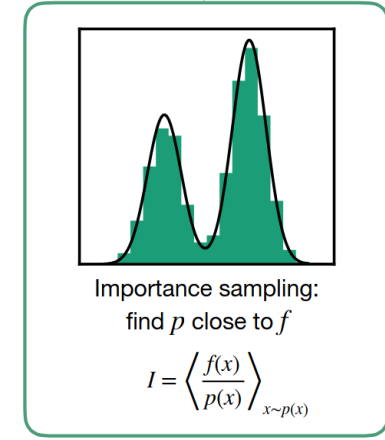
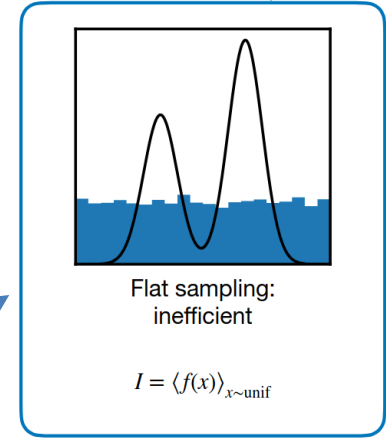


Theory and Event Generators

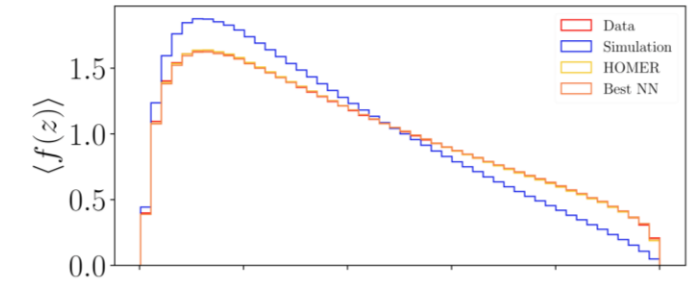
- Normalising Flows (class of invertible network) seem to be king
- Speeding up event generators (~17% ATLAS CPU- [2022 estimate](#))
 - E.g. ML-based Importance Sampling for more efficient phase space integration
 - Plan integration in MadGraph and addition of GPU support
- Hadronization models:
 - ML for tuning existing models (e.g. Lund String model)
 - Or learning a fragmentation model directly from data

Calculate (differential) cross sections

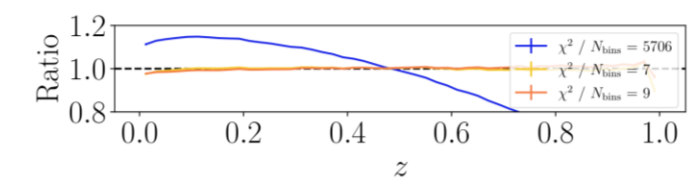
$$d\sigma = \frac{1}{\text{flux}} dx_a dx_b f(x_a) f(x_b) d\Phi_n \langle |M_{\lambda,c,\dots}(p_a, p_b | p_1, \dots, p_n)|^2 \rangle$$



[Talk by R. Winterhalder](#)



[Talk by M. Szewc](#)



Unfolding

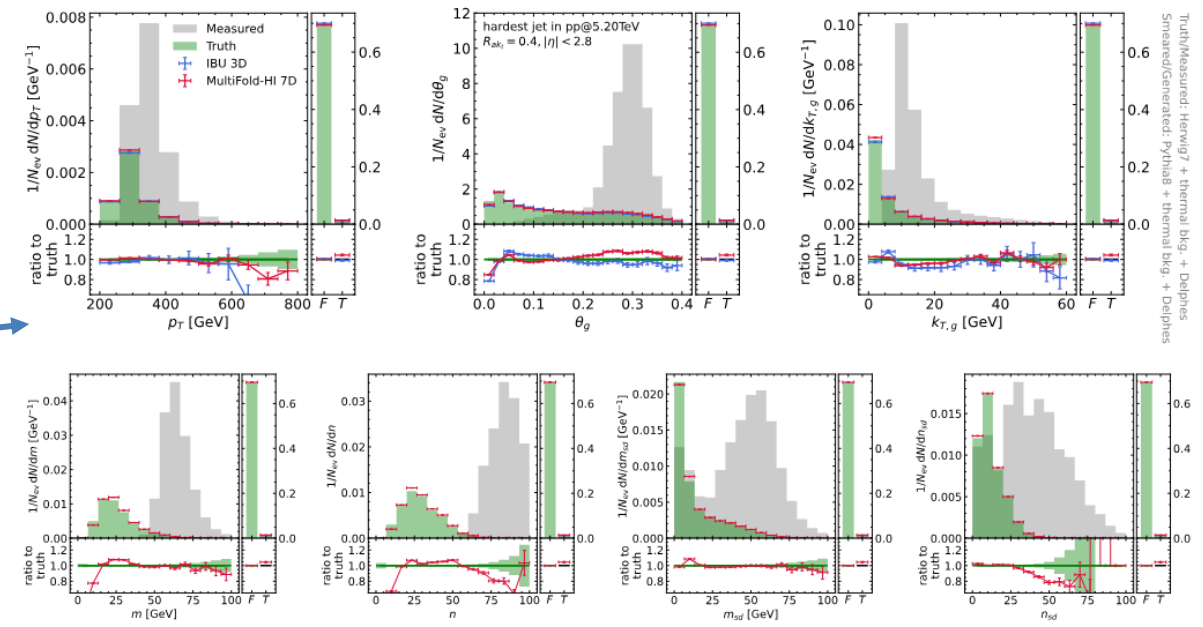
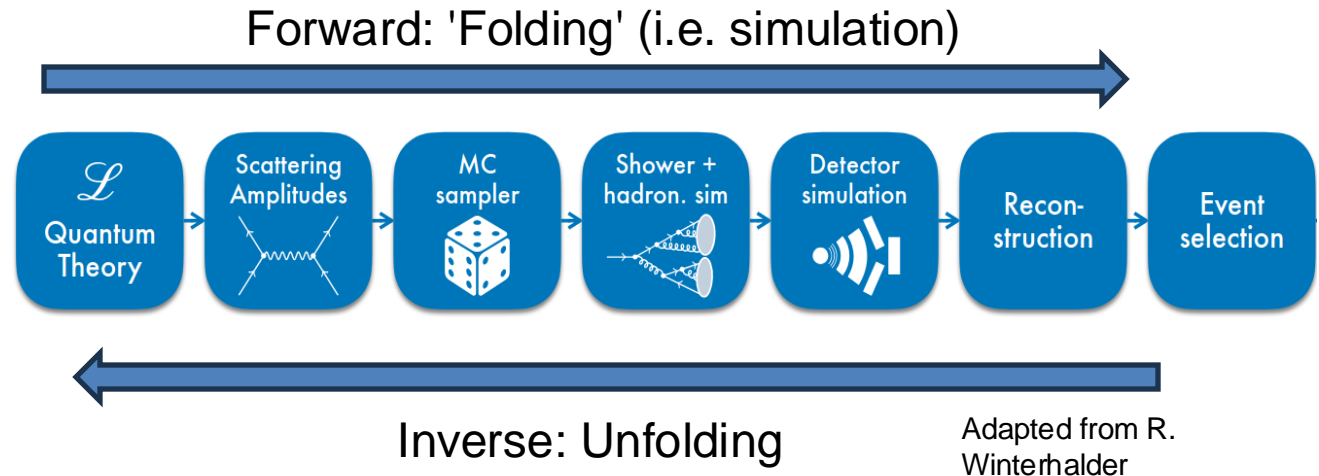
- Unfolding: an inverse problem; removing detector effects from observables- enable comparison to theory

- Limited resolution
- Inefficiency
- Distortions/smearing

- Mapping between distributions with Generative ML/ reweighting procedure

- Things get complicated with Heavy Ions:
 - Can't easily separate underlying event (i.e background) from hard scattering

- ML approaches allow a higher dimensional problem to be tackled and event-wise (unbinned)

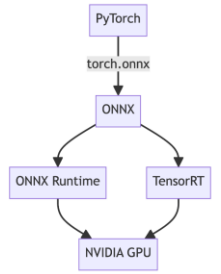


Talk by A. Falcão

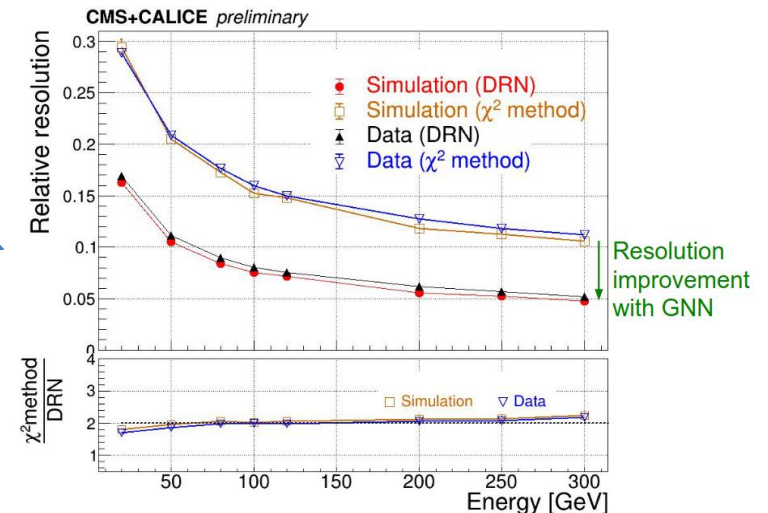
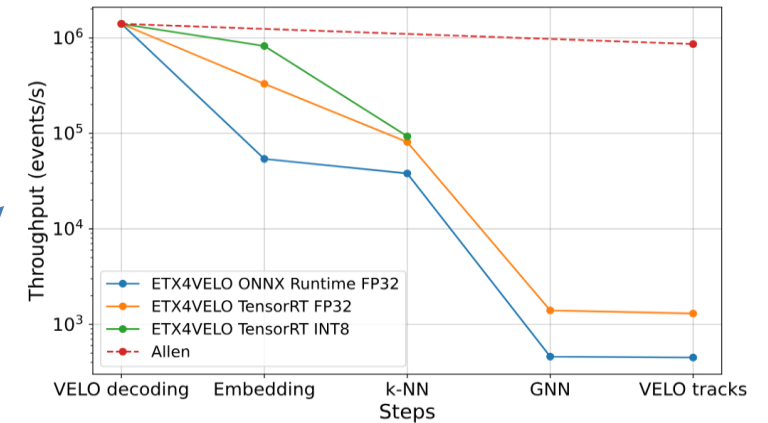


Reconstruction

- Lots of activity: tagging given its own track (4 separate sessions!)
- A few examples:
- Tracking- LHCb have ETX4VELO (Vertex Locator), based on a graph net
 - 'Allen' framework allows event inference in batches on GPU- very high throughput
- Event reconstruction- CMS HGCal endcap
 - Iterative reconstruction pipeline combining classical and ML based algorithms for Particle ID, energy regression etc
 - Some discussion of an 'end-to-end' approach- would bring many challenges (robustness, complexity, interpretability...)



Talk by F. Giasemis



Talk by T. Cuisset



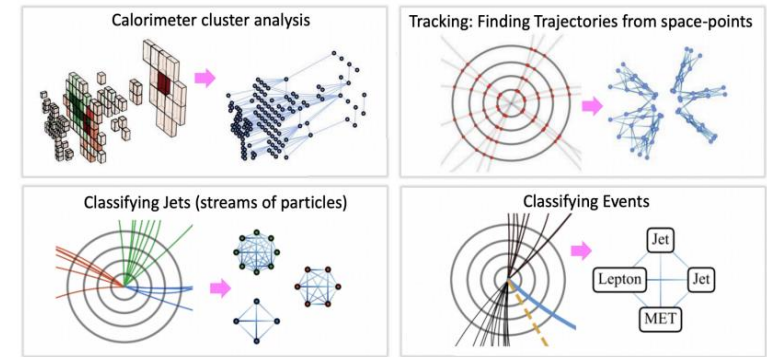
On the Edge

- ML models in an online setting (+ on a specific hardware)
- Eliminate unnecessary information from the model
 - Distillation, pruning, quantization, architecture search
- Models are becoming more and more complex. Harder to keep them fast and accurate
- Looking forward to NGT and hls4ml

Future

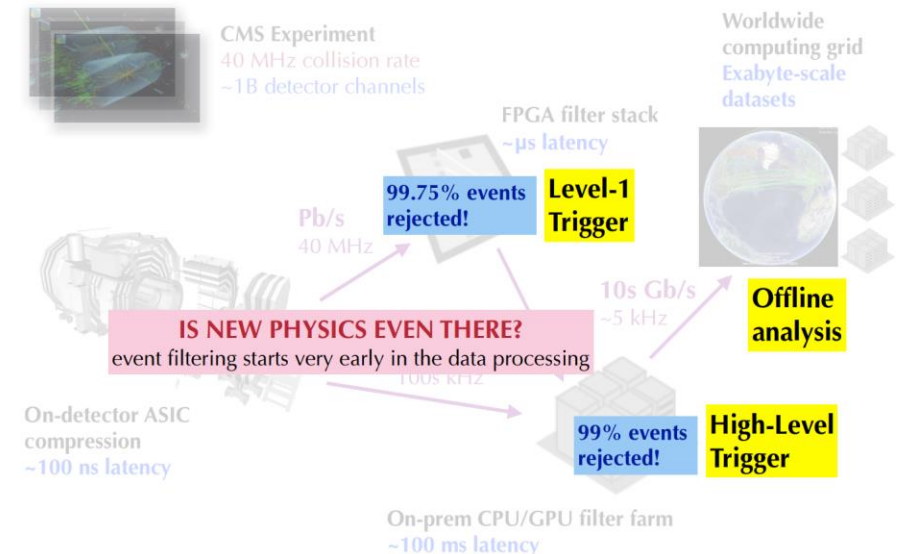
- Extreme-edge: Need for data-driven discovery, beyond experiments' expectations
- Change in setup/distribution, *need for continual learning*

A selection of ML applications, in operation or in development, for online reconstruction
(very much non exhaustive!)



(Summary of various topics from all experiments) [Talk by S. Akar](#)

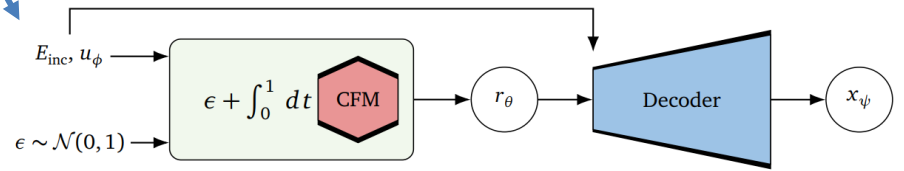
(Edge AI in HEP - highlights) [Talk by J. Ngadiuba](#)



Detector Simulation

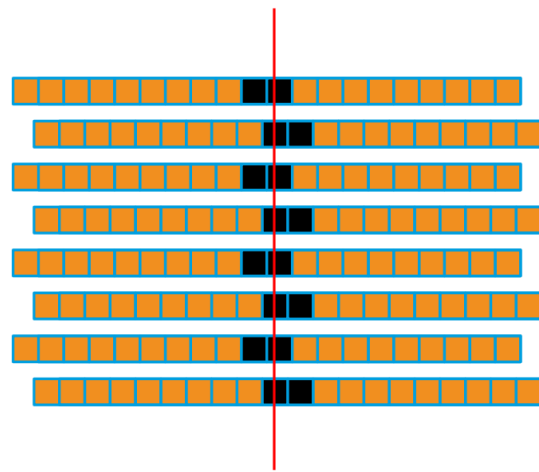
- Presentations from us. It was pretty much our own session!
- Mostly dominated by Flows and Diffusion! Although, a few efforts towards hybrid methods as well
- Community is moving beyond standard image-based ECal setup
 - Considering artifacts due to layer structure
 - High-granularity hadron showers
 - Looking at reconstruction results

Contribution list		Timetable
<div style="display: flex; justify-content: space-between;"> Tue 05/11 Wed 06/11 Thu 07/11 All days </div> <div style="text-align: right;"> Print PDF Full screen Detailed view Filter </div>		
13:00		
14:00	A Library for ML-based Fast Calorimeter Shower Simulation at Future Collider Experiments and Beyond Peter McKeown	
	Towards Detector Agnostic Fast Calorimetry Simulation Salle Séminaires	Piyush Raikwar 14:10 - 14:30
	CaloClouds III: Ultra-Fast Geometry-Independent Highly-Granular Calorimeter Simulation Salle Séminaires	Anatolii Korol 14:30 - 14:50
15:00	Point-Clouds based Diffusion Model on Hadronic Showers Salle Séminaires	Martina Mozzanica 14:50 - 15:10



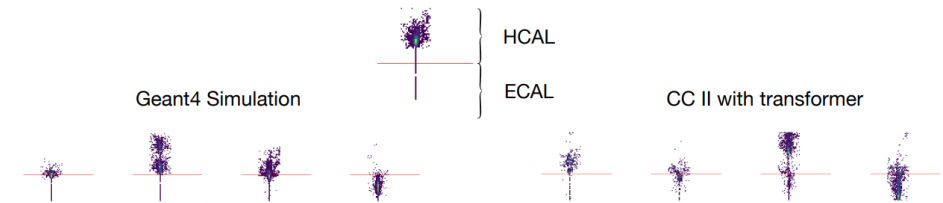
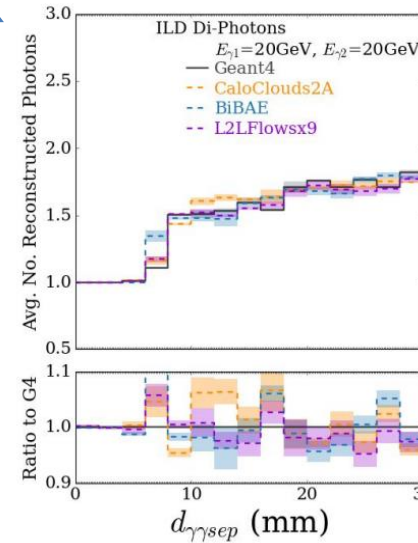
(VAE with flow matching)

[Talk by L. Favaro](#)



(Staggering effect in ILD Ecal and Di-Photons benchmark)

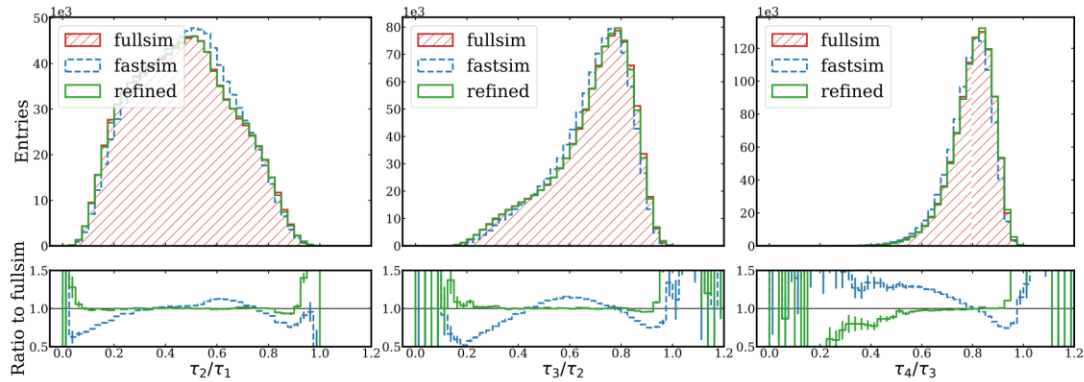
[Talk by A. Korol](#)



(Combined ECal and HCal shower generation) [Talk by M. Mozzanica](#)

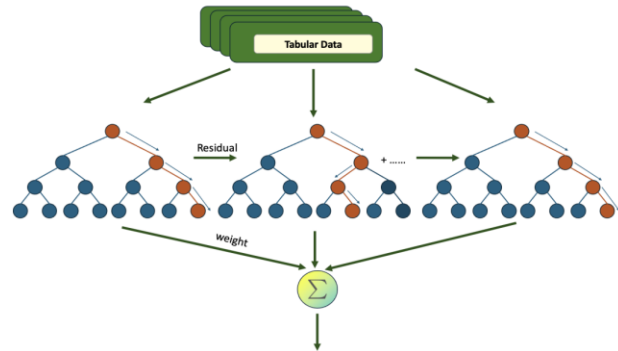


Detector Simulation

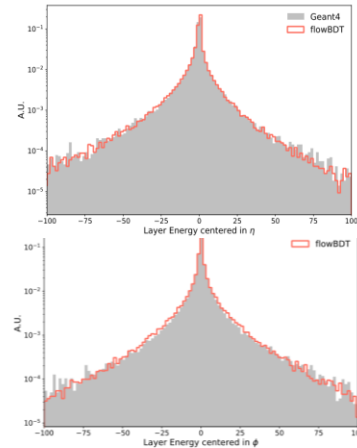


(Refinement of parameterized FastSim)

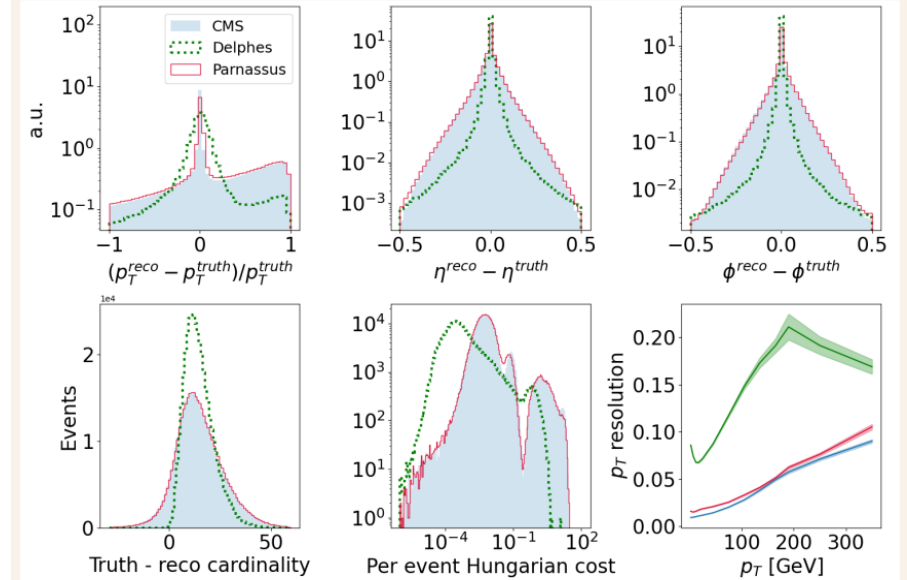
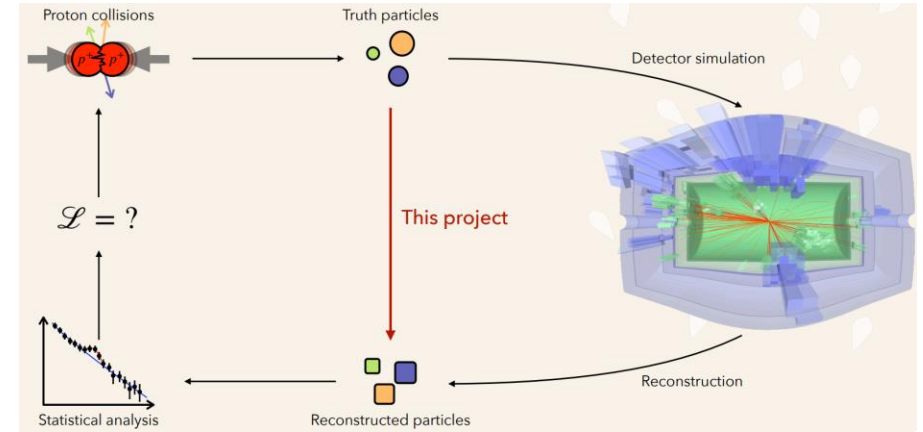
[Talk by L. Stietz](#)



(Shower generation (ATLAS pions, lower granularity) using GBDTs)



[Talk by S. Qian](#)

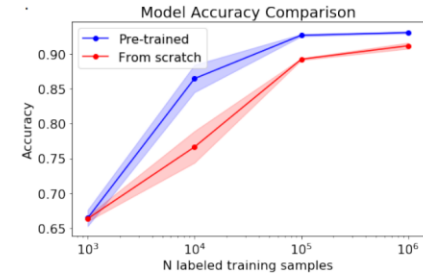
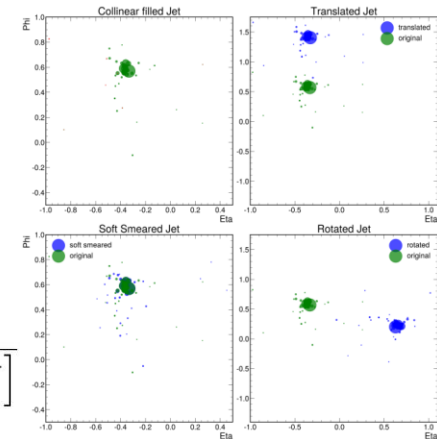
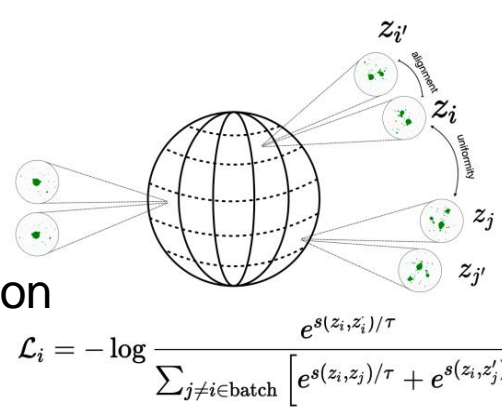


Particle features

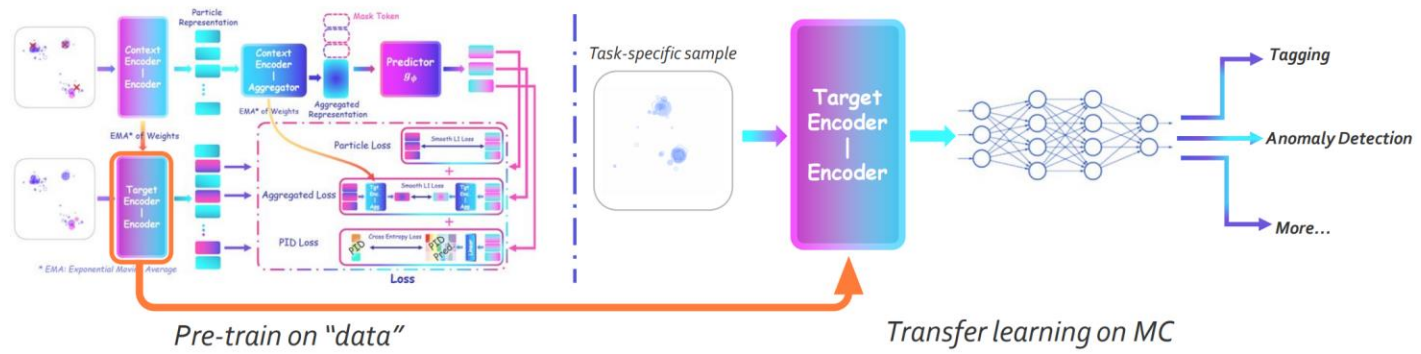
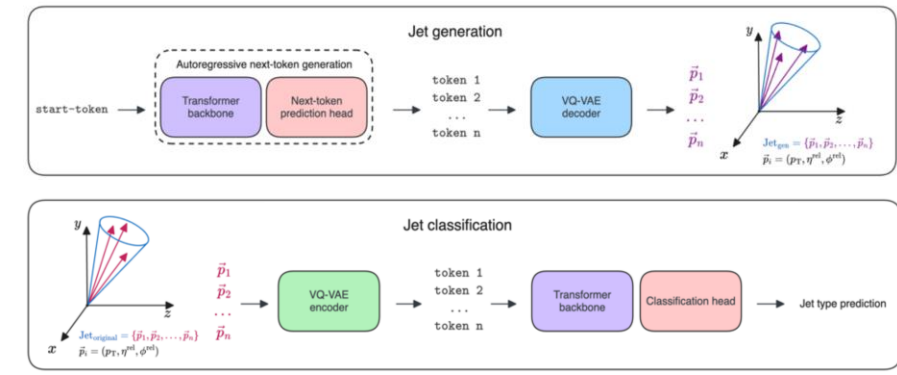
(Sim and Reco in a single step) [Talk by D. Kobylanski](#)

Foundation Models

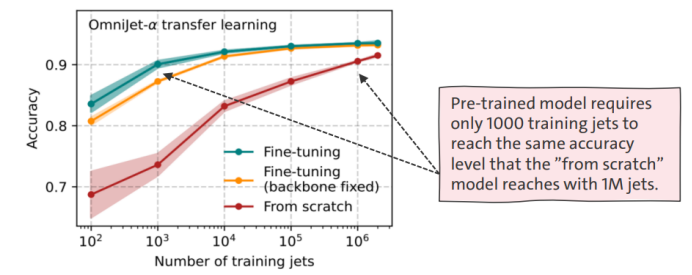
- Several efforts in the direction of "foundation models"
- All of them trained on Jet data
- Same observations – faster adaptation, even on a new task such as tagging, anomaly detection, etc.



(Contrastive learning of Jets. Finetune on top-tagging) [Talk by Z. Zhao](#)



(Training using joint-embedding predictive architecture. Tagging and anomaly detection) [Talk by S. Wang](#)




Pre-trained model requires only 1000 training jets to reach the same accuracy level that the "from scratch" model reaches with 1M jets.

(Next token prediction of Jets. Using the model for Jet generation and tagging) [Talk by A. Hallin](#)



Miscellaneous

- Symbolic AI for scattering amplitudes. [Talk](#)
 - Very interesting talk explaining the challenges encountered and tackling them with simple modifications
- Fair Universe: HiggsML Uncertainty Challenge. [Talk](#)
 - Train an AI model to improve cross section measurement significance
 - Running from September 12 to March 14th. [Competition link](#)
- CaloChallenge wrapped up
 - Final CaloChallenge paper published! [arXiv:2410.21611](#)
 - 59 submissions, 3 datasets, 23 different models
 - Some (very) preliminary discussion on the next CaloChallenge...



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.

[From C. Crause's talk](#)

Summary

- ML (in HEP) is rapidly evolving (as always)
- Increasing number of generative applications- lots of diffusion models
- Lots that is/could be interesting to SFT
 - Increasing interest in ML developments that stretch beyond one experiment
 - Increasing need for 'infrastructure' support/development- we are already involved in many places!
 - Coordinating/supporting community challenges- e.g. SFT will (again) play a central role in the next CaloChallenge



Pizza the Parisian way