

# Fast and Faithful (Calorimeter) Simulations with Normalizing Flows

— AI and Quantum Information Applications in Fundamental Physics, Konjiam, South Korea —

Claudius Krause

Institute for Theoretical Physics, University of Heidelberg  
(and Rutgers, The State University of New Jersey)

February 14, 2023



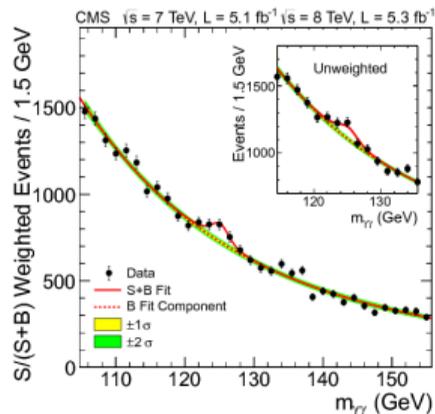
Based on

[2106.05285, 2110.11377, 2210.14245, 2212.06172] and lots [in preparation]

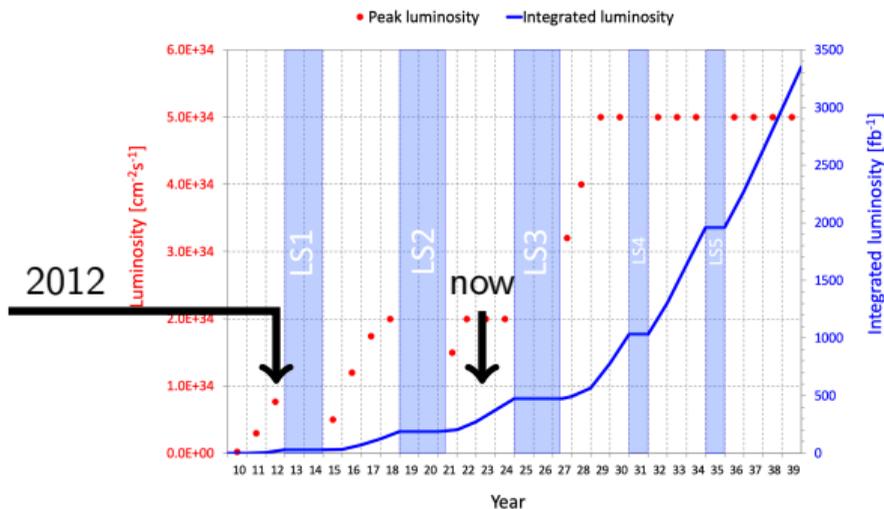
In collaboration with

David Shih, Ian Pang, Matt Buckley, Sascha Diefenbacher, Engin Eren, Frank Gaede, Gregor Kasieczka, Imahn Shekhzadeh, Florian Ernst, Luigi Favaro, Tilman Plehn, Theo Heimel, Ramon Winterhalder, Anja Butter, Joshua Isaacson, Fabio Maltoni, Olivier Mattelaer

# We will have a lot more data in the near future.



CMS Collaboration [arXiv:1207.7235, Phys.Lett.B]

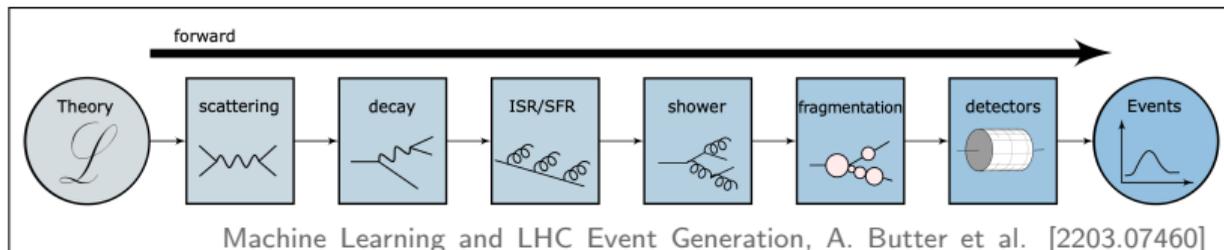


<https://lhc-commissioning.web.cern.ch/schedule/HL-LHC-plots.htm>

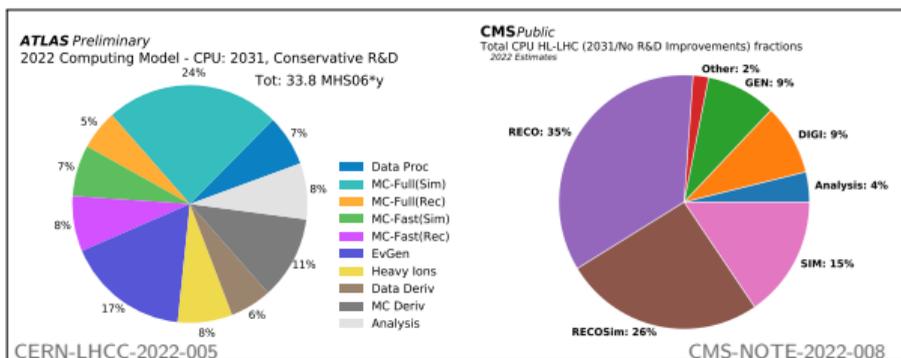
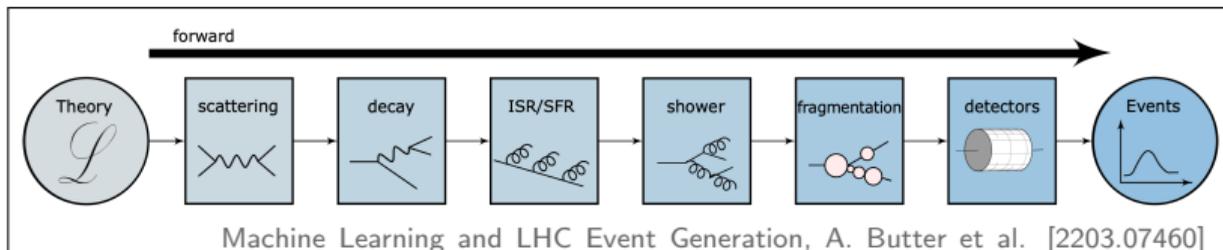
- We will have 20–25 $\times$  more data.

⇒ We want to understand every aspect of it based on 1<sup>st</sup> principles  
(and find New Physics)!

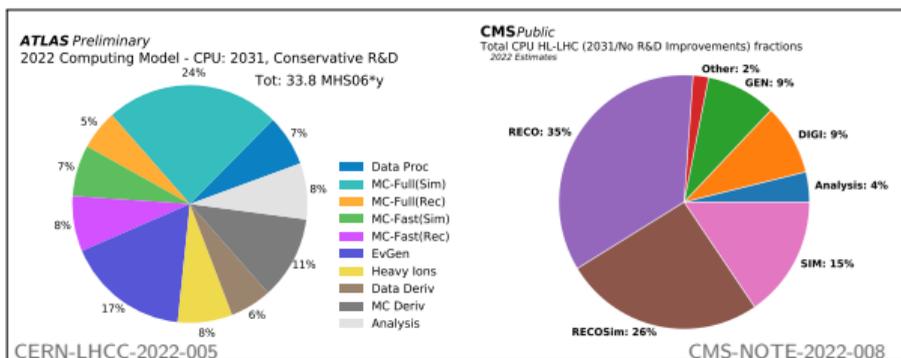
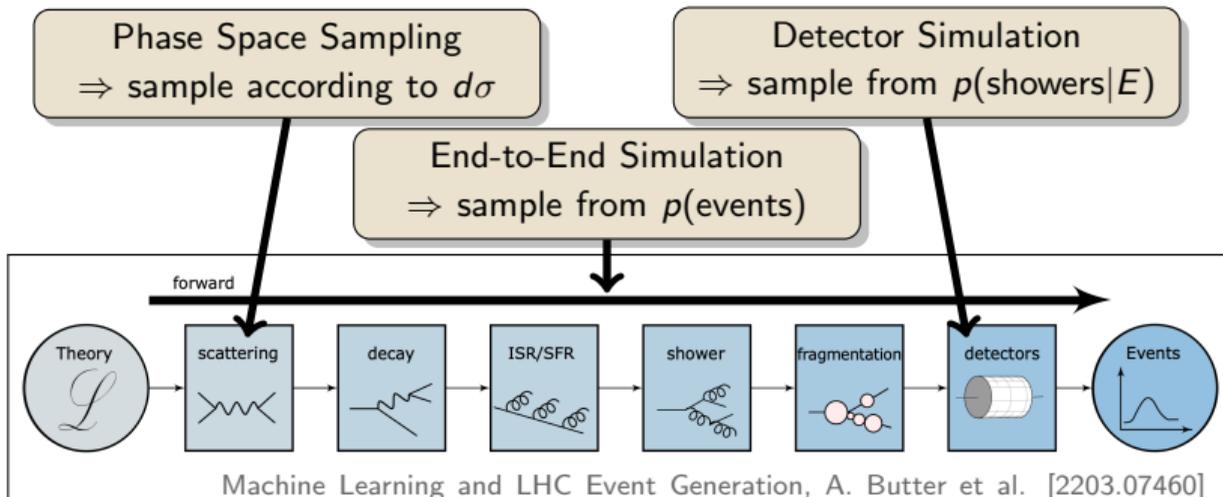
# Simulation bridges Theory and Experiment.



# Simulation bridges Theory and Experiment.



# Simulation bridges Theory and Experiment.



## Deep Generative Models are Random Number Generators.

We have a distribution  $p(x)$  and want to sample (“generate”) new elements that follow it.

given:  $\{x_i\}$

want:  $x \sim p(x)$

- or -

given:  $f(x)$

want:  $x \sim f(x) / \int f(x)$

# Deep Generative Models are Random Number Generators.

We have a distribution  $p(x)$  and want to sample (“generate”) new elements that follow it.

given:  $\{x_i\}$

want:  $x \sim p(x)$

- or -

given:  $f(x)$

want:  $x \sim f(x) / \int f(x)$



“Calorimeter Simulation” via [midjourney.com](#)

- Deep Generative Models have seen a lot of progress in recent years.
- They can generate text, speech, images, . . .
- In HEP, we have seen GANs, VAEs, Normalizing Flows, Diffusion Models, and their derivatives.

## The first models in use were GANs and VAEs.

Variational Autoencoder (VAE): Information is encoded and decoded through a bottleneck.

- Trade-off between generalizability and quality ( $\beta$ -term in the loss)
- Physics is not always compressible.
- Precision not competitive to other generative models.

## The first models in use were GANs and VAEs.

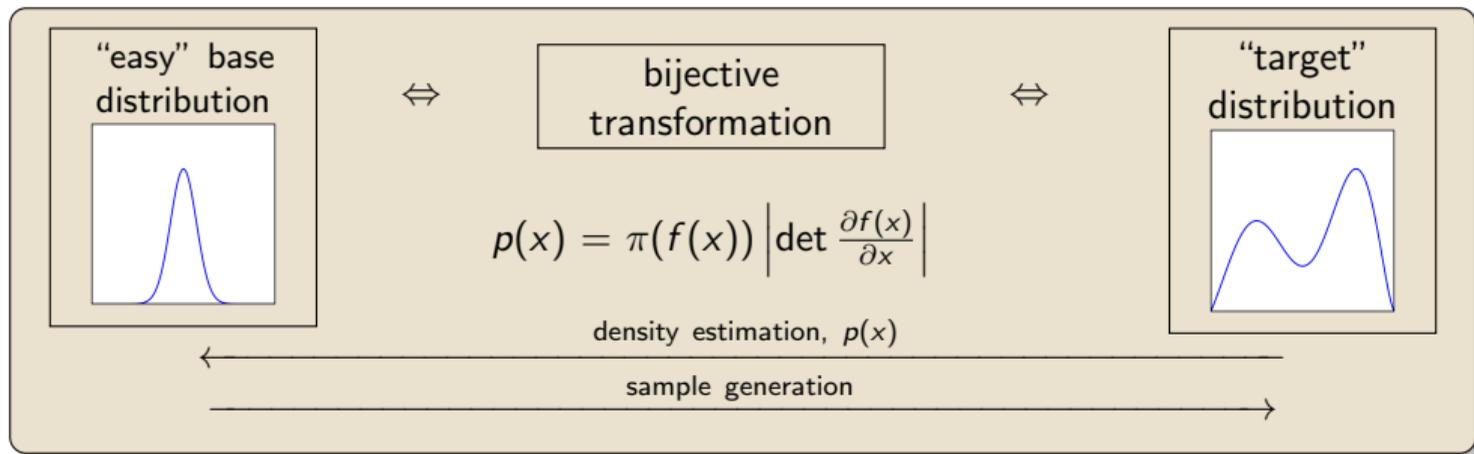
Variational Autoencoder (VAE): Information is encoded and decoded through a bottleneck.

- Trade-off between generalizability and quality ( $\beta$ -term in the loss)
- Physics is not always compressible.
- Precision not competitive to other generative models.

Generative Adversarial Network (GAN): A generator and a critic play a game against each other.

- Delicate optimization on a saddlepoint.
- Suffer from mode-collapse and other artefacts.
- Model selection is difficult.
- Hard to get the distributions to agree at the %-level.

# Normalizing Flows learn a coordinate transformation.



Having access to the log-likelihood (LL) allows several training options:

- ⇒ Based on samples: via maximizing  $\text{LL}(\text{samples})$ .
- ⇒ Based on target function  $f(x)$ : via matching  $p(x)$  to  $f(x)$ .

# How do Normalizing Flows tame Jacobians?

- NFs learn the parameters  $\theta$  of a series of easy transformations.

Dinh et al. [arXiv:1410.8516], Rezende/Mohamed [arXiv:1505.05770]

- Each transformation is 1d & has an analytic Jacobian and inverse.

⇒ We use Rational Quadratic Splines

Durkan et al. [arXiv:1906.04032], Gregory/Delbourgo [IMA J. of Num. An., '82]

- Require a triangular Jacobian for faster evaluation.

⇒ The parameters  $\theta$  depend only on a subset of all other coordinates.

# How do Normalizing Flows tame Jacobians?

- NFs learn the parameters  $\theta$  of a series of easy transformations.

Dinh et al. [arXiv:1410.8516], Rezende/Mohamed [arXiv:1505.05770]

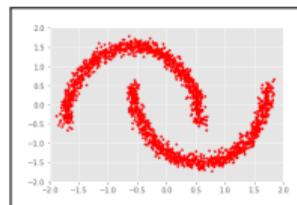
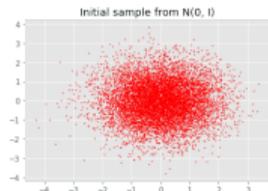
- Each transformation is 1d & has an analytic Jacobian and inverse.

⇒ We use Rational Quadratic Splines

Durkan et al. [arXiv:1906.04032], Gregory/Delbourgo [IMA J. of Num. An., '82]

- Require a triangular Jacobian for faster evaluation.

⇒ The parameters  $\theta$  depend only on a subset of all other coordinates.



<https://engineering.papercup.com/posts/normalizing-flows-part-2/>

# How do Normalizing Flows tame Jacobians?

- NFs learn the parameters  $\theta$  of a series of easy transformations.

Dinh et al. [arXiv:1410.8516], Rezende/Mohamed [arXiv:1505.05770]

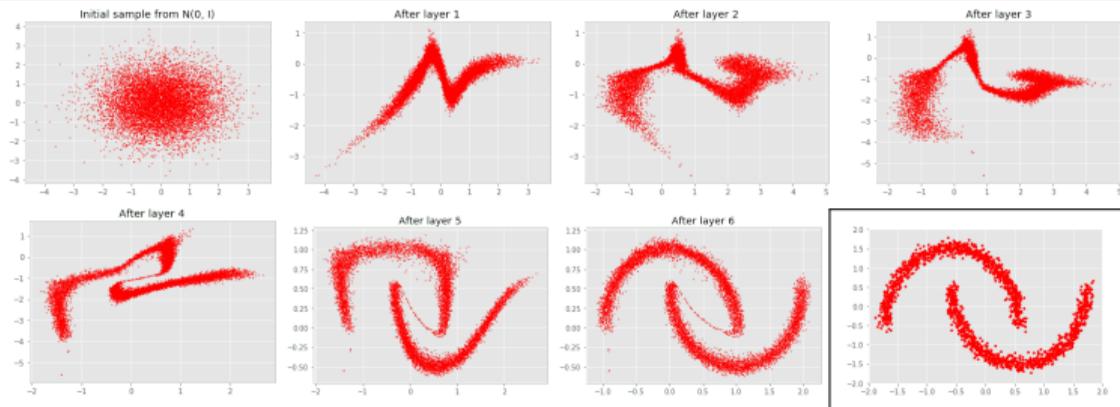
- Each transformation is 1d & has an analytic Jacobian and inverse.

⇒ We use Rational Quadratic Splines

Durkan et al. [arXiv:1906.04032], Gregory/Delbourgo [IMA J. of Num. An., '82]

- Require a triangular Jacobian for faster evaluation.

⇒ The parameters  $\theta$  depend only on a subset of all other coordinates.



<https://engineering.papercup.com/posts/normalizing-flows-part-2/>

# The Bijector is a chain of “easy” transformations.

Each transformation

- must be invertible and have analytical Jacobian
- is chosen to factorize:

$$\vec{C}(\vec{x}; \vec{p}) = (C_1(x_1; p_1), C_2(x_2; p_2), \dots, C_n(x_n; p_n))^T,$$

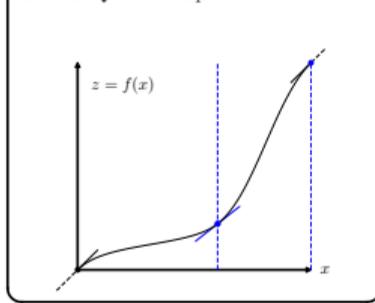
where  $\vec{x}$  are the coordinates to be transformed and  $\vec{p}$  the parameters of the transformation.

Rational Quadratic Splines:

Durkan et al. [arXiv:1906.04032]

Gregory/Delbourgo [IMA Journal of Numerical Analysis, '82]

Rational Quadratic Spline Transformation

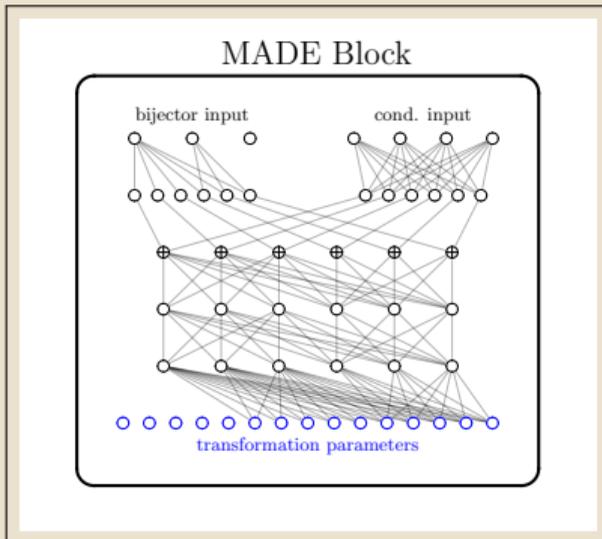


$$C = \frac{a_2\alpha^2 + a_1\alpha + a_0}{b_2\alpha^2 + b_1\alpha + b_0}$$

- numerically easy
- expressive

The NN predicts the bin widths, heights, and derivatives that go in  $a_i$  &  $b_j$ .

# Triangular Jacobians 1: with Autoregressive Blocks



$$\theta_{x_i}(x_{j < i})$$

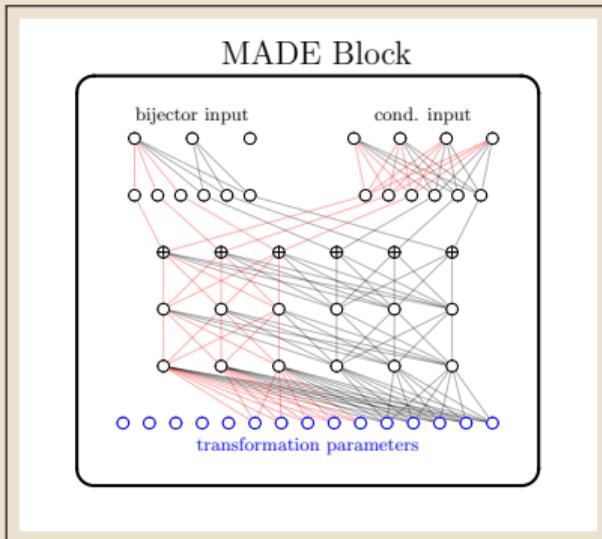
Implementation via masking:

- a single “forward” pass gives all  $\theta_{x_i}(x_{i-1} \dots x_1)$ .  
⇒ very fast
- its “inverse” needs to loop through all dimensions.  
⇒ very slow

Germain/Gregor/Murray/Larochelle [arXiv:1502.03509]

- Masked Autoregressive Flow (MAF) is slow in sampling and fast in inference.  
Papamakarios et al. [arXiv:1705.07057]
- Inverse Autoregressive Flow (IAF) is fast in sampling and slow in inference.  
Kingma et al. [arXiv:1606.04934]

# Triangular Jacobians 1: with Autoregressive Blocks



$$\theta_{x_i}(x_{j < i})$$

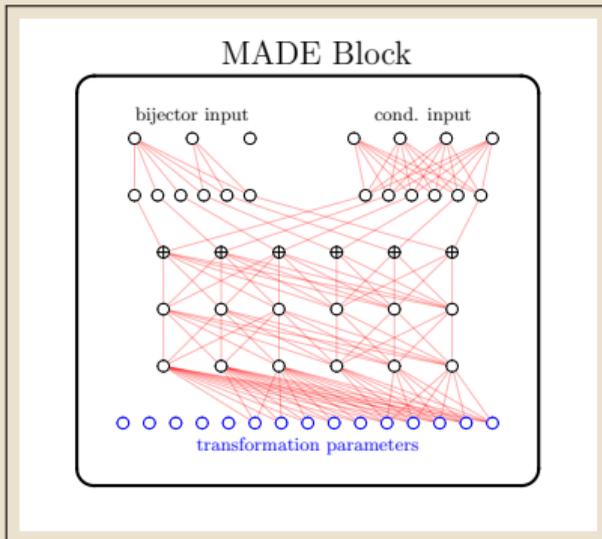
Implementation via masking:

- a single “forward” pass gives all  $\theta_{x_i}(x_{i-1} \dots x_1)$ .  
⇒ very fast
- its “inverse” needs to loop through all dimensions.  
⇒ very slow

Germain/Gregor/Murray/Larochelle [arXiv:1502.03509]

- Masked Autoregressive Flow (MAF) is slow in sampling and fast in inference.  
Papamakarios et al. [arXiv:1705.07057]
- Inverse Autoregressive Flow (IAF) is fast in sampling and slow in inference.  
Kingma et al. [arXiv:1606.04934]

# Triangular Jacobians 1: with Autoregressive Blocks



$$\theta_{x_i}(x_{j < i})$$

Implementation via masking:

- a single “forward” pass gives all  $\theta_{x_i}(x_{i-1} \dots x_1)$ .  
⇒ very fast
- its “inverse” needs to loop through all dimensions.  
⇒ very slow

Germain/Gregor/Murray/Larochelle [arXiv:1502.03509]

- Masked Autoregressive Flow (MAF) is slow in sampling and fast in inference.  
Papamakarios et al. [arXiv:1705.07057]
- Inverse Autoregressive Flow (IAF) is fast in sampling and slow in inference.  
Kingma et al. [arXiv:1606.04934]

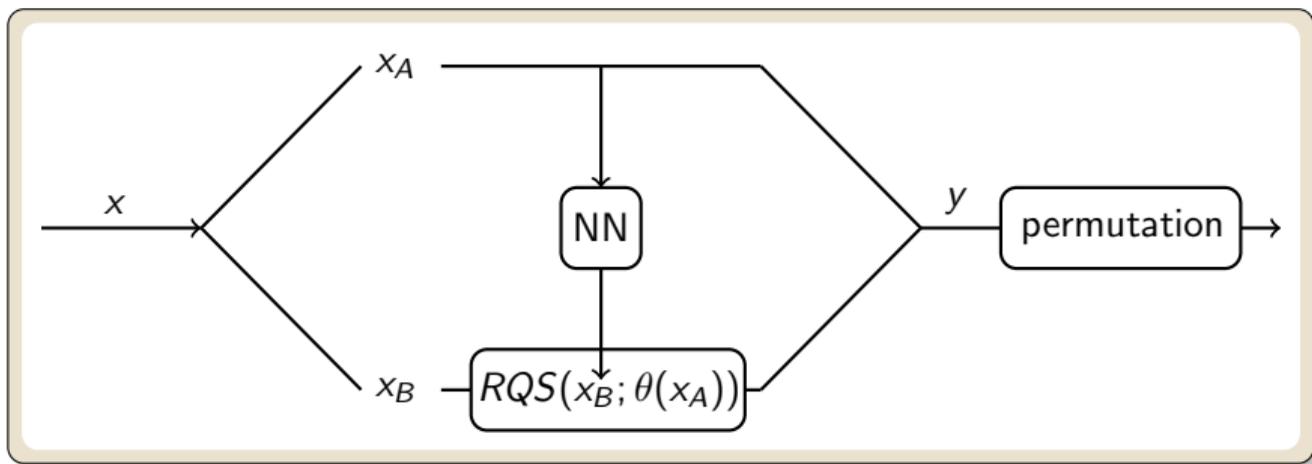
## Triangular Jacobians 2: with Bipartite Blocks

$$\theta_{x \in A}(x \in B) \quad \& \quad \theta_{x \in B}(x \in A)$$

- Coordinates are split in 2 sets, transforming each other.

+ Forward and inverse pass are equally fast.    - Said to be not as expressive.

Dinh et al. [arXiv:1410.8516]



# Normalizing Flows are great!

## Pros and Cons of Normalizing Flows:

- + LL optimization is more stable than saddlepoint optimization of GANs.
- + Do not suffer from mode-collapse.
- + Model selection is straightforward with  $LL(\text{val-set})$ .
- + Flows are versatile (train for one thing, use for another).
- + Empirically: better at learning distributions to the  $\epsilon$ -level
  
- They scale bad with the dimensionality of the problem.
- Some architectures might be slow.
- There are topological constraints.
- Sparse data is hard to learn.

# Evaluating Generative Models is a hard task.

We want to know if  $p_{\text{generated}} = p_{\text{training data}} (= p_{\text{truth}}?)$

- If we have access to  $p(x)$ , we can compute  $f$ -divergences.

↳ Example: KL-divergence  $\int dx p(x) \log \frac{p(x)}{q(x)}$

- Alternatively, we could use Integral Probability Metrics.

↳ Example: Wasserstein distance

- In Computer Science, one uses the Frechét Inception Distance.

- In HEP, we usually look at histograms.

See also: Kansal et al. [arXiv:2211.10295]

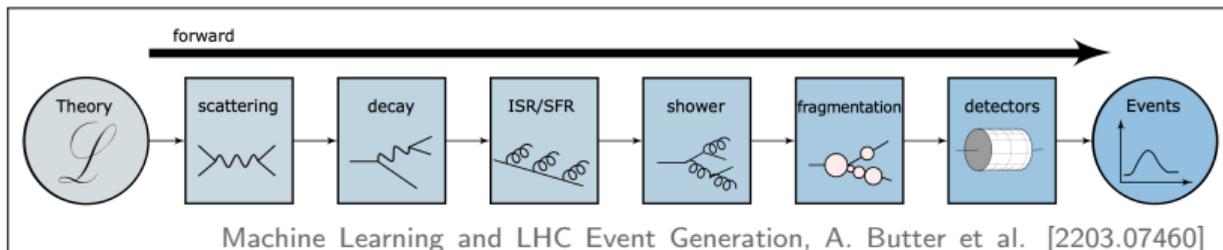
## A Classifier provides the “ultimate metric”.

According to the Neyman-Pearson Lemma we have:

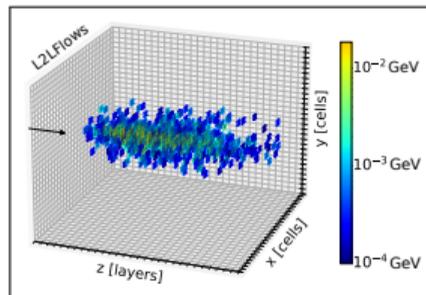
- The likelihood ratio is the most powerful test statistic to distinguish the two samples.
  - A powerful classifier trained to distinguish the samples should therefore learn (something monotonically related to) **this**.
  - If this classifier is confused, we conclude  $p_{\text{GEANT4}}(x) = p_{\text{generated}}(x)$
- ⇒ This captures the full phase space incl. correlations.
- ⇒ However, it is sufficient, but not necessary.
- ? But why wasn't this used before?
- ⇒ Previous deep generative models were separable to almost 100%!

DCTRGAN: Diefenbacher et al. [2009.03796, JINST]

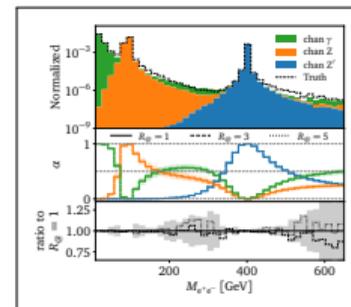
# Fast and Faithful (Calorimeter) Simulations with Normalizing Flows.



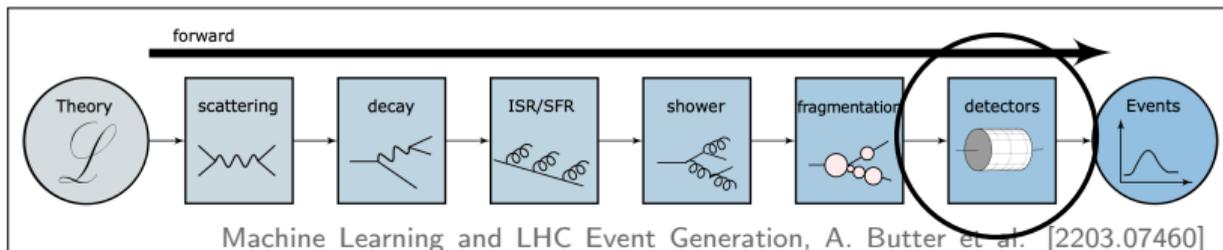
## 1: Calorimeter Simulation



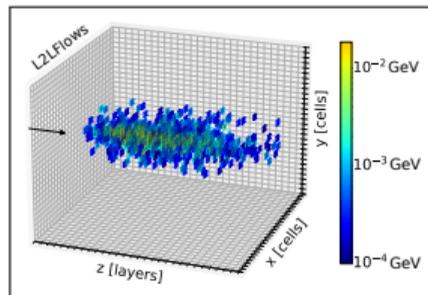
## 2: Phase Space Integration



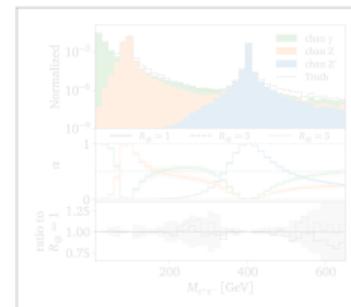
# Fast and Faithful (Calorimeter) Simulations with Normalizing Flows.



## 1: Calorimeter Simulation

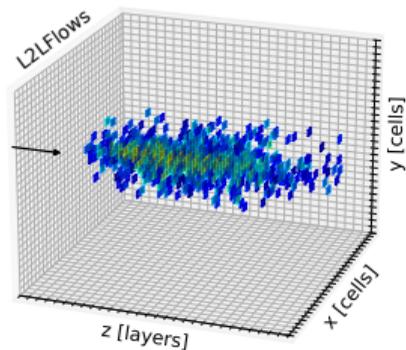
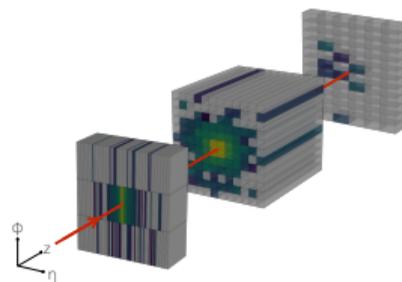


## 2: Phase Space Integration



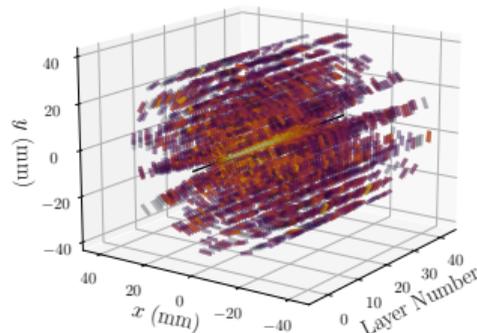
# Different Datasets come with different challenges

## Part I: The CALOGAN Dataset



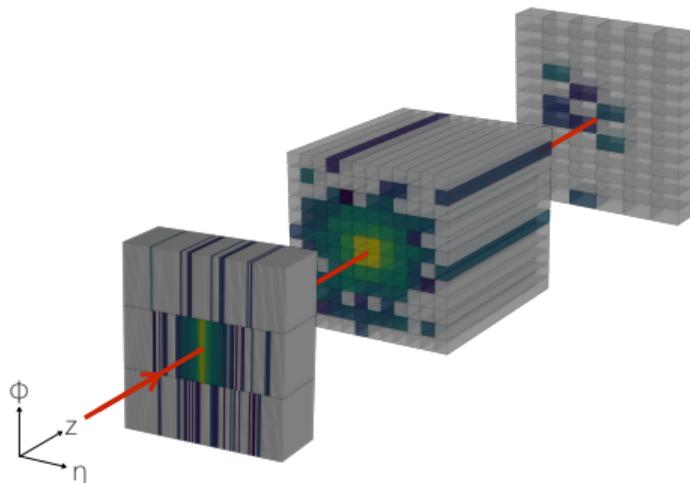
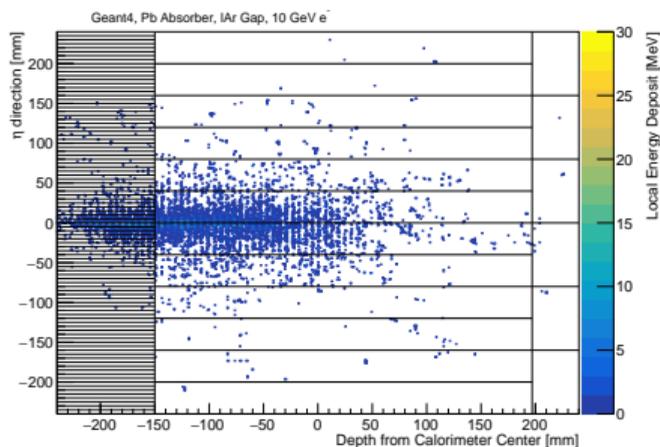
## Part II: The ILD Dataset

## Part III: The CaloChallenge 2022



# I: We use the same calorimeter geometry as CALOGAN

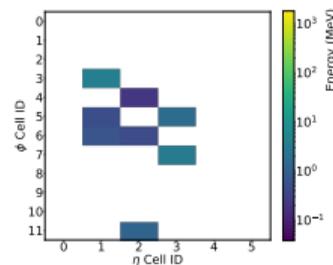
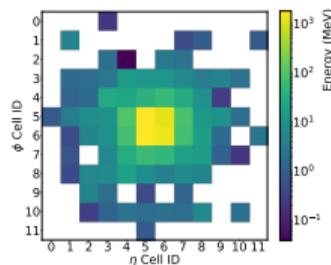
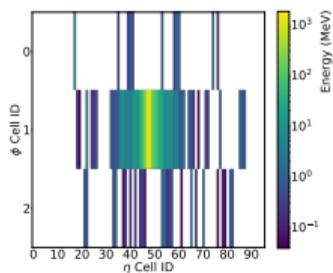
- We consider a toy calorimeter inspired by the ATLAS ECal: flat alternating layers of lead and LAr
- They form three instrumented layers of dimension  $3 \times 96$ ,  $12 \times 12$ , and  $12 \times 6$



CaloGAN: Paganini, de Oliveira, Nachman [1705.02355, PRL; 1712.10321, PRD]

# I: We use the same calorimeter geometry as CALOGAN

- The GEANT4 configuration of CALOGAN is available at <https://github.com/hep-lbdl/CaloGAN>
- We produce our own dataset: available at [DOI: 10.5281/zenodo.5904188]
- Showers of  $e^+$ ,  $\gamma$ , and  $\pi^+$  (100k each)
- All are centered and perpendicular
- $E_{\text{inc}}$  is uniform in  $[1, 100]$  GeV and given in addition to the energy deposits per voxel:



CaloGAN: Paganini, de Oliveira, Nachman [1705.02355, PRL; 1712.10321, PRD]

# I: CALOFlow uses a 2-step approach to learn $p(\vec{\mathcal{I}}|E_{\text{inc}})$ .

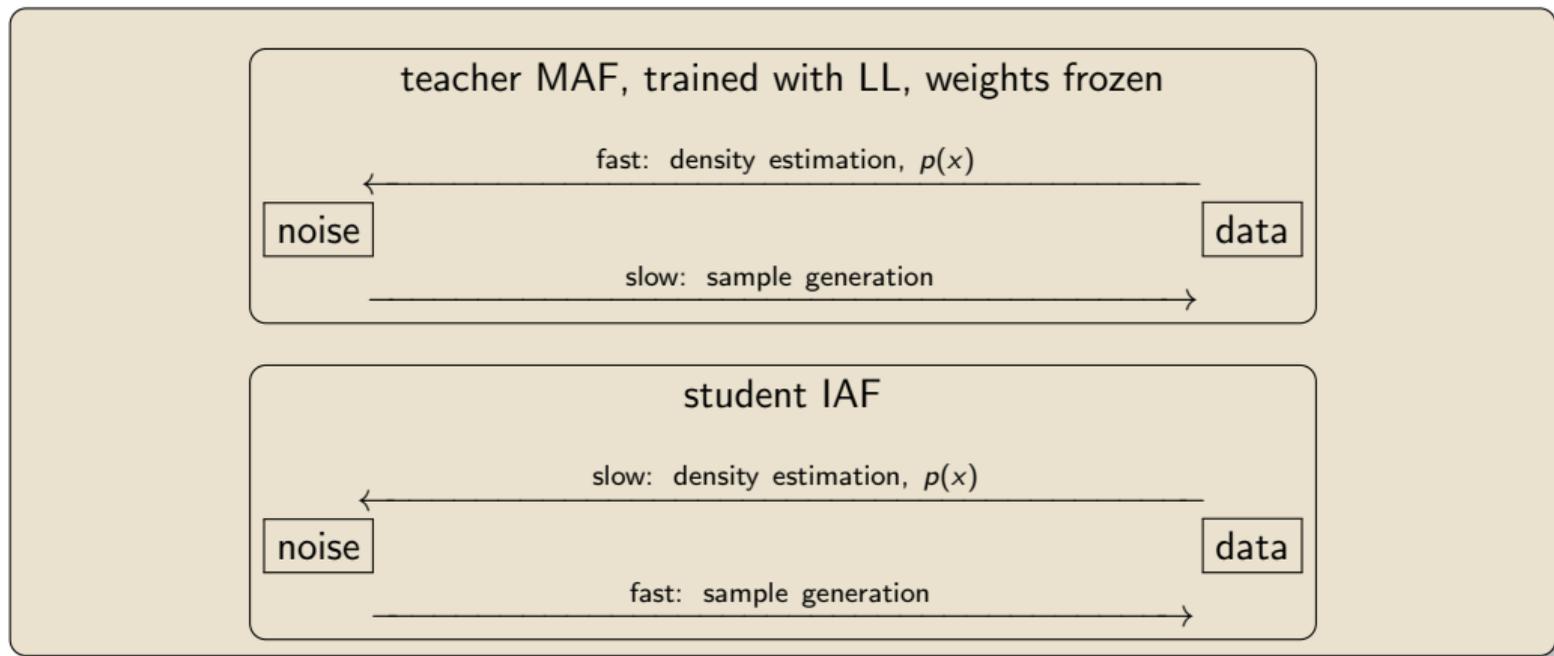
## Flow I

- learns  $p_1(E_0, E_1, E_2|E_{\text{inc}})$
- is a Masked Autoregressive Flow, optimized using the log-likelihood.

## Flow II

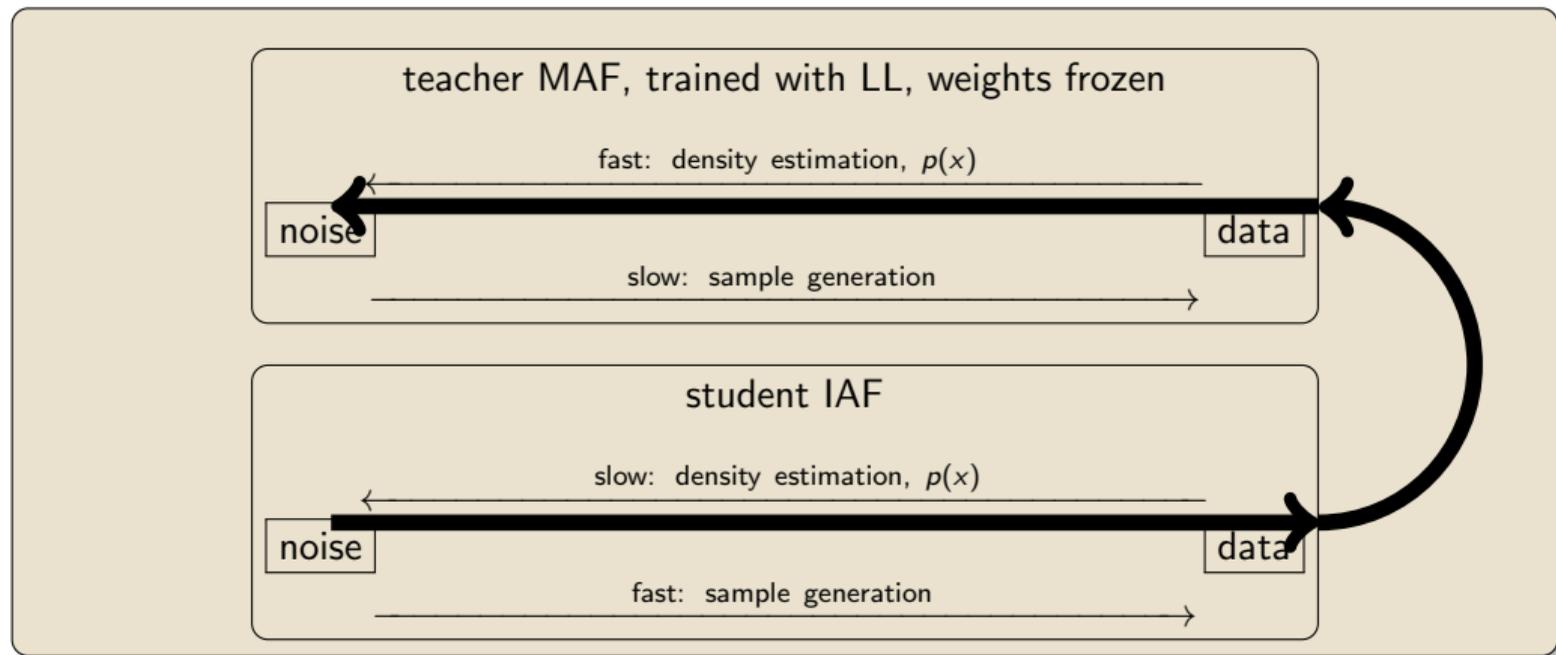
- learns  $p_2(\hat{\mathcal{I}}|E_0, E_1, E_2, E_{\text{inc}})$  of normalized showers
- in CALOFlow v1 (2106.05285 — called “teacher”):
  - Masked Autoregressive Flow trained with log-likelihood
  - Slow in sampling ( $\approx 500\times$  slower than CALOGAN)
- in CALOFlow v2 (2110.11377 — called “student”):
  - Inverse Autoregressive Flow trained with Probability Density Distillation from teacher (log-likelihood prohibitive), van den Oord et al. [1711.10433]  
i.e. matching IAF parameters to frozen MAF
  - Fast in sampling ( $\approx 500\times$  faster than CALOFlow v1)

# I: Probability Density Distillation passes the information from the teacher to the student



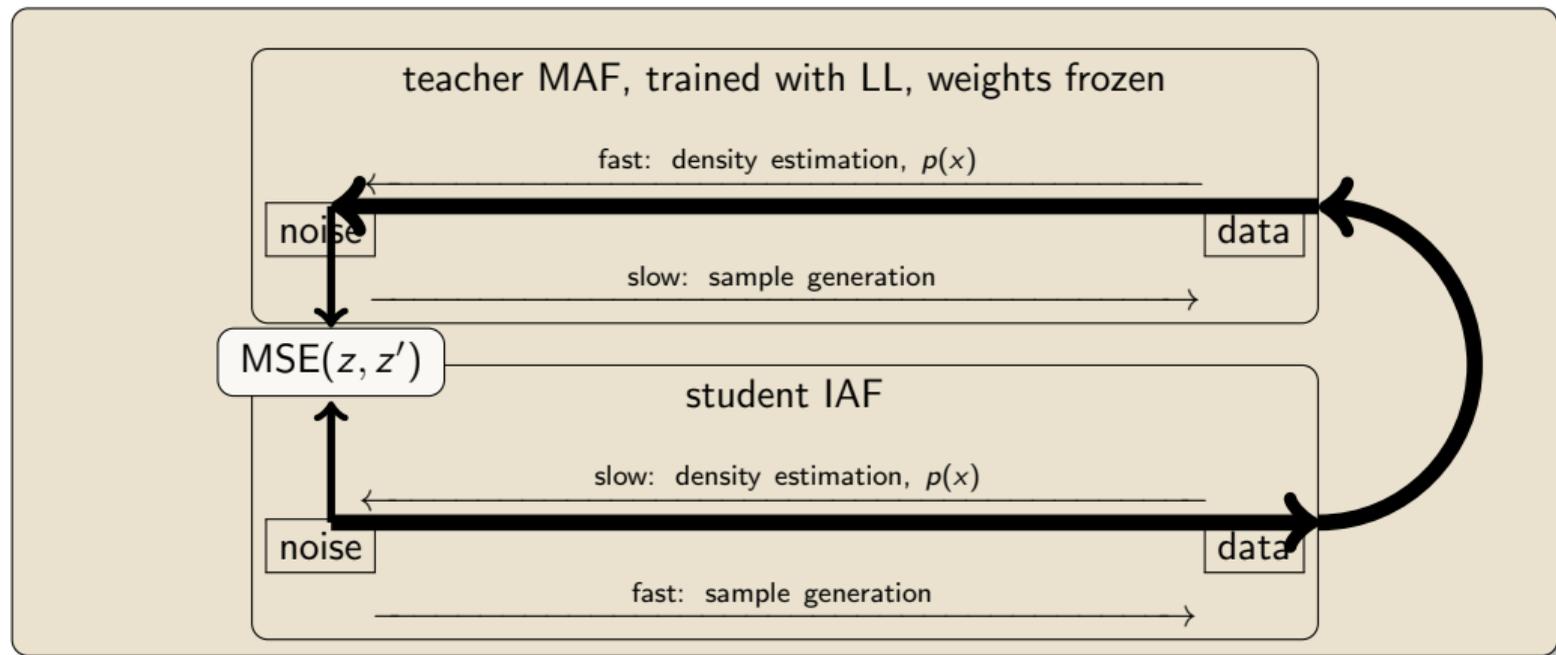
$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) \\ + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



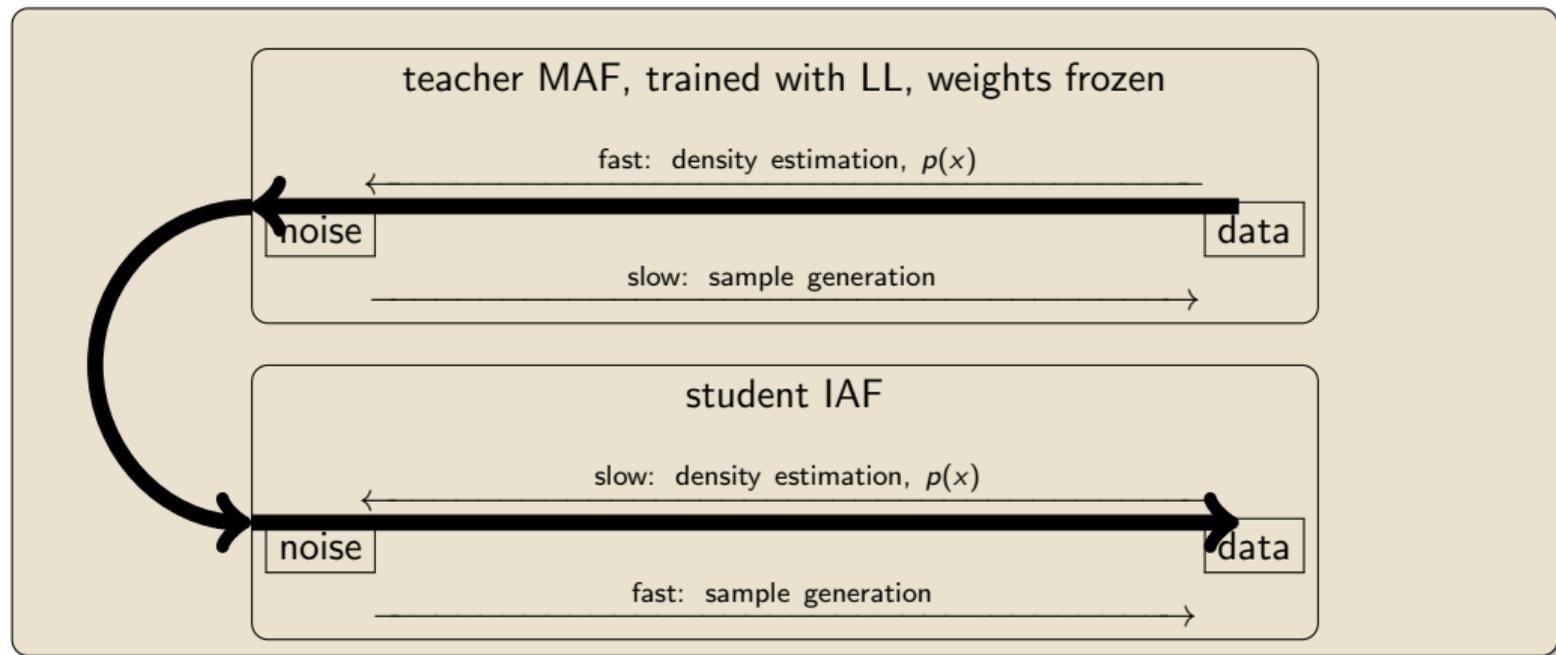
$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) \\ + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



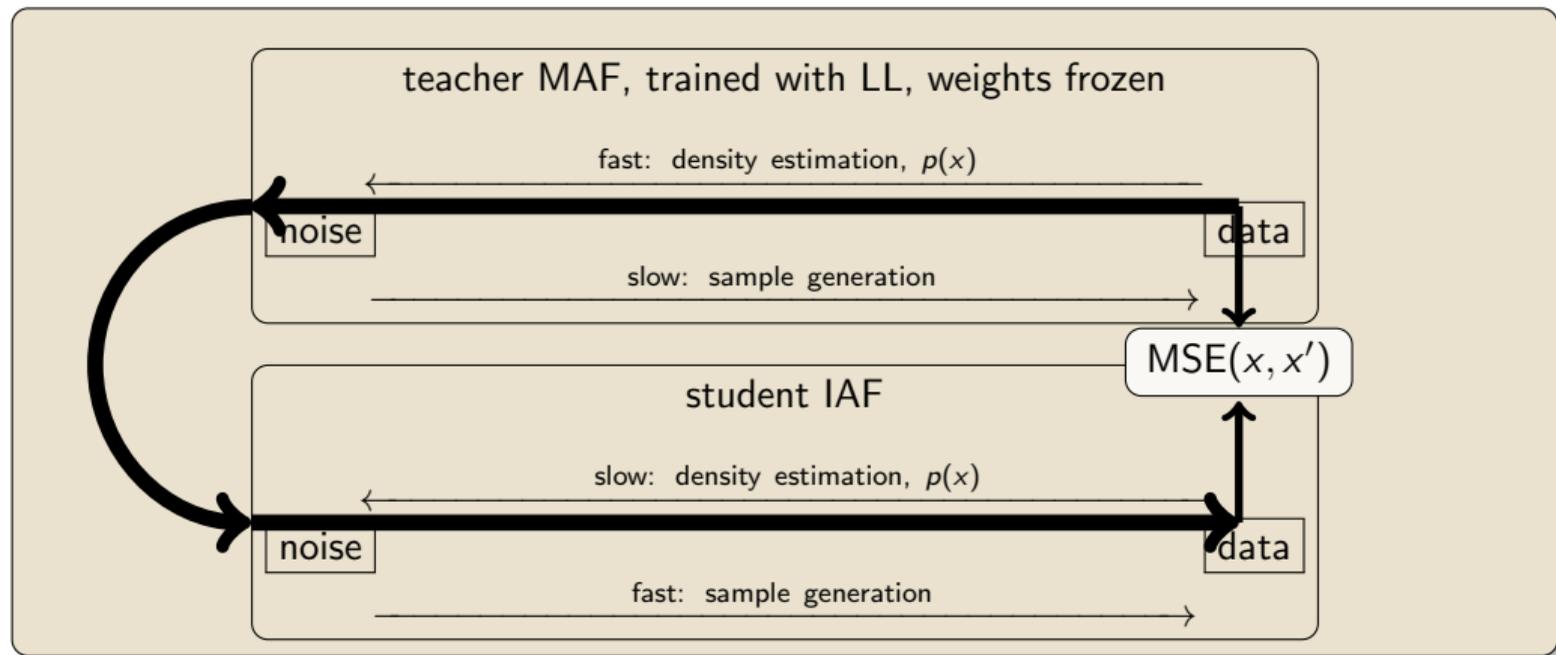
$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) \\ + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



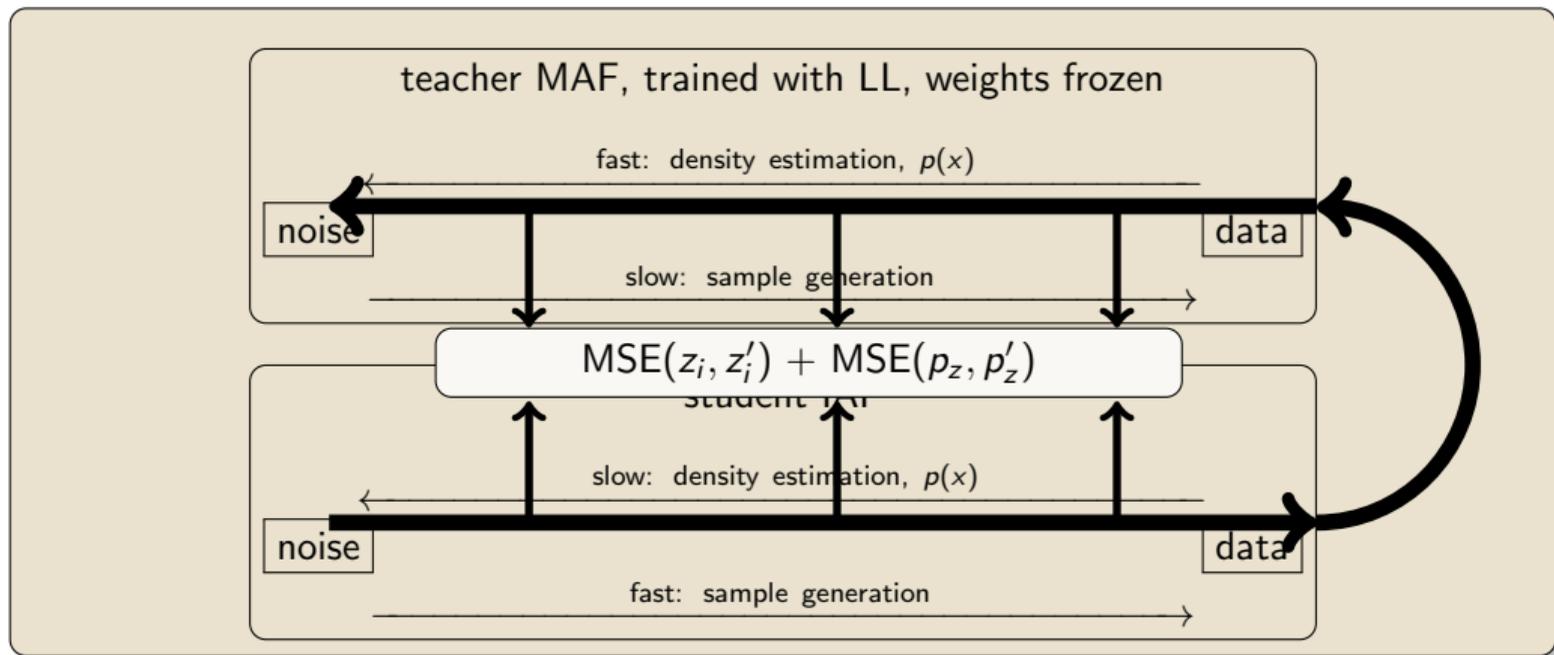
$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) \\ + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



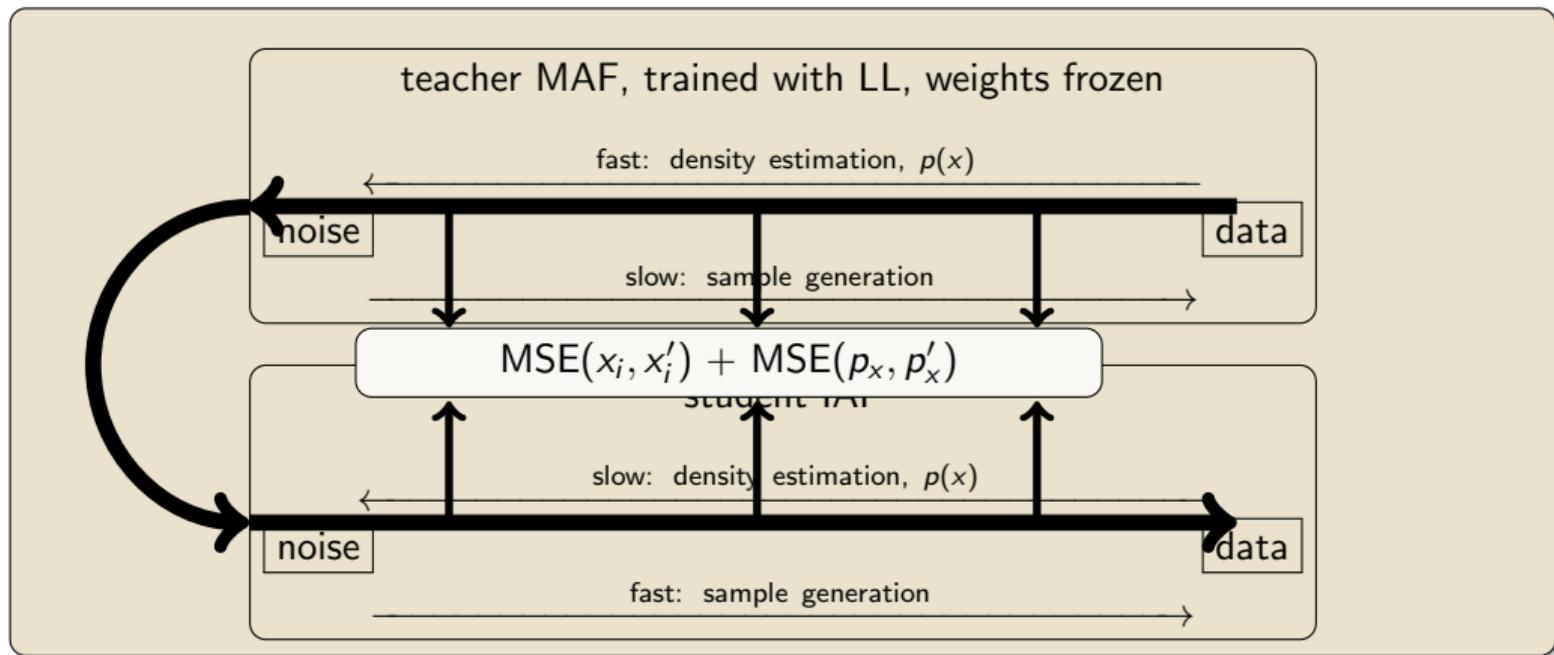
$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) \\ + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



$$\text{Loss} = \text{MSE}(z, z') + \text{MSE}(x, x') + \text{MSE}(z_i, z'_i) + \text{MSE}(x_i, x'_i) + \text{MSE}(p_z, p'_z) + \text{MSE}(p_x, p'_x)$$

# I: Probability Density Distillation passes the information from the teacher to the student



$$\text{Loss} = MSE(z, z') + MSE(x, x') + MSE(z_i, z'_i) \\ + MSE(x_i, x'_i) + MSE(p_z, p'_z) + MSE(p_x, p'_x)$$

# I: CALOFlow passes the “ultimate metric” test.

According to the Neyman-Pearson Lemma we have:  $p_{\text{GEANT4}}(x) = p_{\text{generated}}(x)$  if a classifier cannot distinguish data from generated samples.

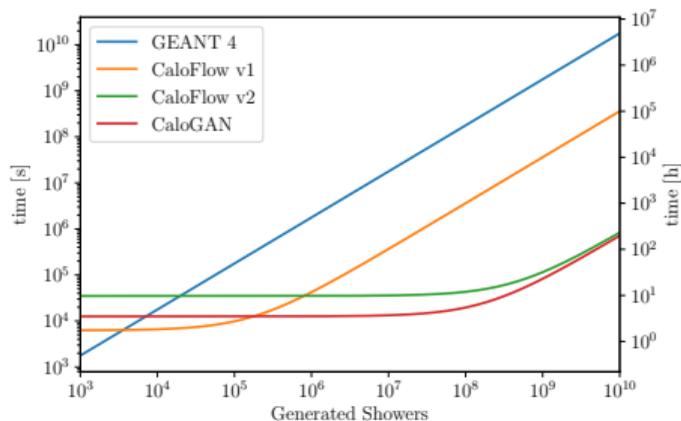
AUC		GEANT4 vs. CALOGAN	GEANT4 vs. (teacher) CALOFlow v1	GEANT4 vs. (student) CALOFlow v2
$e^+$	unnorm.	1.000(0)	0.859(10)	0.786(7)
	norm.	1.000(0)	0.870(2)	0.824(4)
	high-level	1.000(0)	0.795(1)	0.762(3)
$\gamma$	unnorm.	1.000(0)	0.756(48)	0.758(14)
	norm.	1.000(0)	0.796(2)	0.760(3)
	high-level	1.000(0)	0.727(2)	0.739(2)
$\pi^+$	unnorm.	1.000(0)	0.649(3)	0.729(2)
	norm.	1.000(0)	0.755(3)	0.807(1)
	high-level	1.000(0)	0.888(1)	0.893(2)

AUC ( $\in [0.5, 1]$ ): Area Under the ROC Curve, smaller is better, i.e. more confused

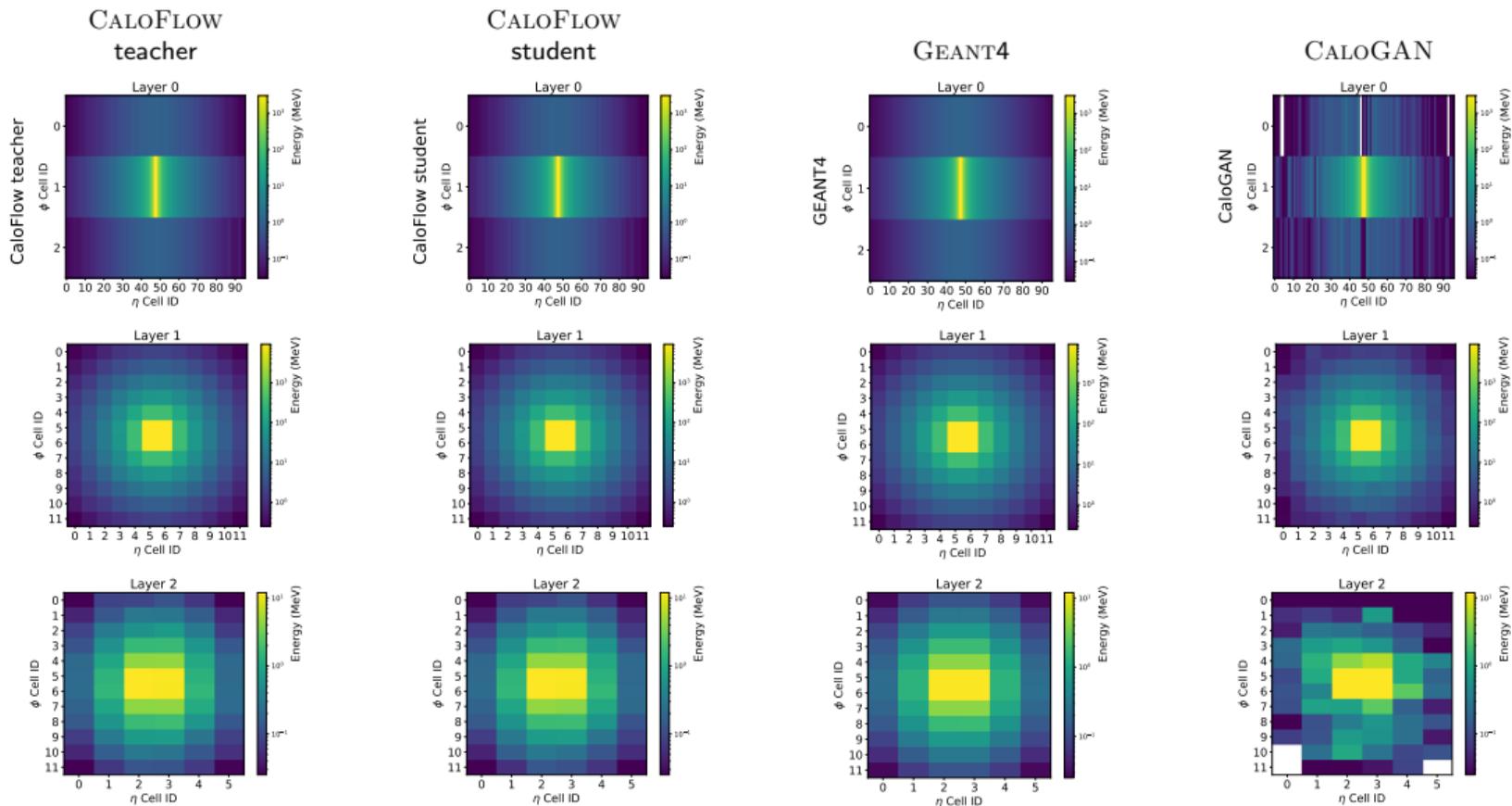
# I: Sampling Speed: The Student beats the Teacher!

	CALOFLOW*		CALOGAN*	GEANT4 <sup>†</sup>
	teacher	student		
training	22+82 min	+ 480 min	210 min	0 min
generation time per shower	36.2 ms	<b>0.08 ms</b>	<b>0.07 ms</b>	1772 ms

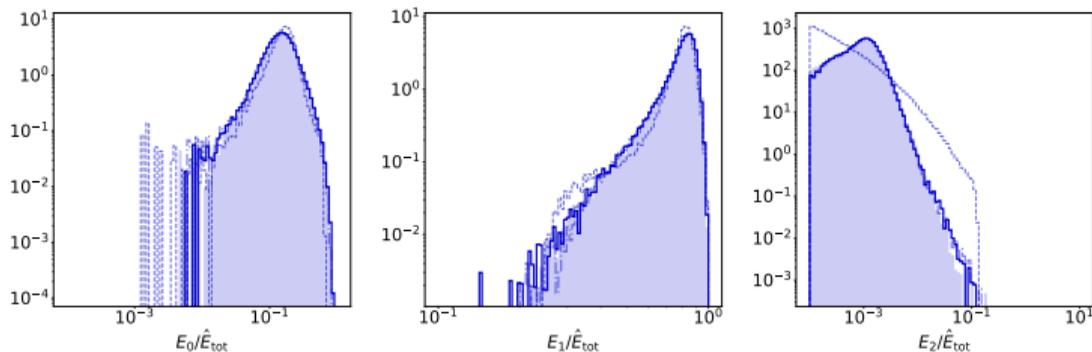
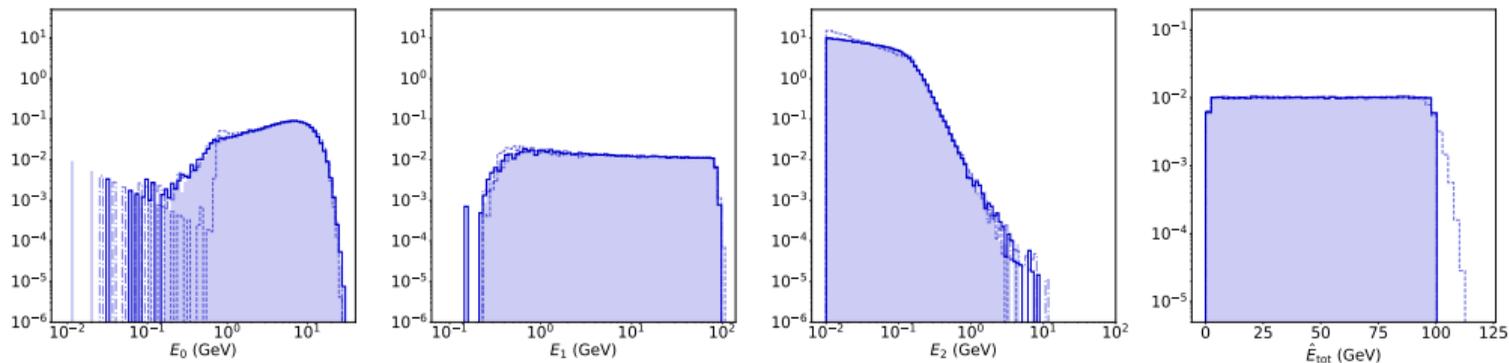
\*: on our TITAN V GPU, †: on the CPU of CaloGAN: Paganini, de Oliveira, Nachman [1712.10321, PRD]



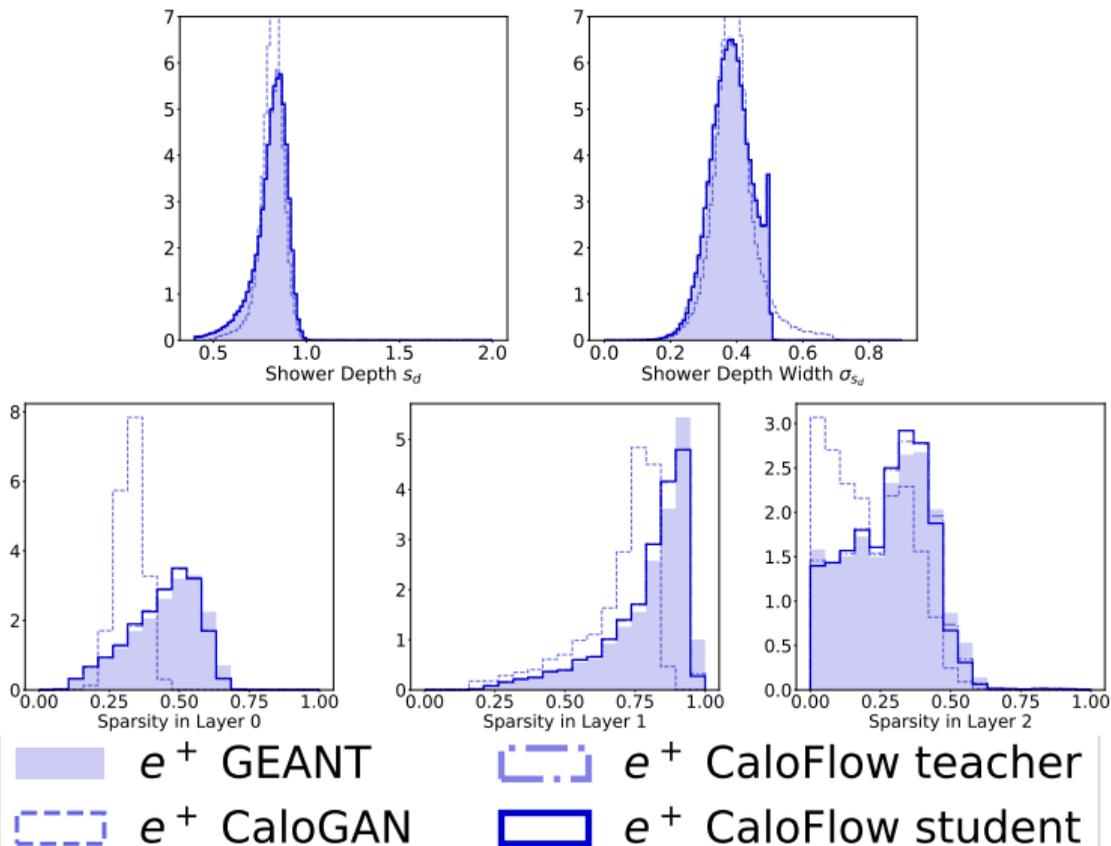
# I: CALOFlow: Comparing Shower Averages: $e^+$



# I: CALOFLOW: histograms: $e^+$



# I: CALOFLOW: histograms: $e^+$



# I: What about other architectures?

Work in progress with Florian Ernst, Luigi Favaro, Tilman Plehn, David Shih:

⇒ We compare the performance of a completely different architecture on the same dataset.

① Replace the autoregressive flow with a coupling layer based one.

② Do it in one step: learn  $p(\hat{\mathcal{I}}_0, \hat{\mathcal{I}}_1, \hat{\mathcal{I}}_2, E_0, E_1, E_2 | E_{\text{inc}})$  (507-dim.)

③ Make it Bayesian to

- ▶ see where the model is uncertain
- ▶ have it regularized
- ▶ have some kind of dropout

Bellagente, Haußmann, Luchmann, Plehn [2104.04543, SciPost]

# I: What about other architectures?

Work in progress with Florian Ernst, Luigi Favaro, Tilman Plehn, David Shih:

⇒ We compare the performance of a completely different architecture on the same dataset.

1 Replace the autoregressive flow with a coupling layer based one.

2 Do it in one step: learn  $p(\hat{\mathcal{I}}_0, \hat{\mathcal{I}}_1, \hat{\mathcal{I}}_2, E_0, E_1, E_2 | E_{\text{inc}})$  (507-dim.)

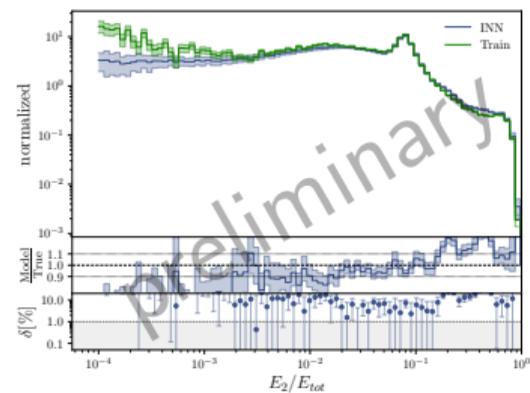
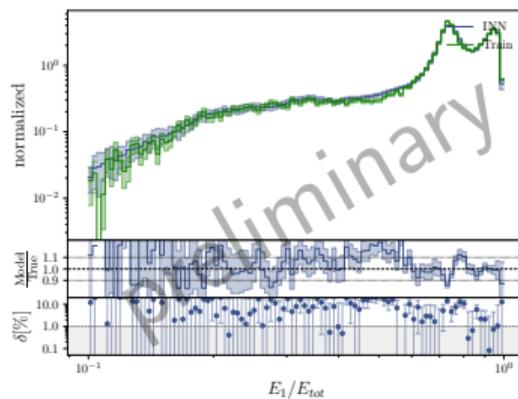
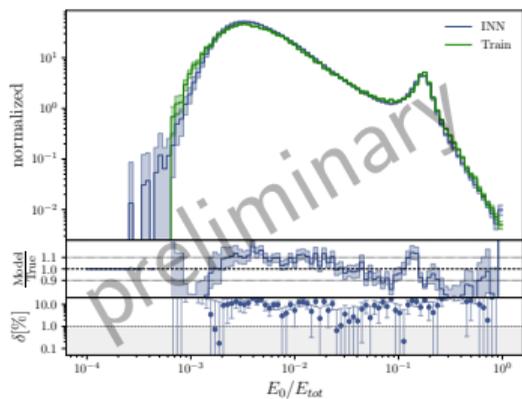
3 Make it Bayesian to

- ▶ see where the model is uncertain
- ▶ have it regularized
- ▶ have some kind of dropout

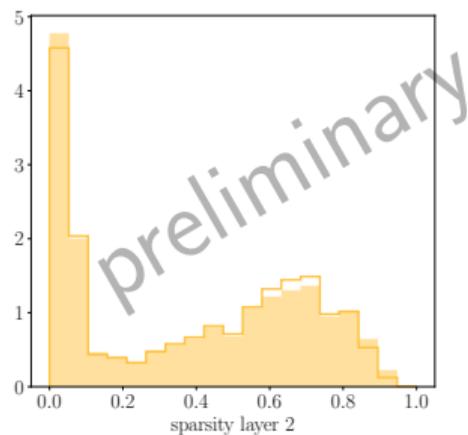
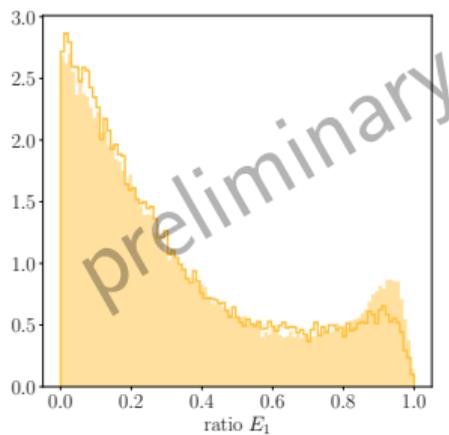
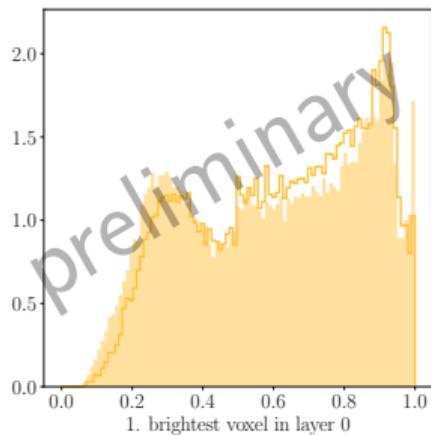
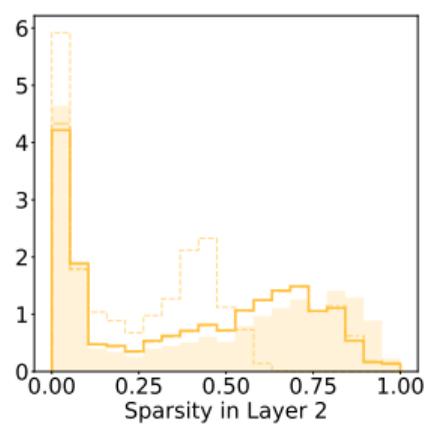
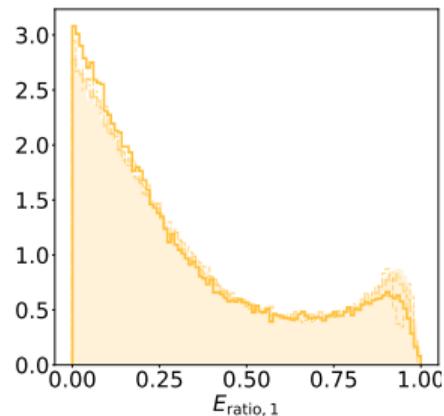
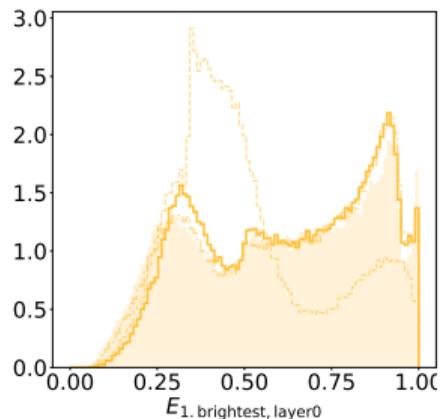
Replace NN weights  $\theta$  with Gaussians and learn their mean and width.

Bellagente, Haußmann, Luchmann, Plehn [2104.04543, SciPost]

# I: Tails are more uncertain.

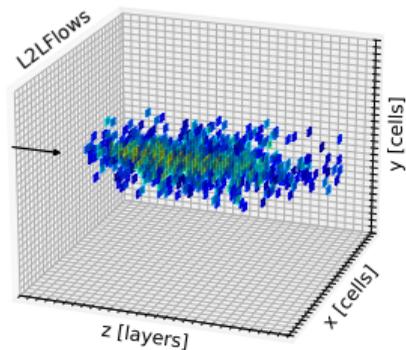
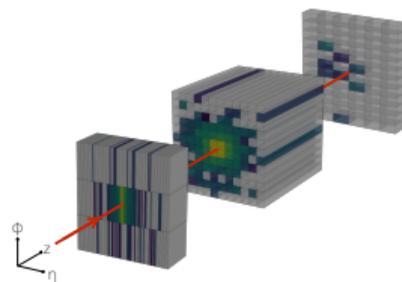


# I: Comparing preliminary $\pi^+$ histograms



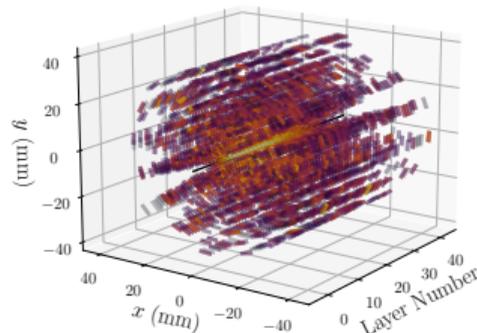
# Different Datasets come with different challenges

## Part I: The CALOGAN Dataset



## Part II: The ILD Dataset

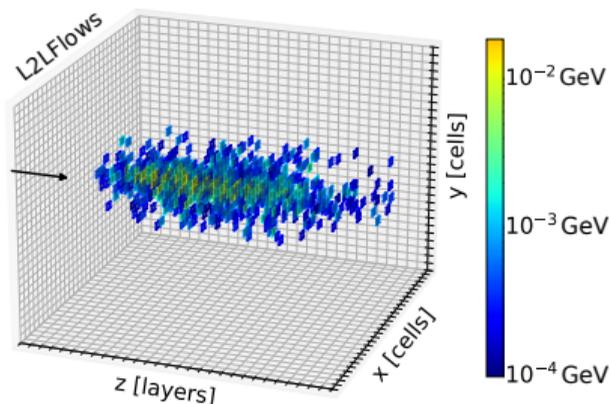
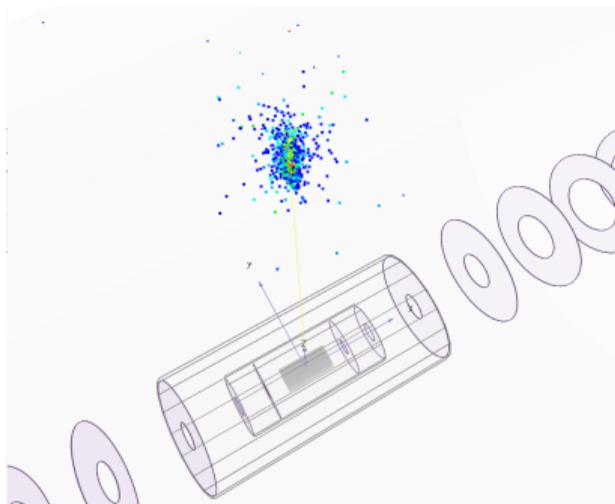
## Part III: The CaloChallenge 2022



## II: How does it scale to more voxels?

Work in Progress with Sascha Diefenbacher, Engin Eren, Frank Gaede, Gregor Kasieczka, Imahn Shekhzadeh, David Shih:

- We consider a the prototype of the ILD: alternating layers of tungsten and silicon
- We consider photon showers with  $E_{\text{inc}} \in [10, 100]$  GeV
- We discretize the data to  $30 \times 30 \times 30$  voxel and focus on the central  $10 \times 10 \times 30$



## II: We also use a 2-step approach to learn $p(\vec{\mathcal{I}}|E_{\text{inc}})$ , like CALOFLOW.

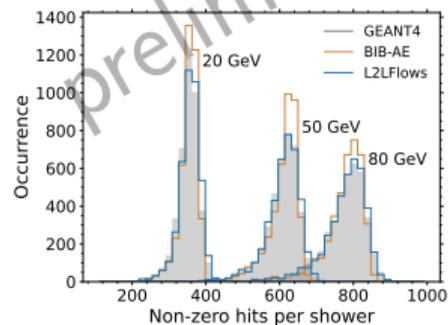
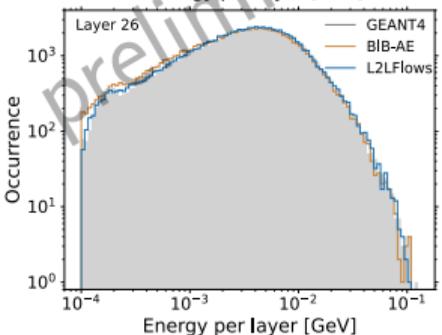
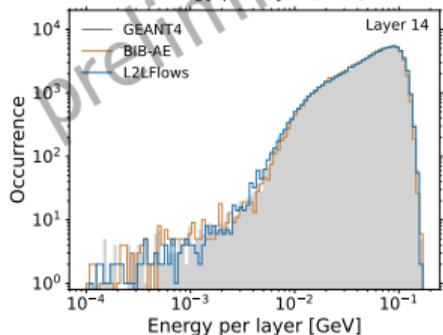
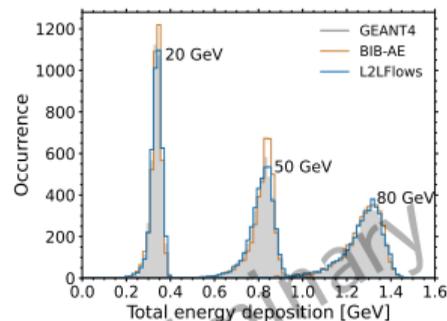
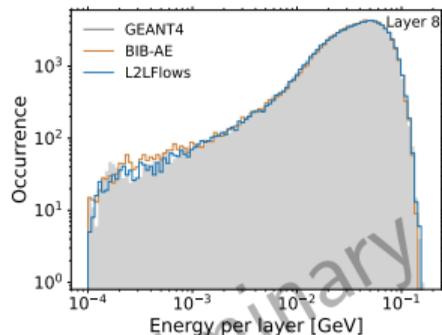
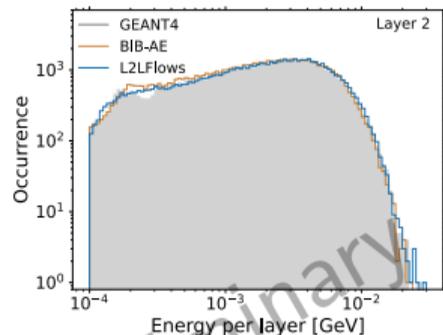
### Step I: Energy Distribution Flow

- learns  $p(E_0, E_1, E_2, \dots, E_{29}|E_{\text{inc}})$
- is a Masked Autoregressive Flow, optimized using the log-likelihood.

### Step II: Causal Flows

- use one Flow for each calorimeter layer (30 in total)
- learn  $p_i(\mathcal{I}_i|\mathcal{I}_{i-n_{\text{cond}}}, \dots, \mathcal{I}_{i-1}, E_0, \dots, E_{29}, E_{\text{inc}})$  of normalized showers.
- use an embedding net to condense the information of  $\mathcal{I}_{i-n_{\text{cond}}}, \dots, \mathcal{I}_{i-1}$ .
- are Masked Autoregressive Flows, optimized using the log-likelihood.
- can be parallelized on 30 machines in training.

## II: L2LFlows: preliminary results



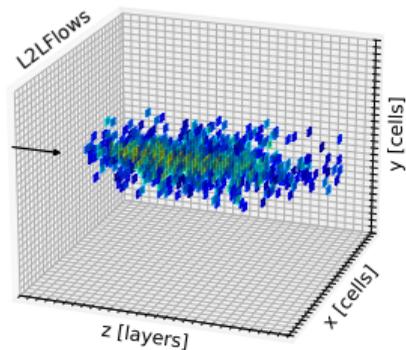
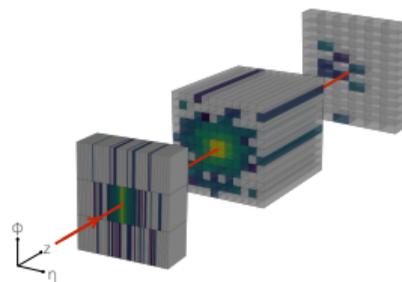
Classifier AUCs:

BIB-AE: 0.995(3)

L2LFlows: 0.852(4)

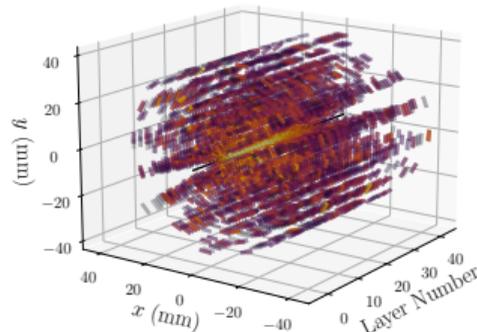
# Different Datasets come with different challenges

## Part I: The CALOGAN Dataset



## Part II: The ILD Dataset

## Part III: The CaloChallenge 2022



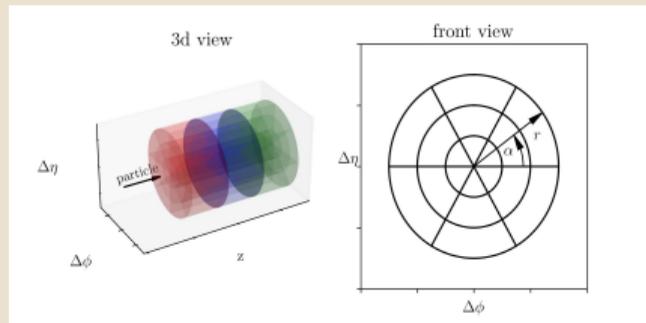
# III: Going the next step: towards deployment in FastSimulation

⇒ need a survey of current approaches on a common dataset!

⇒ Fast Calorimeter Challenge 2022

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

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



● Dataset 1: AtI Fast3 training data

( $\gamma$ : 368,  $\pi$ : 533 voxels)

[2109.02551, Comput.Softw.Big Sci.]

● Dataset 2: simulated detector

( $e^-$ : 6480 voxels)

● Dataset 3: simulated detector

( $e^-$ : 40500 voxels)

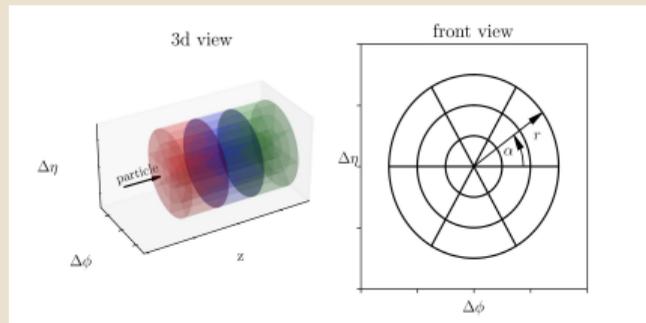
# III: Going the next step: towards deployment in FastSimulation

⇒ need a survey of current approaches on a common dataset!

⇒ Fast Calorimeter Challenge 2022

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

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



- Dataset 1: AtI Fast3 training data ( $\gamma$ : 368,  $\pi$ : 533 voxels)  
↳ Applied CALOFLOW with Ian Pang and David Shih [2109.02551, Comput.Softw.Big Sci.]
- Dataset 2: simulated detector ( $e^-$ : 6480 voxels)  
↳ Work in progress with Matt Buckley, Ian Pang and David Shih
- Dataset 3: simulated detector ( $e^-$ : 40500 voxels)  
↳ Work in progress with Matt Buckley, Ian Pang and David Shih

### III: Reminder: CALOFlow uses a 2-step approach to learn $p(\vec{\mathcal{I}}|E_{\text{inc}})$ .

#### Flow I

- learns  $p_1(E_0, E_1, E_2, \dots | E_{\text{inc}})$
- is a Masked Autoregressive Flow, optimized using the log-likelihood.

#### Flow II

- learns  $p_2(\hat{\mathcal{I}}|E_0, E_1, E_2, \dots, E_{\text{inc}})$  of normalized showers
- in CALOFlow v1 (called “teacher”):

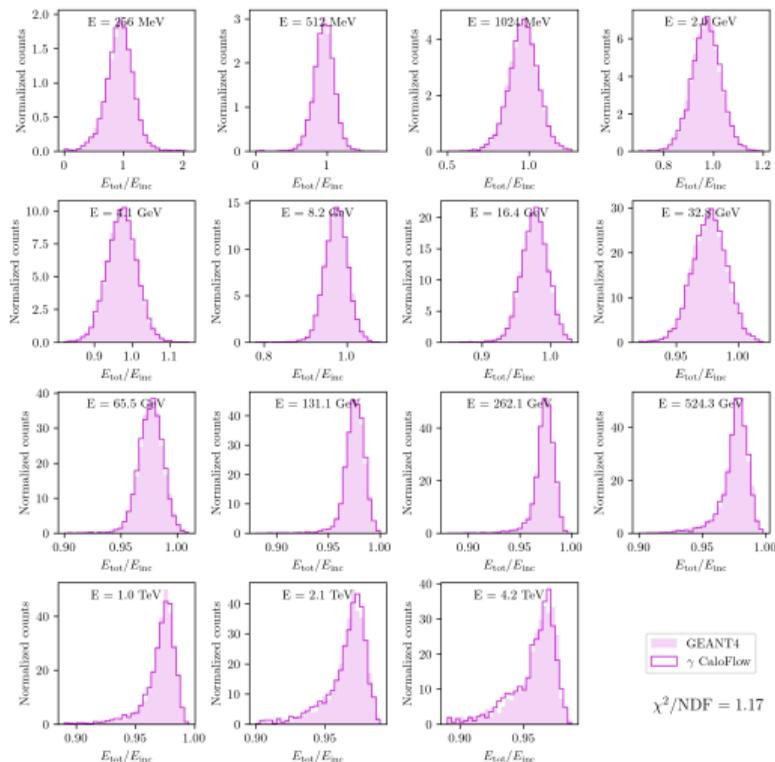
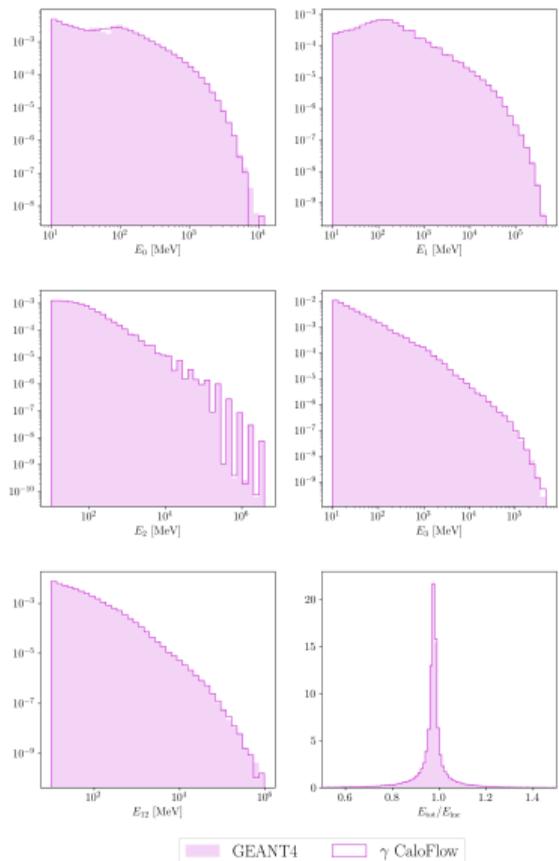
- Masked Autoregressive Flow trained with log-likelihood
- Slow in sampling

- in CALOFlow v2 (called “student”):

- Inverse Autoregressive Flow trained with Probability Density Distillation from teacher (log-likelihood prohibitive), van den Oord et al. [1711.10433]  
i.e. matching IAF parameters to frozen MAF
- Fast in sampling ( $\approx 300 - 500\times$  faster than CALOFlow v1)

# III: CALOFlow on CaloChallenge dataset 1: $\gamma$

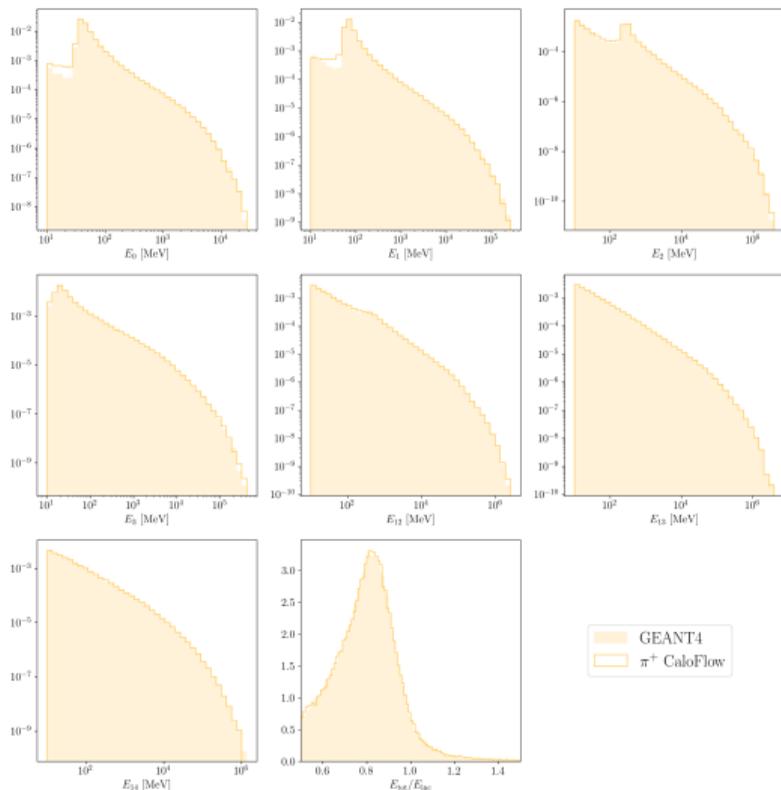
C.K., Ian Pang, David Shih [2210.14245]



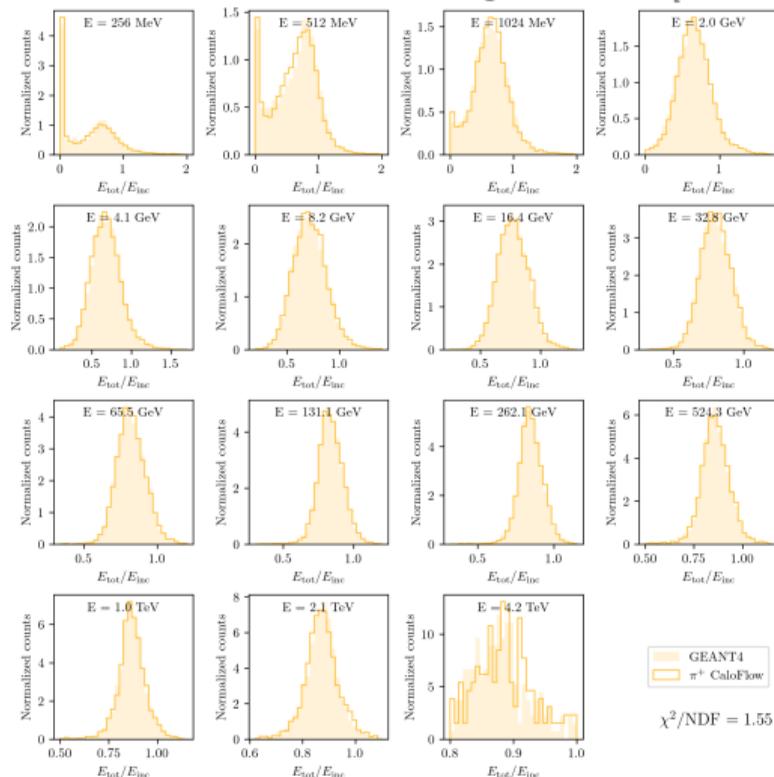
ATLAS  $\chi^2/\text{NDF} = 13.5,$

AtlFast3 [2109.02551, Comput.Softw.Big Sci.]

# III: CALOFLOW on CaloChallenge dataset 1: $\pi^+$



C.K., Ian Pang, David Shih [2210.14245]



ATLAS  $\chi^2/\text{NDF} = 12.7,$

AtlFast3 [2109.02551, Comput.Softw.Big Sci.]

### III: CALOFLOW on CaloChallene dataset 1: numbers

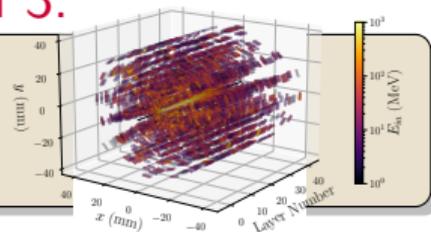
AUC		DNN based classifier	
		GEANT4 vs. CALOFLOW (teacher)	GEANT4 vs. CALOFLOW (student)
$\gamma$	low-level	0.701(3)	0.739(3)
	high-level	0.551(3)	0.556(3)
$\pi^+$	low-level	0.779(1)	0.854(2)
	high-level	0.698(2)	0.726(3)

Particle type	Batch size	Student generation time per event	
		GPU	CPU
$\gamma$	1	57.23 ms	115.88 ms
	10000	0.07 ms	-
$\pi^+$	1	74.09 ms	126.05 ms
	10000	0.09 ms	-

### III: Moving on to datasets 2 and 3.

CaloChallenge datasets 2 and 3 are much bigger:

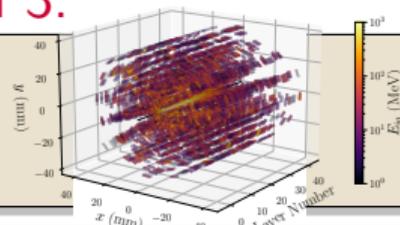
- Dataset 2: 144 voxels in 45 layers  $\rightarrow$  6480 total.
- Dataset 3: 900 voxels in 45 layers  $\rightarrow$  40500 total.



### III: Moving on to datasets 2 and 3.

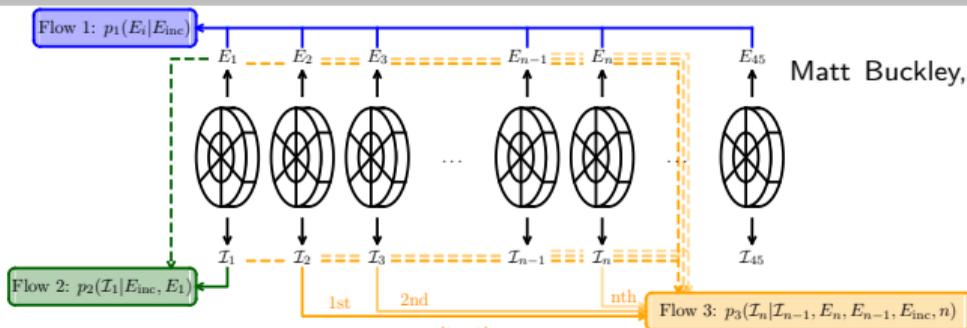
CaloChallenge datasets 2 and 3 are much bigger:

- Dataset 2: 144 voxels in 45 layers  $\rightarrow$  6480 total.
- Dataset 3: 900 voxels in 45 layers  $\rightarrow$  40500 total.



iCALOFlow: Split learning  $p(\vec{\mathcal{I}}|E_{\text{inc}})$  into 3 steps, leveraging the regular detector geometry.

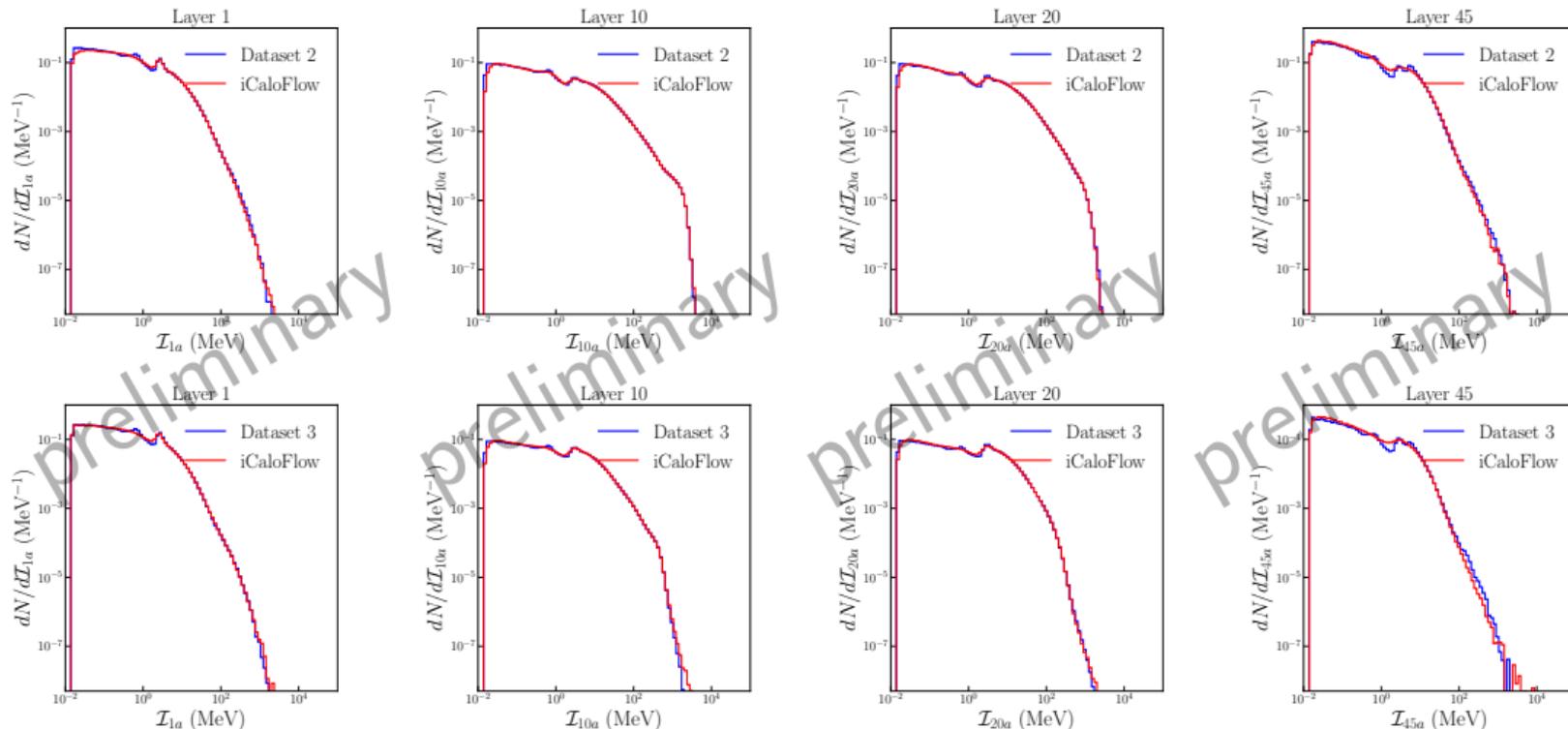
- 1 learns  $p_1(E_1, E_2, E_3, \dots, E_{45}|E_{\text{inc}})$   $\rightarrow$  how energy is distributed among layers.
- 2 learns  $p_2(\mathcal{I}_1|E_1, E_{\text{inc}})$   $\rightarrow$  how the shower in the first layer looks like.
- 3 learns  $p_3(\mathcal{I}_n|\mathcal{I}_{n-1}, n, E_n, E_{n-1}, E_{\text{inc}})$   $\rightarrow$  how the shower in layer  $n$  looks like, given layer  $n - 1$



Work in progress with

Matt Buckley, Ian Pang, David Shih

### III: iCALOFlow: results



Classifier AUCs:	dataset 2, low:	0.797(5)	dataset 2, high:	0.798(3)
	dataset 3, low:	0.911(3)	dataset 3, high:	0.941(1)

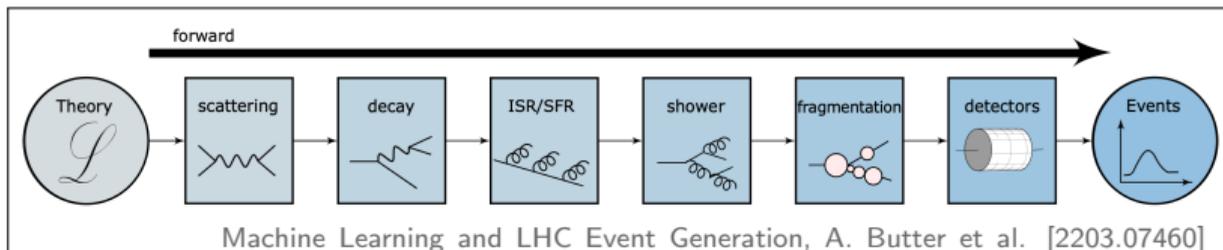
# Conclusions Part 1: Calorimeter Simulations

- ⇒ Normalizing Flows are able to generate high-quality showers, outperforming other generative models.
- ⇒ Training and model-selection is usually more stable.
- ⇒ The naive scaling to higher dimensions requires a lot of compute. But some assumptions on the underlying physics can help reduce the needed resources.

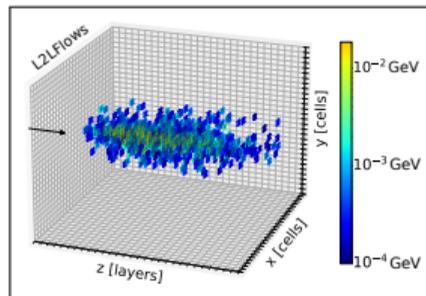
Open issues / points of discussion:

- Is learning at the voxel-level the best choice?
- How can we scale this to a full detector simulation?
- Can we run everything batch-wise?

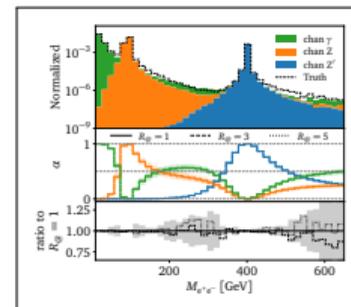
# Fast and Faithful (Calorimeter) Simulations with Normalizing Flows.



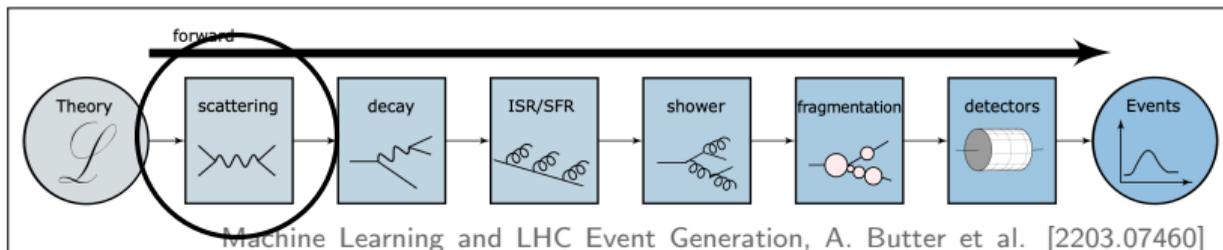
## 1: Calorimeter Simulation



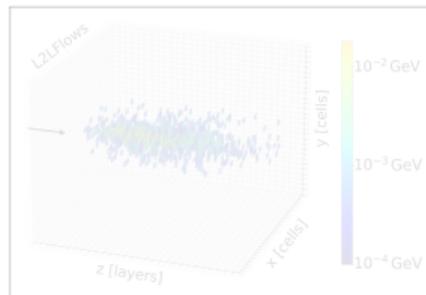
## 2: Phase Space Integration



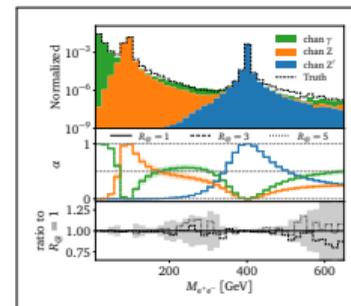
# Fast and Faithful (Calorimeter) Simulations with Normalizing Flows.



## 1: Calorimeter Simulation



## 2: Phase Space Integration



## Phase Space integration uses Importance Sampling.

$$I = \int_0^1 f(\vec{x}) d\vec{x} \quad \xrightarrow{\text{MC}} \quad \frac{1}{N} \sum_i f(\vec{x}_i) \quad \vec{x}_i \dots \text{uniform}, \quad \sigma_{\text{MC}}(I) \sim \frac{1}{\sqrt{N}}$$

$$= \int_0^1 \frac{f(\vec{x})}{q(\vec{x})} q(\vec{x}) d\vec{x} \quad \xrightarrow[\text{importance sampling}]{\text{MC}} \quad \frac{1}{N} \sum_i \frac{f(\vec{x}_i)}{q(\vec{x}_i)} \quad \vec{x}_i \dots q(\vec{x}),$$

In the limit  $q(\vec{x}) \propto f(\vec{x})$ , we get  $\sigma_{\text{IS}}(I) = 0$

We therefore have to find a  $q(\vec{x})$  that approximates the shape of  $f(\vec{x})$ .

$\Rightarrow$  Once found, we can use it for event generation,  
*i.e.* sampling  $p_i, \vartheta_i$ , and  $\varphi_i$  according to  $d\sigma(p_i, \vartheta_i, \varphi_i)$

## Phase Space integration uses Importance Sampling.

$$I = \int_0^1 f(\vec{x}) d\vec{x} \quad \xrightarrow{\text{MC}} \quad \frac{1}{N} \sum_i f(\vec{x}_i) \quad \vec{x}_i \dots \text{uniform}, \quad \sigma_{\text{MC}}(I) \sim \frac{1}{\sqrt{N}}$$

$$= \int_0^1 \frac{f(\vec{x})}{q(\vec{x})} q(\vec{x}) d\vec{x} \quad \xrightarrow[\text{importance sampling}]{\text{MC}} \quad \frac{1}{N} \sum_i \frac{f(\vec{x}_i)}{q(\vec{x}_i)} \quad \vec{x}_i \dots q(\vec{x}),$$

In the limit  $q(\vec{x}) \propto f(\vec{x})$ , we get  $\sigma_{\text{IS}}(I) = 0$

We therefore have to find a  $q(\vec{x})$  that approximates the shape of  $f(\vec{x})$ .

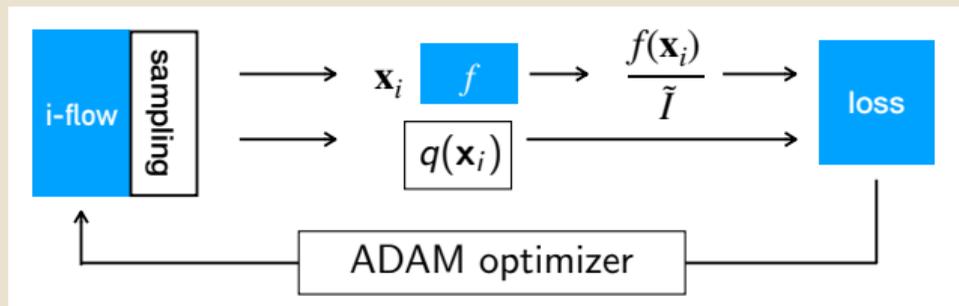
$\Rightarrow$  Once found, we can use it for event generation,  
i.e. sampling  $p_i, \vartheta_i$ , and  $\varphi_i$  according to  $d\sigma(p_i, \vartheta_i, \varphi_i)$

We need both samples  $x$  and their probability  $q(x)$ .

$\Rightarrow$  We use a bipartite, coupling-layer-based Flow.

# i-flow: Numerical Integration with Normalizing Flows.

How it works:



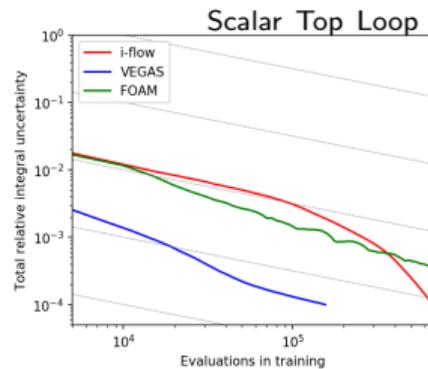
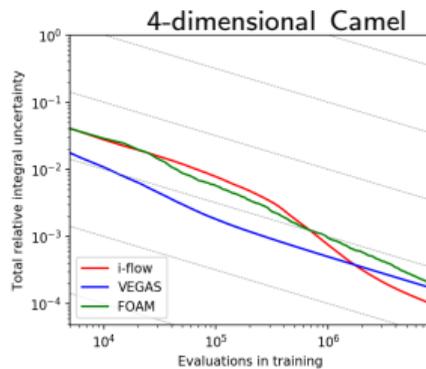
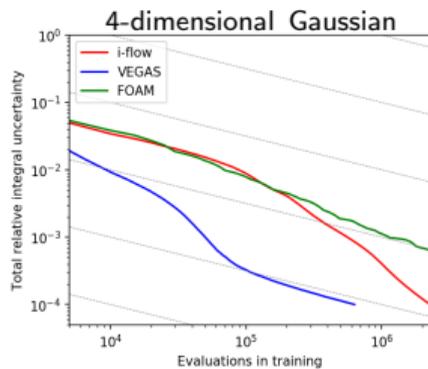
i-flow: C. Gao, J. Isaacson, CK [arXiv:2001.05486, ML:ST]  
[gitlab.com/i-flow/i-flow](https://gitlab.com/i-flow/i-flow)

Statistical Divergences are used as loss functions:

- Kullback-Leibler (KL) divergence:

$$D_{KL} = \int p(x) \log \frac{p(x)}{q(x)} dx \quad \approx \quad \frac{1}{N} \sum \frac{p(x_i)}{q(x_i)} \log \frac{p(x_i)}{q(x_i)}, \quad x_i \dots q(x)$$

i-flow adapts better, but needs more iterations.



## Event Generators need a high-dimensional integrator.

We use Sherpa (**S**imulation of **H**igh-**E**nergy **R**eactions of **P**Articles) to

- compute the matrix element of the process.
- map the unit-hypercube of our integration domain to momenta and angles. To improve efficiency, Sherpa uses a recursive multichannel algorithm.

$$\Rightarrow n_{dim} = \underbrace{3n_{final} - 4}_{\text{kinematics}} + \underbrace{n_{final} - 1}_{\text{multichannel}}$$

Similar statements apply to MadGraph, too.

<https://sherpa.hepforge.org/>

# Event Generators need a high-dimensional integrator.

We use Sherpa (**S**imulation of **H**igh-**E**nergy **R**eactions of **P**Articles) to

- compute the matrix element of the process.
- map the unit-hypercube of our integration domain to momenta and angles. To improve efficiency, Sherpa uses a recursive multichannel algorithm.

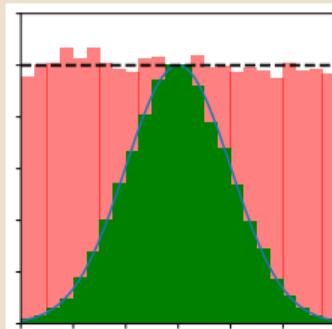
$$\Rightarrow n_{dim} = \underbrace{3n_{final} - 4}_{\text{kinematics}} + \underbrace{n_{final} - 1}_{\text{multichannel}}$$

Similar statements apply to MadGraph, too.

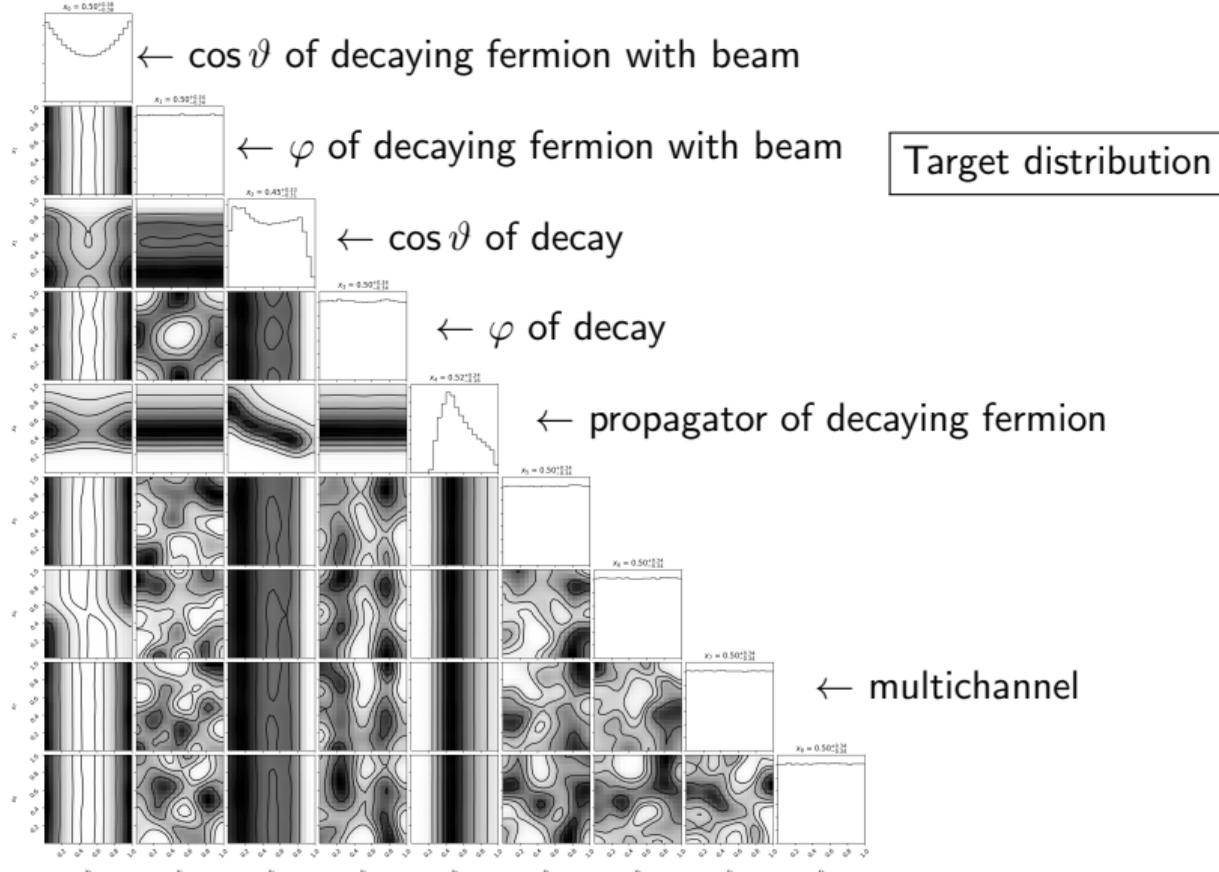
<https://sherpa.hepforge.org/>

Figure of merit: Unweighting efficiency

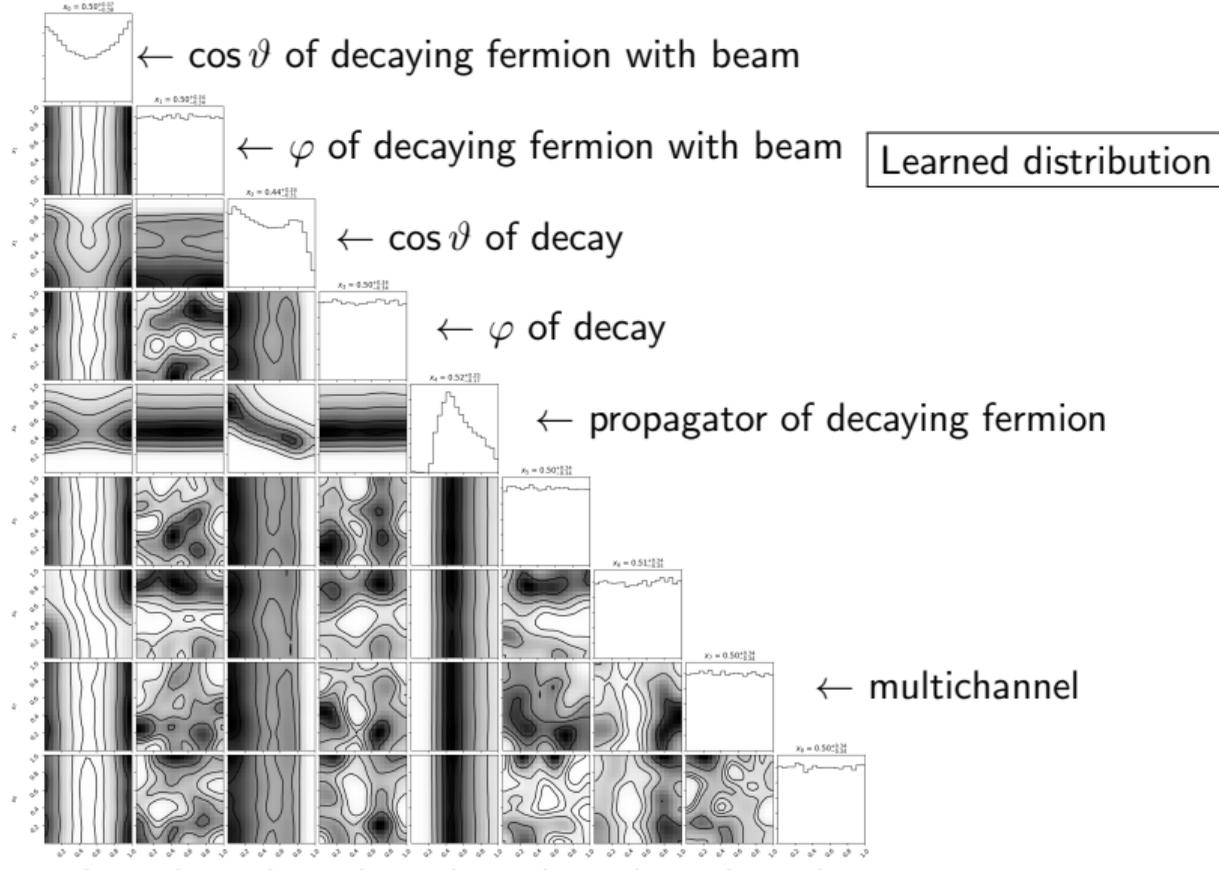
- Unweighting: we need to accept/reject each event with probability  $\frac{f(x_i)}{\max f(x)}$ . The kept events are unweighted and reproduce the shape of  $f(x)$ .
- The unweighting efficiency is the fraction of events that “survives” this procedure.



# An example: $e^+ e^- \rightarrow 3j$ .



# An example: $e^+ e^- \rightarrow 3j$ .



## High Multiplicities are difficult to learn in this setup.

unweighting efficiency $\langle w \rangle / w_{\max}$		LO QCD			
		$n=0$	$n=1$	$n=2$	$n=3$
$W^+ + n$ jets	Sherpa	$2.8 \cdot 10^{-1}$	$3.8 \cdot 10^{-2}$	$7.5 \cdot 10^{-3}$	$1.5 \cdot 10^{-3}$
	i-flow	$6.1 \cdot 10^{-1}$	$1.2 \cdot 10^{-1}$	$1.0 \cdot 10^{-2}$	$1.8 \cdot 10^{-3}$
	Gain	<b>2.2</b>	<b>3.3</b>	<b>1.4</b>	<b>1.2</b>
$W^- + n$ jets	Sherpa	$2.9 \cdot 10^{-1}$	$4.0 \cdot 10^{-2}$	$7.7 \cdot 10^{-3}$	$2.0 \cdot 10^{-3}$
	i-flow	$7.0 \cdot 10^{-1}$	$1.5 \cdot 10^{-1}$	$1.1 \cdot 10^{-2}$	$2.2 \cdot 10^{-3}$
	Gain	<b>2.4</b>	<b>3.3</b>	<b>1.4</b>	<b>1.1</b>
$Z + n$ jets	Sherpa	$3.1 \cdot 10^{-1}$	$3.6 \cdot 10^{-2}$	$1.5 \cdot 10^{-2}$	$4.7 \cdot 10^{-3}$
	i-flow	$3.8 \cdot 10^{-1}$	$1.0 \cdot 10^{-1}$	$1.4 \cdot 10^{-2}$	$2.4 \cdot 10^{-3}$
	Gain	<b>1.2</b>	<b>2.9</b>	<b>0.91</b>	<b>0.51</b>

C. Gao, S. Höche, J. Isaacson, CK, H. Schulz [arXiv:2001.10028, PRD]

## High Multiplicities are difficult to learn in this setup.

unweighting efficiency $\langle w \rangle / w_{\max}$		LO QCD			
		$n=0$	$n=1$	$n=2$	$n=3$
$W^+ + n$ jets	Sherpa	$2.8 \cdot 10^{-1}$	$3.8 \cdot 10^{-2}$	$7.5 \cdot 10^{-3}$	$1.5 \cdot 10^{-3}$
	i-flow	$6.1 \cdot 10^{-1}$	$1.2 \cdot 10^{-1}$	$1.0 \cdot 10^{-2}$	$1.8 \cdot 10^{-3}$
	Gain	<b>2.2</b>	<b>3.3</b>	<b>1.4</b>	<b>1.2</b>
$W^- + n$ jets	Sherpa	$2.9 \cdot 10^{-1}$	$4.0 \cdot 10^{-2}$	$7.7 \cdot 10^{-3}$	$2.0 \cdot 10^{-3}$
	i-flow	$7.0 \cdot 10^{-1}$	$1.5 \cdot 10^{-1}$	$1.1 \cdot 10^{-2}$	$2.2 \cdot 10^{-3}$
	Gain	<b>2.4</b>	<b>3.3</b>	<b>1.4</b>	<b>1.1</b>
$Z + n$ jets	Sherpa	$3.1 \cdot 10^{-1}$	$3.6 \cdot 10^{-2}$	$1.5 \cdot 10^{-2}$	$4.7 \cdot 10^{-3}$
	i-flow	$3.8 \cdot 10^{-1}$	$1.0 \cdot 10^{-1}$	$1.4 \cdot 10^{-2}$	$2.4 \cdot 10^{-3}$
	Gain	<b>1.2</b>	<b>2.9</b>	<b>0.91</b>	<b>0.51</b>

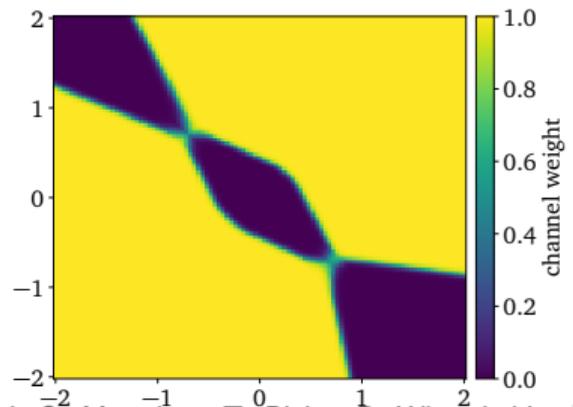
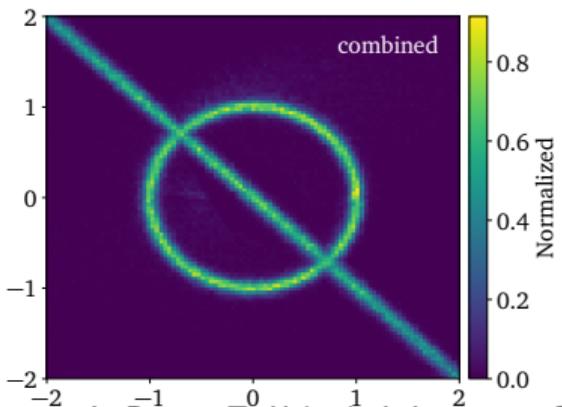
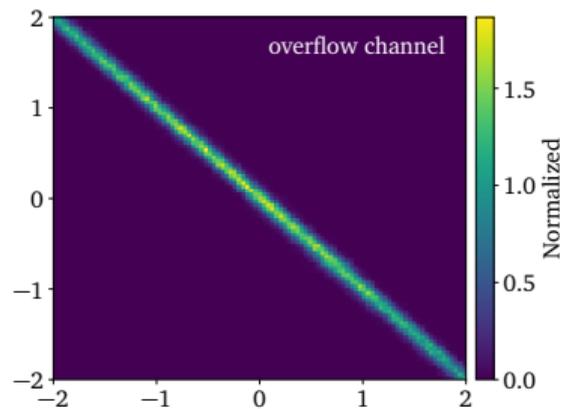
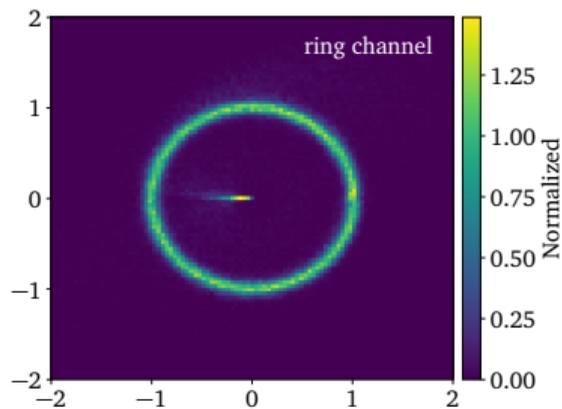
C. Gao, S. Höche, J. Isaacson, CK, H. Schulz [arXiv:2001.10028, PRD]

That cries for improvements:

⇒ MadNIS — Neural Multi-Channel Importance Sampling

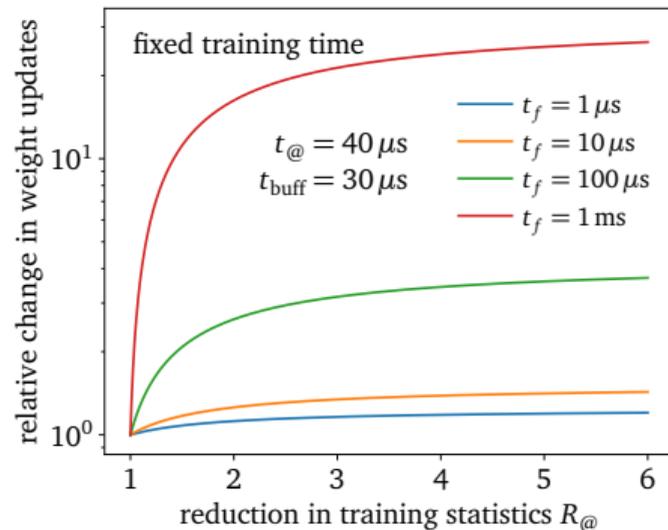
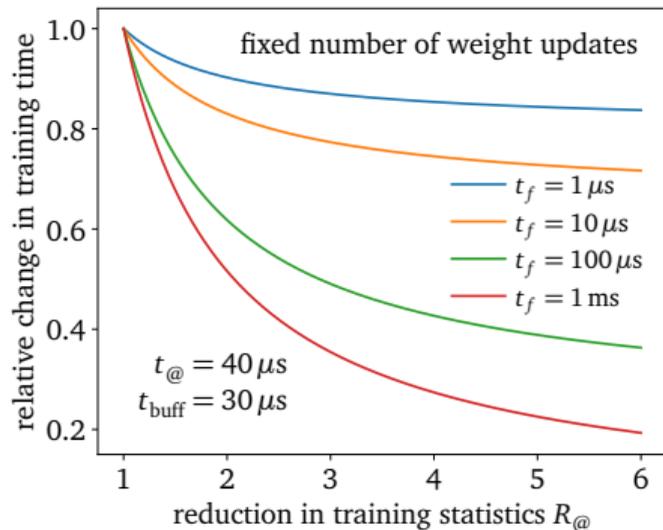
A. Butter, T. Heimgel, J. Isaacson, CK, F. Maltoni, O. Mattelaer, T. Plehn, R. Winterhalder [2212.06172]

# MadNIS learns the channel weights.



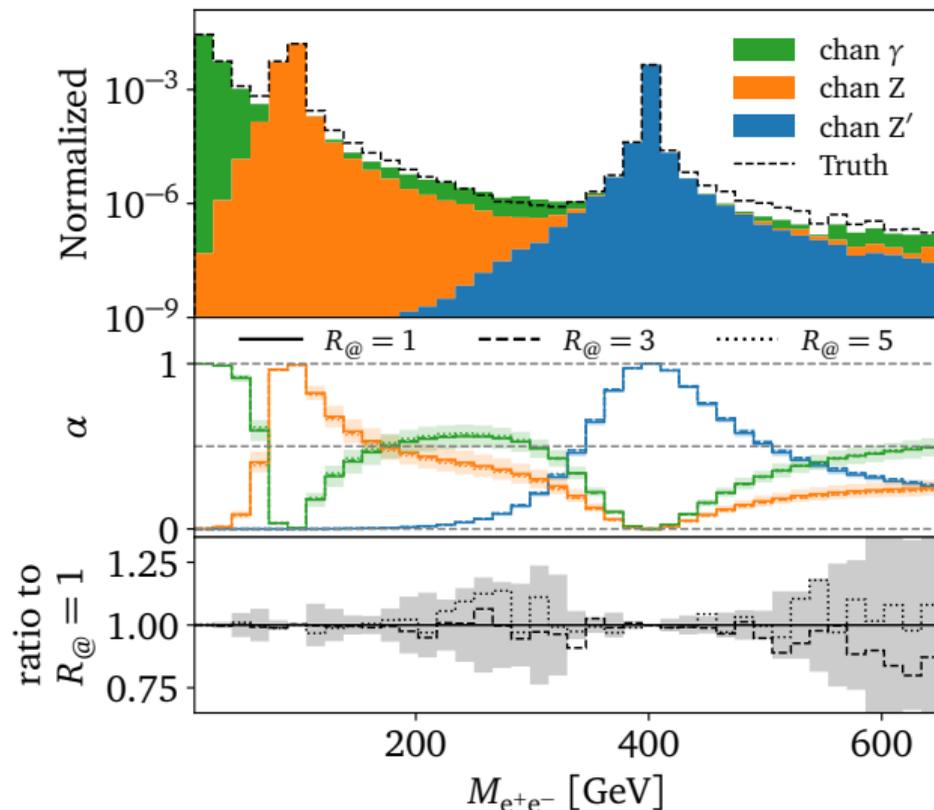
A. Butter, T. Heimes, J. Isaacson, CK, F. Maltoni, O. Mattelaer, T. Plehn, R. Winterhalder [2212.06172]

# MadNIS re-uses expensive matrix elements



A. Butter, T. HeimeI, J. Isaacson, CK, F. Maltoni, O. Mattelaer, T. Plehn, R. Winterhalder [2212.06172]

# MadNIS: putting everything together, the $Z'$ example



A. Butter, T. Heibel, J. Isaacson, CK, F. Maltoni, O. Mattelaer, T. Plehn, R. Winterhalder [2212.06172]

## Conclusions Part 2: Phase Space Integration

- ⇒ Normalizing Flows are perfect for Importance Sampling.
- ⇒ They don't introduce a bias in the result, only increase the uncertainty if not converged.
- ⇒ They can be combined with other parts of MadGraph / Sherpa.

Open issues / points of discussion:

- How is the timing of the full chain?
- What is missing to be production-ready?