

# Use of Normalizing Flows in Particle Physics Simulations

Seminar 2

author: Jan Gavranovič  
advisor: prof. dr. Borut Paul Kerševan

University of Ljubljana  
Faculty of Mathematics and Physics

April 6, 2022

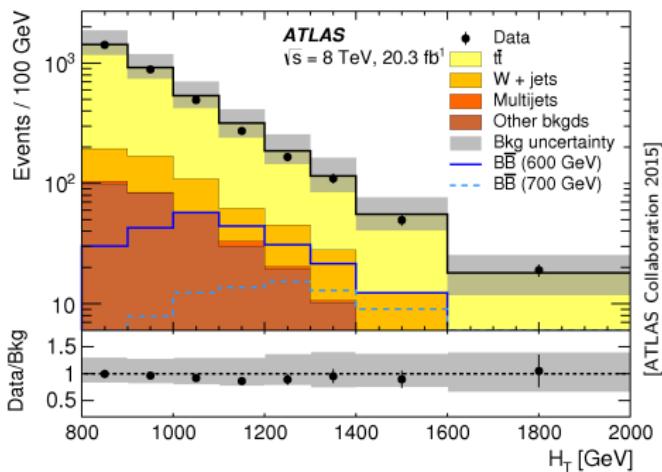
# Overview

1. Computing at the LHC
2. Simulated Higgs boson production
3. Normalizing flows
4. Affine coupling layers
5. ML event generation
6. Conclusion and discussion

# Computing at the LHC

- LHC: proton-proton collisions with center of mass energy of 14 TeV.
- Use of Monte Carlo (MC) simulations.
- Characterize signal hypotheses and study background processes.
- The simulation chain:
  1. **event generation** - simulate  $pp$  interactions using theoretical predictions,
  2. **detector simulation** - particle interactions with the material,
  3. **digitization** - conversion of hits into detector responses,
  4. **reconstruction/analysis**.
- Need large MC statistics for precise predictions  $\Rightarrow$  large computational resources.

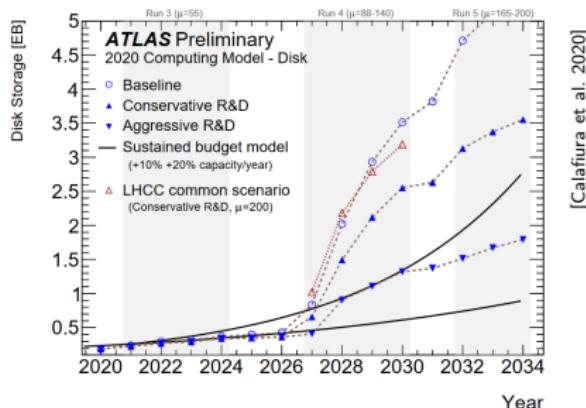
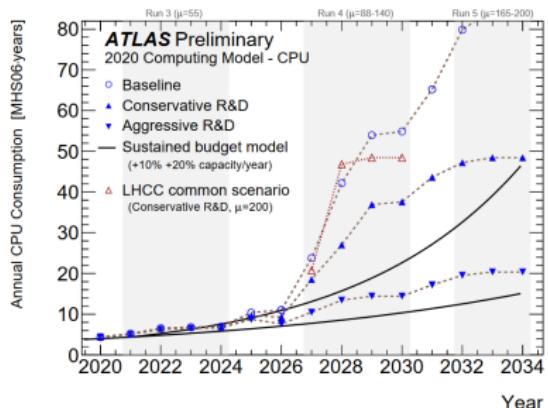
- Example: vector-like  $B$  quarks search.



- High cost of CPU time and disk space  $\Rightarrow$  restricts the discovery process.

# Computing at the LHC

- MC event generation takes up to 10 minutes *per event*.
- Produces  $\sim 1$  MB of data (only  $\sim 1$  kB useful in the final analysis).
- Computing resources provided by the Worldwide LHC Computing Grid.
- Upcoming Run 3 and in the future HL-LHC  $\Rightarrow$  need more computing power.



[Calafiura et al. 2020]

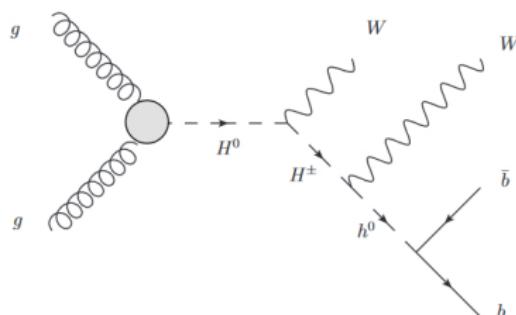
- Promising solution: **use Machine Learning (ML)**.
- **Idea:** create large amounts of events at limited computing cost using a learning algorithm trained on a small set of MC simulated events.

# Simulated Higgs boson production

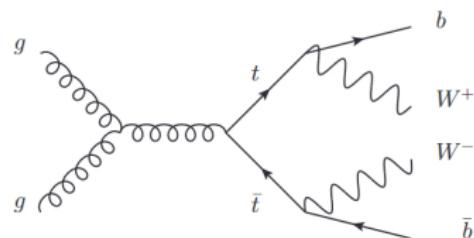
- **Signal:** theoretical Beyond-Standard-Model Higgs boson production

$$gg \rightarrow H^0 \rightarrow W^\mp H^\pm \rightarrow W^\mp W^\pm h^0 \rightarrow W^\mp W^\pm b\bar{b}.$$

- **Background:** identical decay products but distinct kinematic features.



(a) Signal process  $H^0$  and  $H^\pm$ .

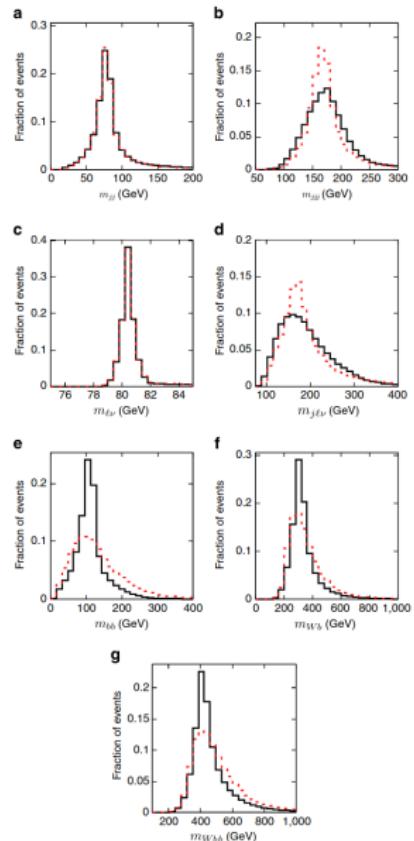


(b) Background top quark process.

- Generated events assuming 8 TeV  $pp$  collisions at the LHC.

# Simulated Higgs boson production

- **Idea:** Teach ML to simulate the final kinematic quantities used in the analysis.
- Observable decay products:
  1. leptons  $\ell$  (electrons or muons),
  2. particle jets  $j$  (streams of particles from quarks or gluons).
- One  $W$  boson  $\rightarrow \ell\nu$ , other  $W$  boson  $\rightarrow jj$
- Decay products:  $\ell\nu b$  and  $jjb$ .
- Low-level features:
  1. jet  $p_T$ ,  $\eta$  and azimuthal angle  $\phi$ ,
  2. jet  $b$ -tag,
  3. lepton  $p_T$ ,  $\eta$  and  $\phi$ ,
  4. missing energy.
- High-level features: $m_{jj}$ ,  $m_{jjj}$ ,  $m_{\ell\nu}$ ,  $m_{j\ell\nu}$ ,  $m_{b\bar{b}}$ ,  $m_{Wb\bar{b}}$ ,  $m_{WWb\bar{b}}$ .
- 18 dimensional feature space for the ML model to learn.



[Baldi et al. 2014]

# Normalizing flows: basics

- Model probability distributions using samples.
- This is **generative modeling** and is a part of unsupervised learning.
- Distribution  $p_x(\mathbf{x})$  from which our data  $\mathbf{x} \in \mathbb{R}^D$  was drawn from?

## Base distribution

Let  $\mathbf{u} \in \mathbb{R}^D$  be a random vector with a *known* probability density function (PDF)  $p_u(\mathbf{u}) : \mathbb{R}^D \rightarrow \mathbb{R}$ . Distribution  $p_u(\mathbf{u})$  is called a *base distribution* and is usually taken to be normal  $\mathcal{N}(\mathbf{0}, \mathbf{I}_{D \times D})$ .

- Express  $\mathbf{x}$  as a transformation  $T$  of a random variable  $\mathbf{u}$ :

$$\mathbf{x} = T(\mathbf{u}), \quad \mathbf{u} \sim p_u(\mathbf{u}).$$

- $T$  must be a *diffeomorphism*  $\Rightarrow$  invertible with differentiable  $T$  and  $T^{-1}$ .
- $p_x(\mathbf{x})$  can then be calculated using the **change of variables formula**:

$$\begin{aligned} p_x(\mathbf{x}) &= p_u(T^{-1}(\mathbf{x})) |\det J_T(T^{-1}(\mathbf{x}))|^{-1} \\ &= p_u(T^{-1}(\mathbf{x})) |\det J_{T^{-1}}(\mathbf{x})|. \end{aligned}$$

# Normalizing flows: basics

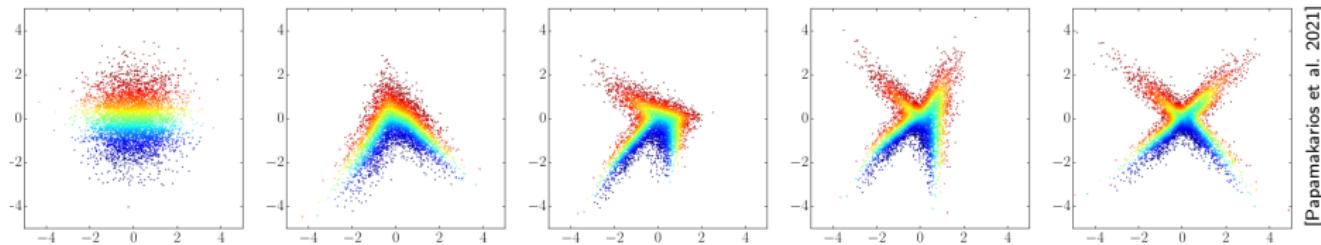
- Invertible and differentiable transformations  $T$  are composable:

$$T = T_K \circ \dots \circ T_1 \quad \text{and} \quad T^{-1} = T_1^{-1} \circ \dots \circ T_K^{-1}.$$

- Construct a **flow** by chaining together many transformations.
- Build a complicated transformation  $T$  from many simpler transformations.

## Normalizing flow

A trajectory of samples from the base (normal) distribution  $p_{\mathbf{u}}(\mathbf{u})$  as they get gradually transformed by each transformation to the target distribution  $p_{\mathbf{x}}(\mathbf{x})$ .



# Normalizing flows: basics

- Forward direction:

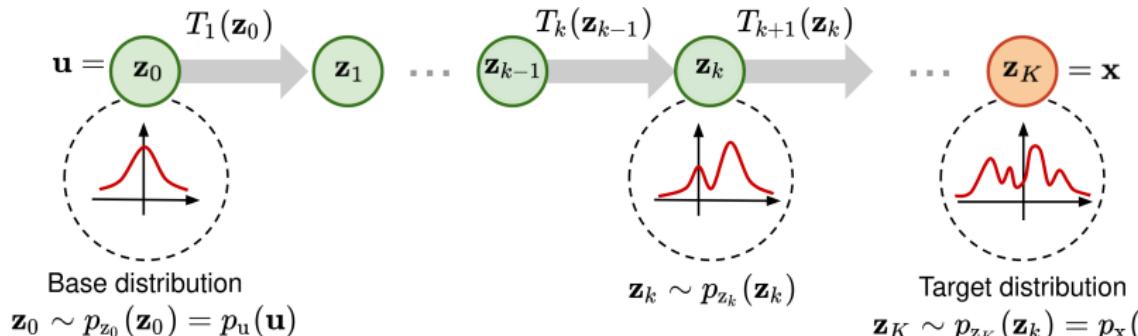
$$\mathbf{z}_k = T_k(\mathbf{z}_{k-1}) \quad \text{for } k = 1, \dots, K \quad \text{with} \quad \mathbf{z}_0 = \mathbf{u}.$$

- Inverse direction:

$$\mathbf{z}_{k-1} = T_k^{-1}(\mathbf{z}_k) \quad \text{for } k = K, \dots, 1 \quad \text{with} \quad \mathbf{z}_K = \mathbf{x}.$$

- The log-determinant of a flow is given by

$$\log |\det J_T(\mathbf{z}_0)| = \log \left| \prod_{k=1}^K \det J_{T_k}(\mathbf{z}_{k-1}) \right| = \sum_{k=1}^K \log |\det J_{T_k}(\mathbf{z}_{k-1})|.$$



[Weng 2018]

# Normalizing flows: learning objective

- Use maximum likelihood estimation.
- Fit a parametric flow model  $p_x(\mathbf{x}; \theta)$  to a target distribution  $p_x^*(\mathbf{x})$ .
- Use average log-likelihood over  $N$  data points

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \log p_x(\mathbf{x}_n; \theta).$$

- Exact log-likelihood of input data is available to us!

## Loss function

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N [\log p_u(T^{-1}(\mathbf{x}_n; \phi); \psi) + \log |\det J_{T^{-1}}(\mathbf{x}_n; \phi)|]$$

Optimized over batches using gradient based methods during training.

# Flow design

## Machine learning flow model $f_\phi$ design

1. is invertible and differentiable,
2. has efficiently computable inverse and Jacobian determinant  $\Rightarrow \mathcal{O}(D)$ .

- General flow implementation form:

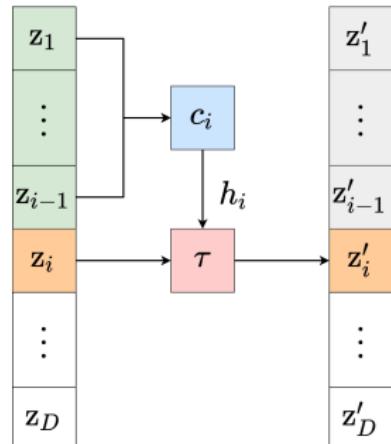
$$z'_i = \tau(z_i; \mathbf{h}_i), \quad z_i = \tau^{-1}(z'_i; \mathbf{h}_i), \quad \mathbf{h}_i = c_i(\mathbf{z}_{<i}).$$

- Autoregressive flow type.
- This makes Jacobian of  $f_\phi$  triangular

$$J_{f_\phi}(\mathbf{z}) = \begin{bmatrix} \frac{\partial \tau}{\partial z_1}(\mathbf{z}_1; \mathbf{h}_1) & \mathbf{0} \\ & \ddots \\ \mathbf{L}(\mathbf{z}) & \frac{\partial \tau}{\partial z_D}(\mathbf{z}_D; \mathbf{h}_D) \end{bmatrix}.$$

- Log-determinant uses only diagonal terms

$$\log |\det J_{f_\phi}(\mathbf{z})| = \sum_{i=1}^D \log \left| \frac{\partial \tau}{\partial z_i} (z_i; \mathbf{h}_i) \right|.$$



[Papamakarios et al. 2021]

# Affine transformer

- Use a simple transformation

$$\tau(z_i; \mathbf{h}_i) = \alpha_i z_i + \beta_i , \quad \mathbf{h}_i = \{\alpha_i, \beta_i\} ,$$

which scales and shifts the input  $\mathbf{z}_{<i}$ .

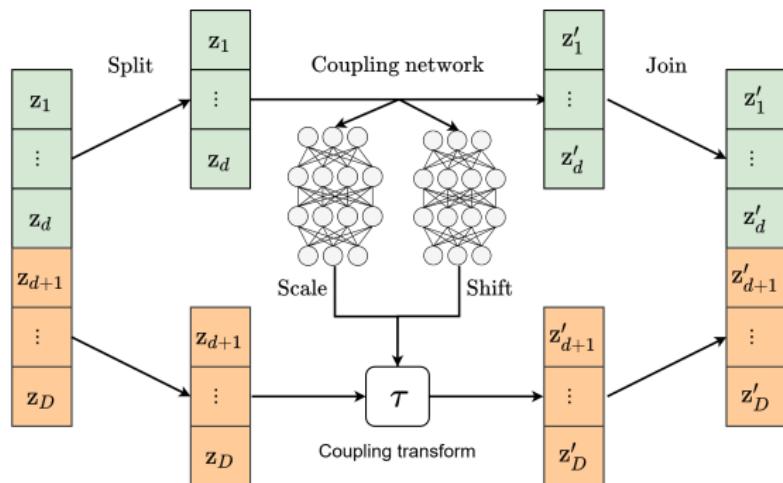
- Invertibility enforced by  $\alpha_i = \exp \tilde{\alpha}_i \Rightarrow \mathbf{h}_i = \{\tilde{\alpha}_i, \beta_i\}$ .
- Jacobian determinant is

$$\log |\det J_{f_\phi}(\mathbf{z})| = \sum_{i=1}^D \log |\alpha_i| = \sum_{i=1}^D \tilde{\alpha}_i .$$

- Scale and shift transformations implemented by two neural networks.
- **Downside:** has limited expressiveness  $\Rightarrow$  stacking multiple affine layers.

# Coupling layer

- A coupling layer splits  $\mathbf{z}$  into two parts at index  $d$  and transforms the second part as a function of the first part.
- Use *real-valued non-volume preserving (realNVP)* model.



[Dinh et al. 2016]

- Forward direction:

$$\mathbf{z}'_{\leq d} = \mathbf{z}_{\leq d}$$

$$\mathbf{h}_{>d} = \{s(\mathbf{z}_{\leq d}), t(\mathbf{z}_{\leq d})\}$$

$$\mathbf{z}'_{>d} = \mathbf{z}_{>d} \odot \exp(s(\mathbf{z}_{\leq d})) + t(\mathbf{z}_{\leq d})$$

- Jacobian is lower triangular

$$J_{f_\phi} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A} & \mathbf{D} \end{bmatrix}$$

with

$$\mathbf{D} = \text{diag}[\exp(s(\mathbf{z}_{\leq d}))] .$$

# Permutation layer

- $\mathbf{z}_{<d}$  remains unchanged when stacking coupling layers into a flow.
- **Solution:** reverse the ordering of output features.
- We can do better! Try to **learn feature permutations**.
- **Idea:** learn a generalization of a permutation matrix  $\mathbf{z}' = \mathbf{W}\mathbf{z}$ .
- **Problem:** need to calculate  $\log|\det\mathbf{W}|$  in a flow  $\Rightarrow \mathcal{O}(D^3)$ .
- Use LU to reduce this to  $\mathcal{O}(D)$ :
  1. initialize  $\mathbf{W}$  as a random rotation matrix

$$\mathbf{W}, \mathbf{R} = \text{QR} [\mathcal{N}(0, \mathbf{I}_{D \times D})] ,$$

2. parameterize  $\mathbf{W}$  by its LU decomposition

$$\mathbf{W} = \mathbf{P}\mathbf{L}(\mathbf{U} + \text{diag}(\mathbf{s})) ,$$

3. the determinant is then

$$\log|\det\mathbf{W}| = \sum_{i=1}^D \log|\mathbf{s}_i| ,$$

4. learning is performed on  $\mathbf{L}$ ,  $\mathbf{U}$  and  $\mathbf{s}$  with  $\mathbf{P}$  fixed.

# Data preprocessing and flow blocks

- Features have very different scales  $\Rightarrow$  poor performance.
- Perform data normalization (feature scaling) before training.
- This was done in three steps:
  1. min-max normalization

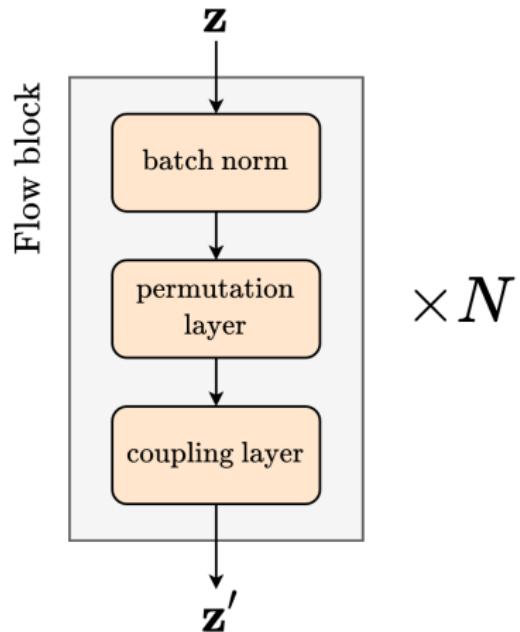
$$x_1 = \frac{x - \min(x)}{\max(x) - \min(x)},$$

2. logit transformation

$$x_2 = \log \frac{x_1}{1 - x_1},$$

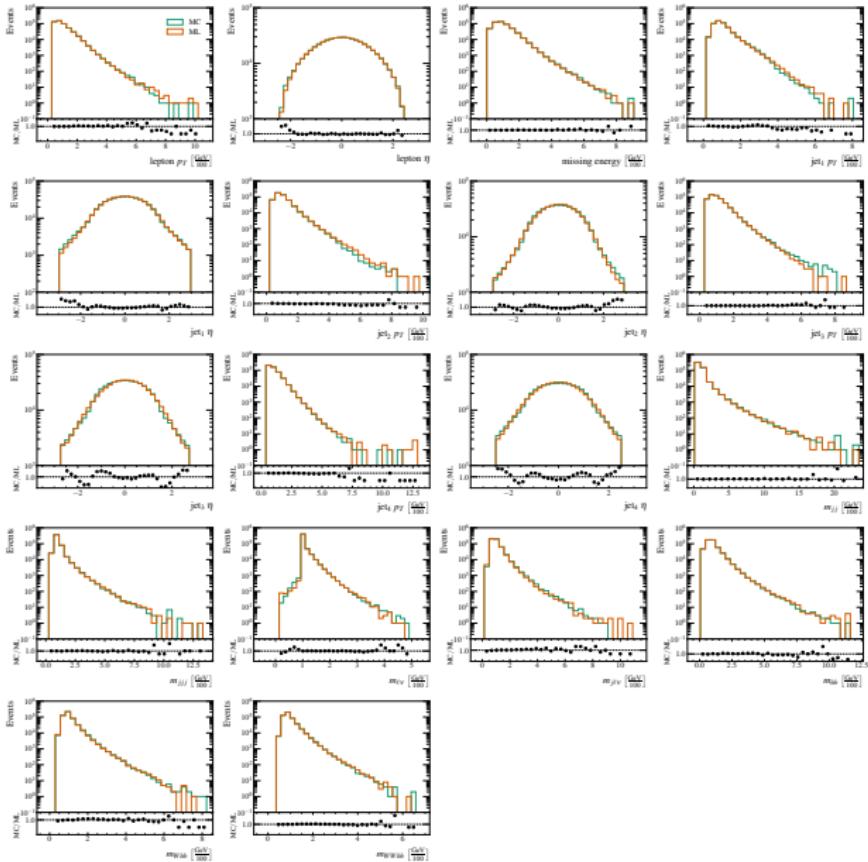
3. standardization

$$x_3 = \frac{x_2 - \mu(x_2)}{\sigma(x_2)}.$$



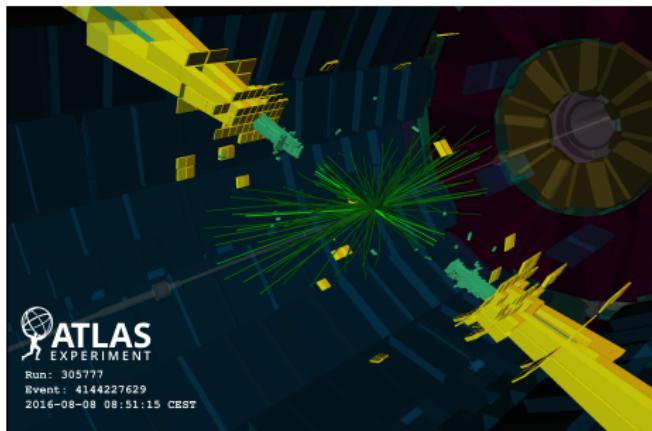
[Kingma et al. 2018]

# ML event generation



# Conclusion

- Machine learning is a rapidly advancing field of study.
- Methods are potentially applicable in particle physics tasks.
- One such method are generative models.
- We have looked at an exact log-likelihood generative approach to event generation called a normalizing flow.
- Increase in the rate of collision at the LHC will result in much higher data amounts and even more complex events to analyze.



- Machine learning will most certainly play an important role in searching for new physics beyond the Standard Model in the future.

# References

- ATLAS Collaboration (June 2015). "Search for vector-like  $B$  quarks in events with one isolated lepton, missing transverse momentum, and jets at  $\sqrt{s} = 8$  TeV with the ATLAS detector". In: *Phys. Rev. D* 91 (11), p. 112011. DOI: 10.1103/PhysRevD.91.112011.
- Calafiura, P et al. (Sept. 2020). *ATLAS HL-LHC Computing Conceptual Design Report*. Tech. rep. Geneva: CERN. URL: <https://cds.cern.ch/record/2729668>.
- Baldi, P., P. Sadowski, and D. Whiteson (July 2014). "Searching for exotic particles in high-energy physics with deep learning". In: *Nature Communications* 5.1. ISSN: 2041-1723. URL: <http://dx.doi.org/10.1038/ncomms5308>.
- Papamakarios, George et al. (2021). "Normalizing flows for probabilistic modeling and inference". In: *Journal of Machine Learning Research* 22.57, pp. 1–64.
- Weng, Lilian (2018). "Flow-based Deep Generative Models". In: Accessed February 28 2022. URL: <https://lilianweng.github.io/posts/2018-10-13-flow-models/>.
- Dinh, Laurent, Jascha Sohl-Dickstein, and Samy Bengio (2016). *Density estimation using Real NVP*. arXiv: 1605.08803 [cs.LG].
- Kingma, Durk P and Prafulla Dhariwal (2018). "Glow: Generative flow with invertible 1x1 convolutions". In: *Advances in neural information processing systems* 31.

# Thank you!



# Backup slides

and some other stuff...

# Experimental cross section and luminosity

Fixed target experiment

- Flux:

$$\Phi_a = \frac{\dot{N}_a}{A} = n_a v_a$$

$n_a$  ... density of particle beam

$v_a$  ... velocity of beam particles

$A$  ... beam area

- Reaction rate:

$$\dot{N} = \Phi_a \cdot N_b \cdot \sigma_b$$

$N_b$  ... target particles within beam area

$\sigma_b$  ... effective area of single scattering center

- Luminosity:

$$L = \Phi_a \cdot N_b$$

For the LHC:

$$N_{a,b} \sim 10^{11}, \quad A \sim 0.0005 \text{ mm}^2,$$

$$U = 27 \text{ km}, \quad n \sim 2800, \quad f \sim 11 \text{ kHz},$$

$$L \sim 10^{34} \text{ cm}^{-2}\text{s}^{-1}, \quad L_{\text{int}} = 189.3 \text{ fb}^{-1}$$

Collider experiment:

- Flux:

$$\Phi_a = \frac{\dot{N}_a}{A} = \frac{N_a \cdot n \cdot \frac{v}{U}}{A} = \frac{N_a \cdot n \cdot f}{A}$$

$N_a$  ... # particles per bunch (beam A)

$N_b$  ... # of particles per bunch (beam B)

$U$  ... circumference of the ring

$n$  ... number of bunches per beam

$v$  ... velocity of beam particles

$f$  ... revolution frequency

- Luminosity:

$$L = f \frac{n N_a N_b}{A} = f \frac{n N_a N_b}{4\pi \sigma_x \sigma_y}$$

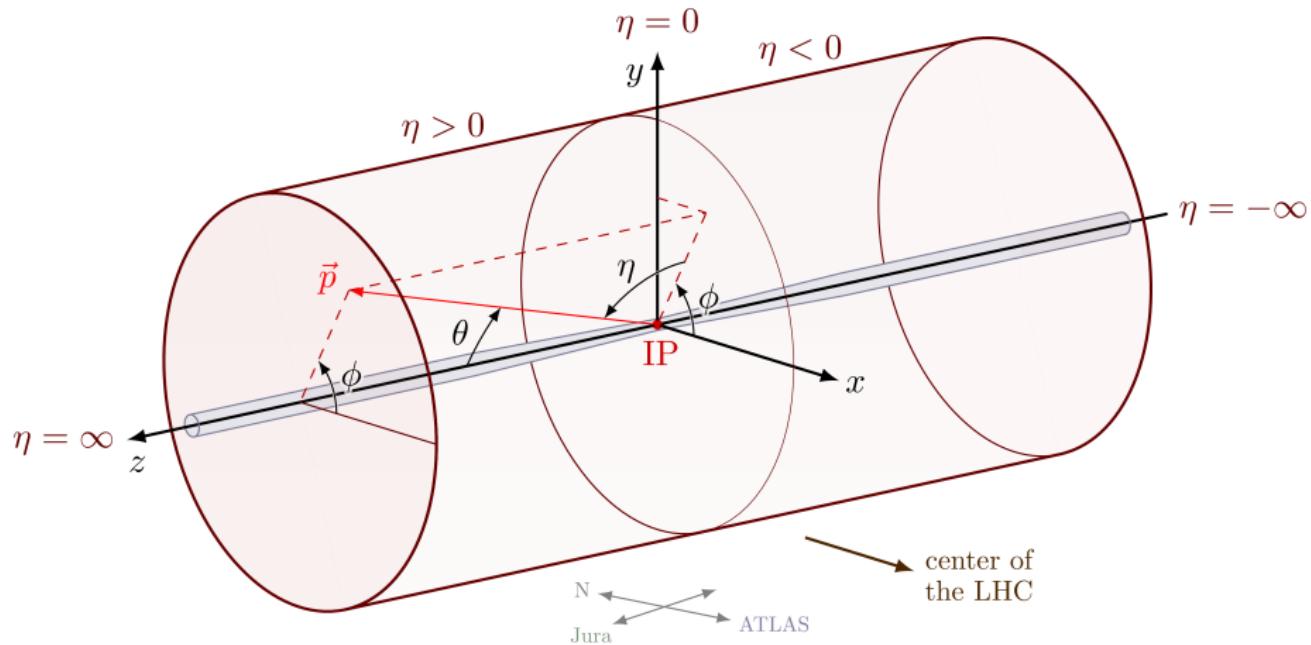
$\sigma_x$  ... std of beam profile in x

$\sigma_y$  ... std of beam profile in y

- Integrated luminosity:

$$\dot{N} \equiv L \cdot \sigma, \quad N = \sigma \cdot \int L dt, \quad \sigma = \frac{N}{L_{\text{int}}}$$

# CMS and ATLAS coordinate system



# Maximum likelihood estimation (MLE)

- Pick the parameters that assign the highest probability to the training data  $\mathcal{D}$ .
- MLE definition:

$$\hat{\theta} \equiv \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta).$$

- Assume independent and identically distributed (**iid**) samples

$$p(\mathcal{D}|\theta) = \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{x}_n, \theta),$$

where  $p(\mathbf{y}_n|\mathbf{x}_n, \theta)$  is a (conditional) likelihood.

- We usually work with the **log likelihood** defined as

$$\mathcal{L}(\theta) \equiv \log p(\mathcal{D}|\theta) = \sum_{n=1}^N \log p(\mathbf{y}_n|\mathbf{x}_n, \theta).$$

- At the heart of all model fitting/machine learning (in one way or another):

$$\hat{\theta} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta).$$

- $\mathcal{L}(\theta)$  is generally some kind of loss function or objective function.

- MLE is thus given by

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \theta) .$$

- Optimization algorithms designed for cost function minimization  $\Rightarrow$  redefine to **negative log-likelihood**:

$$\mathcal{L}(\theta) = \text{NLL}(\theta) = -\log p(\mathcal{D} | \theta) = -\sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \theta) .$$

- For conditional (unsupervised) MLE we finally have

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left[ -\sum_{n=1}^N \log p(\mathbf{y}_n | \theta) \right]$$

or written as an expected value (assuming the law of large numbers)

$$\hat{\theta} = \operatorname{argmin}_{\theta} [-\mathbb{E}_{\mathbf{y} \sim p_{\mathcal{D}}} \log p(\mathbf{y} | \theta)] ,$$

where  $p_{\mathcal{D}}$  is an empirical data distribution.

# Kullback-Leibler (KL) divergence and MLE

- KL divergence for continuous distributions:

$$D_{\text{KL}}(p\|q) \equiv \int p(x) \log \frac{p(x)}{q(x)} dx .$$

- KL divergence for discrete distributions:

$$D_{\text{KL}}(p\|q) \equiv \sum_{k=1}^K p_k \log \frac{p_k}{q_k}$$

and can be rewritten as

$$D_{\text{KL}}(p\|q) = \sum_{k=1}^K p_k \log p_k - \sum_{k=1}^K p_k \log q_k = -H(p) + H(p, q) ,$$

where  $H(p)$  is the entropy of  $p$  and  $H(p, q)$  is the cross-entropy of  $p$  and  $q$ .

- We want to find  $q$  that is as close as possible to  $p$ , measured by KL:

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q) = \operatorname{argmin}_q \left[ \int p(x) \log p(x) dx - \int p(x) \log q(x) dx \right] .$$

- Let  $p$  be the empirical distribution defined as

$$p_{\mathcal{D}}(x) = \frac{1}{N} \sum_{n=1}^N \delta(x - x_n).$$

- We get

$$\begin{aligned} D_{\text{KL}}(p_{\mathcal{D}} \| q) &= - \int p_{\mathcal{D}}(x) \log q(x) dx + C \\ &= - \int \left[ \frac{1}{N} \sum_{n=1}^N \delta(x - x_n) \right] \log q(x) dx + C \\ &= - \frac{1}{N} \sum_{n=1}^N \log q(x_n) + C = \text{NLL}(q) + C, \end{aligned}$$

where  $C$  is a constant independent of  $q$  (entropy  $H$ ).

- Thus minimizing the KL is equivalent to minimizing the NLL which is equivalent to computing the MLE.
- This is the loss function used throughout all of flow training.

# RealNVP coupling layer construction

- $c_i(\mathbf{z}_{*})*$  can be any model with input  $\mathbf{z}_{*}*$  and output  $\mathbf{h}$ .
- A coupling layer splits  $\mathbf{z}$  into two parts at index  $d$  and transforms the second part as a function of the first part:

$$\mathbf{z}'_{\leq d} = \mathbf{z}_{\leq d}$$

$$\mathbf{h}_{>d} = F(\mathbf{z}_{\leq d})$$

$$\mathbf{z}'_{>d} = \tau(\mathbf{z}_{>d}; \mathbf{h}_{>d})$$

- Use *real-valued non-volume preserving (realNVP)* model:

$$\mathbf{z}'_{\leq d} = \mathbf{z}_{\leq d}$$

$$\mathbf{h}_{>d} = \{s(\mathbf{z}_{\leq d}), t(\mathbf{z}_{\leq d})\}$$

$$\mathbf{z}'_{>d} = \mathbf{z}_{>d} \odot \exp(s(\mathbf{z}_{\leq d})) + t(\mathbf{z}_{\leq d})$$

- The inverse

$$\mathbf{z}_{\leq d} = \mathbf{z}'_{\leq d}$$

$$\mathbf{z}_{>d} = (\mathbf{z}'_{>d} - t(\mathbf{z}'_{\leq d})) \odot \exp(-s(\mathbf{z}'_{\leq d}))$$

does not require computing inverses of  $s$  and  $t \Rightarrow$  arbitrarily complex!

- The Jacobian is lower triangular with block like structure.

# RealNVP Jacobian and masking

- Jacobian is lower triangular with block like structure

$$J_{f_\phi} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A} & \mathbf{D} \end{bmatrix} .$$

- Only relevant part is

$$\mathbf{D} = \text{diag}[\exp(s(\mathbf{z}_{\leq d}))] .$$

- Determinant easy to calculate:

$$\log |\det J_{f_\phi}(\mathbf{z})| = \sum_j s(\mathbf{z}_{\leq d})_j .$$

- Splitting and joining can be implemented using binary masks  $\mathbf{b}$ :

$$\mathbf{z}' = \mathbf{b} \odot \mathbf{z} + (1 - \mathbf{b}) \odot (\mathbf{z} \odot \exp(s(\mathbf{b} \odot \mathbf{z})) + t(\mathbf{b} \odot \mathbf{z})) ,$$

$$\mathbf{z} = \mathbf{b} \odot \mathbf{z}' + (1 - \mathbf{b}) \odot (\mathbf{z}' - t(\mathbf{b} \odot \mathbf{z}')) \odot \exp(-s(\mathbf{b} \odot \mathbf{z}')) ,$$

where partitioning is handled implicitly by binary mask construction.

- An example of  $\mathbf{b}$  would be a checkerboard mask

$$b_d = d \bmod 2 \quad \text{where } d = 0, \dots, D-1 .$$

# Invertible $1 \times 1$ convolution

- Cross correlation  $\Leftrightarrow$  convolution:

$$[\mathbf{W} \star \mathbf{X}](k, i, j) = \sum_{u=1}^H \sum_{v=1}^W \sum_{c=1}^{C_{\text{in}}} w_{k,u,v,c} x_{i+u,j+v,c}.$$

- For images  $\mathbf{X}$  with height  $H$ , width  $W$  and number of input channels  $C_{\text{in}}$ .
- A  $1 \times 1$  convolution, i.e.  $H = W = 1$ , is given by

$$[\mathbf{W} \star \mathbf{X}](k, i, j) = \sum_{c=1}^{C_{\text{in}}} w_{k,c,1,1} x_{i,j,c}$$

and changes the number of channels from  $C_{\text{in}}$  to  $C_{\text{out}}$ .

- This can be written as matrix multiplication

$$[\mathbf{W} \star \mathbf{X}](k, l) = \sum_{c=1}^{C_{\text{in}}} w_{k,c} x_{c,l} = \mathbf{W}\mathbf{X}.$$

- Replacing channels with features (images  $\mathbf{X} \rightarrow$  tabular data  $\mathbf{x}$ ) gives us

$$[\mathbf{W} \star \mathbf{x}](k) = \sum_{d=1}^{D_{\text{in}}} w_{k,d} x_d = \mathbf{W}\mathbf{x}.$$

# Batch normalization

- Use batch normalization to further speed up and stabilize flow training.
- Batch norm and its inverse consist of two affine transformations:

$$\text{BN}(\mathbf{z}) = \alpha \odot \frac{\mathbf{z} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta ,$$
$$\text{BN}^{-1}(\mathbf{z}') = \mu + \frac{\mathbf{z}' - \beta}{\alpha} \odot \sqrt{\sigma^2 + \epsilon} .$$

- Updates of batch parameters are done as ( $b$  is the index of the current batch)

$$\mu_{b+1} = \rho \mu_{b-1} + (1 - \rho) \mu_b ,$$

$$\sigma_{b+1}^2 = \rho \sigma_{b-1}^2 + (1 - \rho) \sigma_b^2 ,$$

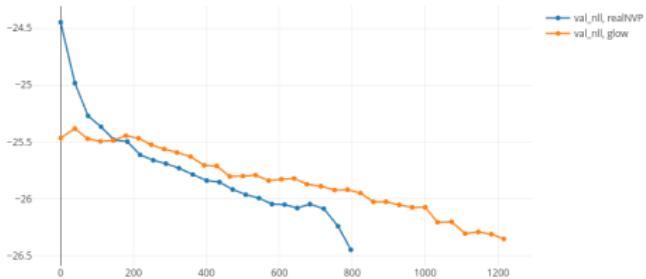
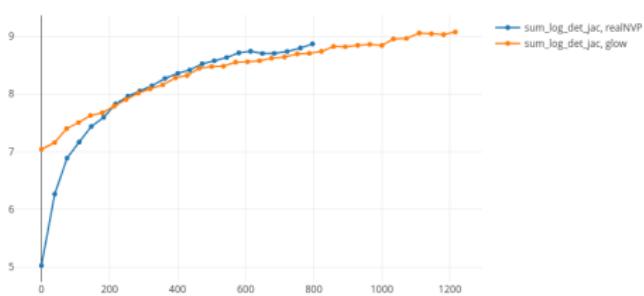
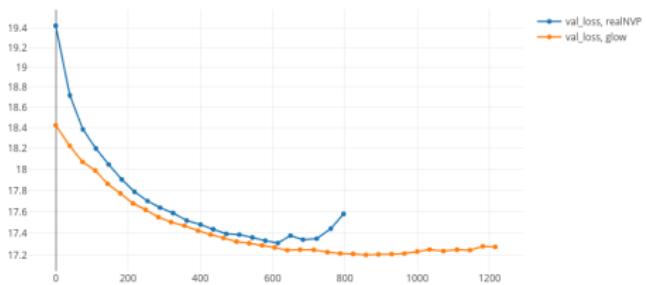
where  $\rho$  is a hyperparameter known as momentum.

- Batch norm is an elementwise operation  $\Rightarrow$  easy to compute determinant:

$$\log |\det J_{\text{BN}}(\mathbf{z})| = \sum_{i=1}^D \log \frac{\alpha_i}{\sqrt{\sigma_i^2 + \epsilon_i}} .$$

- Used as another layer after/before permutation:  $\dots \circ T_k \circ \text{BN} \circ T_{k-1} \circ \dots$

# RealNVP and Glow training curves



# Piecewise-quadratic coupling transform

- Piecewise-polynomial invertible maps as coupling transforms.
- Choose to operate in a unit hypercube:  $[0, 1]^D$ .
- Partitioning (splitting) of input into two parts  $\mathbf{z} = (\mathbf{z}^A, \mathbf{z}^B)$  at  $d$  to ensure computationally tractable Jacobians via separability.
- Each dimension transforms independently:

$$\mathbf{z}'^B = \tau [\mathbf{z}^B; \mathbf{h}(\mathbf{z}^A)] = (\tau_1[\mathbf{z}_1^B; \mathbf{h}(\mathbf{z}^A)], \dots, \tau_d[\mathbf{z}_d^B; \mathbf{h}(\mathbf{z}^A)]) .$$

- Interpret the transformation function  $\tau$  as a CDF and denoting it as  $C$ .
- This can be done due to operation on unit intervals.
- The neural network outputs the unnormalized PDF  $q \Rightarrow$  get  $C$  by integration.
- Piecewise-quadratic coupling transforms  $\Rightarrow$  piecewise-linear PDF.
- Model this PDF by  $K + 1$  vertices and  $K$  bin widths.
- NN outputs a vector of shape  $d(2K + 1)$  that we interpret as a  $d \times (K + 1)$  matrix  $V$  of vertices and  $d \times K$  matrix of bin widths  $W$ .

- The PDF  $q_i$  of  $i$ -th dimension for a data point  $z_i$  is defined as:

$$q_i(z) = \begin{cases} \frac{1}{W_{i1}} (V_{i2} - V_{i1}) z + V_{i1} ; & z < W_{i1} \\ \frac{1}{W_{i2}} (V_{i3} - V_{i2}) (z - W_{i1}) + V_{i2} ; & W_{i1} \leq z < W_{i1} + W_{i2} \\ \vdots & \vdots \\ \frac{1}{W_{iK}} (V_{i(K+1)} - V_{iK}) \left( z - \sum_{k=1}^{K-1} W_{ik} \right) + \underbrace{V_{iK}}_{\text{offset}} ; & \sum_{k=1}^{K-1} W_{ik} \leq z < 1 \end{cases}$$

slope
shift to origin

and is a type of linear interpolation function given by

$$\text{lerp}(a, b, t) = a + (b - a)t .$$

- Integrating we get the CDF (trivial once drawn on a piece of paper)

$$C_i(z_i^B, W, V) = \int_0^{z_i^B} q_i(z) dz = \frac{(\alpha_i^B)^2}{2} (V_{ib+1} - V_{ib}) W_{ib} + \alpha_i^B V_{ib} W_{ib} + \sum_{k=1}^{b-1} \frac{1}{2} (V_{ik+1} + V_{ik}) W_{ik} , \quad \alpha_i^B = \frac{1}{W_{ib}} (z_i^B - \sum_{k=1}^{b-1} W_{ik}) ,$$

where  $b \in [1, K]$  is the index of  $z_i^B$  data point bin and  $\alpha_i^B$  is the relative position of  $z_i^B$  in  $b$ .

- Note that  $W$  is normalized by the softmax function and  $V$  by the area  $\sum_{k=1}^{b-1} \frac{1}{2}(V_{ik+1} + V_{ik})W_{ik}$  where by  $V$  we mean  $V = \text{softplus}(V) + \epsilon$ .
- Solving the quadratic CDF equation for  $\alpha_i^B$  we get the inverse

$$\alpha_i^B = \frac{-V_{ib}}{V_{ib+1} - V_{ib}} + \sqrt{\frac{V_{ib}^2}{(V_{ib+1} - V_{ib})^2} - 2\beta_i^B}$$

where

$$\beta_i^B = -\frac{1}{W_{ib}(V_{ib+1} - V_{ib})} \left[ C_i(z_i^B, W, V) - \sum_{k=1}^{b-1} \frac{1}{2}(V_{ik+1} + V_{ik})W_{ik} \right].$$

- $z_i^B$  can thus be calculated from

$$z_i^B = \alpha_i^B W_{ib} + \sum_{k=1}^{b-1} W_{ik}.$$

- The log-determinant (for the forward direction) is given by

$$\log \left| \det \left( \frac{\partial C_i(z_i^B, W, V)}{\partial z_i^B} \right) \right| = \log [\alpha_i^B(V_{ib+1} - V_{ib})W_{ib} + V_{ib}W_{ib}] .$$

# Shift & scale vs piecewise-quadratic

