# Machine Learning for HEP

## Lecture III — Generative models for the LHC

Dall-E

UNIVERSITÀ DEGLI STUDI DI MILANO

BND Graduate School — Blankenberge 2024
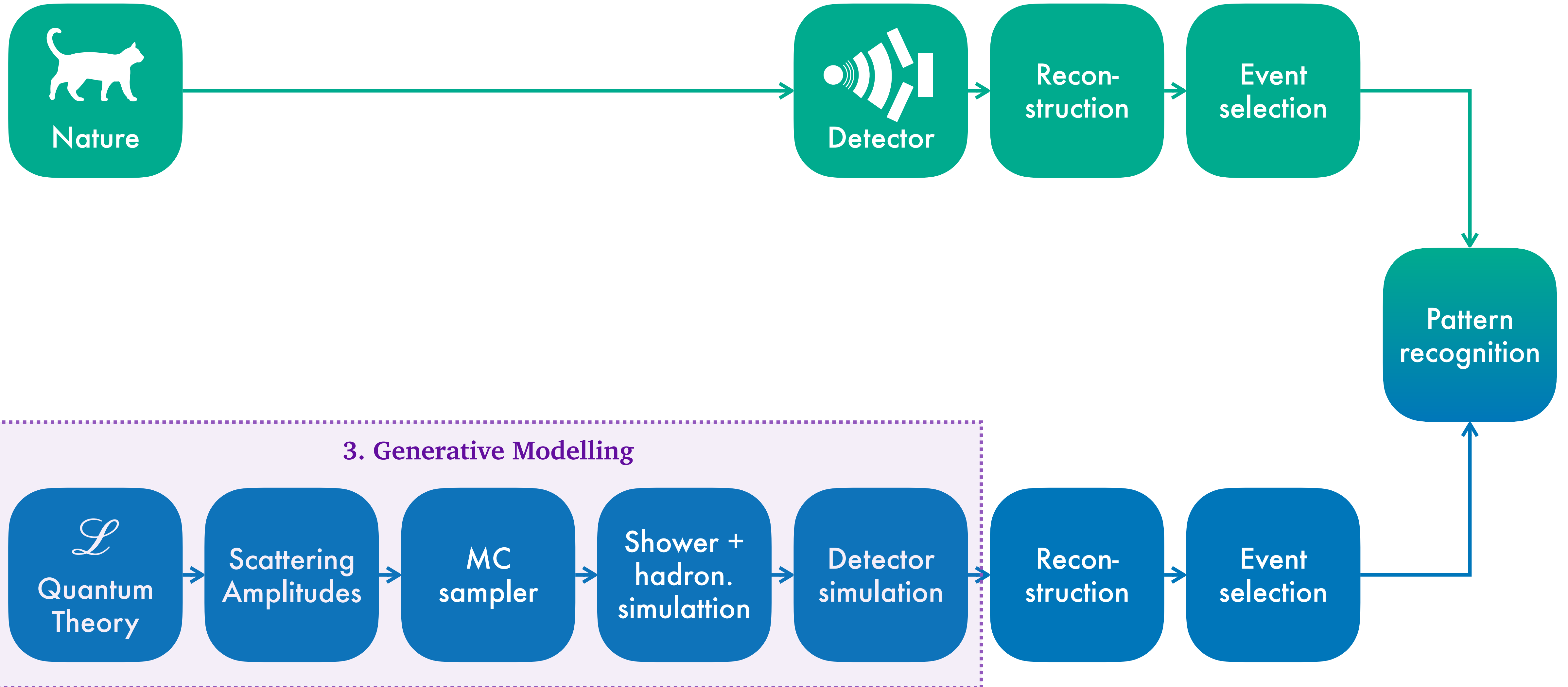
Ramon Winterhalder

# LHC analysis + ML

# Lecture III

## Generative Models for the LHC

**Text prompt**

*Realistic photography of Max Planck sitting at his laptop being super happy while coding*

Text analysis
(Typically **Transformer**)

Image generation
(Typically **Diffusion**)

Generated with Midjourney

# Generative Models

**GAN**

GAN Art (2018)
→ sold for $432,500

**Diffusion Models**

State-of-the-art
image generation

**Transformer**

ChatGPT

State-of-the-art
text generation

We have: $\boxed{p_{\text{truth}} \equiv p_{\text{data}}(x)}$ $\longrightarrow$ We want to generate new samples $\boxed{x \sim p_{\omega}(x) \simeq p_{\text{data}}(x)}$

The distribution $p_{\text{truth}}$ is usually given as:

- **explicit** as function (e.g. $\mathrm{d}\sigma \propto$ differential cross-section)

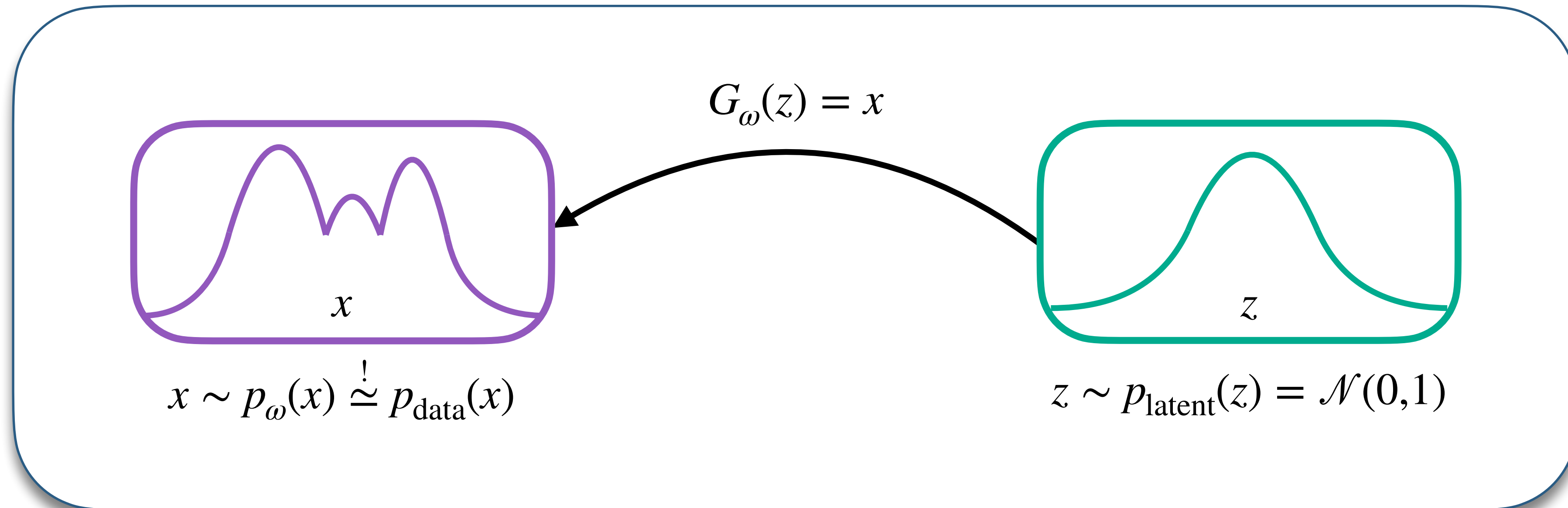- **implicit** via a set of training data $\boxed{\{x\} \sim p_{\text{data}}(x)}$

In **particle physics:**

- Event generation
- Calorimeter simulation
- Unfolding
- MEM (transfer function)

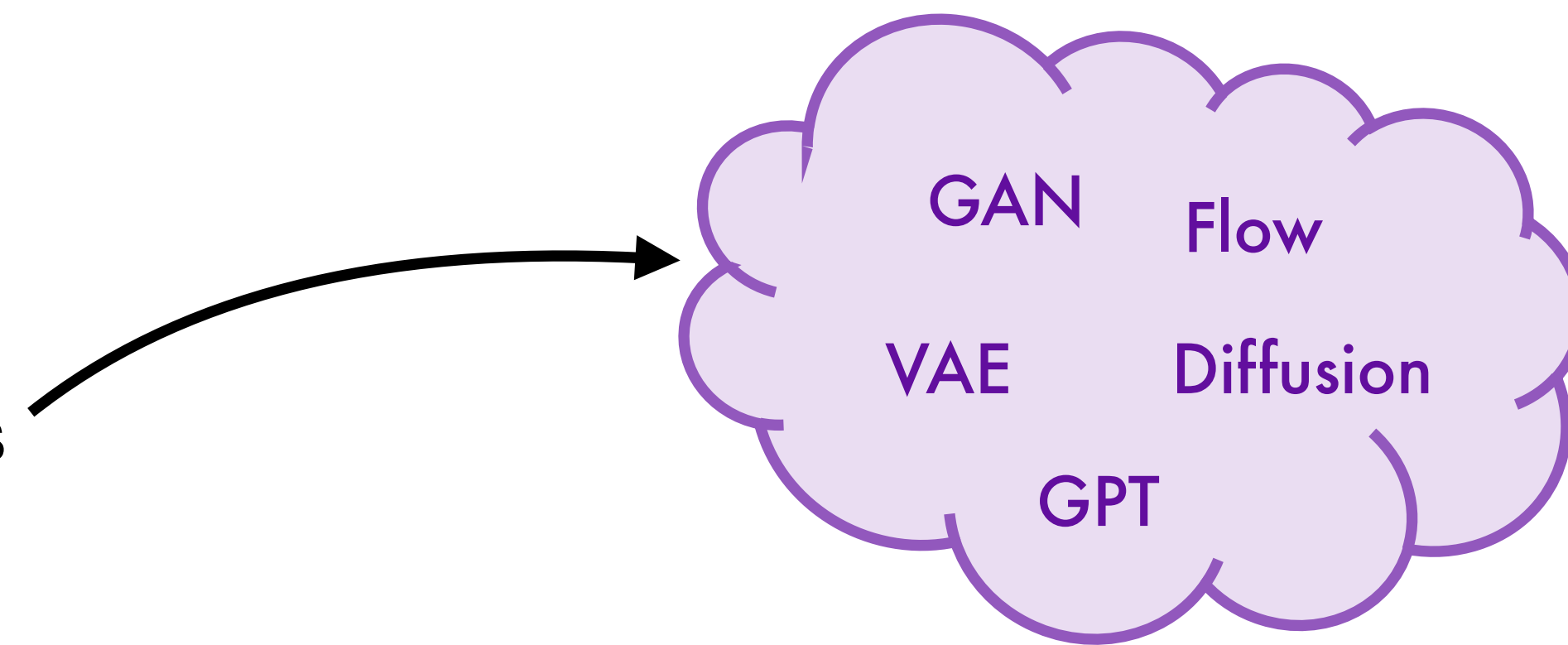→ this is a stochastic (random) process (RNG)

→ needs "random" input

$$G_\omega(z) = x$$

$$x \sim p_\omega(x) \stackrel{!}{\simeq} p_{\text{data}}(x)$$

$$z \sim p_{\text{latent}}(z) = \mathcal{N}(0,1)$$

→ How to **construct** and **train** $G_\omega(z)$?

→ **Multiple types** of generative models

GAN  Flow

VAE  Diffusion

GPT

# Types of deep generative models

# Deep generative models

# Deep generative models

**β-VAE**

**Hierarchical VAE**

**②**

**Diffusion Probabilistic Model**

## Variational Autoencoder

## Diffusion Model

**VQ-VAE**

Score-matching Model

**Conditional Flow Matching**

SurVAE

**①**

**③**

Wasserstein GAN

**Normalizing Flow**

Continuous NFs

## Generative Adversarial Network

## Maximum-likelihood Models

LS-GAN

Relativistic GAN

Autoregressive Transfomer (GPT)

# Part I

Generative Adversarial Networks

Synthetic samples

Generator $G_\omega(z)$

Real samples

Discriminator $D_\varphi(x)$

Real or Fake?

**Synthetic samples**

**Real samples**

Generator $G_\omega(z)$

Discriminator $D_\varphi(x)$

**Real** or **Fake**?

$$\mathscr{L}_D = -\left\langle \log D_\varphi(x) \right\rangle_{x \sim p_\text{data}} - \left\langle \log(1 - D_\varphi(x)) \right\rangle_{x \sim p_{\hat{\omega}}}$$

$$= -\left\langle \log D_\varphi(x) \right\rangle_{x \sim p_\text{data}} - \left\langle \log(1 - D_\varphi(G_{\hat{\omega}}(z))) \right\rangle_{z \sim p_\text{latent}}$$

**Iterative Training**

$$\mathscr{L}_G = \left\langle \log(1 - D_{\hat{\varphi}}(x)) \right\rangle_{x \sim p_\omega}$$

$$= \left\langle \log(1 - D_{\hat{\varphi}}(G_\omega(z))) \right\rangle_{z \sim p_\text{latent}}$$

## Problems with GANs
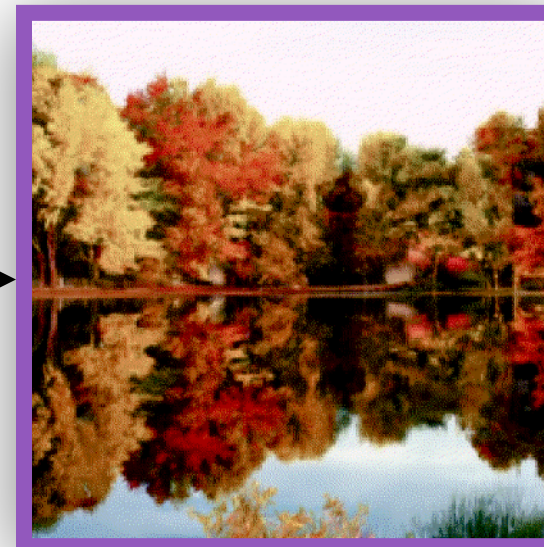
⊖ Min-max training unstable
→ vanishing gradients

⊖ Metric for sucess?
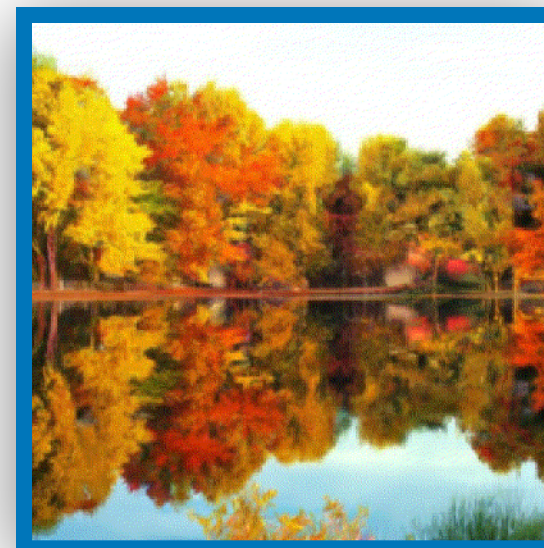→ loss only shows competition!

⊖ Mode collapse

## Ways to improve

⊕ Wasserstein GAN **[1701.07875, 1704.00028]**

⊕⊕ Gradient penalty **[1705.09367, 1801.04406]**

⊕⊕ Spectral normalization **[1802.05957]**

⊕ Other losses
– Least square **[1611.04076]**
– Relativistic GAN **[1807.00734]**
– …

Synthetic
samples

Generator
$G_\omega(z)$

Real
samples

Discriminator
$D_\varphi(x)$

Real
or
Fake?

Iterative
Training

$$\mathscr{L}_G = \left\langle \log(1 - D_{\hat\varphi}(x)) \right\rangle_{x \sim p_\omega}$$

$$= \left\langle \log(1 - D_{\hat\varphi}(G_\omega(z))) \right\rangle_{z \sim p_{\text{latent}}}$$

$$\left\langle \quad (x)) \right\rangle_{x \sim p_{\hat\omega}}$$

$$\left\langle \quad (G_{\hat\omega}(z)) \right\rangle_{z \sim p_{\text{latent}}}$$

# Part II

**Diffusion Models**
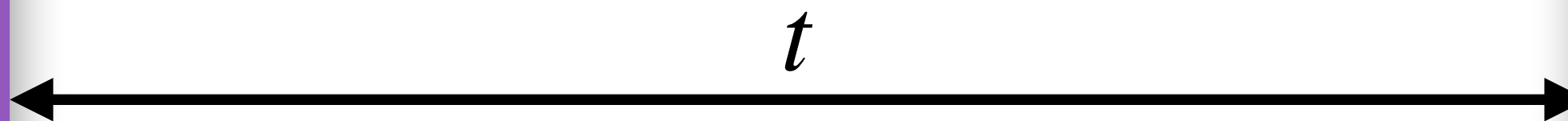
Midjourney

Dall-E

# Recently…

Stable Diffusion

Midjourney AI

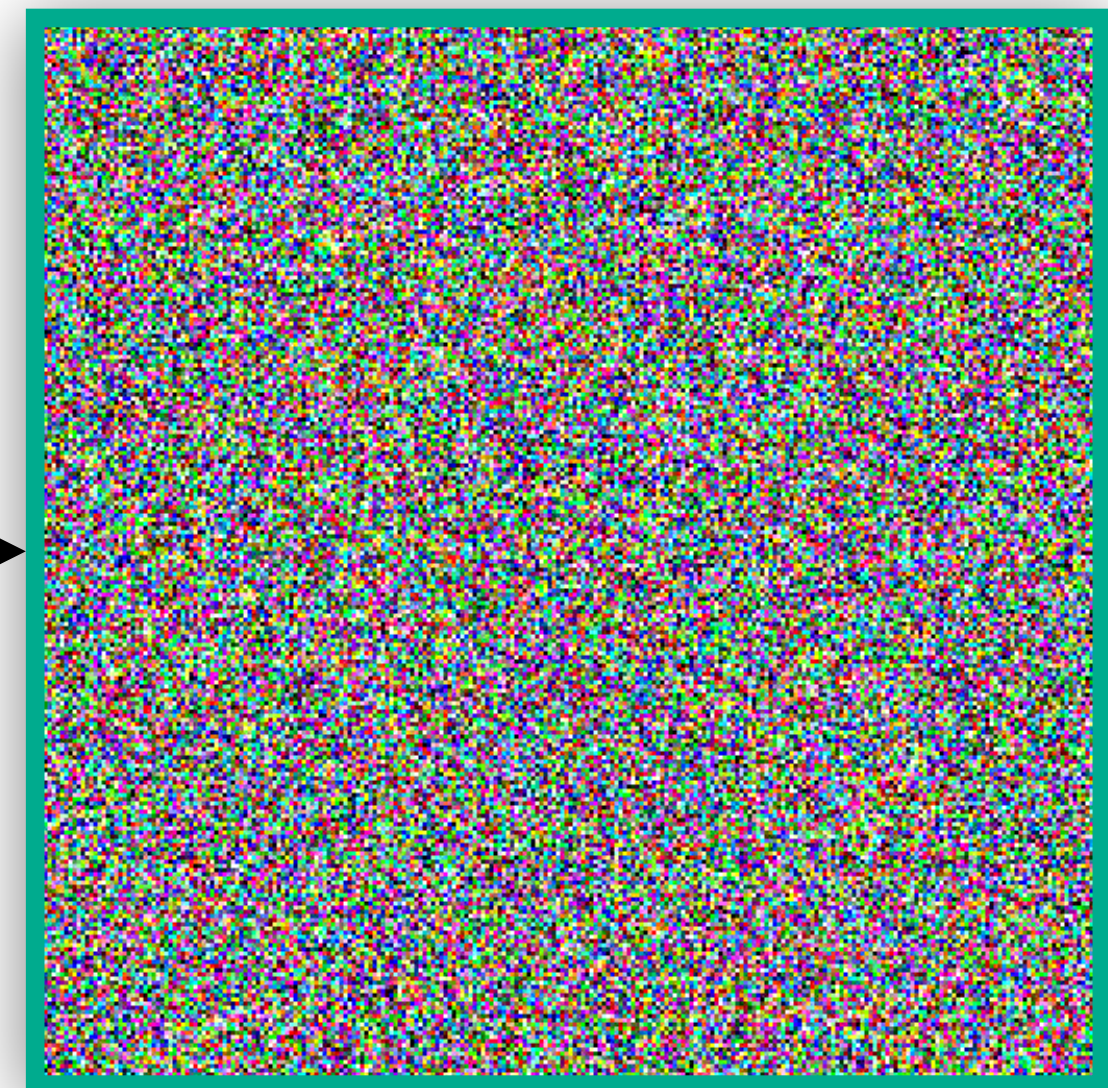Define mapping as **time-dependent** diffusion process



$t$

$$x \sim p_\omega(x) \overset{!}{\simeq} p_{\text{data}}(x)$$

**Parametrization of time**

- discrete: $t = 0, 1, \ldots, T$
- continuous: $t \in [0,1]$

$$z \sim p_{\text{latent}}(z) = \mathcal{N}(0,1)$$

→ Gradually **add noise** to data samples to transform them to gaussians

← Gradually **remove noise** from gaussians to obtain data samples

**Bayes' Theorem**

$$p(z \,|\, x) = \frac{p(x \,|\, z)\, p(z)}{p(x)}$$

← Prior

← Evidence

**Problem: evidence intractable**

$$p(x) = \int \mathrm{d}z\, p(x \,|\, z)\, p(z)$$

$$p_\omega(x \,|\, z)$$

$X$     $Z$

$$p(z \,|\, x)$$

Likelihood:   $p_\omega(x \,|\, z)$   ← Decoder

**Variational inference**

Posterior:   $p(z \,|\, x)$   ← Want!

Encoder:   $q_\phi(z \,|\, x) = \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))$

**Match!**

[1312.6114, 1906.02691]

**Minimize KL**

$$\text{KL}(q_\phi(z|x), p(z|x)) = \int \mathrm{d}z \, q_\phi(z|x) \, \log\left(\frac{q_\phi(z|x)}{p(z|x)}\right)$$

**Maximize ELBO**



ELBO ≠ Elmo

**Minimize KL**

$$\mathrm{KL}(q_\phi(z|x), p(z|x)) = \int \mathrm{d}z\, q_\phi(z|x)\, \log\left(\frac{q_\phi(z|x)}{p(z|x)}\right)$$

**Maximize ELBO**

$$\log p(x) = \int \mathrm{d}z\, q_\varphi(z|x)\, \log p(x)$$

Bayes' theorem

$$= \int \mathrm{d}z\, q_\varphi(z|x)\, \log \frac{p_\omega(x|z)\, p(z)}{p(z|x)}$$

$$= \int \mathrm{d}z\, q_\varphi(z|x) \left[ \log p_\omega(x|z) - \log \frac{q_\phi(z|x)}{p(z)} + \log \frac{q_\phi(z|x)}{p(z|x)} \right]$$

$$= \left\langle \log p_\omega(x|z) \right\rangle_{q_\phi(z|x)} - \mathrm{KL}(q_\phi(z|x), p(z)) + \mathrm{KL}(q_\phi(z|x), p(z|x))$$

$$\geq \left\langle \log p_\omega(x|z) \right\rangle_{q_\phi(z|x)} - \mathrm{KL}(q_\phi(z|x), p(z))$$

$\geq 0$

ELBO

**SurVAE [2007.02731]**

$$\log p(x) = \int \mathrm{d}z\, q_\varphi(z|x) \left[ \log p(z) + \log \frac{p_\omega(x|z)}{q_\phi(z|x)} + \log \frac{q_\phi(z|x)}{p(z|x)} \right]$$

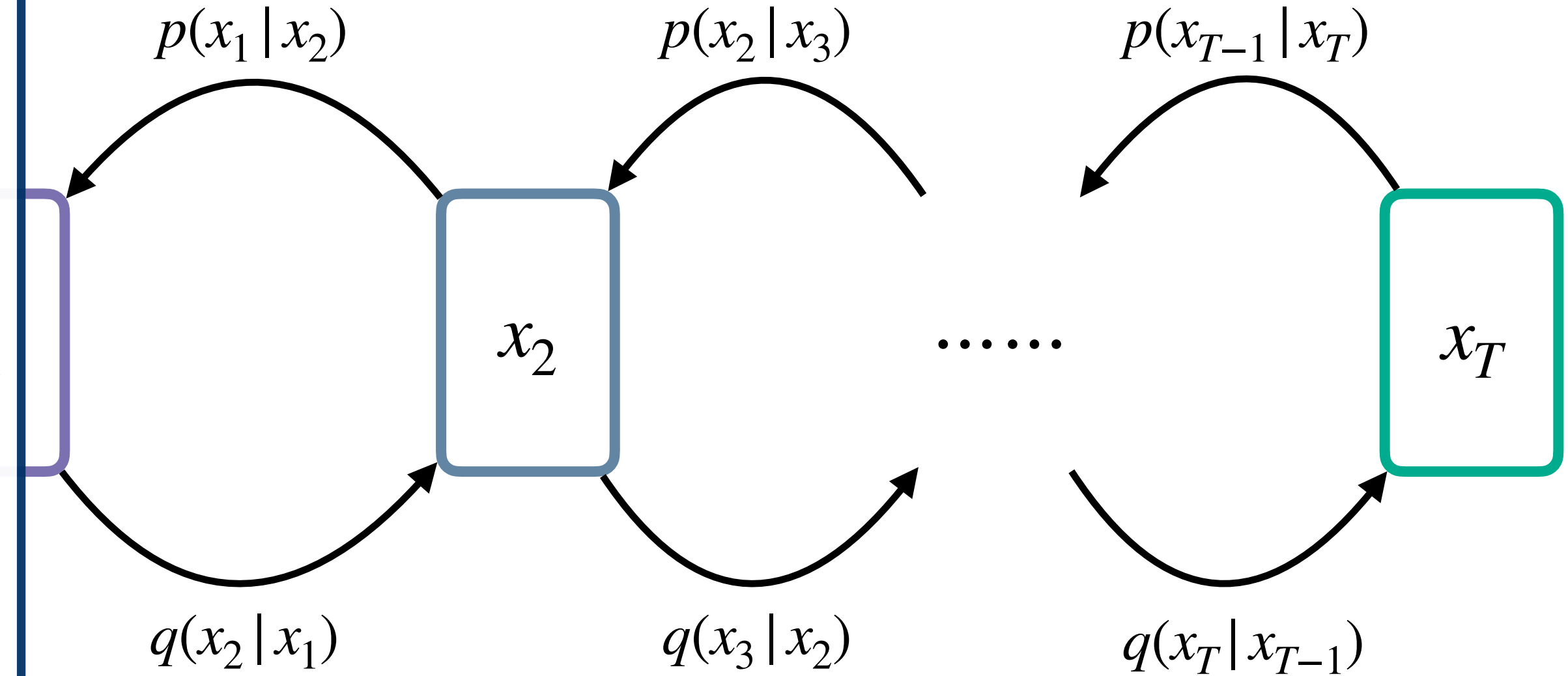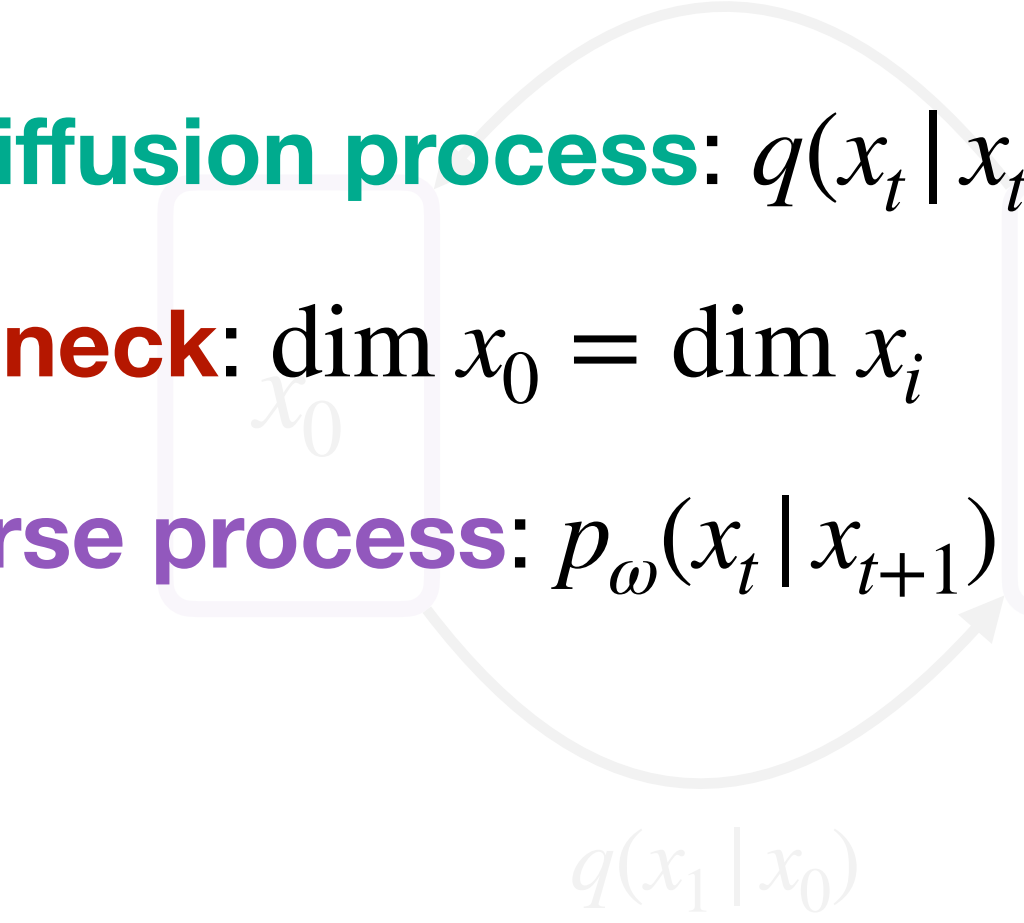$$p(x_0|x_1) \qquad p(x_1|x_2) \qquad p(x_2|x_3) \qquad p(x_{T-1}|x_T)$$

$$x_0 \qquad x_1 \qquad x_2 \qquad \ldots\ldots \qquad x_T$$

$$q(x_1|x_0) \qquad q(x_2|x_1) \qquad q(x_3|x_2) \qquad q(x_T|x_{T-1})$$

$$x_0 \qquad x_1 \qquad x_2 \qquad x_T = z$$

## Hierachical VAEs = DDPMs

- explicit **diffusion process**: $q(x_t | x_{t-1})$

- no **bottleneck**: $\dim x_0 = \dim x_i$

- find **reverse process**: $p_\omega(x_t | x_{t+1})$

$$p(x_1 | x_2) \qquad p(x_2 | x_3) \qquad p(x_{T-1} | x_T)$$

$$x_2 \qquad \ldots\ldots \qquad x_T$$

$$q(x_2 | x_1) \qquad q(x_3 | x_2) \qquad q(x_T | x_{T-1})$$

## Key facts - Diffusion Model

- $\oplus$ State-of-the-art in precision

- $\oplus$ Fast and stable training

- $\ominus$ Slow evaluation

$$x_0 \qquad x_1 \qquad x_2 \qquad x_T = z$$

[2006.11239]

# Part III

## Normalizing Flows

[1505.05770, 1908.09257]

forward $G$

density estimation

$x$

$z$

$p(x)$

generation/sampling

$p_{\text{latent}}(z)$

Inverse $G^{-1} \equiv \overline{G}$

Conservation of probability: $\boxed{p(x)\,\mathrm{d}x = p_{\text{latent}}(z)\,\mathrm{d}z}$ with $\boxed{z = G_\omega(x) \quad x = \overline{G}_\omega(z)}$

Change-of-variables formula: $\boxed{p_\omega(x) = p_{\text{latent}}(z = G_\omega(x)) \cdot \left| \dfrac{\partial G_\omega(x)}{\partial x} \right|}$

forward $G$

density estimation

$x$

$z$

generation/sampling

$p(x)$

$p_{\text{latent}}(z)$

Inverse $G^{-1} \equiv \overline{G}$

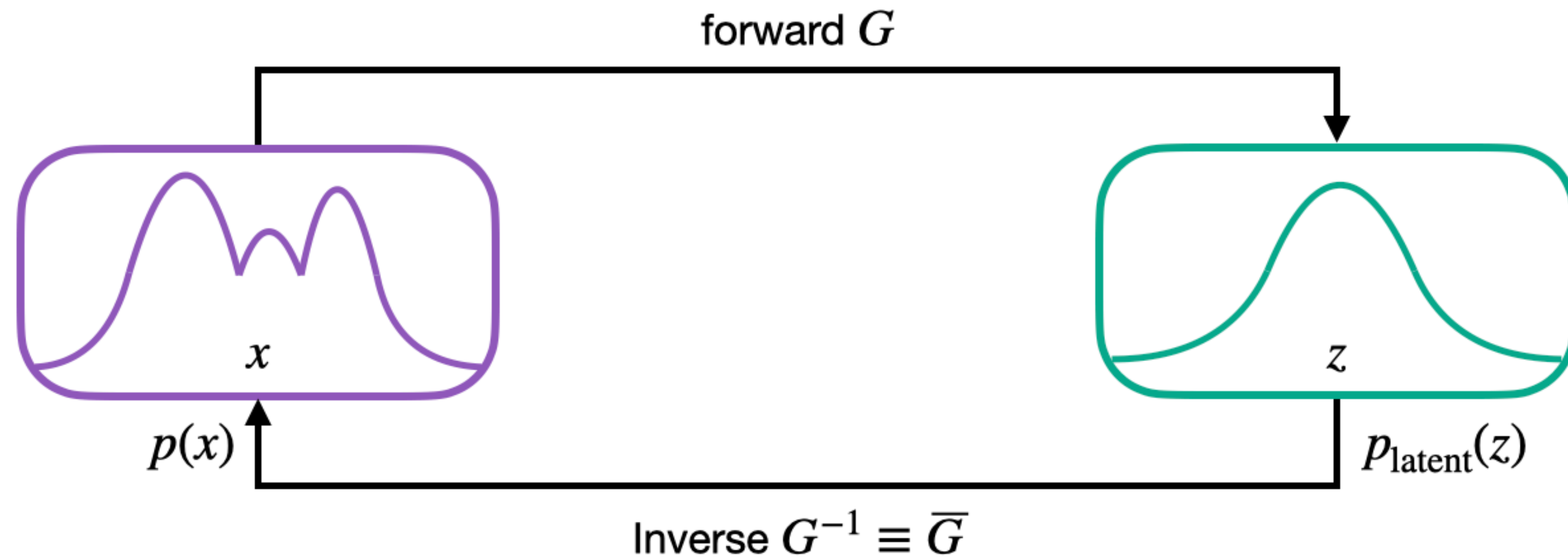Conservation of probability: $\boxed{p(x)\,\mathrm{d}x = p_{\text{latent}}(z)\,\mathrm{d}z}$  with  $\boxed{z = G_\omega(x) \quad x = \overline{G}_\omega(z)}$

Change-of-variables formula: $\boxed{\log p_\omega(x) = \log p_{\text{latent}}(z = G_\omega(x)) + \log \left| \dfrac{\partial G_\omega(x)}{\partial x} \right|}$

[1505.05770, 1908.09257]

# How to train it?

forward $G$

$x$

$p(x)$

$z$

$p_{\text{latent}}(z)$

Inverse $G^{-1} \equiv \overline{G}$

$$\log p_\omega(x) = \log p_{\text{latent}}(z = G_\omega(x)) + \log \left| \frac{\partial G_\omega(x)}{\partial x} \right|$$

$\longrightarrow$ Match $p_\omega(x)$ with $p_{\text{data}}(x)$

Kullback-Leibler divergence:

$$\mathrm{KL}(p_{\text{data}}(x) \,|\, p_\omega(x)) = \int \mathrm{d}x \, p_{\text{data}}(x) \, \log \frac{p_{\text{data}}(x)}{p_\omega(x)}$$

$$= - \int \mathrm{d}x \, p_{\text{data}}(x) \, \log p_\omega(x) + \int \mathrm{d}x \, p_{\text{data}}(x) \, \log p_{\text{data}}(x)$$

No $\omega$ dependence

Negative log-likelihood loss:

$$\mathscr{L}_{\text{NLL}} = - \int \mathrm{d}x \, p_{\text{data}}(x) \, \log p_\omega(x) = \left\langle -\log p_\omega(x) \right\rangle_{x \sim p_{\text{data}}}$$
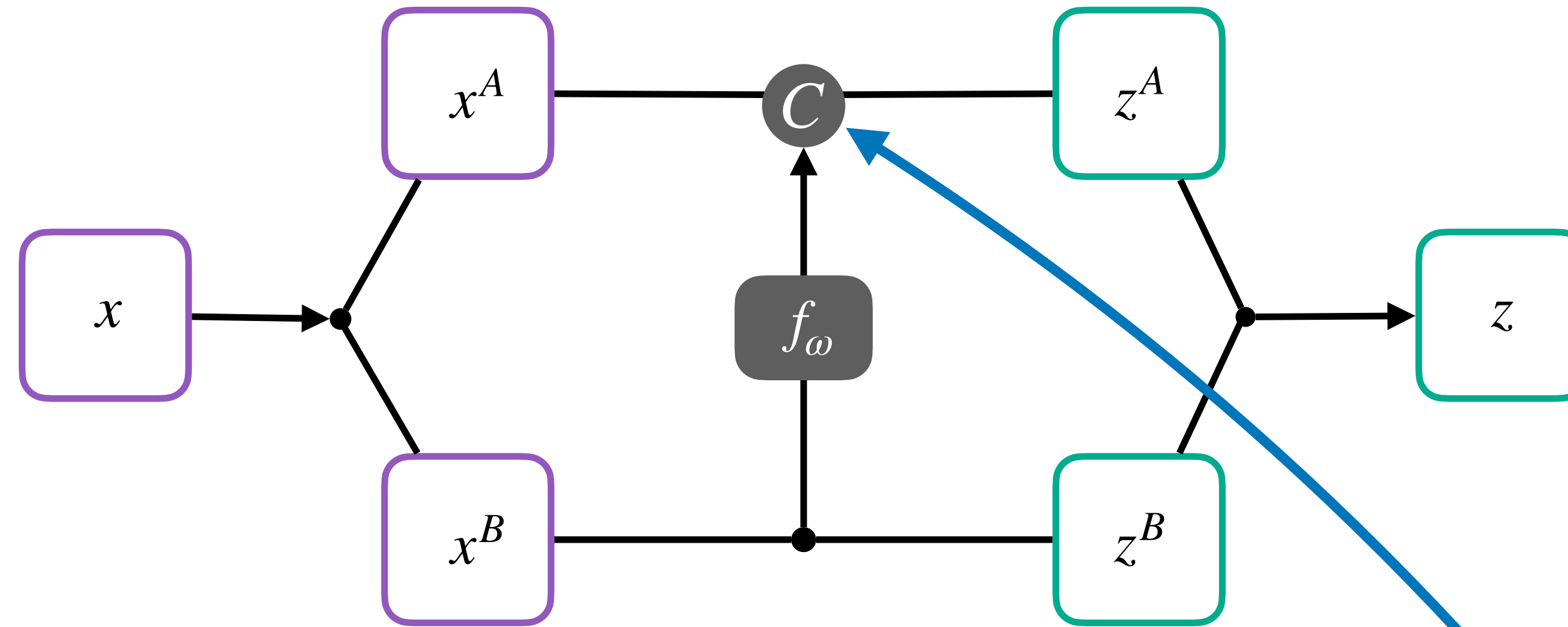
$$\log p_\omega(x) = \log p_{\text{latent}}(z = G_\omega(x)) + \log \left| \frac{\partial G_\omega(x)}{\partial x} \right|$$

⟶ Requires tractable Jacobian!

In general: $g_\omega(x) = \left| \dfrac{\partial G_\omega(x)}{\partial x} \right|$ is $d \times d$ matrix ⟶ Scales with $\mathcal{O}(d^3)$ ☹

Solution: **Autoregressive transformations** $\quad z = \begin{pmatrix} z_1 \\ \vdots \\ z_d \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

$$
\begin{aligned}
z_1 &\equiv z_1(x_1) \\
z_2 &\equiv z_2(x_1, x_2) \\
&\vdots \\
z_d &\equiv z_d(x_1, x_2, \ldots, x_d)
\end{aligned}
$$

⟶ $J_{ij}(x) = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_2}{\partial x_1} & \cdots & \frac{\partial z_d}{\partial x_1} \\ 0 & \frac{\partial z_2}{\partial x_2} & \cdots & \frac{\partial z_d}{\partial x_1} \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & \frac{\partial z_d}{\partial x_d} \end{pmatrix}$ ⟶ $\det J = \prod_i J_{ii} \sim \mathcal{O}(d)$ ☺

Forward pass:
$$z^A = C(x^A; f_\omega(x^B))$$
$$z^B = x^B$$

$$J_{ij}(x) = \begin{pmatrix} \frac{\partial C}{\partial x^A} & \frac{\partial C}{\partial f_\omega} \frac{\partial f_\omega}{\partial x^B} \\ 0 & I_m \end{pmatrix}$$

Inverse pass:
$$x^A = C^{-1}(z^A; f_\omega(z^B))$$
$$x^B = z^B$$

**What is the function $C$?**

Forward pass:
$$z^A = C(x^A; f_\omega(x^B))$$
$$z^B = x^B$$

Inverse pass:
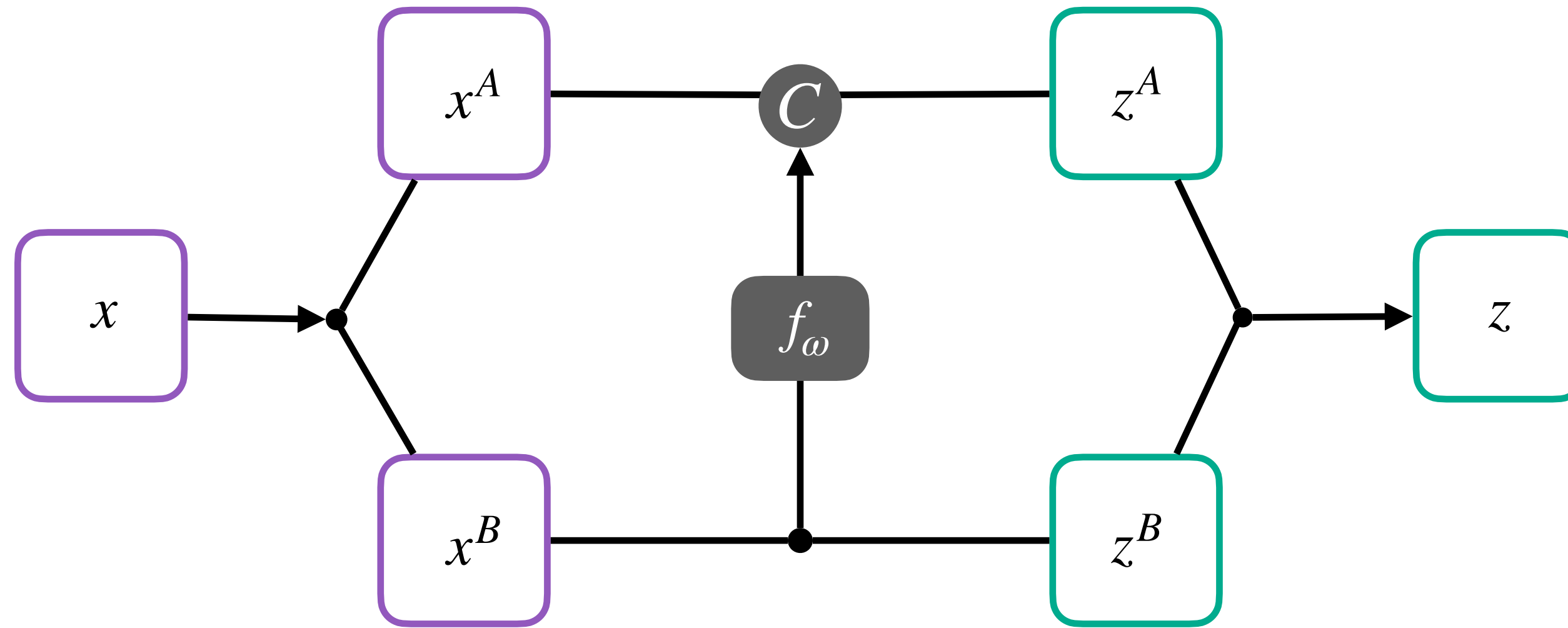$$x^A = C^{-1}(z^A; f_\omega(z^B))$$
$$x^B = z^B$$

$$J_{ij}(x) = \begin{pmatrix} \frac{\partial C}{\partial x^A} & \frac{\partial C}{\partial f_\omega} \frac{\partial f_\omega}{\partial x^B} \\ 0 & I_m \end{pmatrix}$$

Affine
[1605.08803]

$$C^A = \alpha_\omega(x^B) \cdot x^A + \mu_\omega(x^B)$$

**parametrized by NN**

Forward pass:
$$z^A = C(x^A; f_\omega(x^B))$$
$$z^B = x^B$$

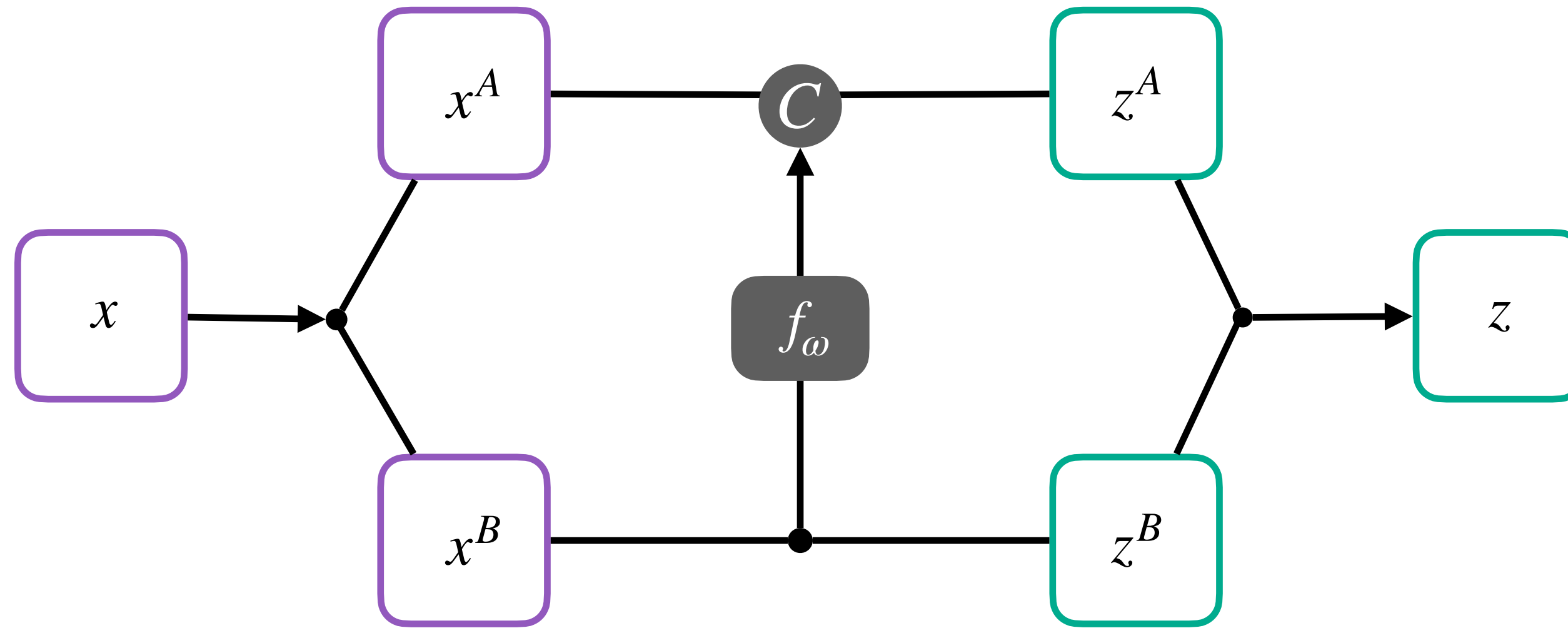Inverse pass:
$$x^A = C^{-1}(z^A; f_\omega(z^B))$$
$$x^B = z^B$$

$$J_{ij}(x) = \begin{pmatrix} \frac{\partial C}{\partial x^A} & \frac{\partial C}{\partial f_\omega}\frac{\partial f_\omega}{\partial x^B} \\ 0 & I_m \end{pmatrix}$$

Affine
[1605.08803]
$$C^A = \alpha_\omega(x^B) \cdot x^A + \mu_\omega(x^B)$$

Quadratic
[1808.03856]
$$C = a_\omega x^2 + b_\omega x + c_\omega$$

Rational quadratic
[1906.04032]
$$C = \frac{a_\omega x^2 + b_\omega x + c_\omega}{d_\omega x^2 + e_\omega x + f_\omega}$$

## Key facts - Normalizing Flows

⊕ Fast training and evaluation

⊕ Tractable and fast likelihoods

⊖ Reduced flexibility and expressivity

Application
dependent

## Key facts - Diffusion Model

⊕ State-of-the-art in precision

⊕ Fast and stable training

⊖ Slow evaluation

$x^A$ — $C$ — $z^A$

$f_\omega$

$z$

$x^B$ — $z^B$

$$J_{ij}(x) = \begin{pmatrix} \frac{\partial C}{\partial x^A} & \frac{\partial C}{\partial f_\omega} \frac{\partial f_\omega}{\partial x^B} \\ 0 & I_m \end{pmatrix}$$
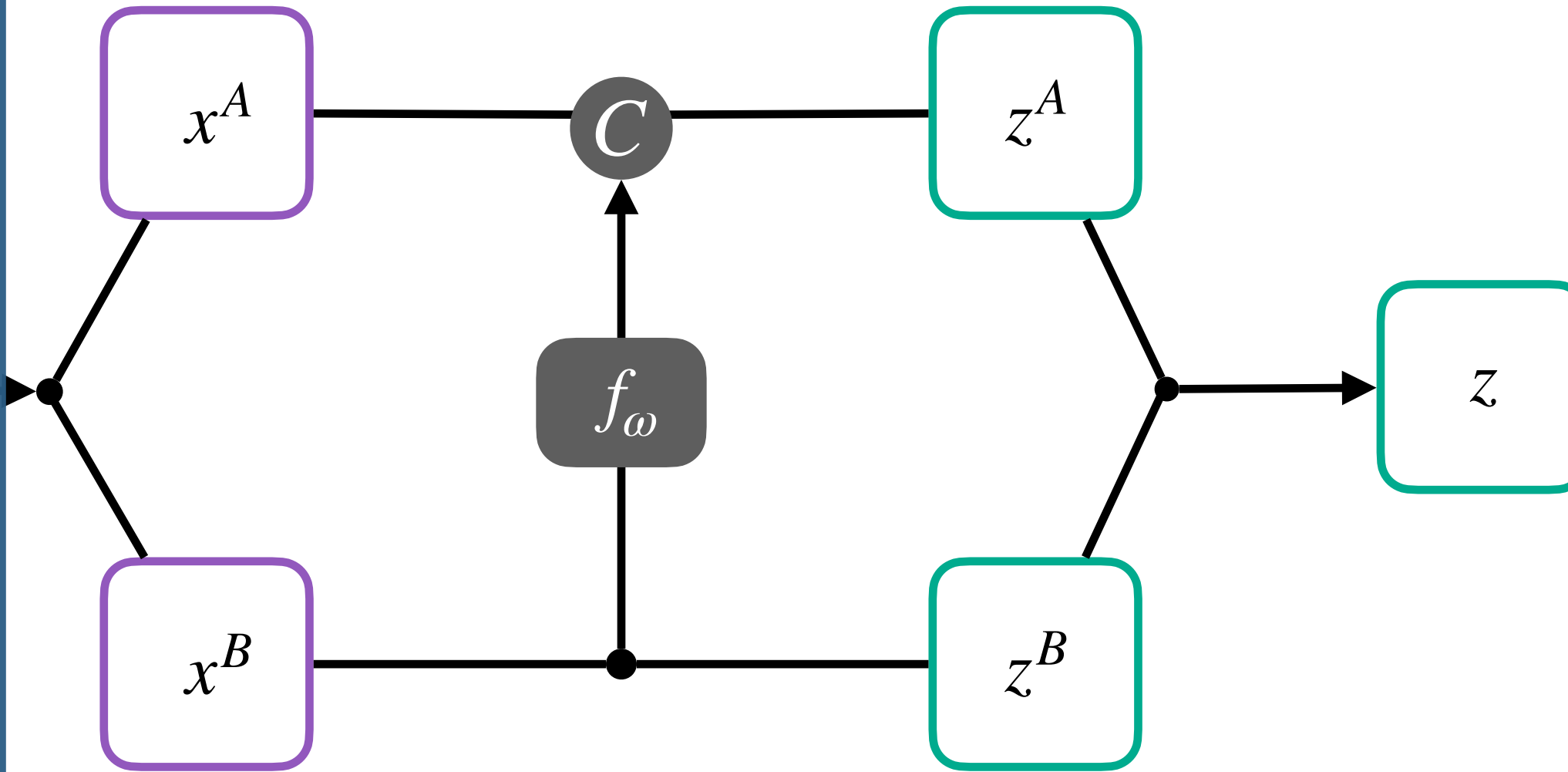
Affine [1605.08803]

$$C^A = \alpha_\omega(x^B) \cdot x^A + \mu_\omega(x^B)$$

Quadratic [1808.03856]

$$C = a_\omega x^2 + b_\omega x + c_\omega$$

Rational quadratic [1906.04032]

$$C = \frac{a_\omega x^2 + b_\omega x + c_\omega}{d_\omega x^2 + e_\omega x + f_\omega}$$

# Questions?

# Machine Learning

## Generative Models

### VAE
- OTUS [2101.08944]
- Jet Simulation [2203.00520]
- ELSA [2305.07696]

### NF
- Precision Generation [2110.13632]
- MEM [2210.00019]
- MADNIS [2212.06172,..]
- CaloFlow I-IV [2106.05285,…]

### Diffusion Models
- PC-JeDi [2303.05376]
- FPCD [2304.01266]
- DDPM & CFM [2305.10475]

### GAN
- EPiC-GAN [2301.08128]
- How to GAN [1907.03764]
- CaloGAN [1712.10321]

### Transformer
- JetGPT [2305.10475]

## Supervised Learning

### Classification
- Point Clouds [2102.05073]
- PELICAN [2211.00454]
- Energy Flow Networks [1810.05165]
- Landscape of Top tagger [1902.09914]
- Bayesian Tagger [1904.10004]

### Regression
- MadMiner [1907.10621,…]
- NNPDF [2109.02653]
- Matrix Elements [2206.14831]
- Symbolic regression [2109.10414]

## Unsupervised Learning

### Clustering
- Jet Clustering [2008.06064]
- 3D Pixel Clustering [2007.03083]

### Anomaly Detection
- Normalized AE [2206.14225]
- CATHODE [2109.00546]
- CWoLA Hunting [1902.02634]
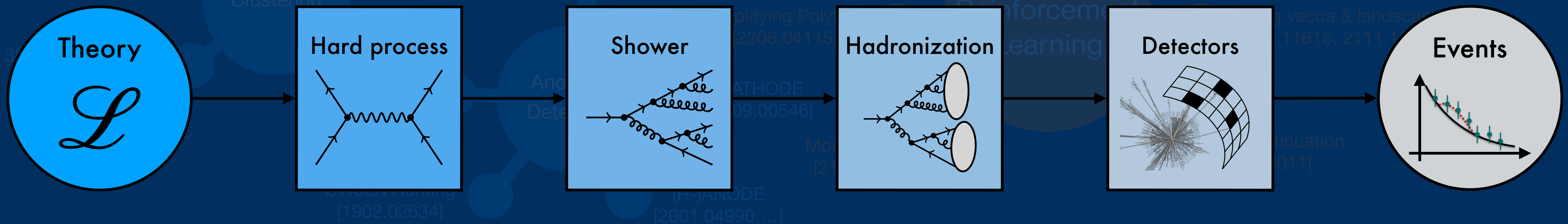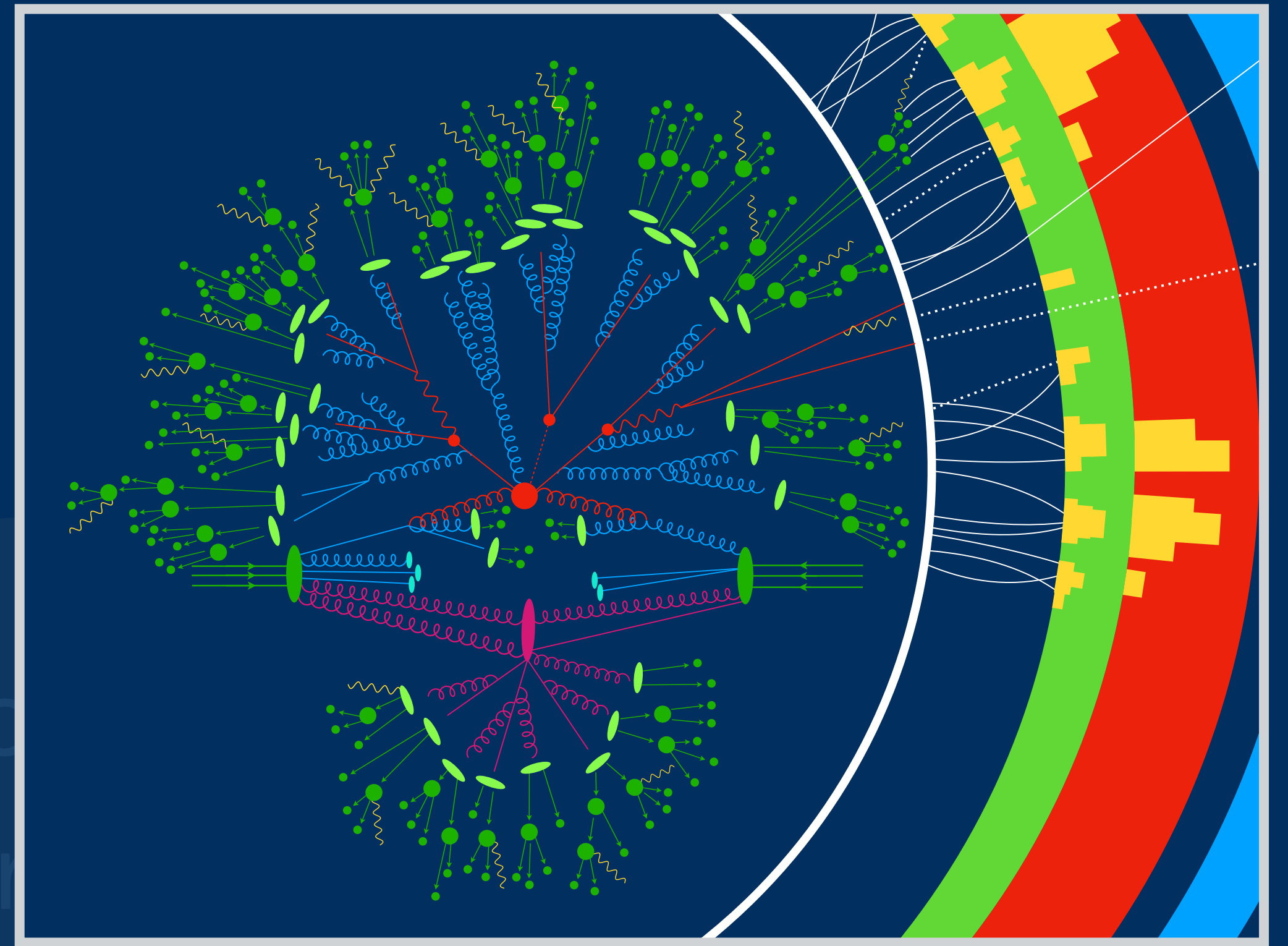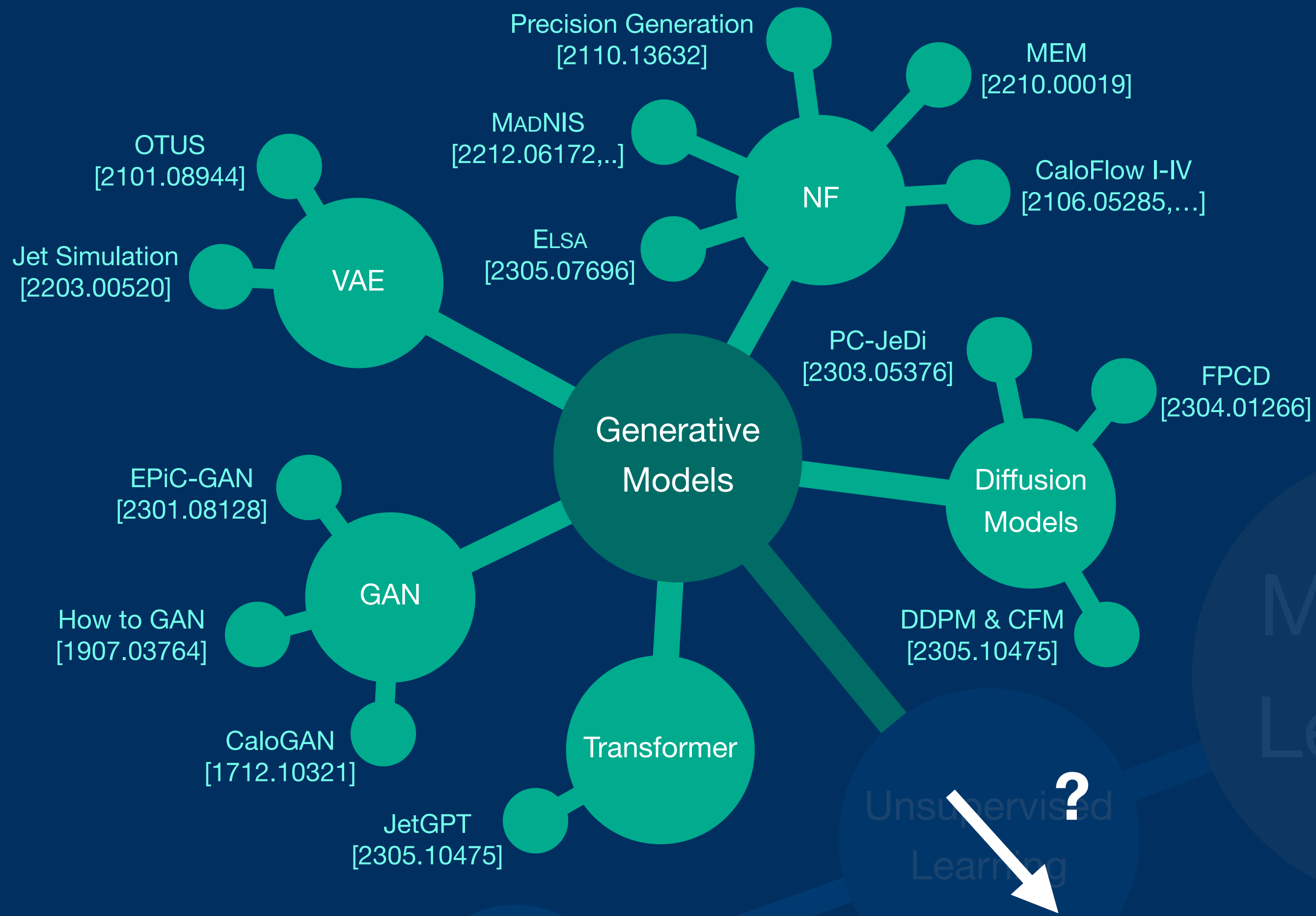- (R-)ANODE [2001.04990,…]

## Reinforcement Learning
- Flavor structure [2304.14176]
- Simplifying Polylogs [2206.04115]
- String vacua & landscape [1903.11616, 2111.11466]
- Model Building [2103.04759]
- Analytic continuation [2112.13011]

## Generative Models

**VAE**
- OTUS [2101.08944]
- Jet Simulation [2203.00520]

**NF**
- Precision Generation [2110.13632]
- MaDNIS [2212.06172,..]
- ELSA [2305.07696]
- MEM [2210.00019]
- CaloFlow I-IV [2106.05285,...]

**GAN**
- EPiC-GAN [2301.08128]
- How to GAN [1907.03764]
- CaloGAN [1712.10321]

**Transformer**
- JetGPT [2305.10475]

**Diffusion Models**
- PC-JeDi [2303.05376]
- FPCD [2304.01266]
- DDPM & CFM [2305.10475]

?

Theory $\mathcal{L}$ → Hard process → Shower → Hadronization → Detectors → Events

**Phase-space integration**
**Event generation**

Theory $\mathcal{L}$ → Hard process → Shower → Hadronization → Detectors → Events

**Importance sampling**
BDT [1707.00028], NN [1810.11509, 2009.07819]
NF [2001.05486, 2001.05478, 2001.10028, 2005.12719,
2112.09145, 2212.06172, 2311.01548]
Chili [2302.10449]

**Surrogate regression**
Full weight [2109.11964],
Matrix element [1912.11055, 2002.07516,
2006.16273, 2106.09474, 2107.06625, 2109.11964,
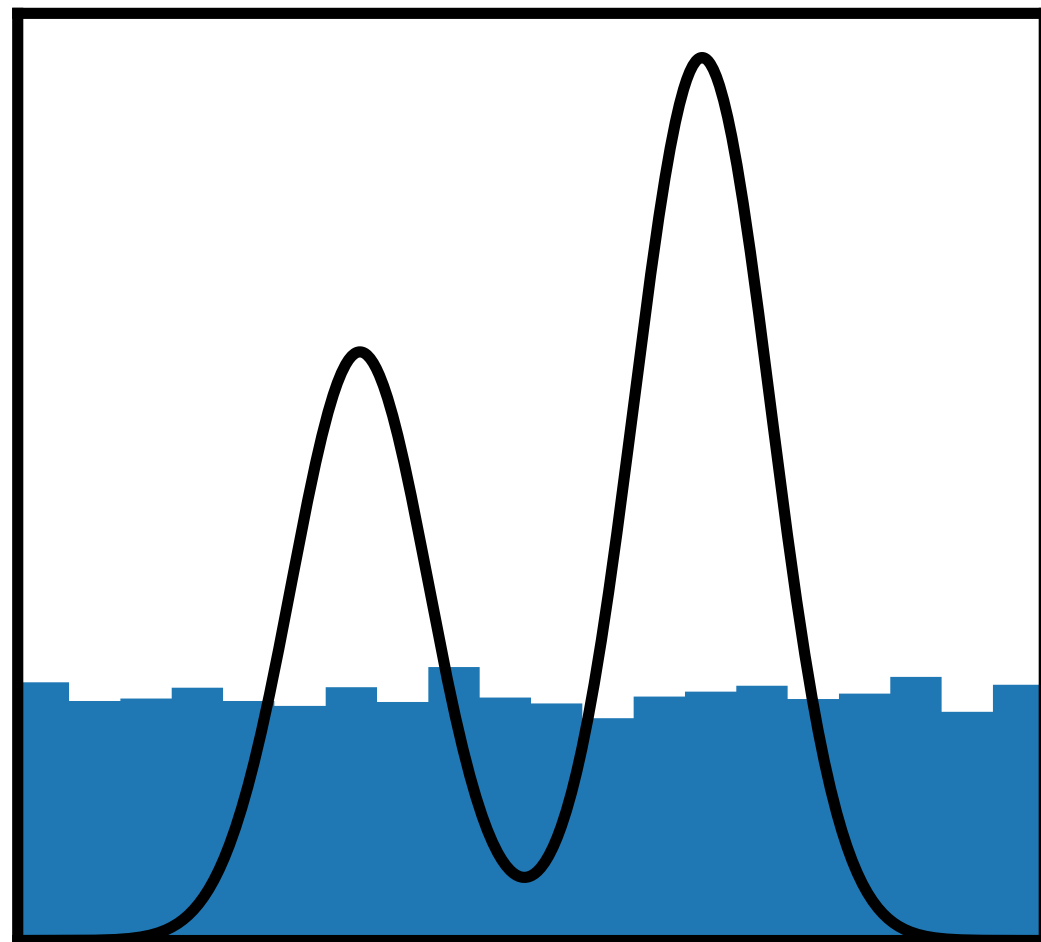2206.14831, 2301.13562, 2302.04005, 2306.07726]

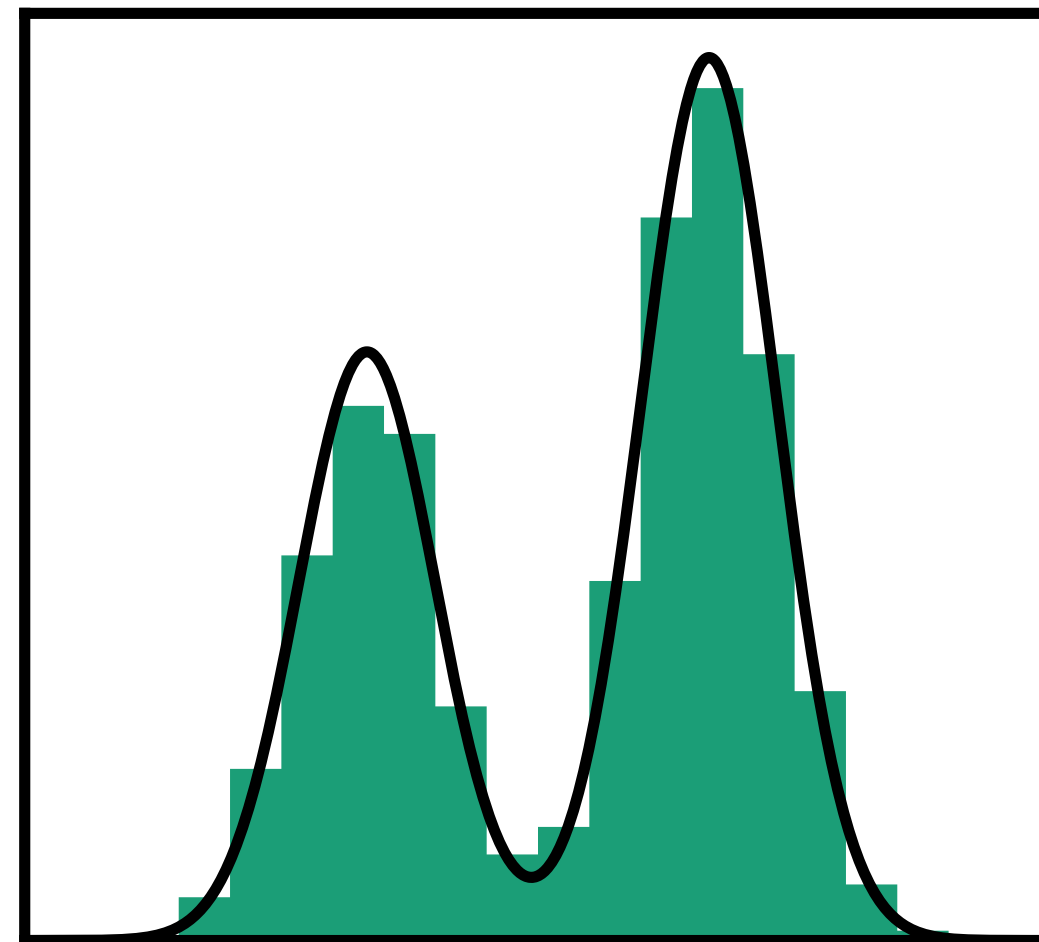# Example I

**Neural importance sampling with MadNIS**

Heimel, Huetsch, Maltoni, Mattelaer, Plehn, RW [2311.01548]
Heimel, RW, Butter, Isaacson, Krause, Maltoni, Mattelaer, Plehn [2212.06172]

Calculate (differential) cross sections

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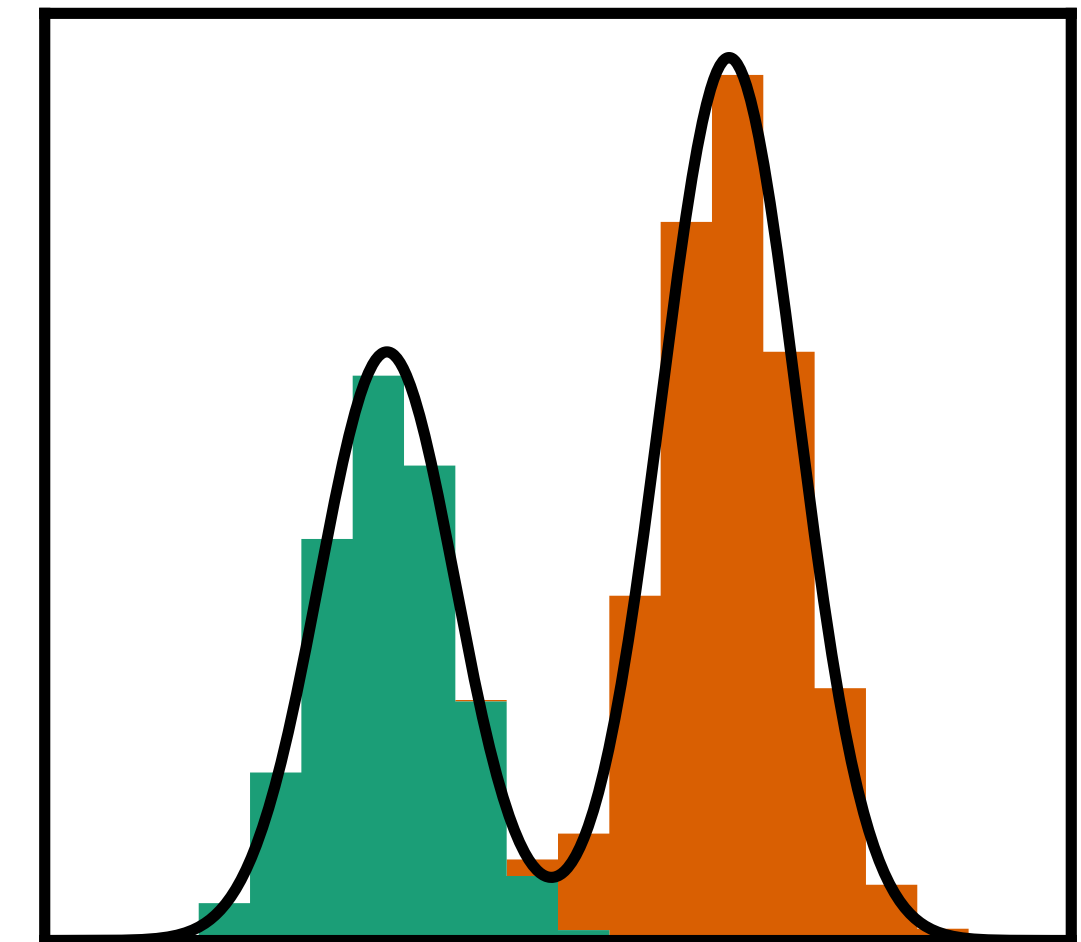


Flat sampling:
inefficient

$$I = \left\langle f(x)\right\rangle_{x\sim\mathrm{unif}}$$

Importance sampling:
find $p$ close to $f$

$$I = \left\langle\frac{f(x)}{p(x)}\right\rangle_{x\sim p(x)}$$

Multi-channel:
one map for each channel

$$I = \sum_i\left\langle\alpha_i(x)\frac{f(x)}{p_i(x)}\right\rangle_{x\sim p_i(x)}$$

Calculate (differential) cross sections

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

**Sum over channels**

MadGraph: build channels
from Feynman diagrams

**Integrand**

MadGraph: $\mathrm{d}\sigma/\mathrm{d}x$

$$I = \sum_i\left\langle\alpha_i(x)\,\frac{f(x)}{p_i(x)}\right\rangle_{x\sim p_i(x)}$$

**Channel weights**

MadGraph: $\alpha_i^{\text{MG}}(x)\sim|M_i|^2$

**Channel mappings**

MadGraph: use amplitude structure, …
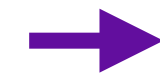Analytic mappings + refine with **VEGAS**

Factorize probability

$$p(x) = p(x_1) \cdots p(x_n)$$

Fit bins with equal probability
and varying width



⊕ Computationally cheap

⊖ High-dim and rich peaking functions
→ **slow convergence**

⊖ Peaks not aligned with grid axes
→ **phantom peaks**



[G. P. Lepage, 1978]

Calculate (differential) cross sections

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

**Sum over channels**

MadGraph: build channels
from Feynman diagrams

**Integrand**

MadGraph: $\mathrm{d}\sigma/\mathrm{d}x$

$$I = \sum_i \left\langle \alpha_i(x)\,\frac{f(x)}{p_i(x)} \right\rangle_{x \sim p_i(x)}$$

**Channel weights**

MadGraph: $\alpha_i^{\mathrm{MG}}(x) \sim |M_i|^2$

**Channel mappings**

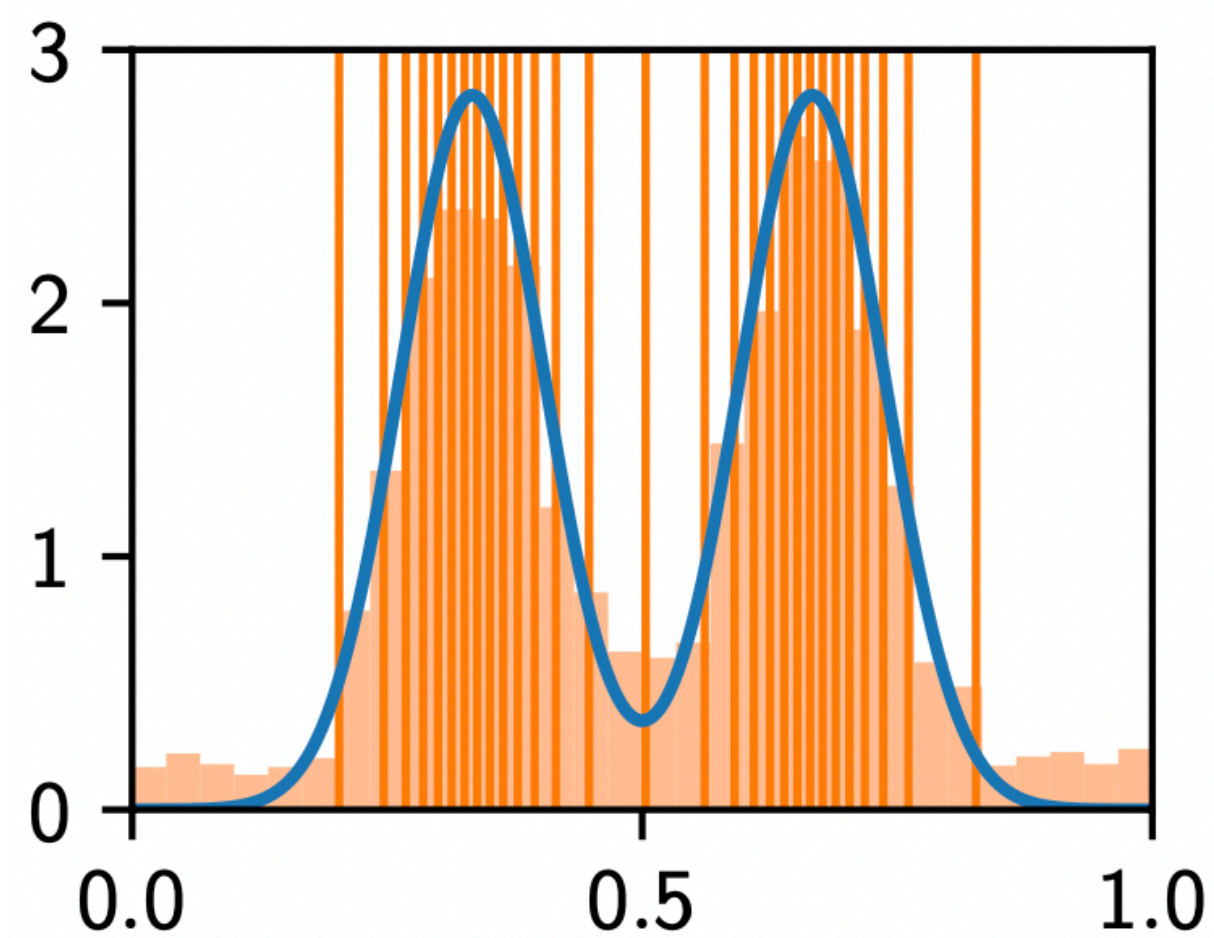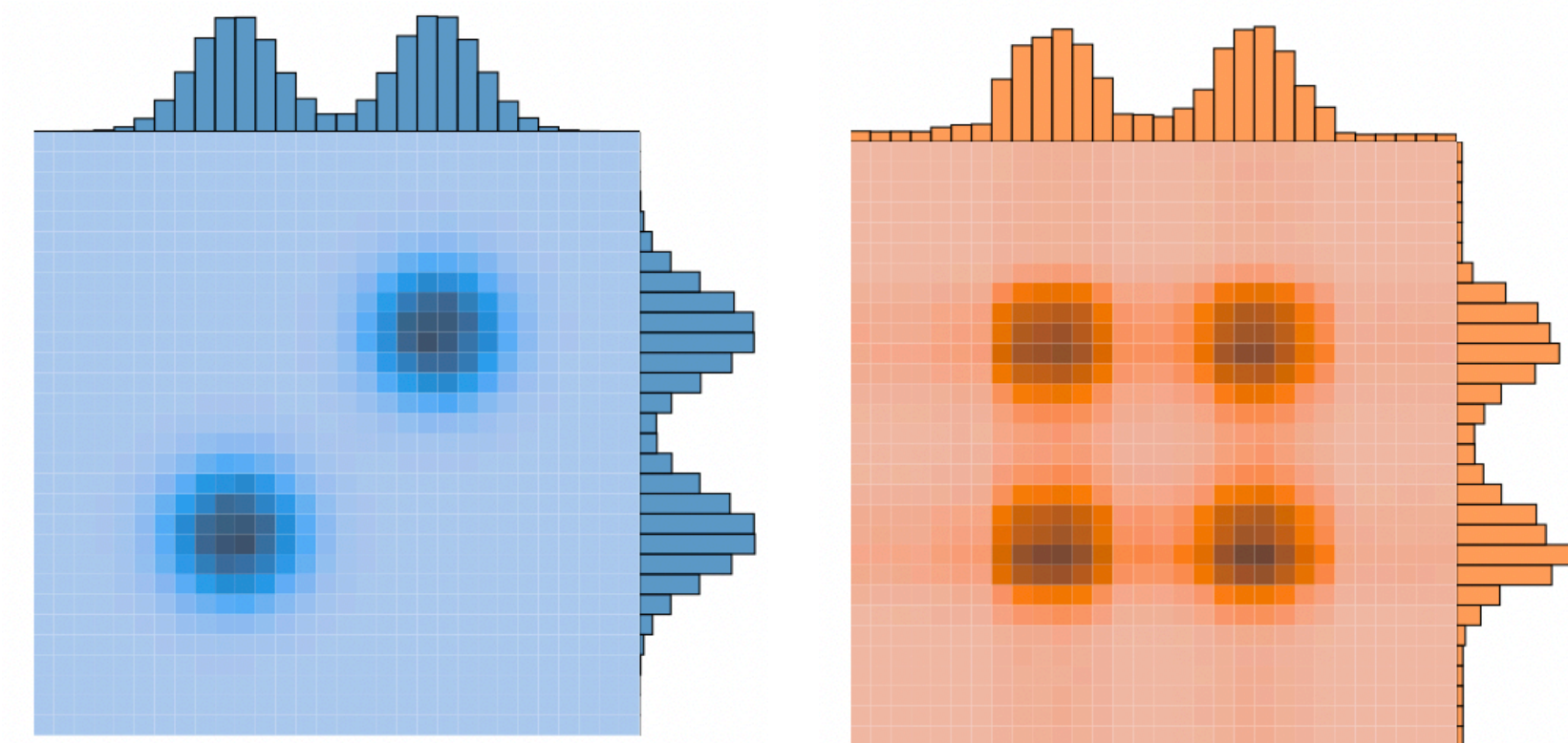MadGraph: use amplitude structure, …
Analytic mappings + refine with VEGAS

Calculate (differential) cross sections

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

**Sum over channels**

MadGraph: build channels
from Feynman diagrams

**Integrand**

MadGraph: $\mathrm{d}\sigma/\mathrm{d}x$

$$I = \sum_i \left\langle\, \alpha_i(x)\, \frac{f(x)}{p_i^\omega(x)}\, \right\rangle_{x \sim p_i^\omega(x)}$$

**Channel weights**

MadGraph: $\alpha_i^{\mathrm{MG}}(x) \sim |M_i|^2$

**Learned channel mappings**

MadGraph: use amplitude structure, …
Analytic mappings + ~~refine with VEGAS~~

refine with **NF**

Calculate (differential) cross sections

$$\mathrm{d}\sigma = \frac{1}{\text{flux}} \mathrm{d}x_a \mathrm{d}x_b \, f(x_a)f(x_b) \, \mathrm{d}\Phi_n \left\langle \, |M_{\lambda,c,\ldots}(p_a, p_b \mid p_1, \ldots, p_n)|^2 \, \right\rangle$$

**Sum over channels**

MadGraph: build channels
from Feynman diagrams

**Integrand**

MadGraph: $\mathrm{d}\sigma/\mathrm{d}x$

$$I = \sum_i \left\langle \alpha_i^{\xi}(x) \frac{f(x)}{p_i^{\omega}(x)} \right\rangle_{x \sim p_i^{\omega}(x)}$$

**Learned Channel weights**

MadGraph: $\alpha_i^{\mathrm{MG}}(x) \sim |M_i|^2$

$$\alpha_i(x) \to \alpha_i^{\xi}(x) = \alpha_i^{\mathrm{MG}}(x) \cdot K_i^{\xi}(x)$$
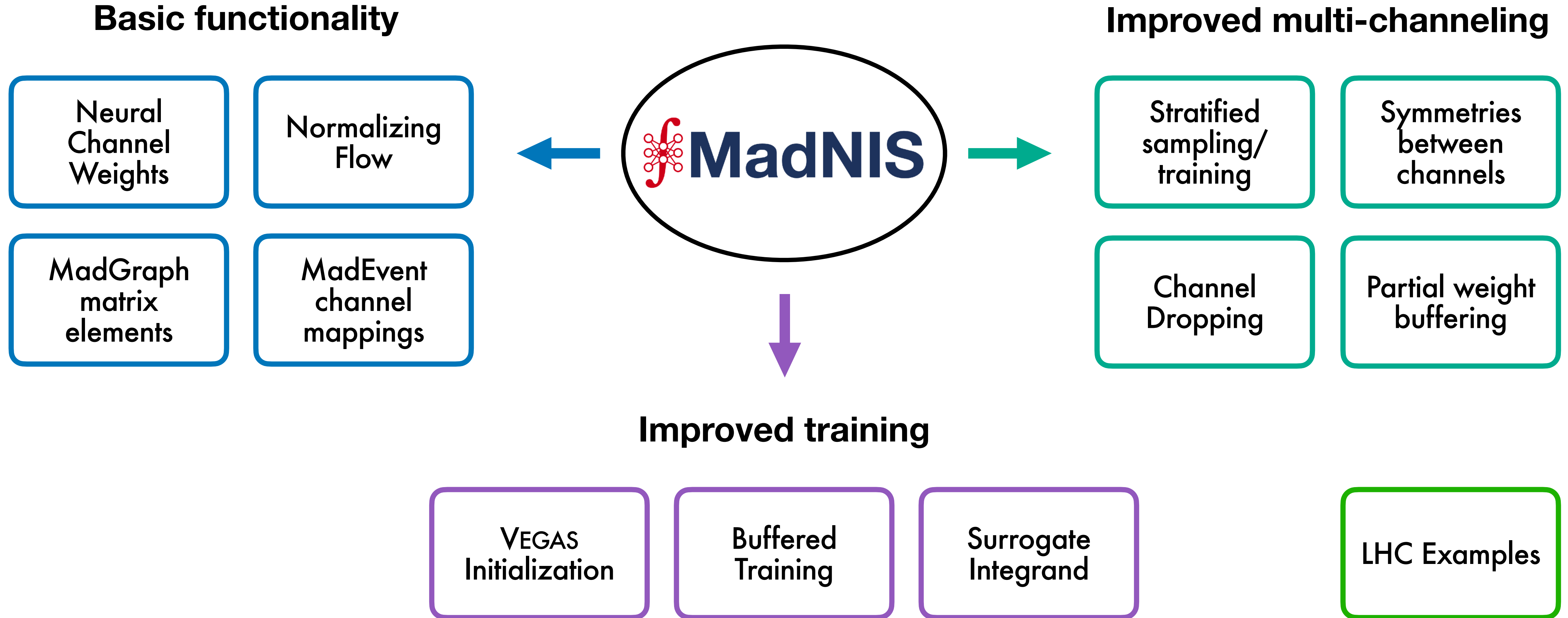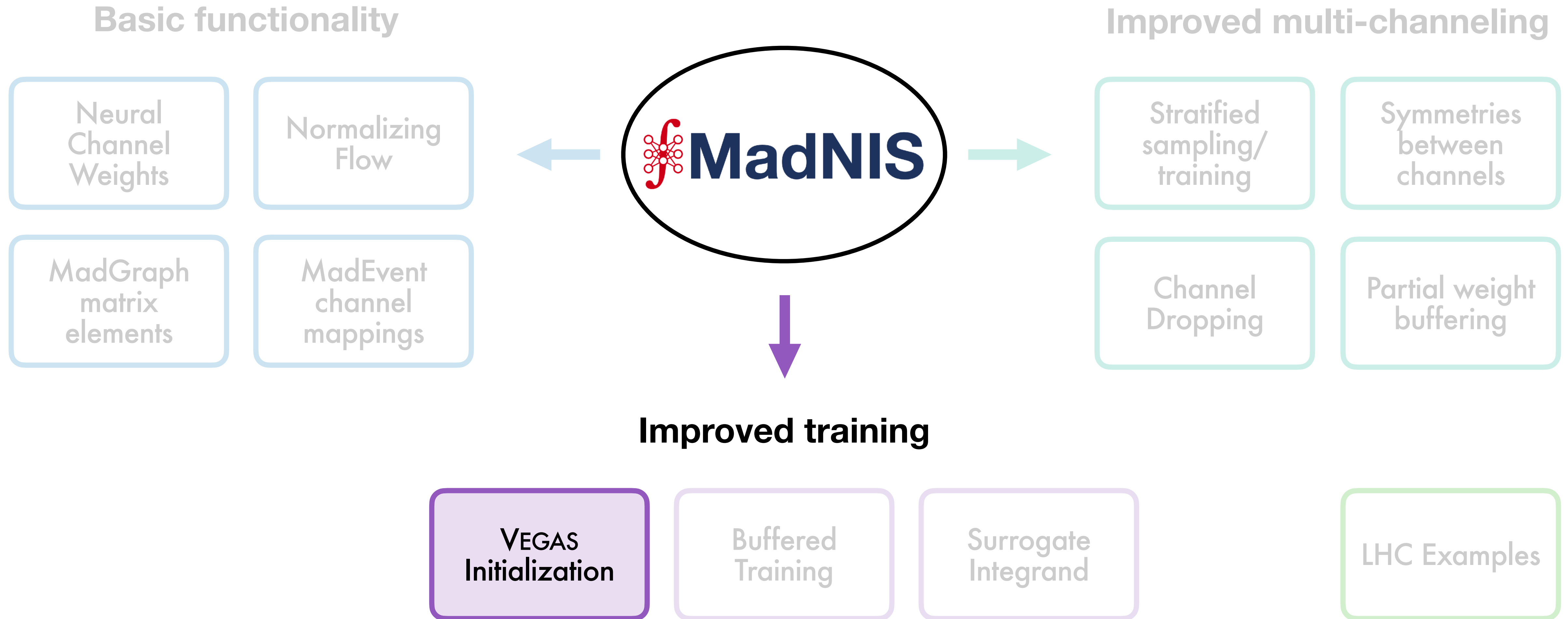
parametrize with **NN**

**Learned channel mappings**

MadGraph: use amplitude structure, …
Analytic mappings + ~~refine with V~~ᴇɢᴀs~~

refine with **NF**

**Basic functionality**

Neural Channel Weights

Normalizing Flow

MadGraph matrix elements

MadEvent channel mappings

$\int$MadNIS

**Improved multi-channeling**

Stratified sampling/ training

Symmetries between channels

Channel Dropping

Partial weight buffering

**Improved training**

VEGAS Initialization

Buffered Training

Surrogate Integrand

LHC Examples

**Basic functionality**

**Improved multi-channeling**



| Neural Channel Weights | Normalizing Flow |
|---|---|
| MadGraph matrix elements | MadEvent channel mappings |

**MadNIS**

| Stratified sampling/ training | Symmetries between channels |
|---|---|
| Channel Dropping | Partial weight buffering |

**Improved training**

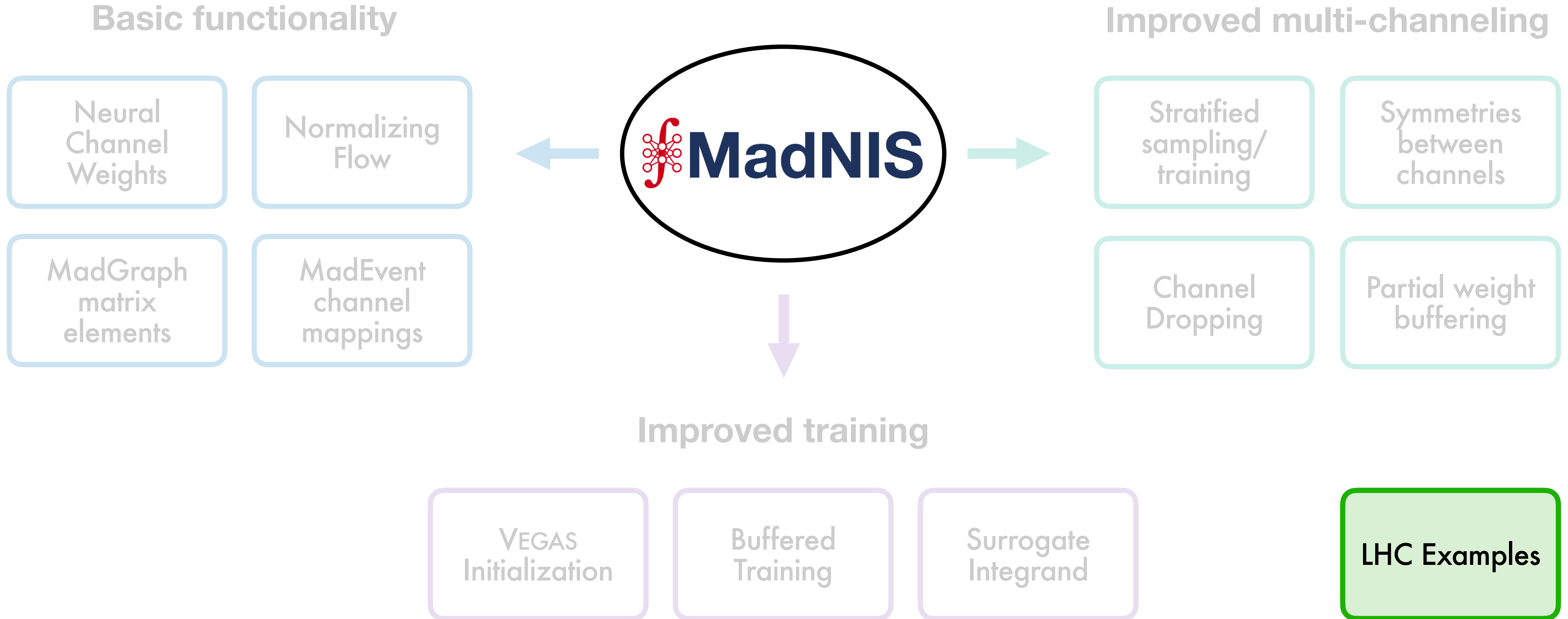| VEGAS Initialization | Buffered Training | Surrogate Integrand |
|---|---|---|

LHC Examples

|  | VEGAS | Flow |
|---|---|---|
| Training | Fast | Slow |
| Correlations | No | Yes |

Combine advantages:

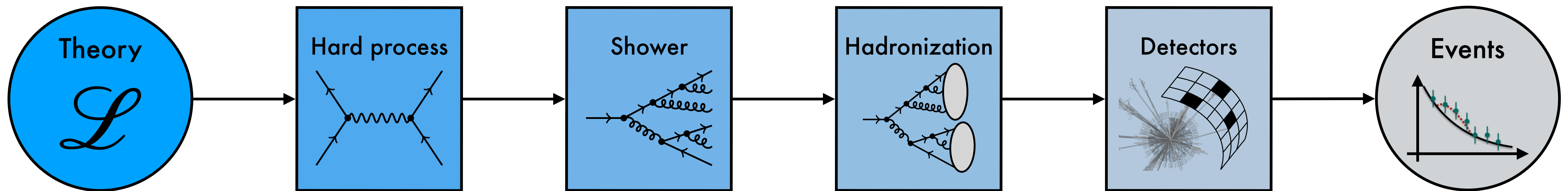Pre-trained VEGAS grid as starting point for flow training

VEGAS grid → Bin reduction → Initialization

Concat

$y_1$ — $y_2$

Subnet → RQ-Spline

RQ-Spline ← Subnet

$z_1$ ↔ $z_2$

Split

**Basic functionality**

| Neural Channel Weights | Normalizing Flow |
| --- | --- |

| MadGraph matrix elements | MadEvent channel mappings |

**Improved multi-channeling**

| Stratified sampling/ training | Symmetries between channels |
| --- | --- |

| Channel Dropping | Partial weight buffering |

**Improved training**

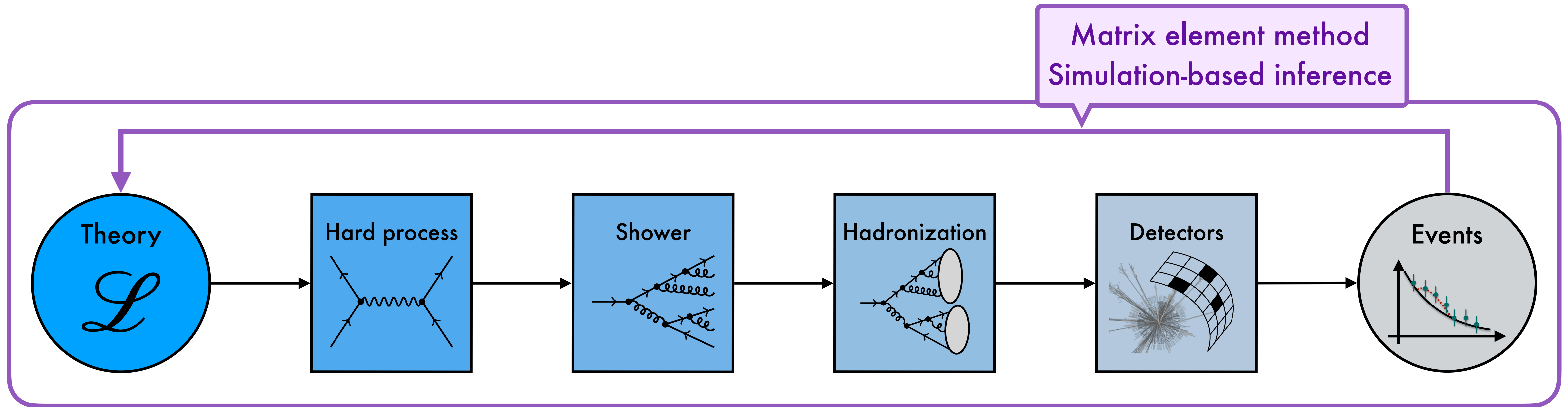| VEGAS Initialization | Buffered Training | Surrogate Integrand |
| --- | --- | --- |

LHC Examples

1. excellent results with all improvements

2. Larger improvements for processes with large interference terms

Matrix element method
Simulation-based inference

Theory $\mathcal{L}$ → Hard process → Shower → Hadronization → Detectors → Events

Simulation-based inference
[1506.02169, 1601.07913, 1805.00013, 1805.00020,
1805.12244, 1907.10621, 2101.07263, 2210.01680,
2305.10500, 2308.05704,…]

Matrix element method (MEM)
[hep-ex/9808029, hep-ex/0406031, hep-ex/0605118,
1003.1316, 1007.3300, 1010.2263, 1211.3011, 1304.6414,
1502.02485, 1511.05980, 1511.06170, 1512.03429,
1606.03107, 1710.10699, 1712.03266, 1805.08555,
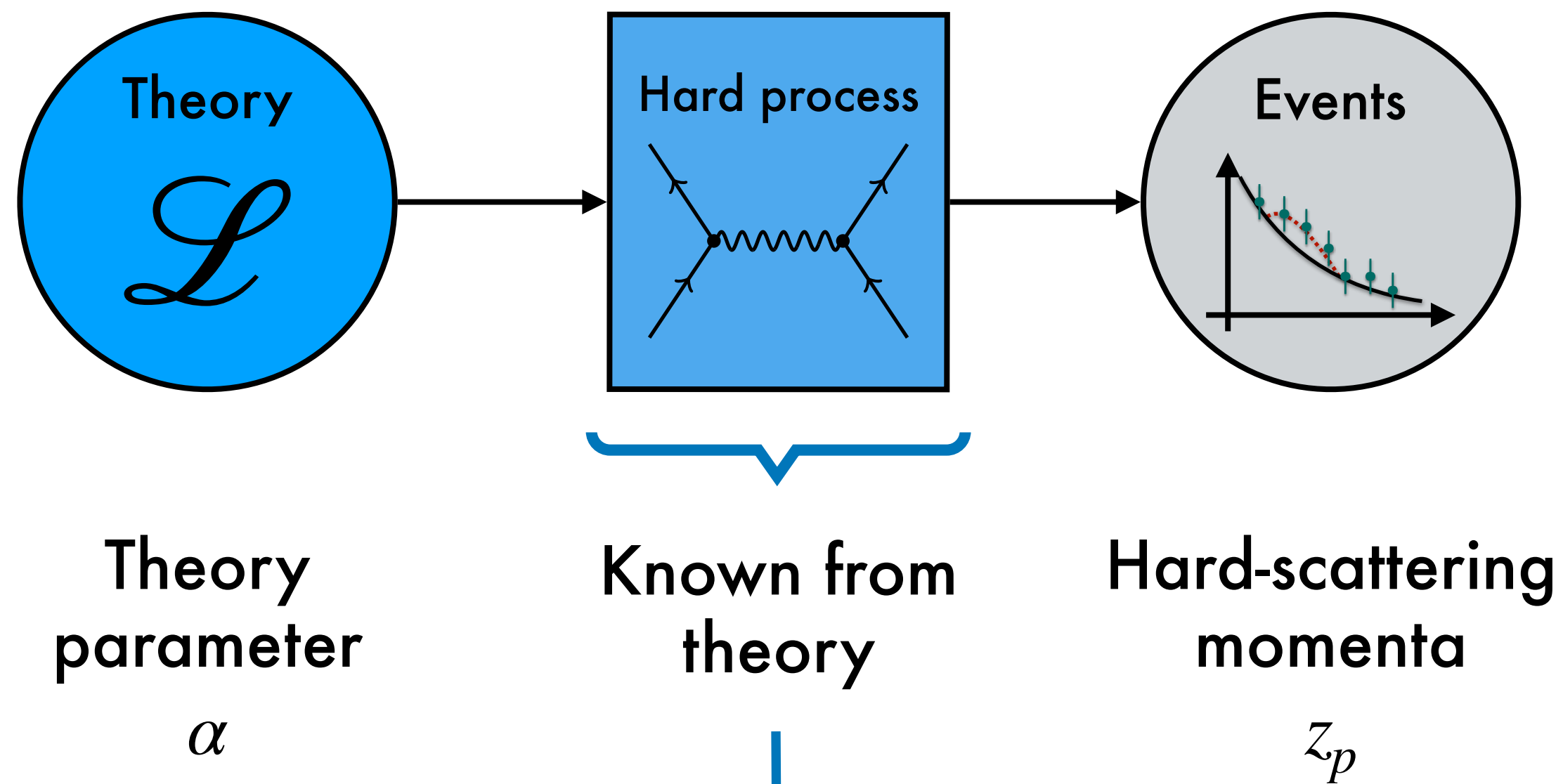2008.10949, 2210.00019, 2310.07752,…..]

# Example II

## Matrix Element Method

Heimel, Huetsch, RW, Plehn, Butter [2310.07752]

Butter, Heimel, Martini, Peitzsch, Plehn [2210.00019]

Theory

$\mathcal{L}$

Hard process

Events

Theory
parameter
$\alpha$

Known from
theory

Hard-scattering
momenta
$z_p$

Likelihood from differential cross section

$$p(z_p \mid \alpha) = \frac{1}{\sigma(\alpha)} \frac{\mathrm{d}\sigma(\alpha)}{\mathrm{d}z_p}$$
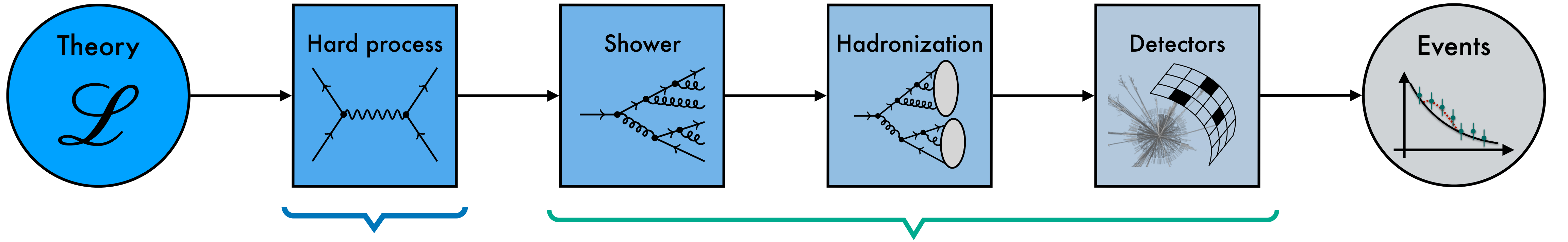
**Classical analysis**

⊖ hand-crafted observables

⊖ binned data

→ not all information is used ☹

**Matrix Element Method (MEM)**

⊕ based on first principles

⊕ estimates uncertainties reliably

⊕ optimal use of information

→ perfect for processes with few events ☺

Theory
parameter
$\alpha$

Known from
theory

Likelihood intractable

Reconstructed
momenta
$x$

**MEM master formula:** $$p(x \mid \alpha) = \int \mathrm{d}z_p \; p(z_p \mid \alpha) \; p(x \mid z_p) \; \epsilon(z_p)$$

Integrate out hard momenta

Acceptance function

$$p(x \mid \alpha) = \int \mathrm{d}z_p \quad p(z_p \mid \alpha) \quad p(x \mid z_p) \quad \epsilon(z_p)$$

**Efficient MC integration**

importance sampling:
**Normalizing Flow**

$$z_p \sim p(z_p \mid x, \alpha)$$

**Theory knowledge**

differential
cross-section

$$\frac{1}{\sigma(\alpha)} \frac{\mathrm{d}\sigma(\alpha)}{\mathrm{d}z_p}$$

**Transfer function**

density estimation:
**Normalizing Flow**

solve combinatorics:
**Transformer**

**Transfermer**

**Acceptance function**

learn with simple
**Classifier network**

**Single Higgs production** with anomalous
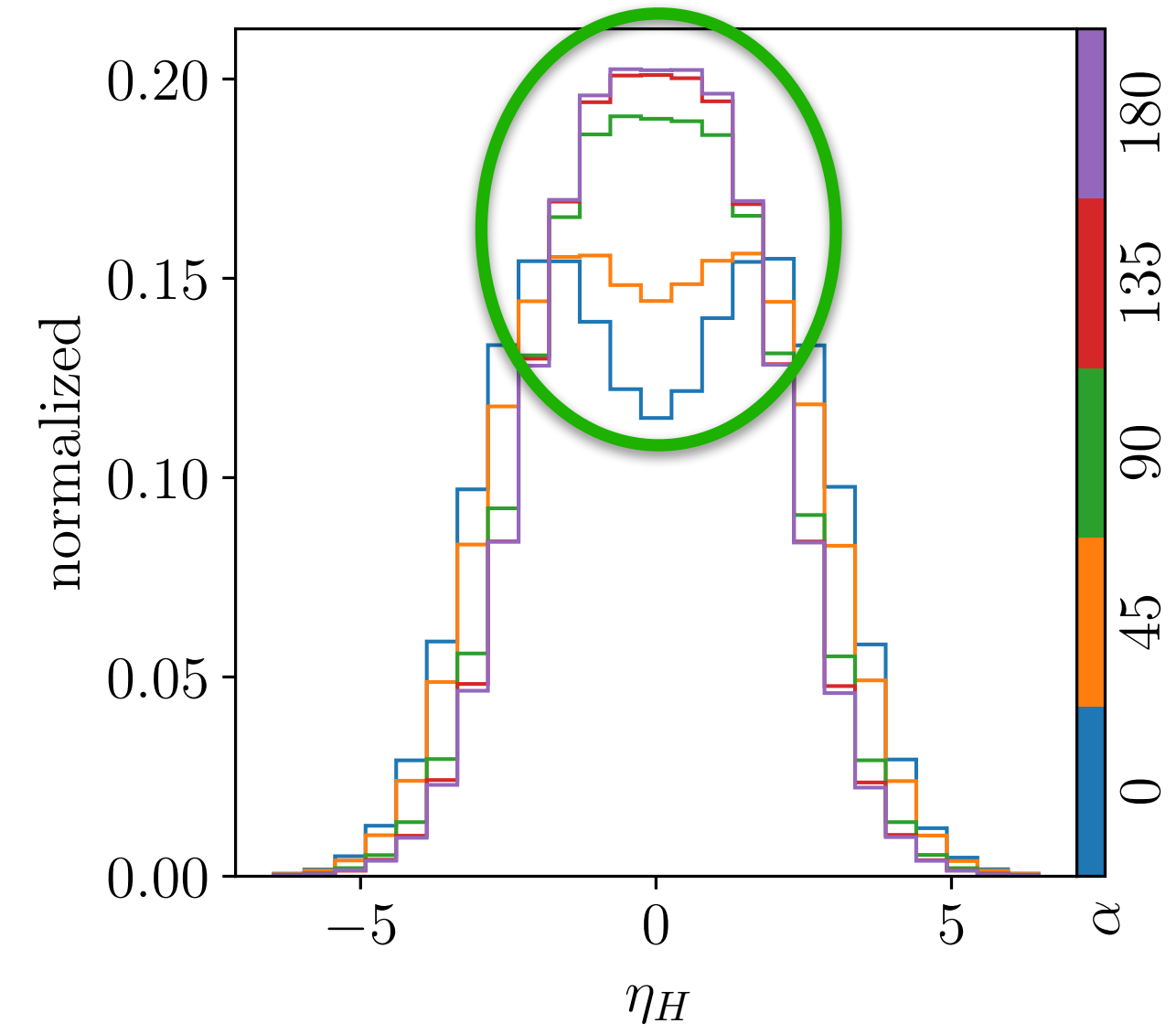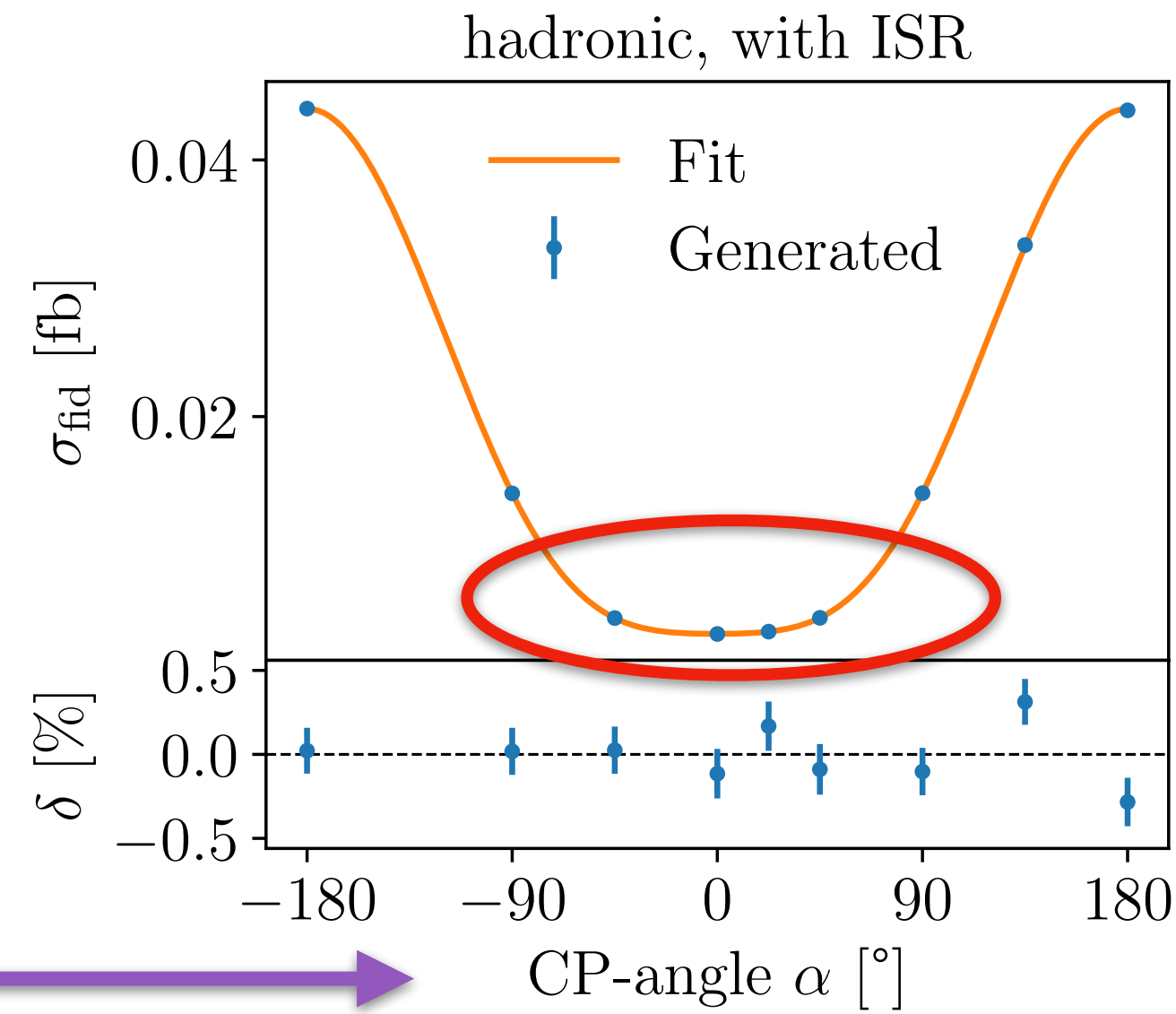non-CP-conserving Higgs coupling

↓

**Hadronic decay of top + ISR**

$tHq \rightarrow (bjj) \, (\gamma\gamma) \, j + QCD \, jets$

↓

$$\mathscr{L}_{t\bar{t}H} = -\frac{y_t}{\sqrt{2}}\left[\cos\alpha \, \bar{t}t + \frac{2}{3}i\sin\alpha \, \bar{t}\gamma_5 t\right]H$$
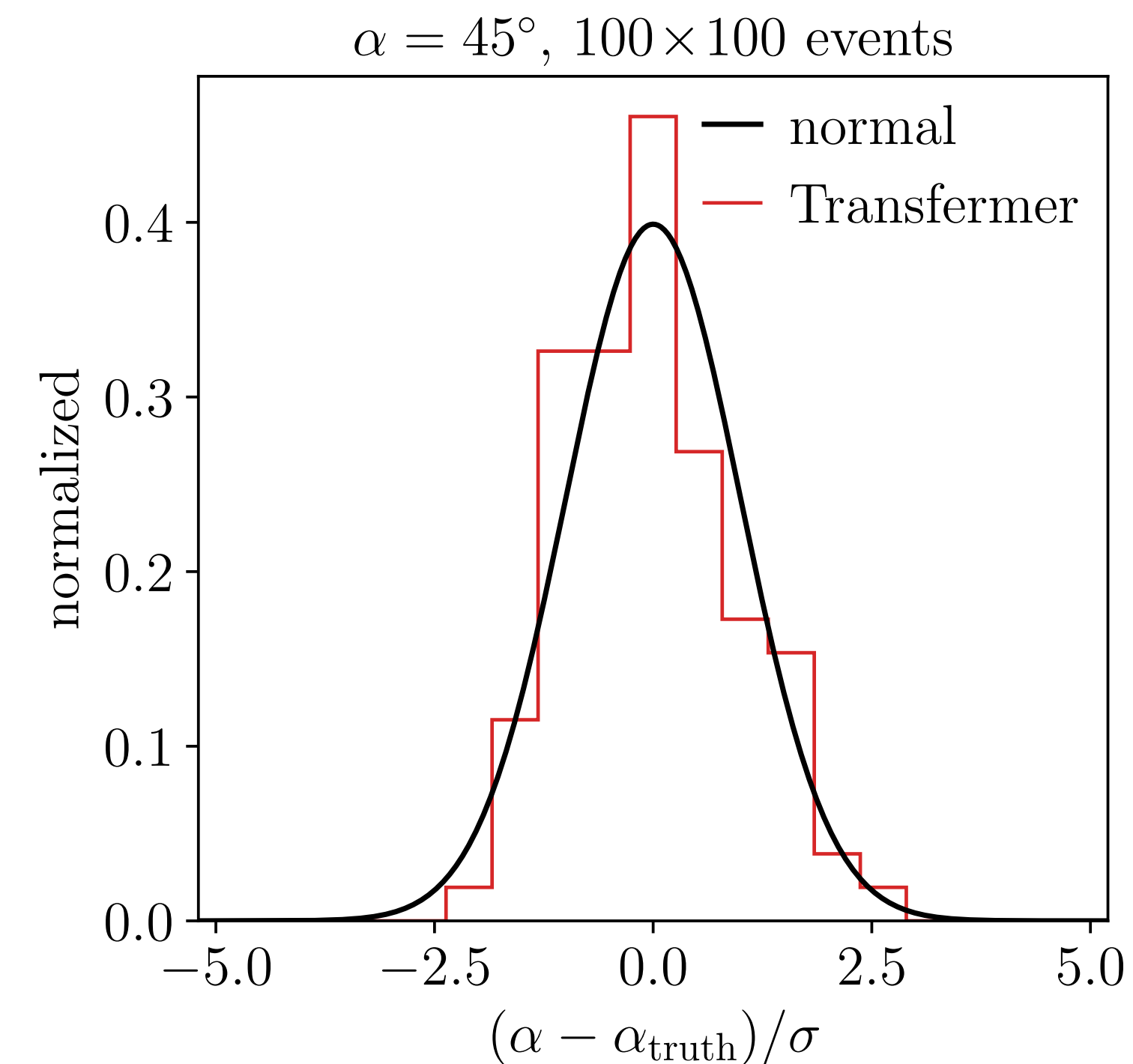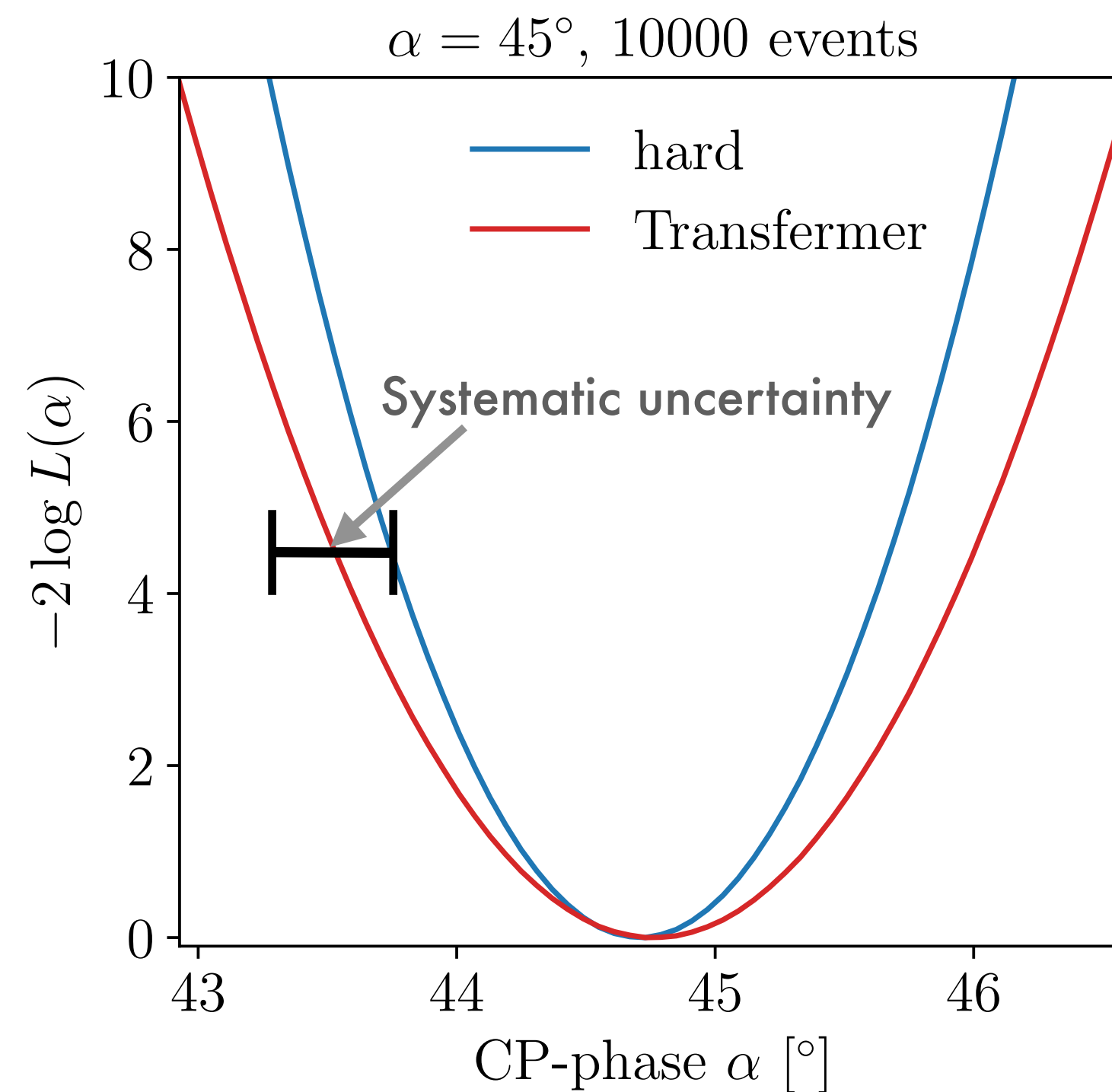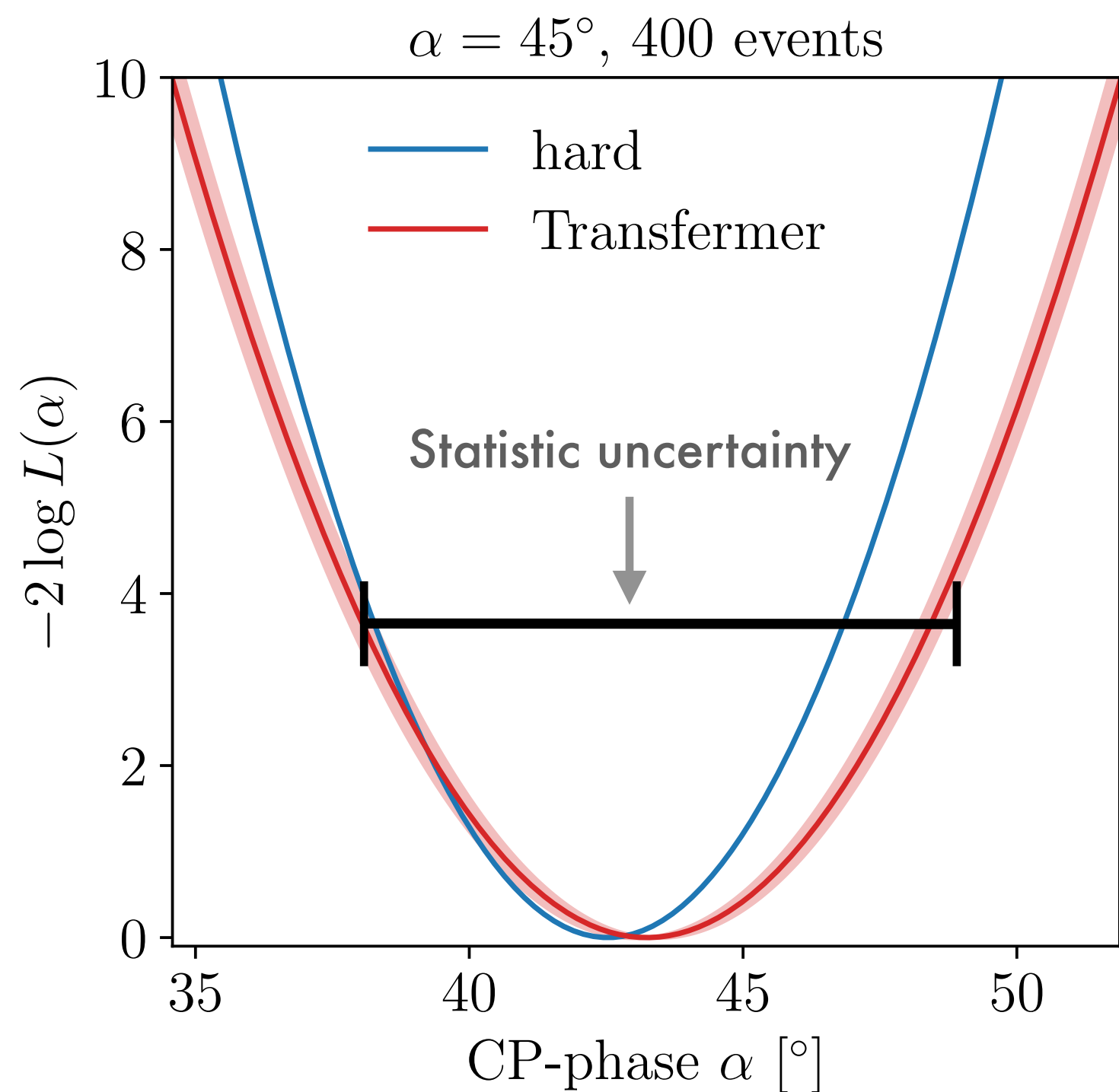
Anomalous coupling
with **CP-angle** $\alpha$ →



hadronic, with ISR

⊖ low total cross section
⊖ low variation of rate

⊕ kinematics sensitive

→ ideal use case for **MEM**

→ **well-calibrated likelihoods**, both for low and high event count

→ Uncertainty bands: **MC integration error** &
  **syst. error from limited training statistics** (Bayesian NN)