

The background of the slide is a vibrant cosmic scene. It features several large, glowing galaxies in shades of blue, purple, and white. Interspersed among these are numerous smaller, colorful nebulae and star clusters, creating a rich, multi-colored starfield. The overall effect is that of a deep-space exploration or a view of the universe's vastness.

# **Search for New Physics with Machine Learning**

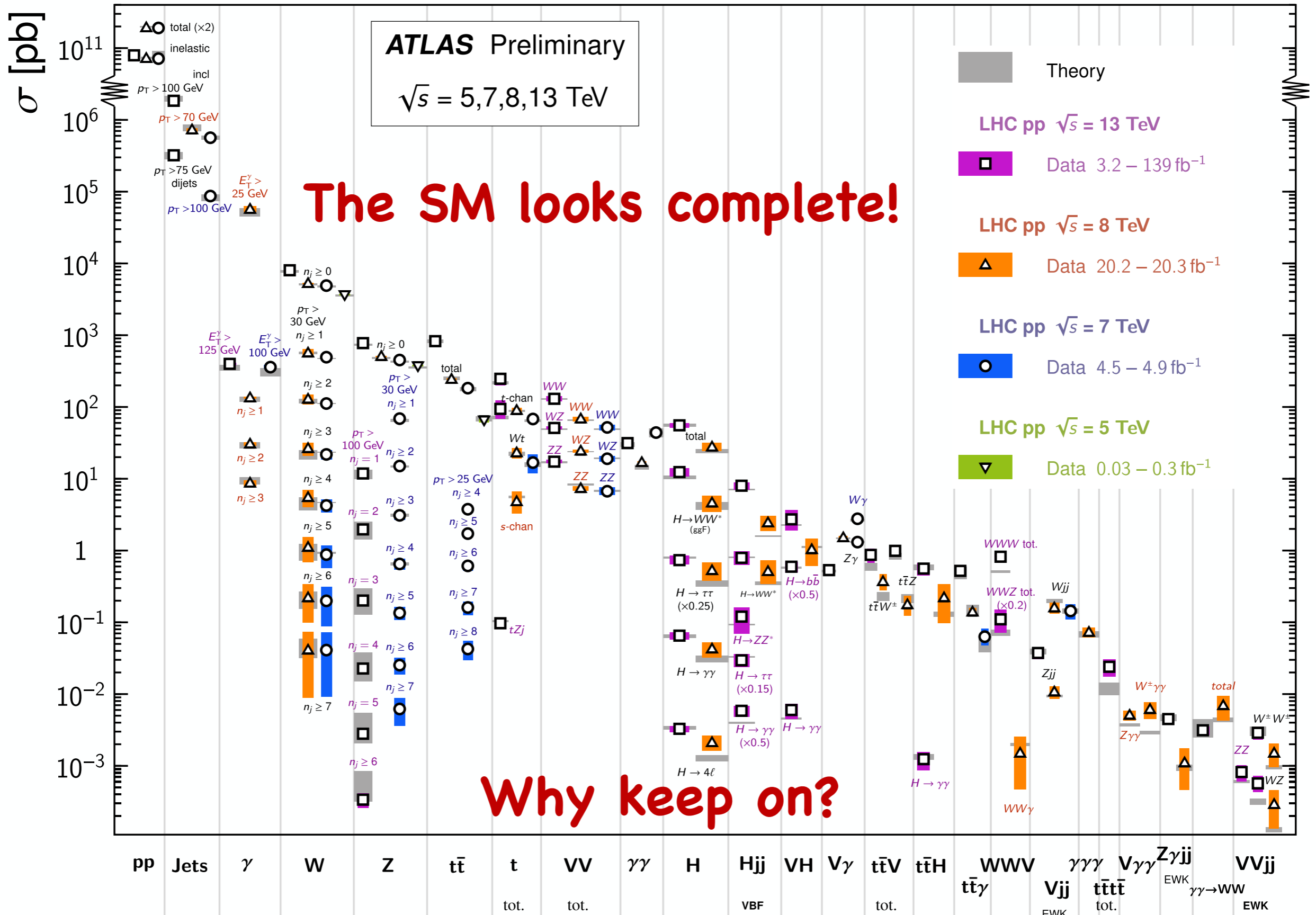
concepts, applications and recent progress

**Roman Pasechnik**  
Lund University

# Impressive performance of the Standard Model

## Standard Model Production Cross Section Measurements

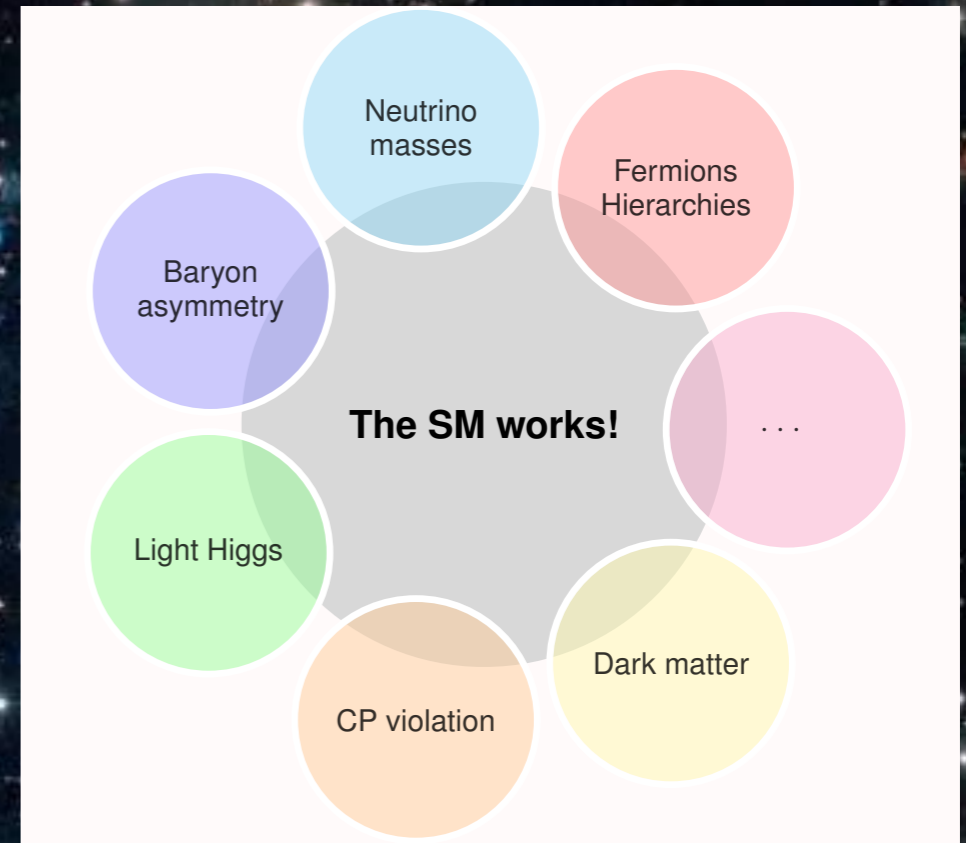
Status: February 2022



# Open questions

Unexplained phenomena:

- Dark Matter
- Matter-antimatter asymmetry
- Dark Energy



Unsatisfactory structure of the SM:

- Hierarchy problems (Higgs, flavour)
- Naturalness
- Quantum Gravity etc

Plethora of Beyond the SM theories to test against data

# Goals of Particle Physics

- Explore physics at **the highest energy scale** (TeV scale at the LHC)
  - search of new Higgs scalars ( = Higgs “partners”)
  - search for low-energy traces of supersymmetry (SUSY)
  - investigate various scenarios of physics beyond the SM
- **Precision measurements** of SM processes:
  - test nontrivial predictions of the SM, including very rare processes
  - search for deviations from the SM sensitive to new physics at high scales
  - improve precision of the SM parameter measurements
  - study the Higgs boson and the EWSB sector
  - study QCD dynamics and parton content of the proton
  - QCD/MC tools development
  - exploit the SM measurements as “Standard Candles” to tune and test detector performance

# Emerging anomalies in collider data

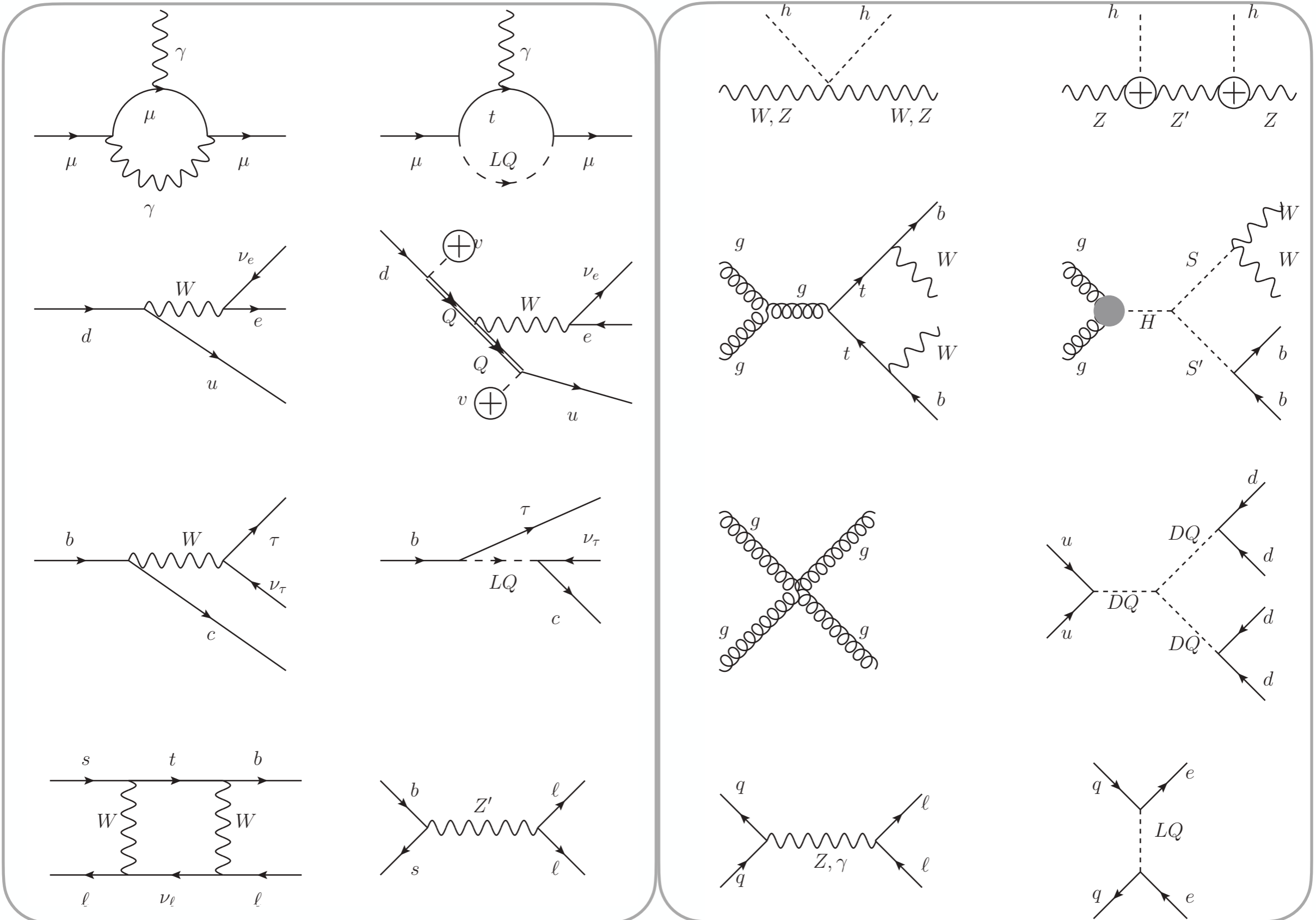
Crivellin, Mellado, Nature Reviews Physics 6, 294 (2024)

**Standard Model**

**New Physics**

**Standard Model**

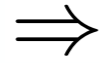
**New Physics**



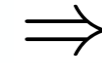
# Why ML is relevant for Particle Physics

Karagiorgi, Kasieczka, Kravitz, Nachman, Shih, arXiv:2112.03769

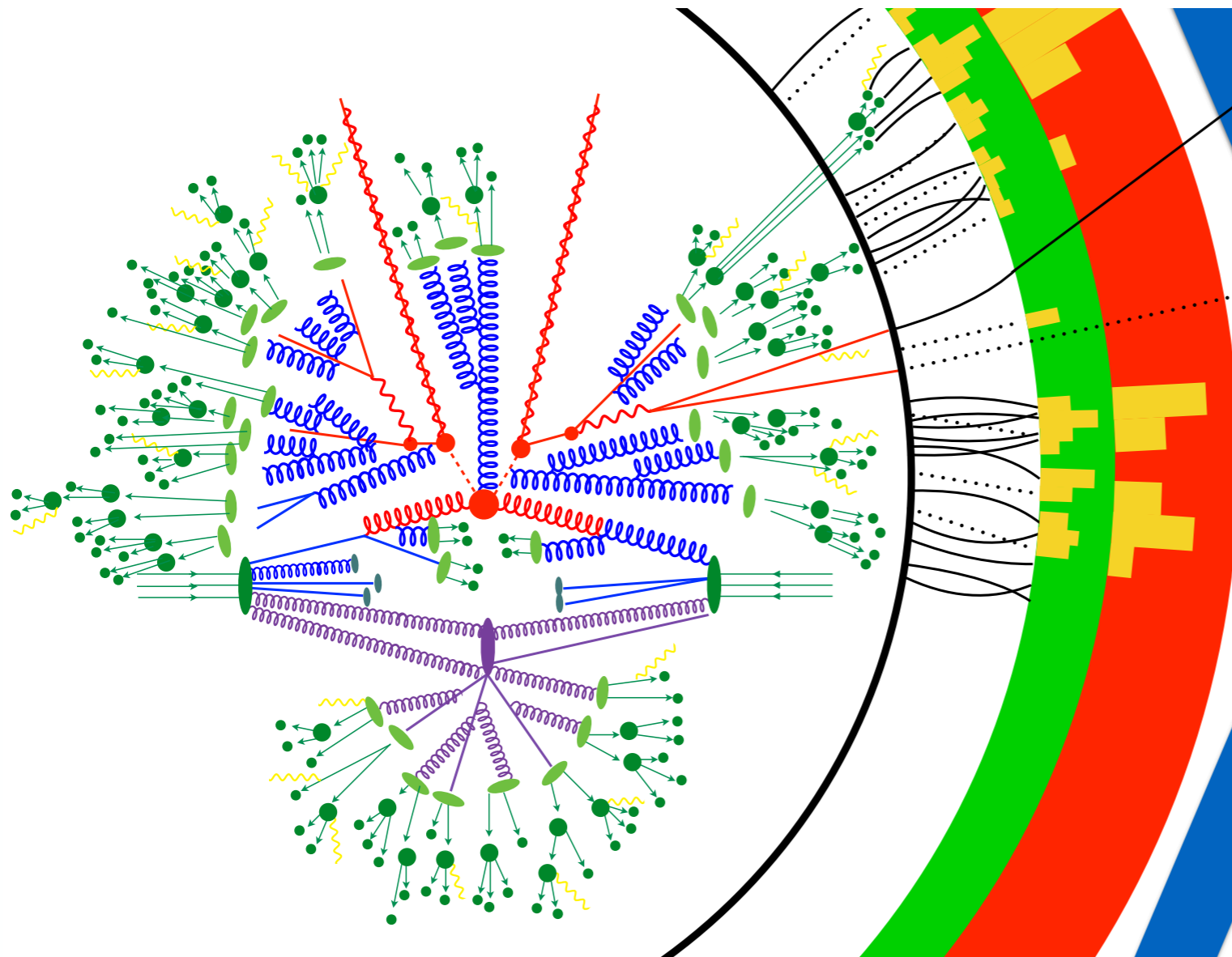
Microscopic physics  
(particle interactions,  
scattering...)



Large-distance phenomena  
(hadronisation, propagation, jets...)



Detector response



Particle physics  
experiments produce enormous  
datasets of high complexity!

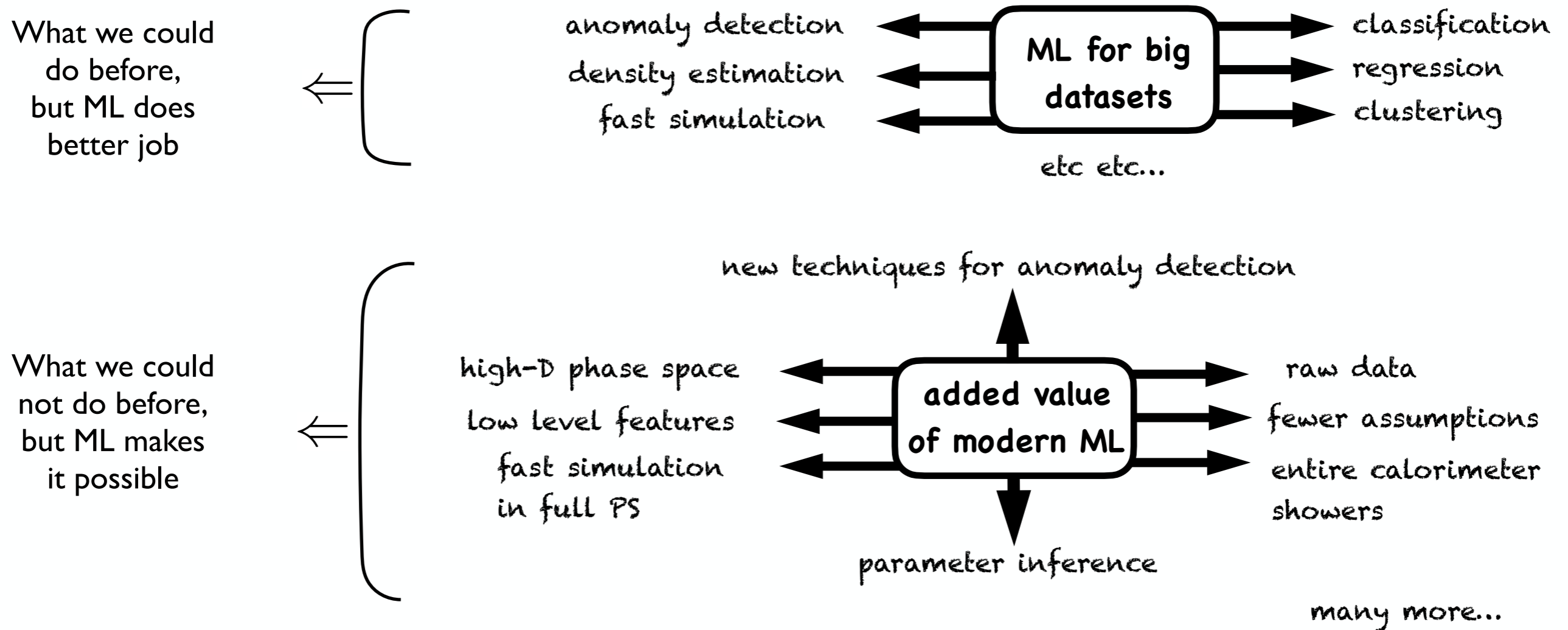
ML is a powerful new tool that  
enables us to get more physics  
out of these big datasets

We live in a very special point of history  
similar to invention of a telescope

ML would enable us to see features we could not see before — “telescope” for Big Data!

# What ML can do for particle physicists

Mehta et al, Phys. Rept. 810, 1-124 (2019)

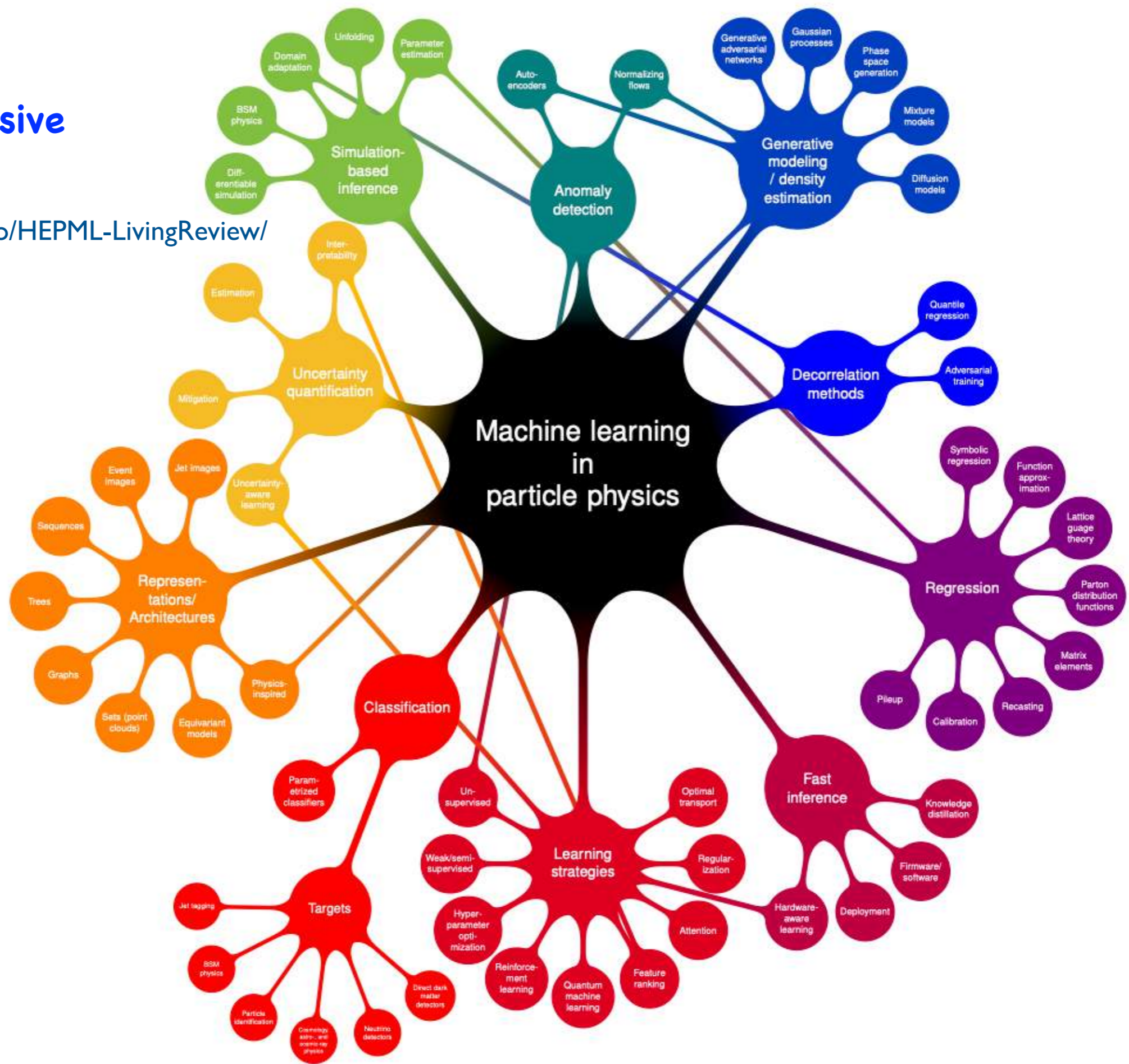


## Biggest advantages of ML

- ◆ greatly enhanced sensitivity/precision [x10-100]
- ◆ accelerated simulation [e.g. fast simulation]
- ◆ accelerated/efficient extraction of physics [inference]
- ◆ ML is cross-disciplinary [same methods can be used in many fields]
- ◆ ML is not only for experimentalists - theorists use "simulated data" a lot!

For comprehensive review, see

<https://iml-wg.github.io/HEPML-LivingReview/>

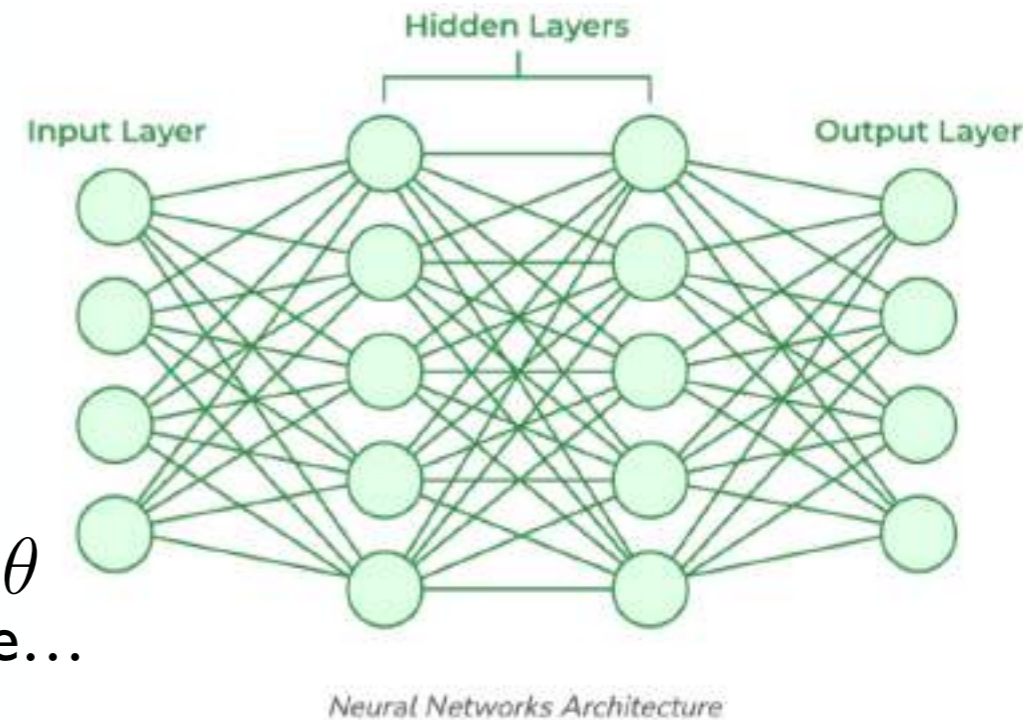


# Learning from data

## ML = sophisticated curve fitting!

- ◆ the data instances  $\vec{x}_i \in R^d \quad i = 1 \dots N$  are considered to be drawn i.i.d. from some data distribution  $p_{\text{data}}(x)$
- ◆ often we are interested to **“learn” a function**  $f(x, \theta)$  from the data for some parameters (e.g. weights, biases)  $\theta$   
 $\Rightarrow$  Neural Network, or made up of NNs, or smth else...

Example: feed-forward NN (MLP)



- ◆ the learning process (e.g. NN training) is optimisation of **an objective** (loss function)

$$L(\theta) = \sum_{\text{data}} \mathcal{L}[f(x, \theta)]$$

- ◆ **supervised learning** [regression or classification] — we want  $f(x, \theta)$  to take a specific form e.g. binary classification  $f(\vec{x}_i, \theta) = y_i$  for truth labels  $y_i = 1$  or  $0$  (QCD jets vs top jets)
- ◆ the objective is to get as close to the truth labels as possible e.g. minimise **the mean squared error (MSE) loss**

$$\mathcal{L} = (f(\vec{x}_i, \theta) - y_i)^2$$

- ◆ Truth labels often come from simulation, or when categorisation of the data is obvious e.g. hand-labelled data (cat vs dog, natural images etc)

# Less than supervised

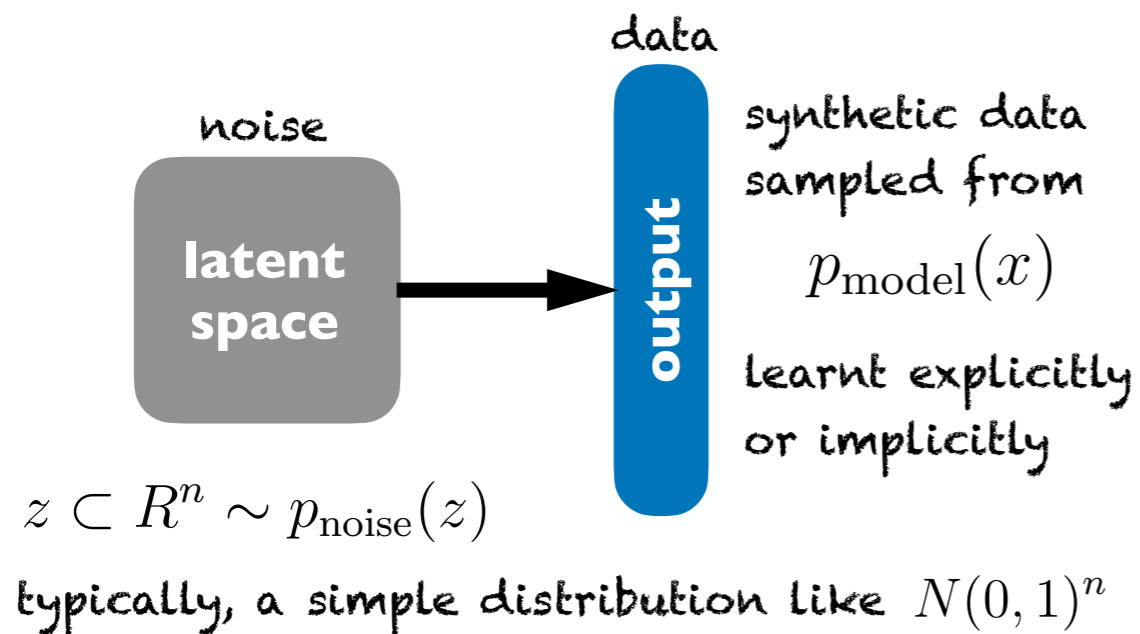
Data often come from experiments (e.g. LHC) without truth labels (non human interpretable) while simulations are not perfect -> What to do?

- ◆ in Particle Physics the data is very complex and not clearly separable into categories  
⇒ large overlap in their distributions
- ◆ **Unsupervised learning** — no labels available at all — data-driven/simulation-free approach!  
⇒ uncovering hidden patterns, structures or relationships in the data  
⇒ Examples: clustering similar data points, reducing dimensionality for visualisation
- ◆ “Noisy” labels: **weakly supervised learning** — a **powerful tool in New Physics searches**:  
⇒ data-derived labels that correlate with a given category [e.g. signal S vs background B] but may not be EXACTLY in that category  
⇒ Splitting the data into “signal region” [S-enriched] + “control region” [B-enriched], a generative model is applied for anomaly detection
- ◆ A mix of labelled and unlabelled data: **semi-supervised learning**  
⇒ simulation+data to mitigate the simulation effects, or when parts cannot be labeled
- ◆ Data-driven methods to learn the objective: **self-supervised learning** — by using symmetries or deleting parts of the data and trying to fill that in [relevant in Large Language Models]  
⇒ useful for learning embeddings (e.g. using data-derived labels on jets related or not related by rotation, one learns a jet representation encoding the symmetry)

# Generative ML

Can we generate more samples that follow the same distribution as the data?

- ◆ We want to learn the data probability distribution  $p_{\text{data}}(x)$  [density estimation] and then sample from it — often, a very difficult task!  
⇒ we can learn to sample from  $p_{\text{data}}(x)$  without actually learning this function
- ◆ **Generative modelling** — to learn  $p_{\text{model}}(x)$  as close as possible to  $p_{\text{data}}(x)$  and then sample to generate (dream up) synthetic data capturing underlying patterns of the original dataset



## Methods of generative ML

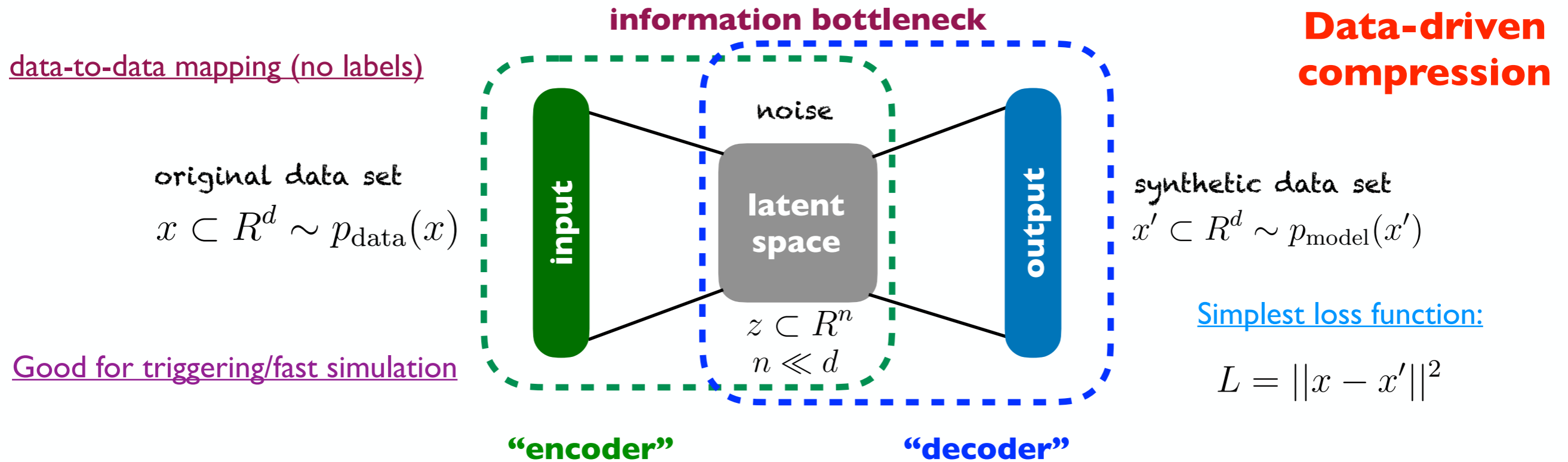
- ⇒ Generative Adversarial Networks (GANs)
- ⇒ Variational Autoencoders (VAE)
- ⇒ Normalising Flows
- ⇒ Diffusion Models

## Common applications in Particle Physics:

- ◆ fast simulation e.g. “surrogate” modelling (training on few samples to generate more), phase space sampling (integrations)
- ◆ **anomaly detection** — outliers (clean up data), group anomalies (bumps from NP)
- ◆ **simulation based inference** e.g. EFT fits

# Autoencoders

- ◆ AE maps data back to itself through reduced latent space trying to figure out a latent representation of the data that captures its essential features



- ◆ AE trained on “normal” events as an “anomaly detector” [only outliers, not overdensity]
  - ⇒ an outlier causes the loss to have a “cluster” at large values
  - ⇒ anomalies may be separated in the latent space — compression enhances clustering
  - BUT no guarantee that AE learns sensible latent space [no probabilistic interpretation]
- ◆ Simulation-based AE for New Physics search Farina, Nakai, Shih, PRD 101, 075021 (2020)  
Heimel et al, SciPost Phys. 6, 030 (2019)
  - ⇒ take simulated jet images as data: QCD jets as B and NP jets as S (anomaly)
  - ⇒ if S is rare, AE separates the anomaly well (in the tail of the loss)
- ◆ AE as a “complexity detector” Finke et al, JHEP 06, 161 (2021)
  - ⇒ train on QCD jets, finds top jets, BUT train on top jets, does not find QCD jets

# Variational Autoencoders

Can we enforce the latent space to have a suitable probabilistic interpretation?

◆ VAE as a latent variable model:  $z \sim p(z)$  (the “prior”) while  $x \sim p_\theta(x|z)$  we get a set  $\{x_i, z_i\}$  drawn from  $p(x, z)$  — by integrating out  $z$  we get data distribution  $p_\theta(x)$

◆ To determine the conditional probability  $p_\theta(x|z)$  — learn it by **maximising the maximum likelihood estimation (MLE)** w.r.t NN parameters  $\theta$

$$\text{MLE} = \sum_{x \sim p_{\text{data}}(x)} \log p_\theta(x) \quad \text{with Bayesian evidence} \quad p_\theta(x) = \int p_\theta(x|z)p(z)dz \quad \text{hard!}$$

“encoder” “decoder”

◆ **Variational “posterior”**  $r_\psi(z|x)$  — still samples z-space but differently depending on x

◆ Utilising MC sampling of the integral from the posterior and applying Jensen inequality:

$$\log p_\theta(x) \geq \sum_{z \sim r_\psi(z|x)} \log \frac{p_\theta(x|z)p(z)}{r_\psi(z|x)} = \log p_\theta(x) - \text{KL}(r_\psi(z|x) || p_\theta(z|x)) \quad \text{“evidence lower bound” (ELBO)}$$

Kullback-Leibler divergence  
 true posterior for a given  $\theta$

◆ Taking normal distribution  $p_\theta(x|z) = N(\mu_\theta(z), \beta)$

$$\text{ELBO} = E_{r_\psi(z|x)} \left[ -\frac{\|x - \mu_\theta(z)\|^2}{2\beta^2} \right] - \text{KL}(r_\psi(z|x) || p(z))$$

**Maximised!**

reconstruction error of vanilla AE

VAE as a “regularised” vanilla AE with a “smoothing” KL-term (posterior tends to the prior)

# Normalising Flows

Papamakarios et al,  
J. ML Res., 22(57) 1, 2021

## Invertible map between the data and latent spaces

- ◆ We can achieve both — get the latent space  $z$  and do density estimation (DE)

$$\begin{array}{ccc} z = f(x, \theta) & \longleftrightarrow & x = f^{-1}(z, \theta) \\ p(z) & \longleftrightarrow & p_{\theta}(x) = p(z) \left| \det \frac{df}{dz} \right| \end{array}$$

Optimisation problem:  
to fit parameters  $\theta$   
to the data

- ◆ We directly optimise the negative-log likelihood [tends to perform better than VAE]

$$L = - \sum_x \log p_{\theta}(x)$$

best model that generalises best to the  
unseen test set

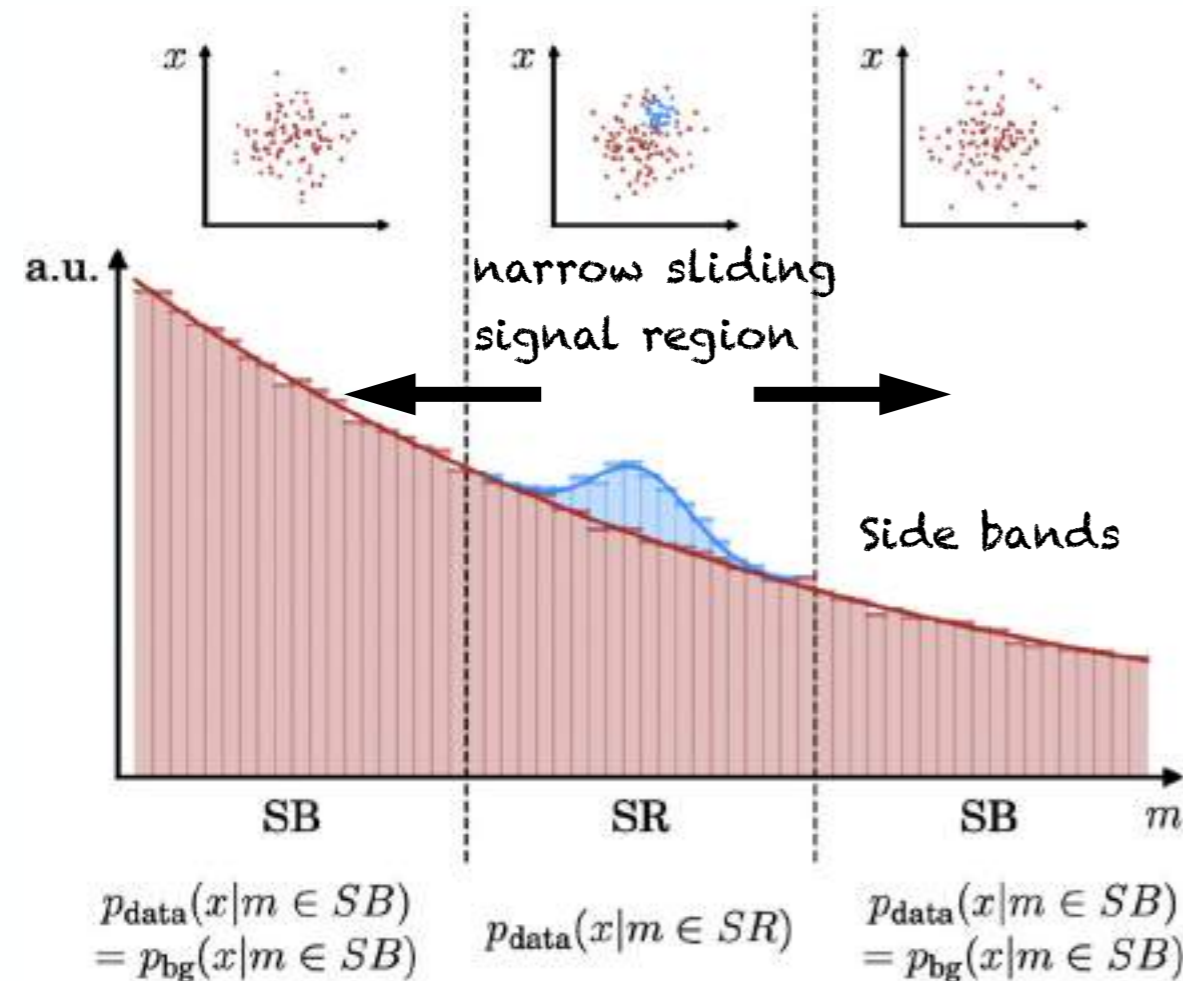
- ◆ Trade off: more ambitious/challenging, bad scaling with dimensions;  
 $\Rightarrow$  Jacobian of  $d \times d$  matrix takes  $O(d^3)$  operations repeated many times
- ◆ A small family of invertible functions with upper-triangular matrices:  $O(d)$  operations!  
 $\Rightarrow$  “autoregressive transformation”  $z_1 = f_1(x_1), \quad z_2 = f_2(x_1, x_2), \quad \dots \quad z_n = f_n(x_1 \dots x_n)$
- ◆ To gain on expressivity [= having enough parameters to learn the transform], one can chain multiple such transforms, permute between them —  
 $\Rightarrow$  **Masked Autoregressive Flow (MAF)**:  $z_i = \alpha_i(x_1, \dots, x_{i-1}, \theta) z_i + \mu_i(x_1 \dots x_{i-1}, \theta)$  where the coefficients can be thought as outputs of one big NN (slow sampling)  
 $\Rightarrow$  **Inverse Autoregressive Flow (IAF)**: inverse algorithm (slow training)
- ◆ “**Density distillation**”: first train the MAF then distill it into an IAF [fit IAF to MAF]

# Resonant anomaly detection: SIC

Golling et al, EPJC 84, 241(2024)

## How we use ML in searches for new phenomena in Particle Physics?

- ◆ Assume  $S$  is localised (resonant) in some feature (typically, invariant mass) and  $B$  is smooth



- ◆ **Inclusive bump hunt** — standard technique for new particle searches at colliders:
  - ⇒ split the distribution into SR and SBs
  - ⇒ smooth interpolation provides fully data-driven  $B$  in the SR
  - ⇒ discovery significance via Poisson  $\sigma = \frac{S}{\sqrt{B}}$
- ◆ How ML can enhance the bump?
  - ⇒ **multivariate bump hunt**: looking for correlated excesses in other features  $x$  (e.g. jet substructure, missing energy etc)

- ◆ **Anomaly Score**  $R(x)$  — large for  $S$ , and small for  $B$  — must be uncorrelated with the mass
- ◆ **Significance Improvement Characteristic (SIC)** — how much the significance is improved by a cut on  $R(x)$  based upon how many  $S, B$  events survived:

$$\text{SIC} = \frac{\epsilon_S}{\sqrt{\epsilon_B}} \quad \text{for cut efficiencies } \epsilon_{S,B}$$

We want a ML technique that produces large SIC in a data-driven way

**Proof of concept: inject a small signal to simulated  $B$ , see how it shows up in SIC**

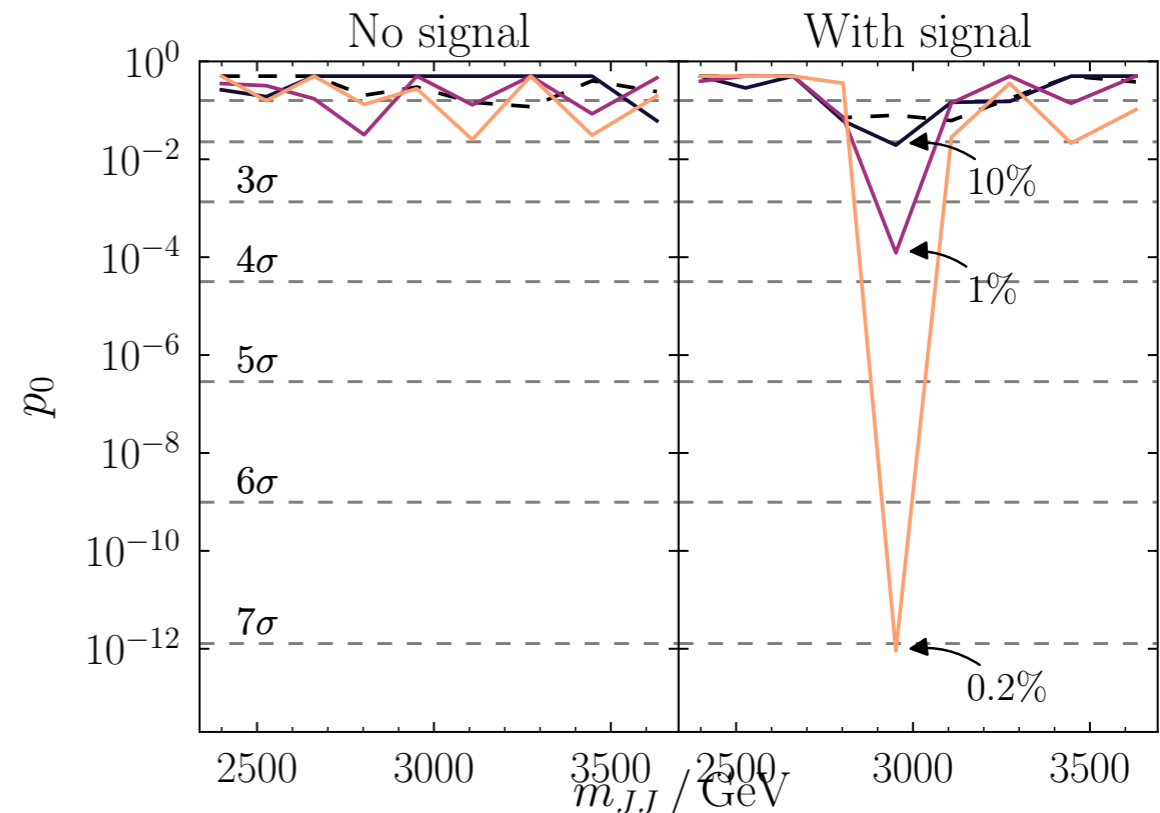
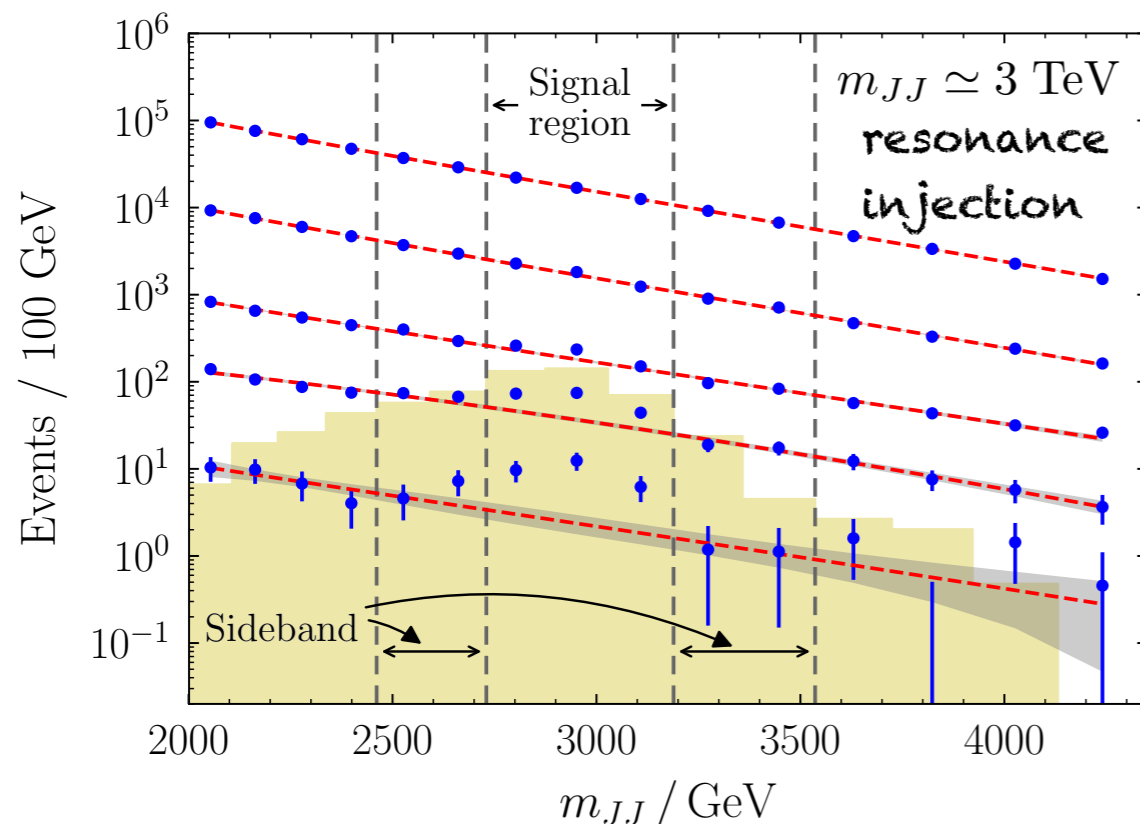
# CWoLa Hunting

Collins, Howe, Nachman, PRD 99, 014038 (2019)  
ATLAS search in real data: 2005.02983

How can we evaluate anomaly score and improve discovery significance?

- ◆ For any S, the S-to-B likelihood ratio  $p_S(x)/p_B(x)$  is an **optimal S/B classifier!**  
⇒ how do we learn approximations to  $p_{S,B}(x)$  using ML?
- ◆ **“Classification W/o Labels” (CWoLa)** — train a classifier to learn SB vs SR to approximate the ideal anomaly score  $R_{\text{ideal}} = p_{\text{data}}(x)/p_B(x)$  fully correlated with the ratio  $p_S(x)/p_B(x)$   
⇒ if features  $x$  and the mass are uncorrelated in B, train a binary classifier to learn  $R(x)$  using noisy labels SR (label 1) and SB (label 0) — weak supervision!  
⇒ output: probability that a given  $x$  comes from SR

Bayes' theorem: 
$$p(\text{SR}|x) = \frac{p_{\text{SR}}(x)}{p_{\text{SR}}(x) + p_{\text{SB}}(x)} = \frac{R(x)}{1 + R(x)} \quad \longrightarrow \quad R(x)$$



# Anomaly detectors with correlation

## What happens if features are correlated with mass?

- ◆ CWoLa stops working — becomes just a B-to-B classifier!
- ◆ “Anomaly detection with density estimation” (ANODE) Nachman, Shih, PRD 101, 075042 (2020)  
⇒ no classifier, train conditional DEs (Normalising Flows) to learn SR and SB densities conditioned on mass
- ◆ “Classifying Anomalies thorough outer DE” (CATHODE) Hallin et al, PRD 106, 055006 (2022)  
combine best of CWoLa and ANODE: learn SB density in ANODE, interpolate into SR, train a classifier on SR data in CWoLa — great performance!
- ◆ CURTAINS method Raine, Klein, Sengupta, Golling, Front. Big Data 6, 899345 (2023)  
instead of Normalising Flows, do invertible NN that learns to map SB to SB using optimal transport loss
- ◆ LaCATHODE method Hallin et al, PRD 107, 114012 (2023)  
using CATHODE in latent space
- ◆ Methods using simulation with re-weighting:
  - SALAD Andreassen, Nachman, Shih, PRD 101, 095004 (2020)
  - FETA Golling, Klein, Mastandrea, Nachman, PRD 107, 096025 (2023)
  - SA-CWoLa Benkendorfer, Pottier, Nachman, PRD 104, 035003 (2021)

# Anomaly detectors: performance

Golling et al, EPJC 84, 241(2024)

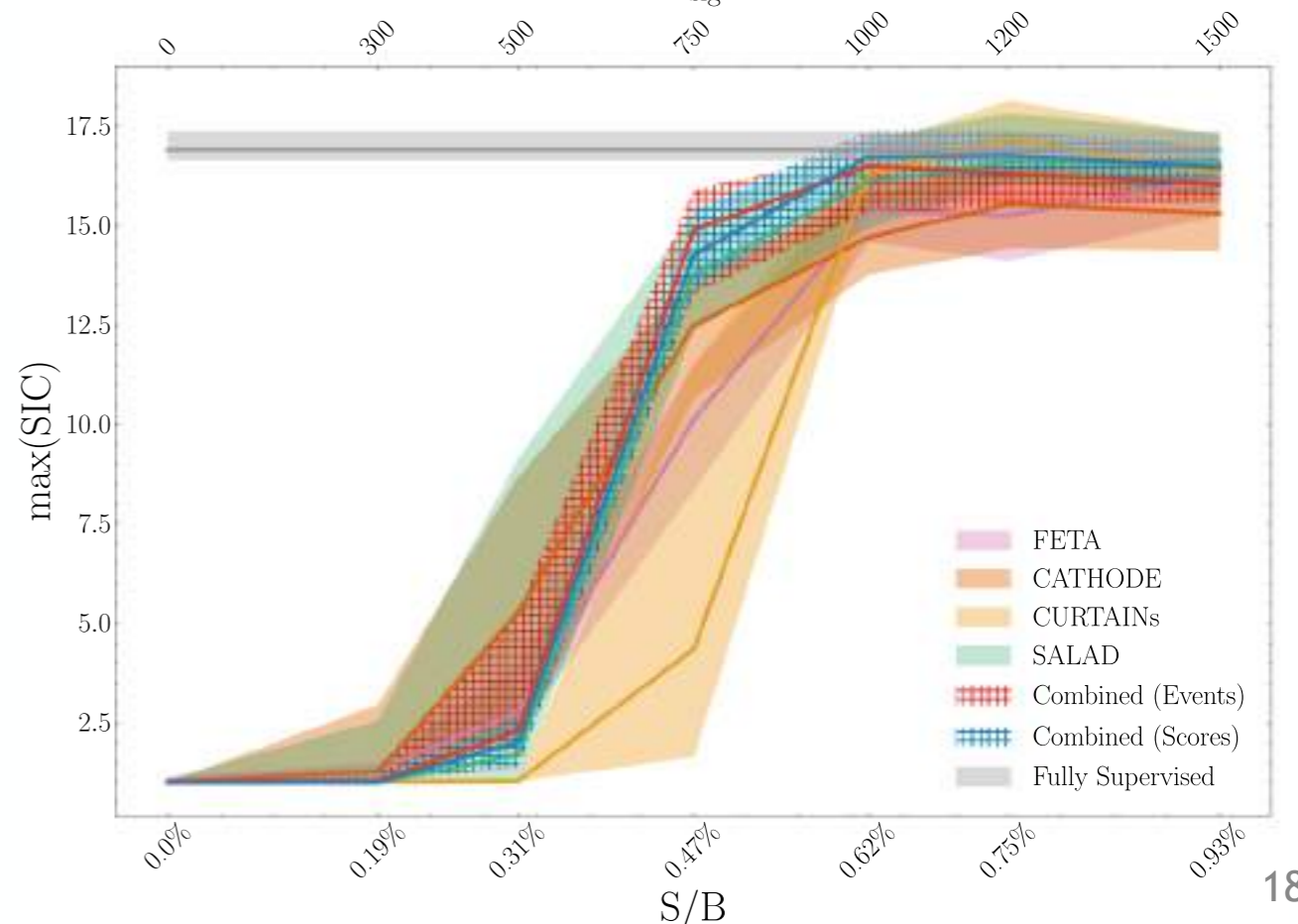
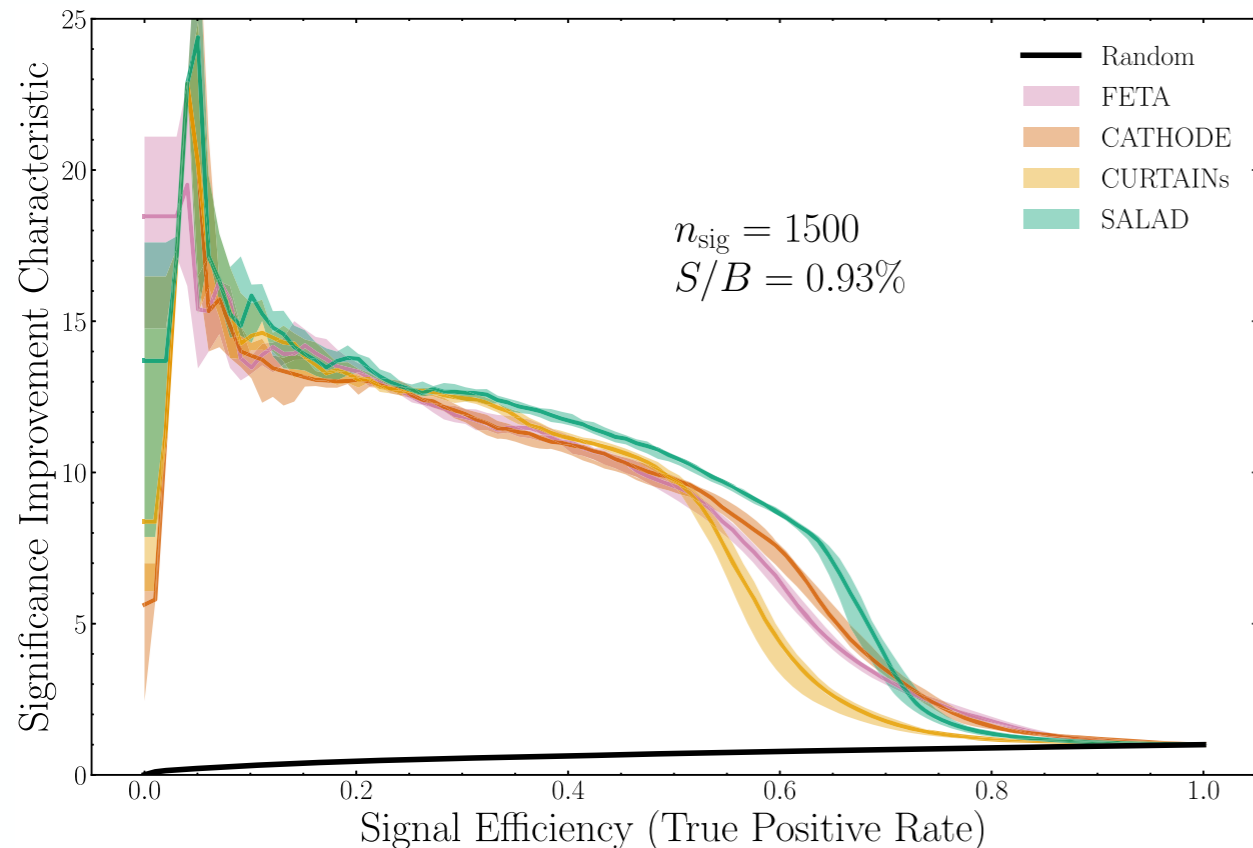
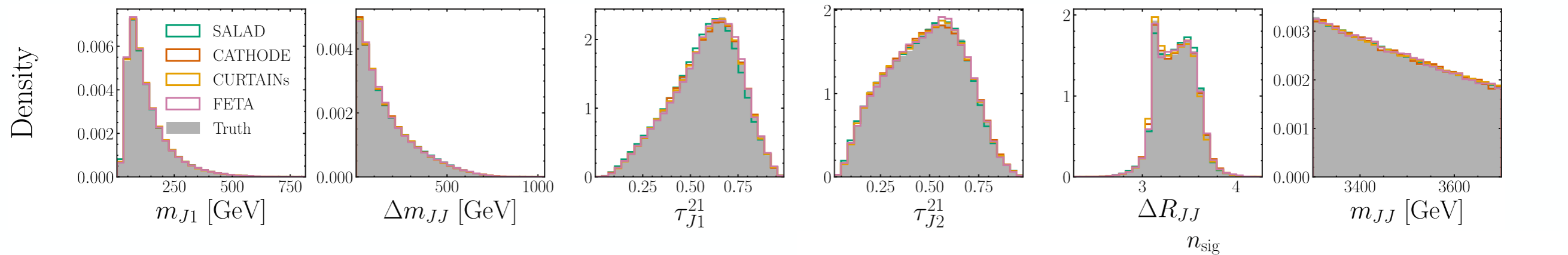
## ◆ ML model building: LHC Olympics R&D data set - fully labeled

Pythia + Delphes: 13 TeV pp, leading jet  $p_T > 1.2$  TeV

SM background: QCD di-jets

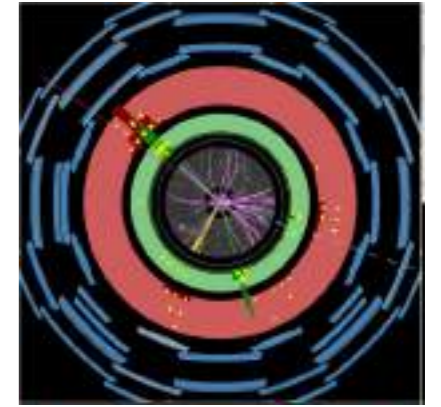
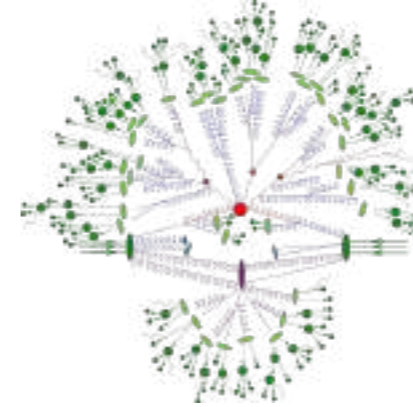
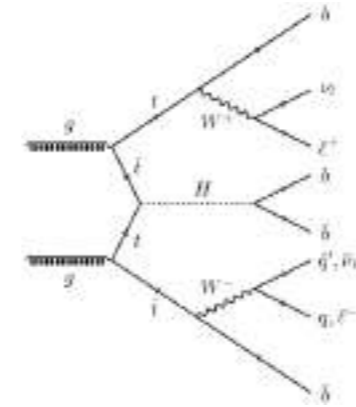
R&D signal:  $pp \rightarrow Z' \rightarrow X(\rightarrow jj)+Y(\rightarrow jj)$       $m_{Z'} = 3.5$  TeV      $m_Y = 100$  GeV      $m_X = 500$  GeV

## Synthetic data in 5D feature space + di-jet mass (SM B only)



# Simulation based inference

## Data generation



O(20) Fundamental physics parameters  $\theta$

O(10) particles

O(100) particles

O(10<sup>8</sup>) detector elements

$$p(\text{100M detector elements} \mid \text{SM / BSM parameters}) ?$$

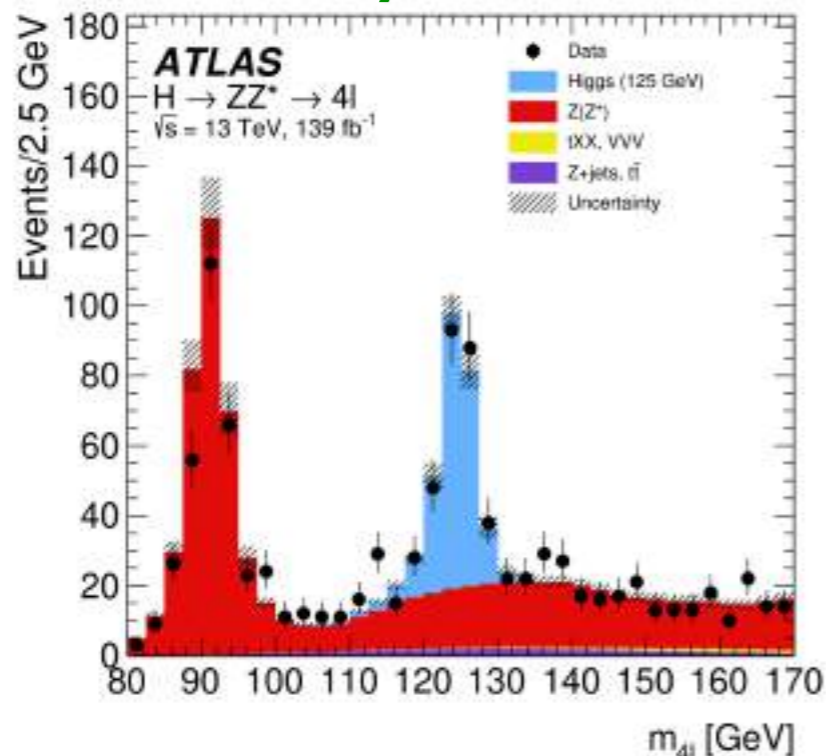
$$p(x|\theta) = \int dz p(x|z_h)p(z_h|z_p)p(z_p|\theta)$$

unobserved random processes: mathematically intractable!

100M detector elements

SM / BSM parameters

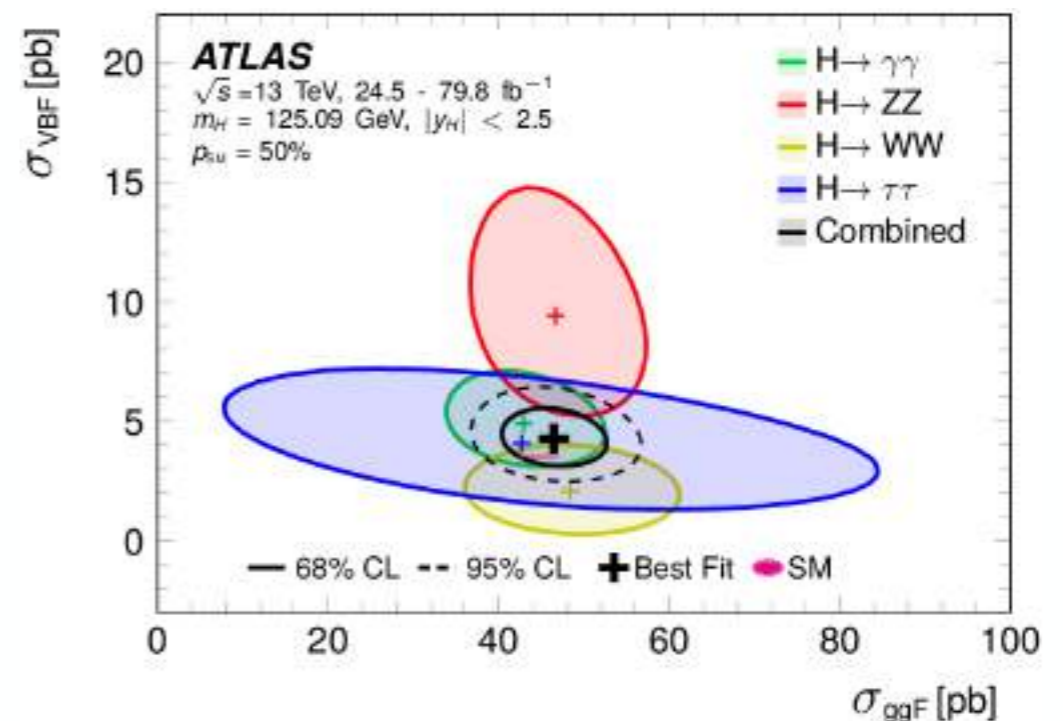
## Density estimation



in 1D we do this all the time!



## Likelihood for inference



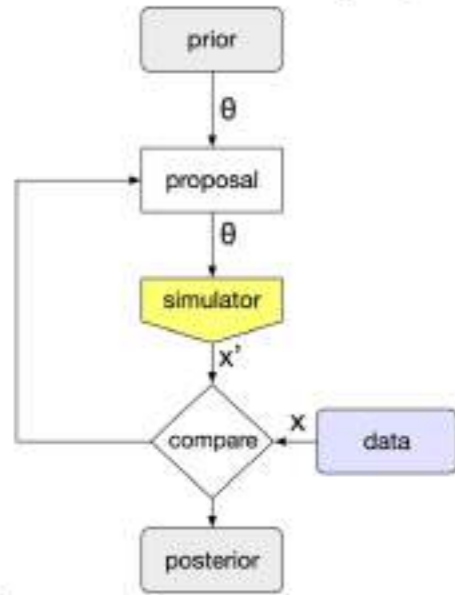
This is extremely challenging in large d's!

# Simulation based inference with ML

## Summary of different approaches:

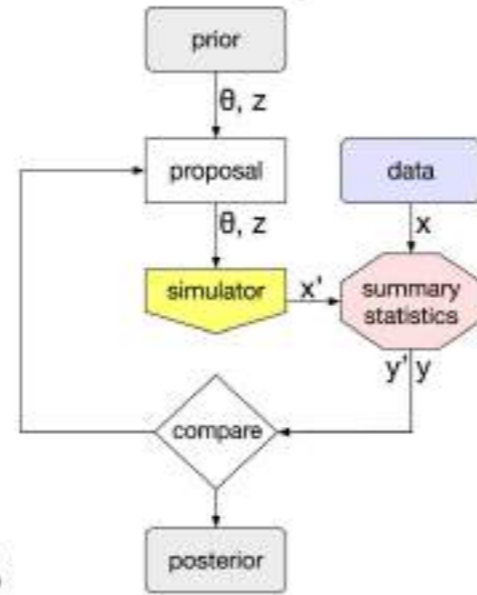
Cranmer, Brehmer, Louppe, PNAS 117, 30055 (2020)

**Approximate Bayesian Computation with Monte Carlo sampling**



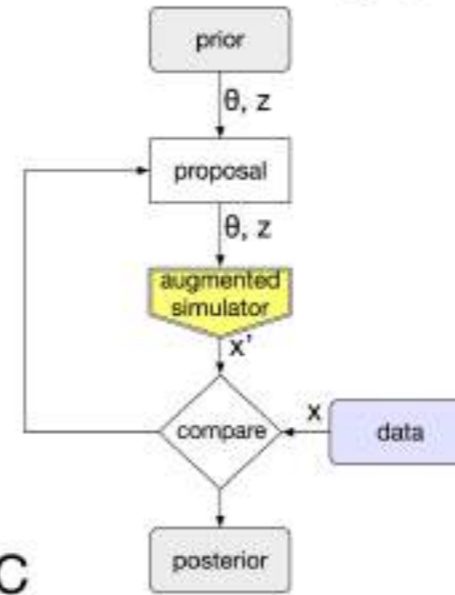
A

**Approximate Bayesian Computation with learned summary statistics**



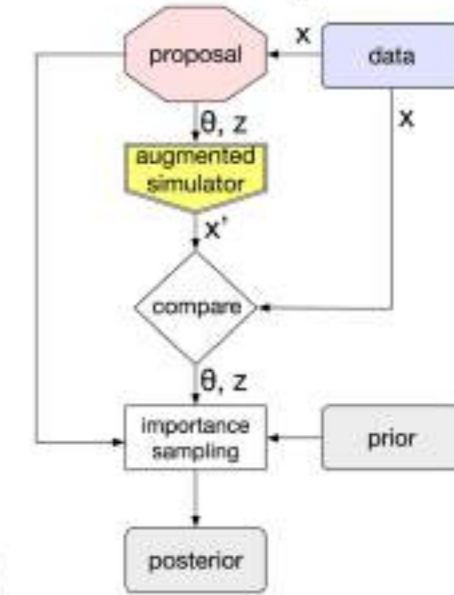
B

**Probabilistic Programming with Monte Carlo sampling**



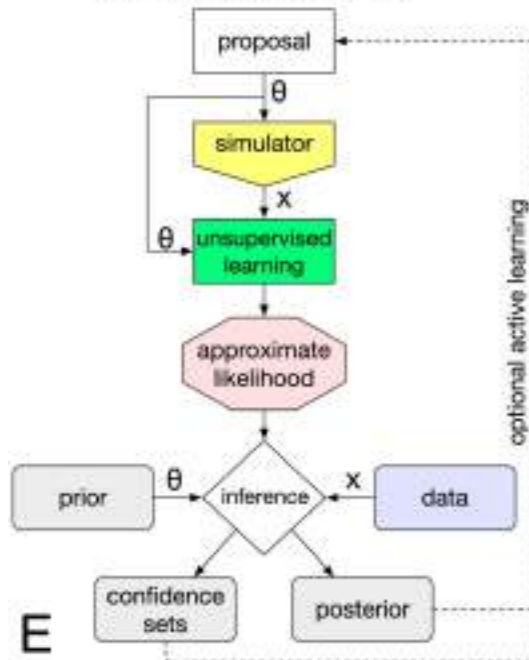
C

**Probabilistic Programming with Inference Compilation**



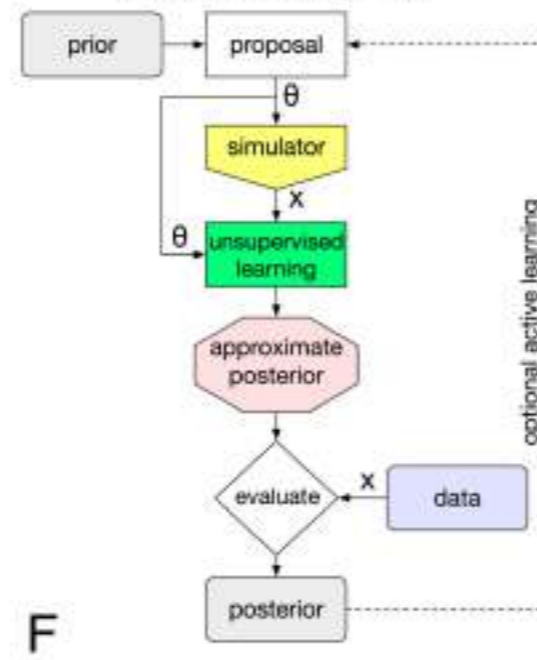
D

**Amortized likelihood**



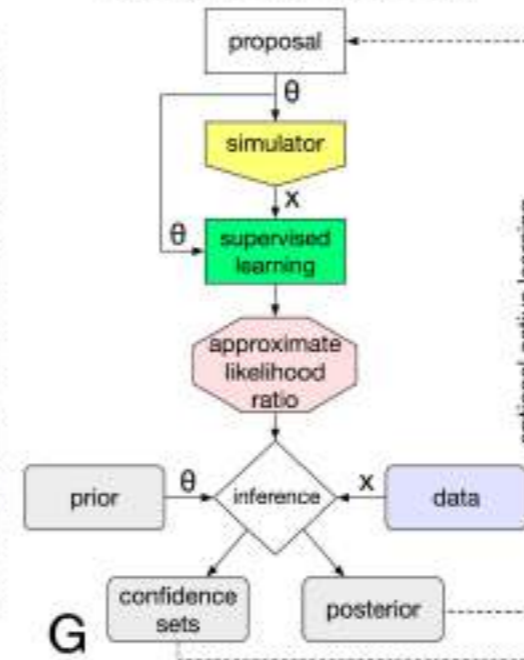
E

**Amortized posterior**



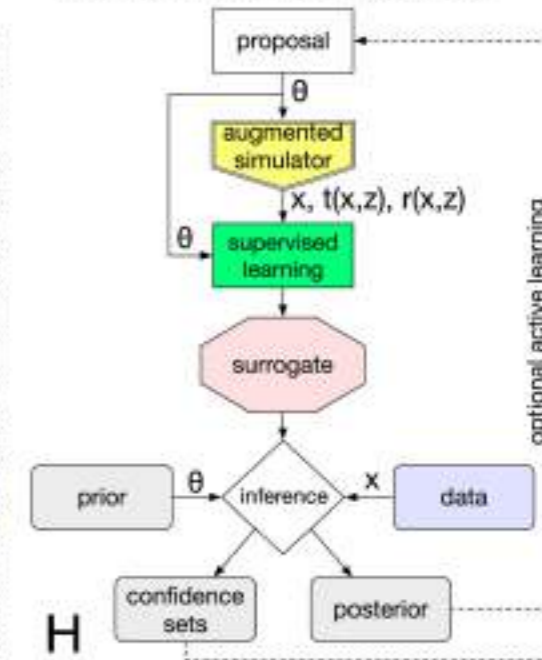
F

**Amortized likelihood ratio**



G

**Amortized surrogates trained with augmented data**



H

Use simulator to train NN to approximate likelihood ratio (using classifier), likelihood or posterior (using Normalising Flow), then use NN to do parameter inference on observed data

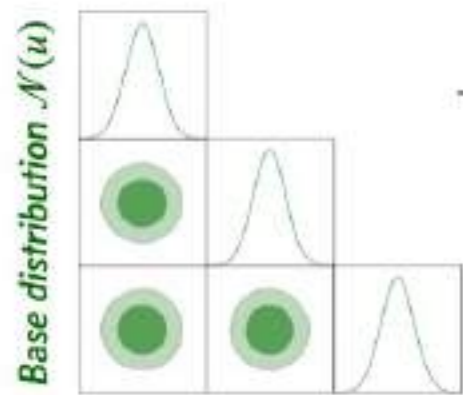
# Conditional neural posterior DE

Getting posterior distribution of model parameters for a given (simulated) data

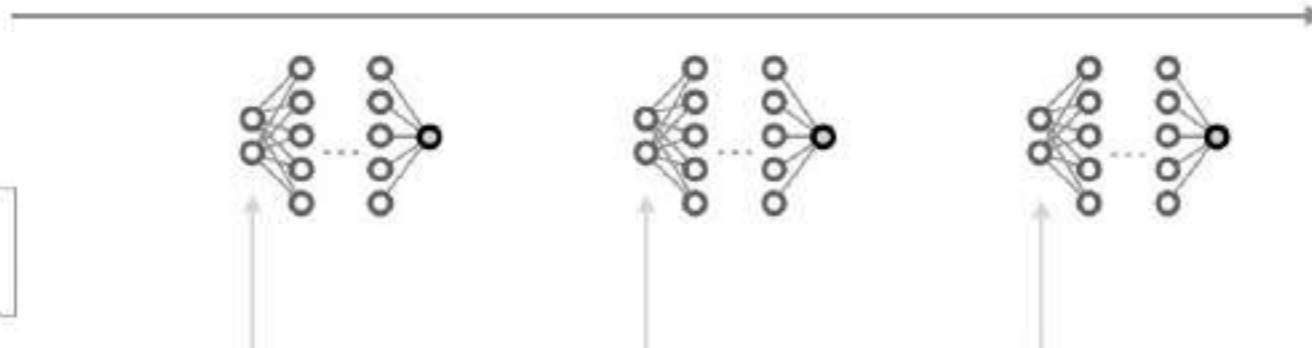
Posterior conditioned on  $x$

random noise

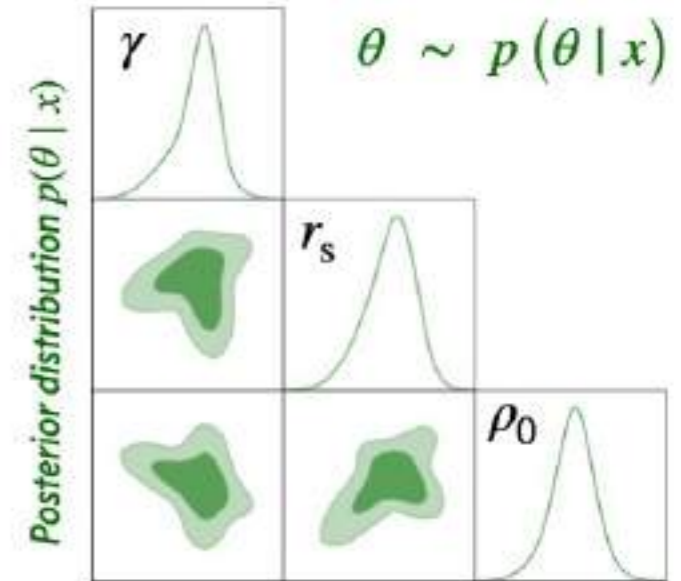
Normalising Flow



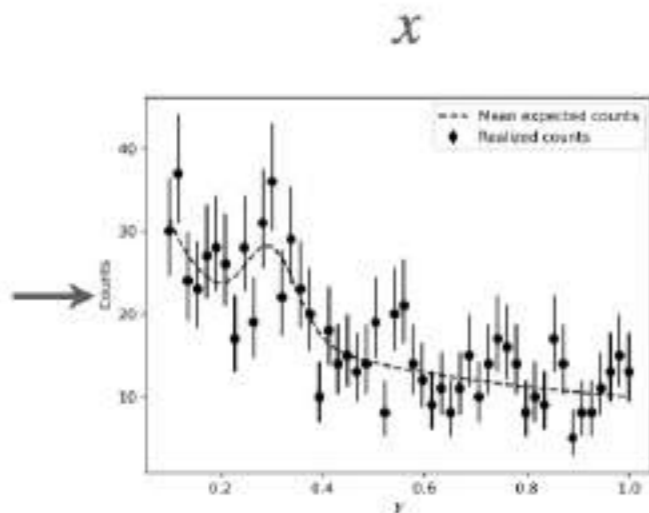
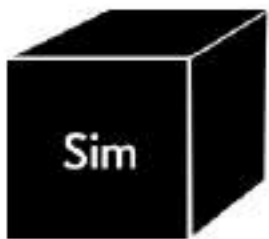
Conditional flow transformation  $\theta = f_\phi(u)$



Feeding some info on  $x$  at each layer

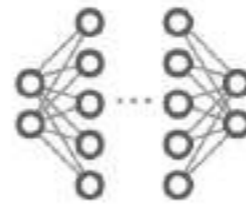


$\theta \sim p(\theta)$



a NN evaluates a smaller summary of data  $x$

$s_\phi(x)$



Feature extractor



Conditioning

Finally, evaluate the model on the measured data to get inference (e.g. confidence intervals)

Optimized simultaneously:

- Feature extractor  $s_\phi$
- Flow transformation  $f_\phi$

Physics-informed NN architectures may offer new ways to gain on efficiency of training and precision

# Summary

- ◆ ML is a new exciting research field that bridges Data Science, theoretical and experimental physics together
- ◆ ML is revolutionising research in Particle Physics and other data-intensive frontier fields offering a huge improvement in precision and computational efficiency, and new applications
- ◆ ML enables to hunt for subtle features in large and complex datasets and potentially to infer the best fundamental physics model combining vast amounts of data from different measurements and theoretical constraints
- ◆ While New Physics remains elusive, ML offers new opportunities for future discoveries