# Motivation

➔ **Background modelling** is one of the main challenges in particle physics analyses
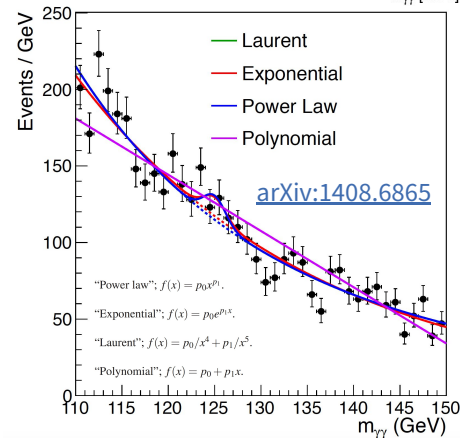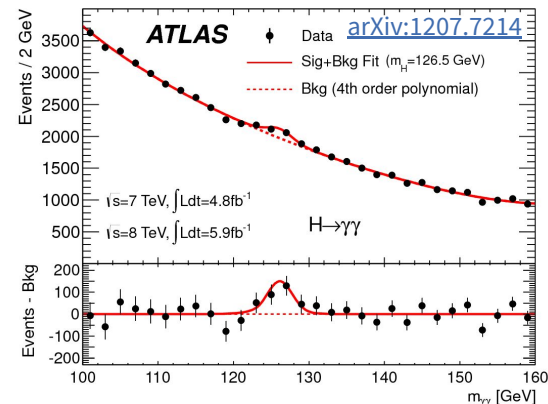
➔ Common techniques of background modelling:

◆ **MC simulation:**
- Not always possible to model background with sufficient accuracy
- Often computationally costly to produce large samples → significant statistical uncertainties

◆ **Parametric Models:**
- Does the true shape belong to the family of curves parametrised by the chosen function?
  - Taken into account as **"spurious signal"** systematic uncertainty [1, 2]
  - Discrete **profiling of an ensemble of parametric forms** [3, 4, 5]

[arXiv:1207.7214 (1), arXiv:2007.07830 (2), arXiv:1408.6865 (3), arXiv:2002.06398 (4), arXiv:2009.04363 (5)]
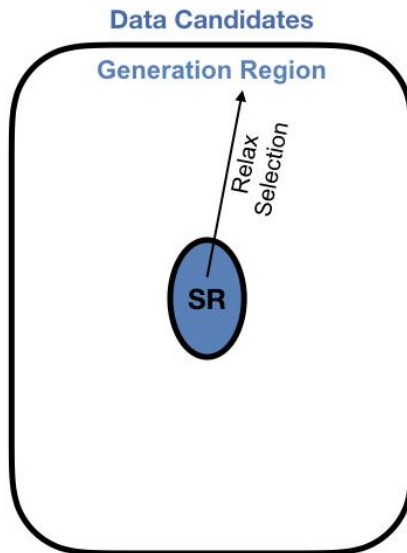
# Background Modelling

➔  Non-parametric data-driven background modelling:

# Background Modelling

➜ Non-parametric data-driven background modelling:

    1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)

# Background Modelling

➜ Non-parametric data-driven background modelling:
   1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
   2. Obtain conditional PDF of relevant variables ($x_1$, $x_2$,…, $x_n$)

**Data Candidates**

**Generation Region**

*Relax Selection*

**SR**

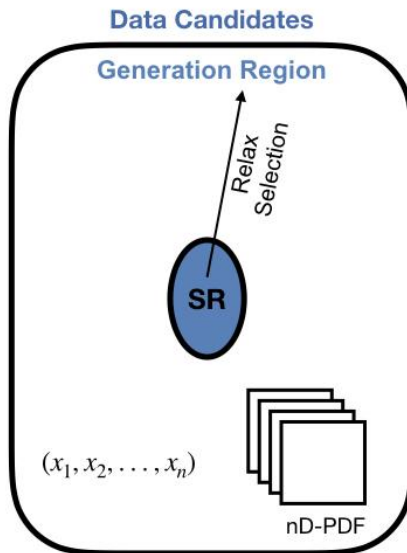$(x_1, x_2, \ldots, x_n)$

nD-PDF

# Background Modelling

➜ Non-parametric data-driven background modelling:
1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
2. Obtain conditional PDF of relevant variables ($x_1$, $x_2$,..., $x_n$)
3. Generate sample of pseudo-candidates

# Background Modelling

➔ Non-parametric data-driven background modelling:

1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
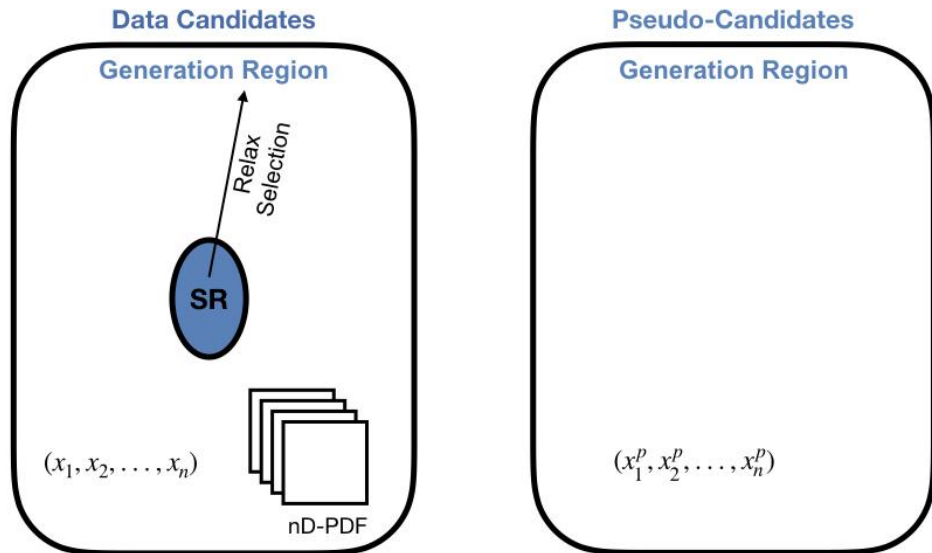2. Obtain conditional PDF of relevant variables ($x_1$, $x_2$,…, $x_n$)
3. Generate sample of pseudo-candidates
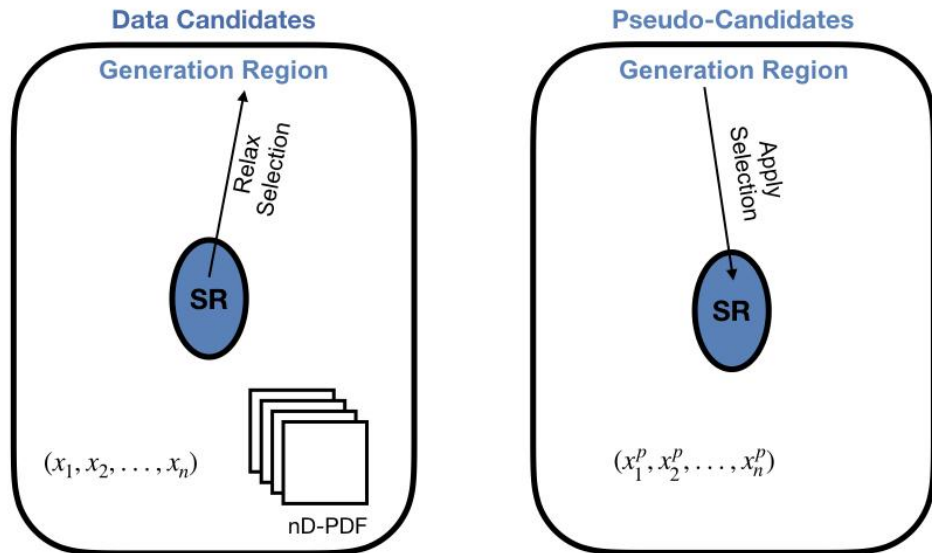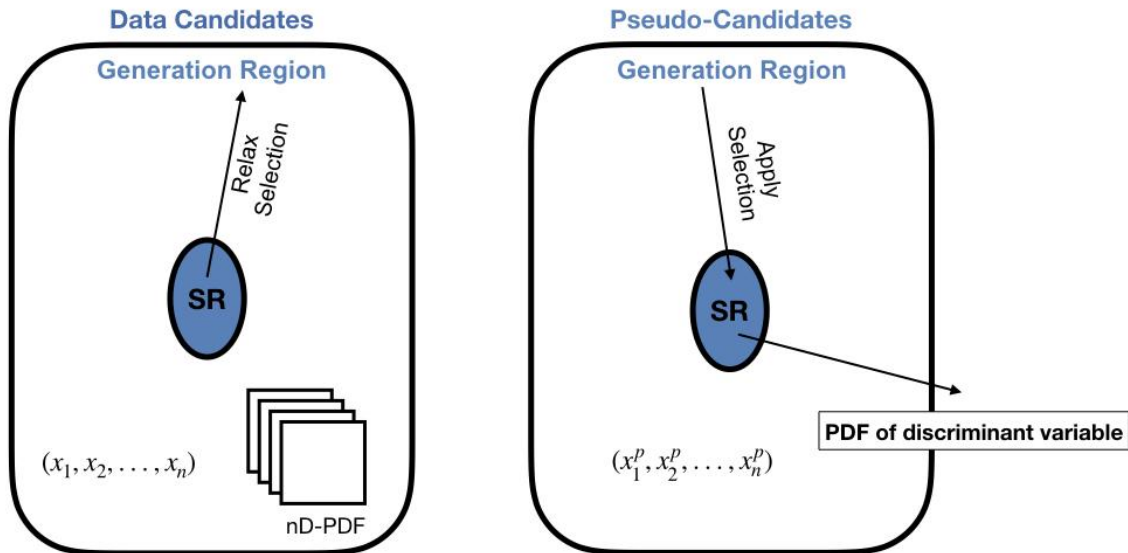4. Apply **Signal Region** requirements to pseudo-candidates sample

# Background Modelling

➔ Non-parametric data-driven background modelling:

1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
2. Obtain conditional PDF of relevant variables $(x_1, x_2, \ldots, x_n)$
3. Generate sample of pseudo-candidates
4. Apply **Signal Region** requirements to pseudo-candidates sample - obtain PDF of **discriminant variable** for statistical analysis

**Data Candidates**

Generation Region

Relax Selection

**SR**

$(x_1, x_2, \ldots, x_n)$

nD-PDF

**Pseudo-Candidates**

Generation Region

Apply Selection

**SR**

$(x_1^p, x_2^p, \ldots, x_n^p)$

PDF of discriminant variable

# Background Modelling

➜ Non-parametric data-driven background modelling:

1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
2. Obtain conditional PDF of relevant variables ($x_1$, $x_2$,…, $x_n$)
3. Generate sample of pseudo-candidates
4. Apply **Signal Region** requirements to pseudo-candidates sample - obtain PDF of **discriminant variable** for statistical analysis
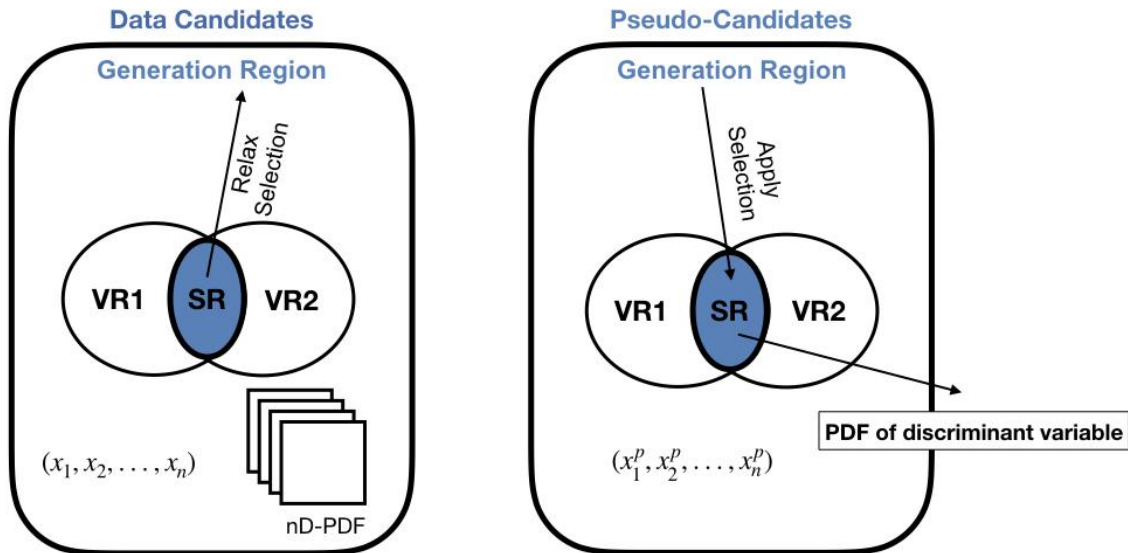   - Intermediate **Validation Regions** to check method

# Background Modelling

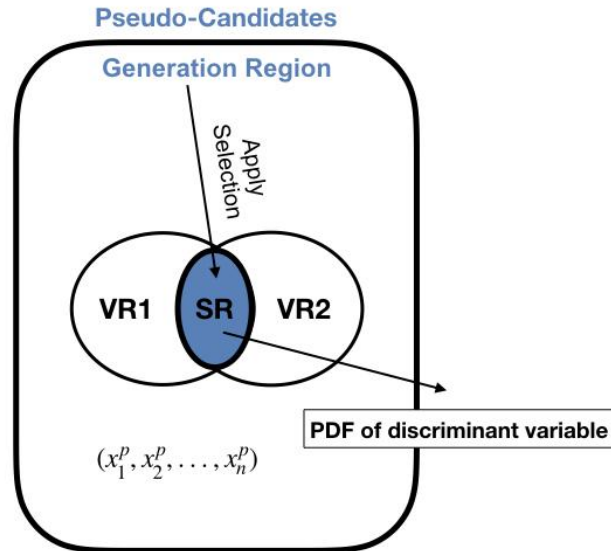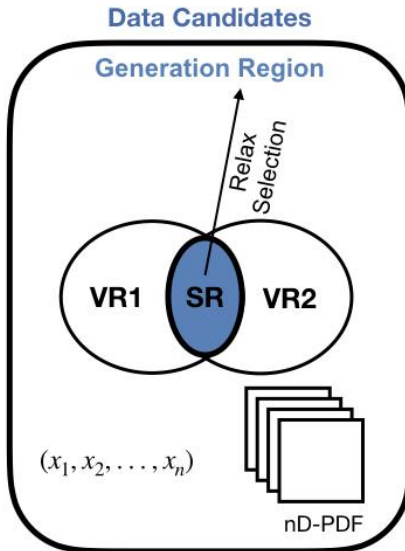➔ Non-parametric data-driven background modelling:

1. Obtain sample of data events enriched in background by relaxing event selection requirements (**Generation Region**)
2. Obtain conditional PDF of relevant variables ($x_1$, $x_2$,…, $x_n$)
3. Generate sample of pseudo-candidates
4. Apply **Signal Region** requirements to pseudo-candidates sample - obtain PDF of **discriminant variable** for statistical analysis
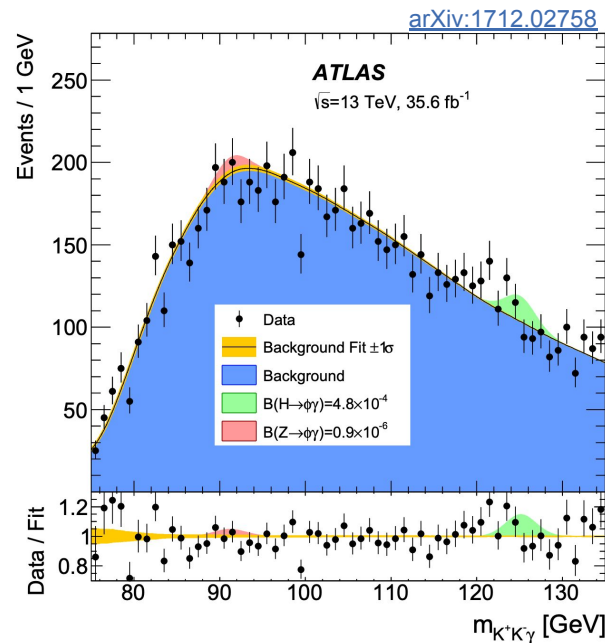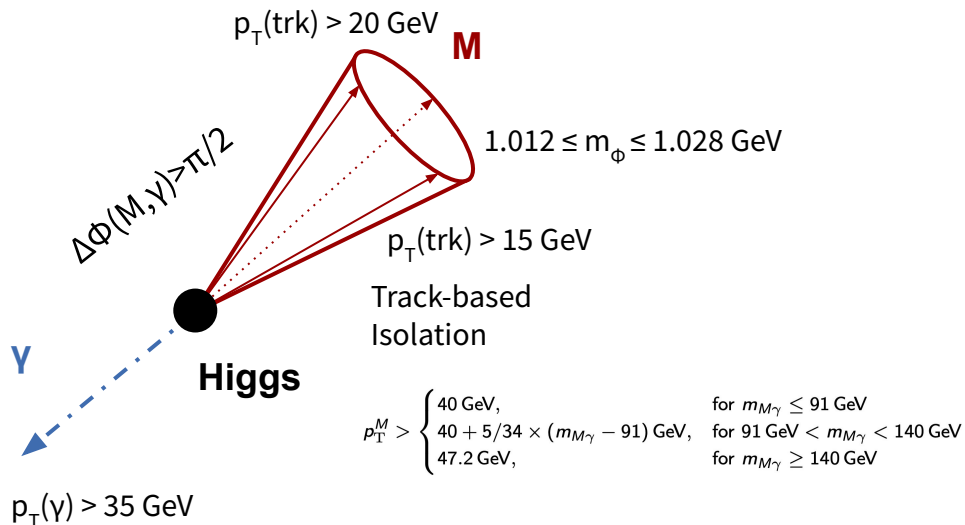   - Intermediate **Validation Regions** to check method

**Today will present 2 implementations of the method**

UNIVERSITY OF BIRMINGHAM

# Ancestral Sampling

# Case Study: H→Φγ

➔ H→φ(K⁺K⁻)γ has potential to probe **Higgs coupling to strange quark**

   ◆ **Distinct experimental signature**: pair of collimated high-$p_T$ isolated tracks recoiling against isolated photon

   ◆ Main background : **photon + jet** and **dijet**

     ● difficult to model accurately using MC - ideal use case for method

     ● **photon + jet MC sample** used to exemplify model application

**ATLAS**
$\sqrt{s}$=13 TeV, 35.6 fb$^{-1}$

Events / 1 GeV

Data
Background Fit ±1σ
Background
B(H→φγ)=4.8×10$^{-4}$
B(Z→φγ)=0.9×10$^{-6}$

Data / Fit

$m_{K^+K^-\gamma}$ [GeV]

$p_T$(trk) > 20 GeV  **M**

$1.012 \leq m_\phi \leq 1.028$ GeV

$\Delta\Phi(M,\gamma)>\pi/2$

$p_T$(trk) > 15 GeV

Track-based Isolation

**γ**  **Higgs**

$$p_T^M > \begin{cases} 40\,\text{GeV}, & \text{for } m_{M\gamma} \leq 91\,\text{GeV} \\ 40 + 5/34 \times (m_{M\gamma} - 91)\,\text{GeV}, & \text{for } 91\,\text{GeV} < m_{M\gamma} < 140\,\text{GeV} \\ 47.2\,\text{GeV}, & \text{for } m_{M\gamma} \geq 140\,\text{GeV} \end{cases}$$

$p_T$(γ) > 35 GeV

# Building the model for H→Φγ

1. Relax $p_T(M)$ and Iso(M) requirements

| Region | $p_T(M)$ cut | Iso(M) cut |
|--------|--------------|------------|
| **GR** | x | x |
| **VR1** | ✓ | x |
| **VR2** | x | ✓ |
| **SR** | ✓ | ✓ |

UNIVERSITY OF BIRMINGHAM
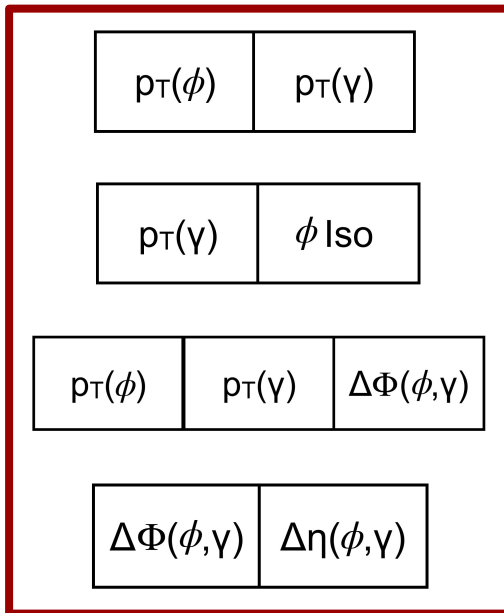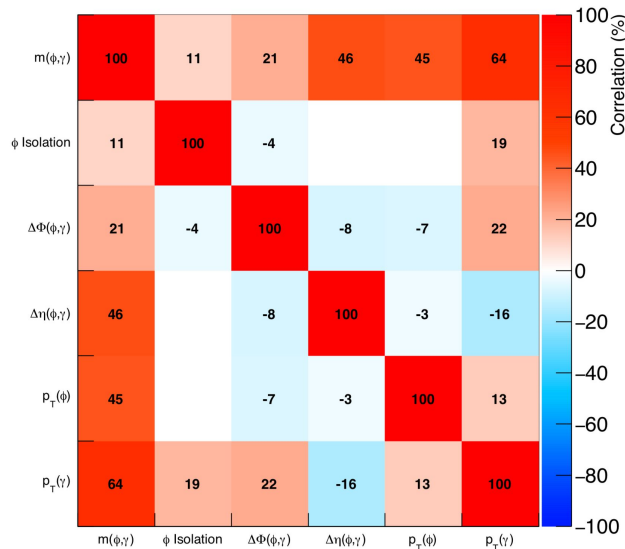
# Building the model for H→Φγ

**Relax Selection**

**Obtain Conditional PDFs**

**Generate pseudo candidates**

**Apply Selection**

2. Build PDFs of relevant **kinematic** and **isolation** variables in generation region - we need the φ and γ 4-momentum vectors to ultimately obtain **m(φγ)** + Iso(φ)

♦ 1D, 2D and 3D **histograms** to be sampled from in generation step

♦ only most important correlations are explicitly described



| $p_T(\phi)$ | $p_T(\gamma)$ |

| $p_T(\gamma)$ | $\phi$ Iso |

| $p_T(\phi)$ | $p_T(\gamma)$ | $\Delta\Phi(\phi,\gamma)$ |

| $\Delta\Phi(\phi,\gamma)$ | $\Delta\eta(\phi,\gamma)$ |

UNIVERSITY OF BIRMINGHAM

# Building the model for H→Φγ

# Building the model for H→Φγ

<table>
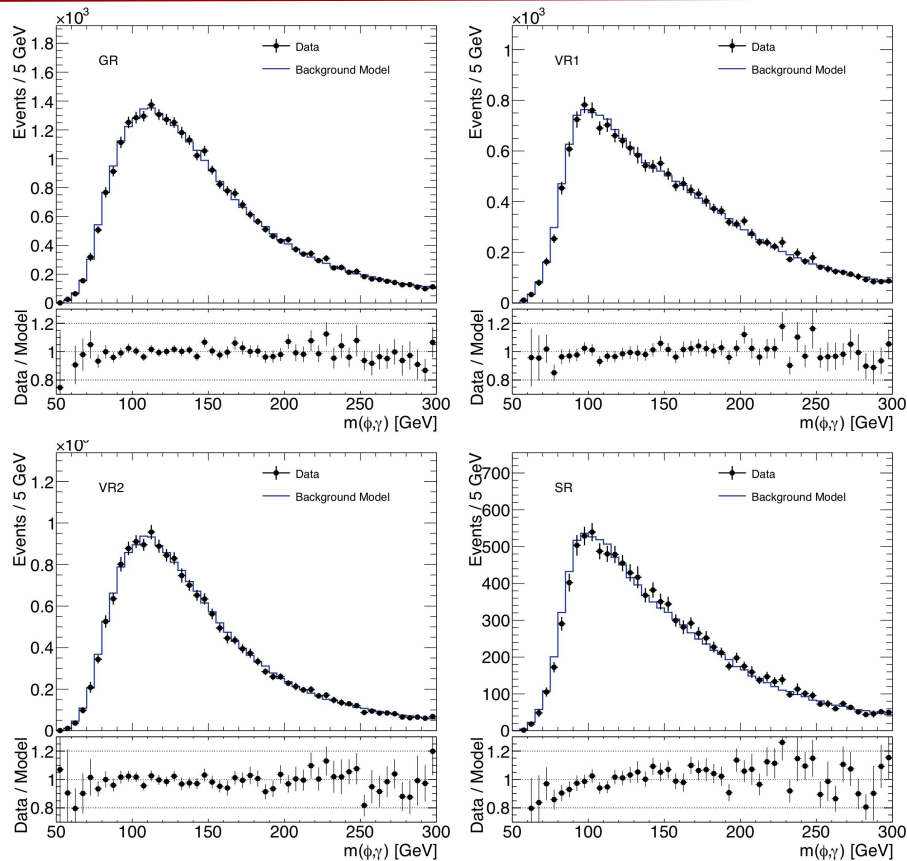<tr><td>Relax Selection</td></tr>
<tr><td>Obtain Conditional PDFs</td></tr>
<tr><td>Generate pseudo candidates</td></tr>
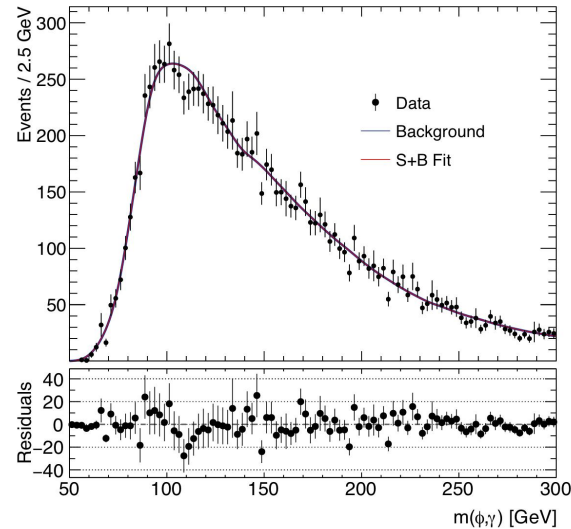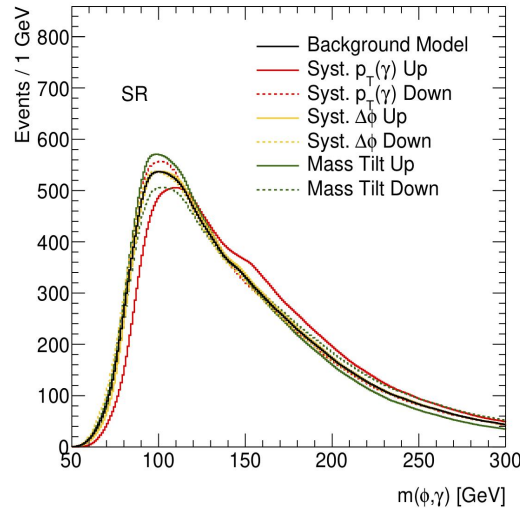<tr><td>**Apply Selection**</td></tr>
</table>

4. Apply $p_T(M)$ and Iso(M) requirements to sample of pseudo-candidates
   - obtain PDF of **m(φγ)** for statistical analysis in Signal and Validation Regions

# Implementation in Statistical Analysis

➔ **Systematic uncertainties** are provided through variations of the nominal PDFs

　◆　selected to capture different modes of potential deformations of the background shape

➔ Maximum likelihood fit to invariant mass

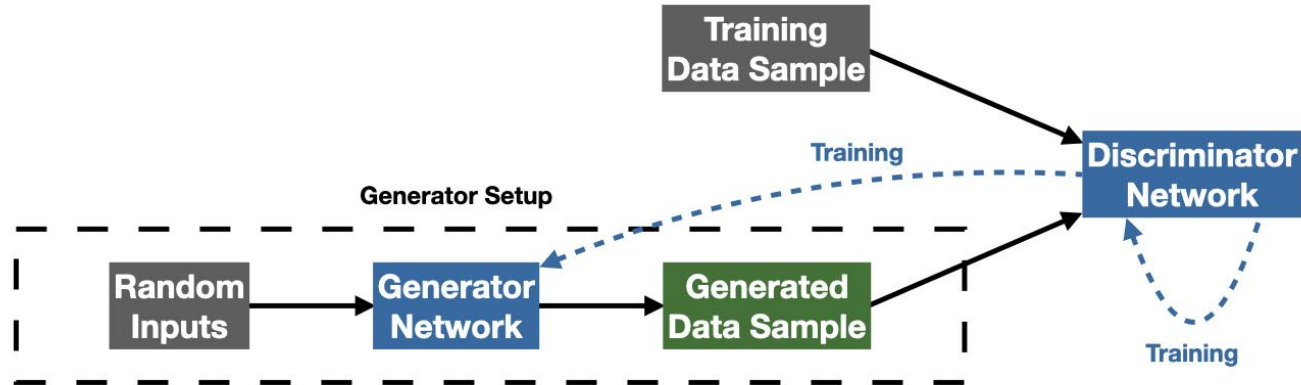　◆　each variation controlled by a nuisance parameter - directly constrained by data in fit



| Parameter | Value | Uncertainty ($\pm 1\sigma$) |
|---|---|---|
| $\mu_{\text{signal}}$ | $-0.03$ | $\pm 0.55$ |
| $\mu_{\text{bkgd}}$ | $1.01$ | $\pm 0.01$ |
| Shape: $p_{\text{T}}(\gamma)$ shift | $0.27$ | $\pm 0.15$ |
| Shape: $\Delta\Phi(\phi,\gamma)$ tilt | $0.26$ | $\pm 0.43$ |
| Shape: $m(\phi,\gamma)$ tilt | $0.11$ | $\pm 0.24$ |

UNIVERSITY OF BIRMINGHAM

# Conditional Generative Adversarial Networks

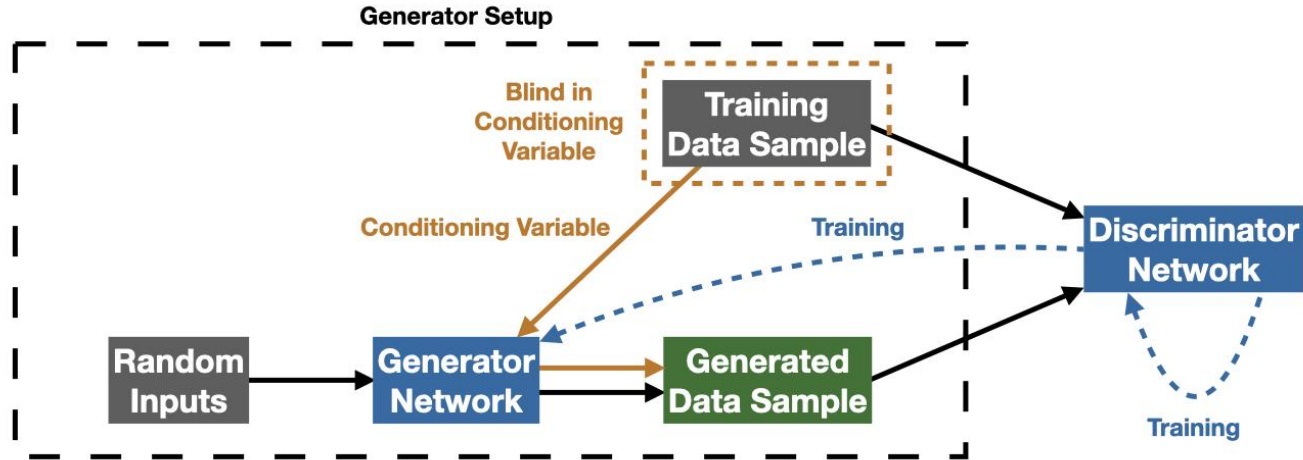# Generative Adversarial Networks

➔ Challenge for ancestral sampling:
  ◆ application in multivariate analyses
  ◆ signal region blinding
➔ Generalisation of method: use **GANs trained on data** to produce background model
  ◆ **Generator** - learns generative model from data sample
  ◆ **Discriminator** - simultaneously trained to discriminate the generator output from data

# Conditional Generative Adversarial Networks

➜ Possible **signal contamination** in training data:
  ◆ **Condition** GAN (cGAN) on a blinding variable, allowing **SR to be blinded during training** - cGAN extrapolates prediction into SR
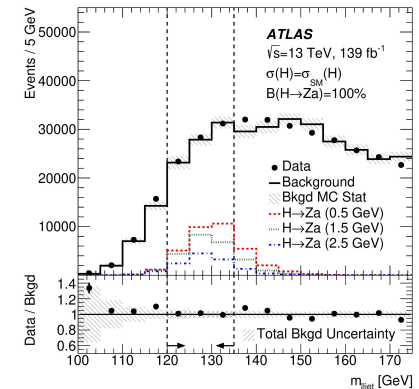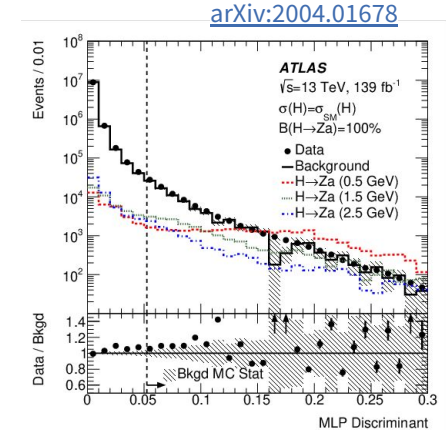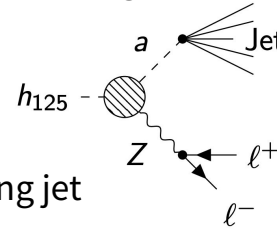
# Case Study: H→Za

➔ Light pseudo-scalars produced in Higgs decays feature in BSM theories like two-Higgs-doublet model and the 2HDM with additional scalar singlet

➔ Search for **H→Z(ll)+a**, with a→hadrons:
  ◆ Main background: **Z + jets**
  ◆ background discrimination relies on **MVA** techniques, using jet substructure variables
  ◆ ideal case study for implementation of background modelling using cGANs
    ● systematics arising by limited stats in MC simulation (arXiv:2004.01678)
    ● use of MVA techniques makes it impractical to use ancestral sampling

➔ **Z + jets MC sample** used to exemplify model application

arXiv:2004.01678

UNIVERSITY OF BIRMINGHAM

# Building the model for H→Za

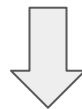| |
|---|
| **Relax Selection** |
| Obtain Conditional PDFs |
| Generate pseudo candidates |
| Apply Selection |

1. Remove MLP-based selection
   ◆ **& blind signal region** to avoid signal contamination

   **Use $m_{\mu\mu j}$ as blinding variable**

   ⬇

   123 GeV ≤ $m_{\mu\mu j}$ ≤135 GeV blinded

# Building the model for H→Za

<table>
<tr><td>
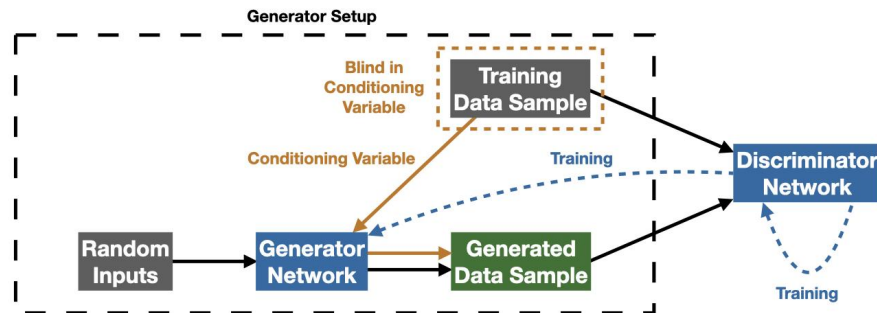
Relax
Selection

**Obtain
Conditional
PDFs**

Generate
pseudo
candidates

Apply
Selection

</td><td>

2. cGans trained using **blinded data**

◆ learn generative model of the conditional probability distribution of the data, given value of blinding variable

◆ Use **ensemble** of cGANs and take average:

● 100 cGANs trained, 5 best based on $\chi^2$ metric kept for analysis

Generator and discriminator:
● 5 layers x 256 hidden nodes with leaky ReLU activation function
● binary cross entropy loss function and L2 regularisation



</td></tr>
</table>

# Building the model for H→Za

<table>
<tr><td>Relax Selection</td></tr>
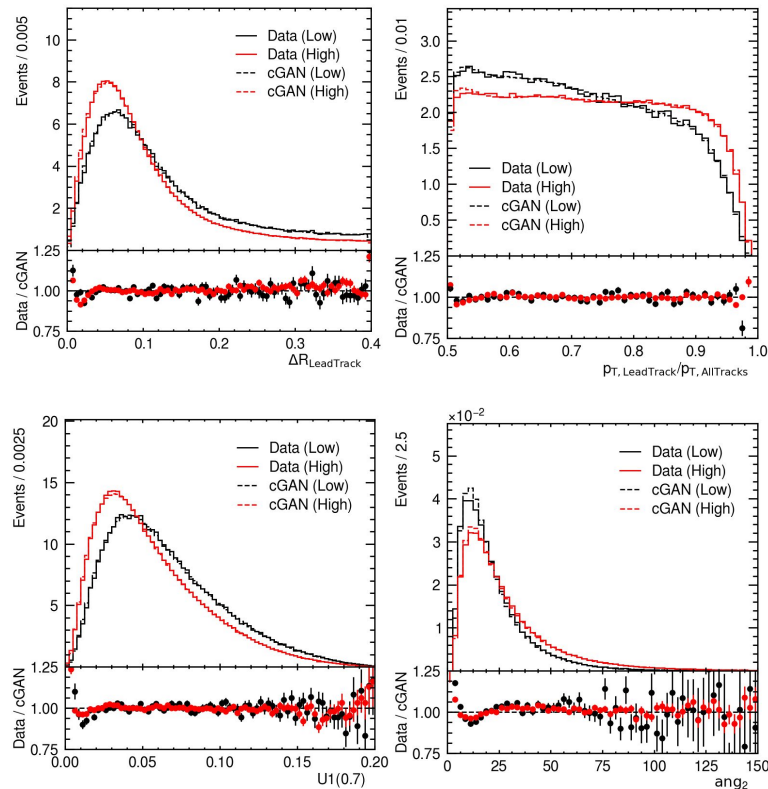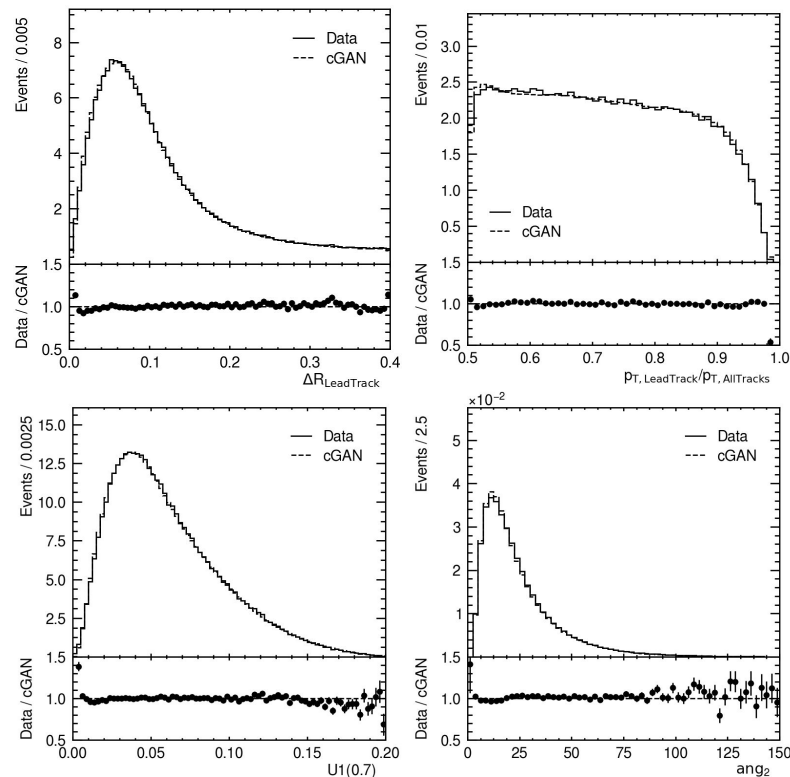<tr><td>Obtain Conditional PDFs</td></tr>
<tr><td>**Generate pseudo candidates**</td></tr>
<tr><td>Apply Selection</td></tr>
</table>

3. Input inclusive distribution of the conditioning variable into cGAN:
   - ◆ cGAN **extrapolates** the conditional generative model into signal region
   - ◆ obtain prediction of **MLP input variables**

$m_{\mu\mu j}$ **sidebands**

# Building the model for H→Za

| |
|---|
| Relax Selection |
| Obtain Conditional PDFs |
| **Generate pseudo candidates** |
| Apply Selection |

3. Input inclusive distribution of the conditioning variable into cGAN:
   - ◆ cGAN **extrapolates** the conditional generative model into signal region
   - ◆ obtain prediction of **MLP input variables**



m$_{\mu\mu j}$ SR

# Building the model for H→Za

<table>
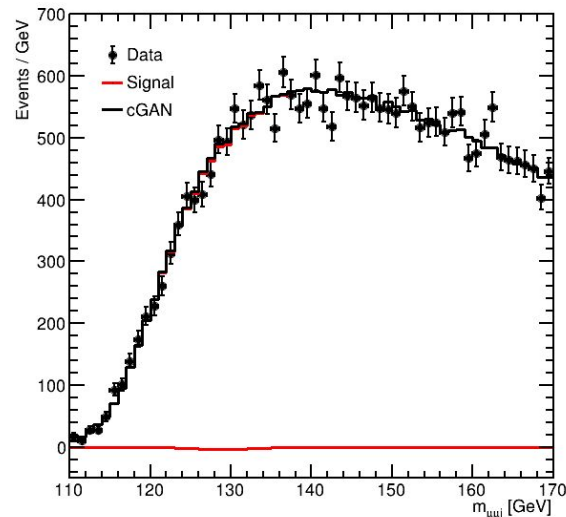<tr><td>Relax Selection</td></tr>
<tr><td>Obtain Conditional PDFs</td></tr>
<tr><td>Generate pseudo candidates</td></tr>
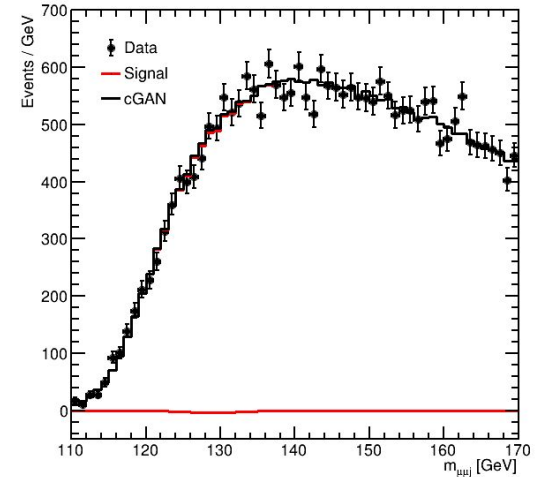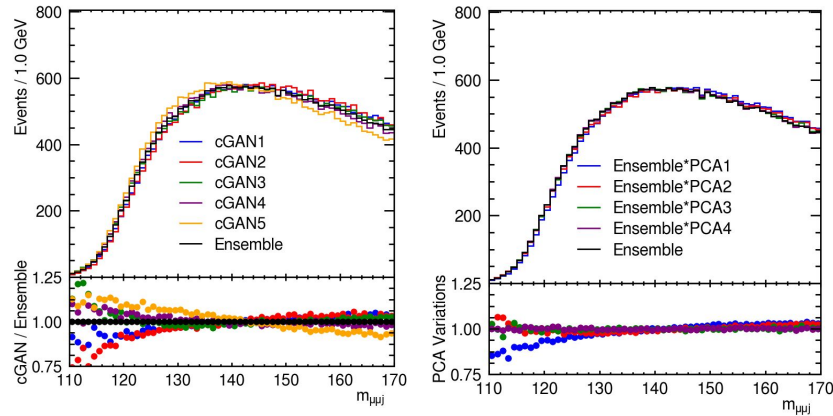<tr><td>**Apply Selection**</td></tr>
</table>

4. Apply MLP selection to pseudo-candidates sample
   ◆ obtain PDF of $m_{\mu\mu j}$ in SR for statistical analysis

UNIVERSITY OF BIRMINGHAM

# Implementation in Statistical Analysis

➔ **Systematic uncertainties** are provided through shape variations:
  ◆ Differences between ensemble and individual cGANs
  ◆ Principal component analysis performed to orthogonalise differences
➔ Maximum likelihood fit to Higgs invariant mass
  ◆ each variation controlled by a nuisance parameter - directly constrained by data in fit



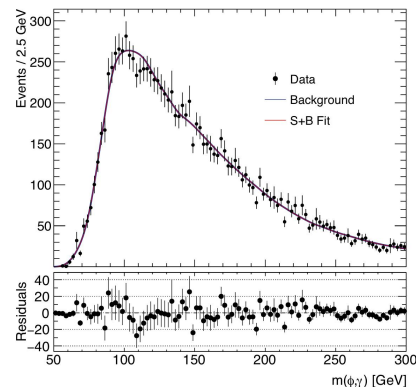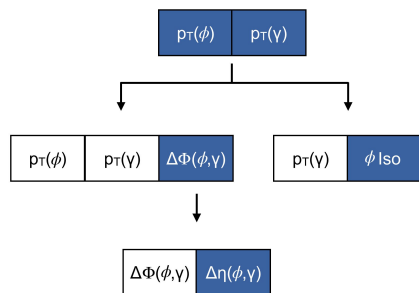| Parameter | Value | Uncertainty ($\pm 1\sigma$) |
|---|---|---|
| $\mu_{\text{signal}}$ | $-0.003$ | $\pm 0.010$ |
| $\mu_{\text{bkgd}}$ | $1.001$ | $\pm 0.008$ |
| Shape uncertainty 1 | $-0.36$ | $\pm 0.27$ |
| Shape uncertainty 2 | $-0.31$ | $\pm 0.52$ |

UNIVERSITY OF BIRMINGHAM

# Summary

➔ A novel **non-parametric**, **data-driven** background modelling technique was presented
  ◆ Addresses typical shortcomings of often employed background modelling techniques
  ◆ Dataset from a **relaxed event selection** to create a model based on **conditional probabilities**
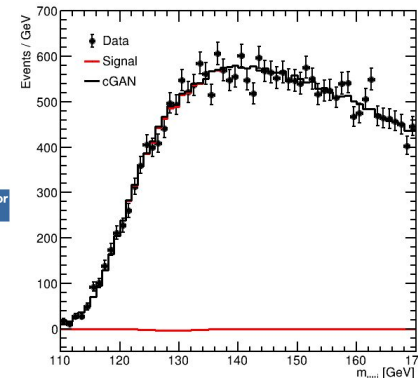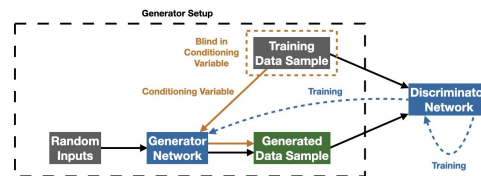  ◆ Two distinct ways of building the conditional PDF:

**Ancestral sampling**

**Conditional Generative Adversarial Networks**

- Sample from histograms of relevant variables in data, built with respect to most important correlations
- Already used in multiple analysis! [Phys. Rev. Lett. 114 (2015) 121801, Phys. Rev. Lett. 117, 111802 (2016), JHEP 07 (2018) 127, Phys. Lett. B 786 (2018) 134]

- Generalisation of ancestral sampling
- Use GANs trained on data to produce background model
- **Condition** GAN (cGAN) on a blinding variable, allowing **SR to be blinded during training**
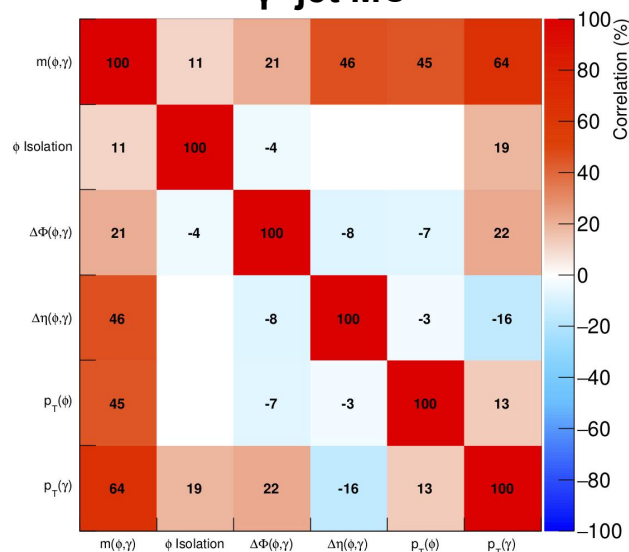
# BACK-UP

# Building the model for H→Φγ

# Ancestral Sampling

➔ Background model is **robust under signal contamination**

◆ Signal Injection tests to evaluate robustness