

# **Data driven background estimation in HEP** using Generative Adversarial Networks

#### CHEP 2023 - May 9th, 2023

Victor Lohezic (victor.lohezic@cern.ch) Fabrice Couderc, Julie Malclès, Özgür Şahin **IRFU - CEA Saclay** 









Eur. Phys. J. C 83, 256 (2023) https://doi.org/10.1140/epjc/ s10052-023-11347-8



## Introduction

GAN based data-driven technique to estimate background processes with a misidentified object in collider events. We will showcase this technique for the  $\gamma$  + Jets background process of the H  $\rightarrow \gamma \gamma$ analysis.

In the H  $\rightarrow \gamma \gamma$  analysis, dominant backgrounds are :  $\gamma \gamma$  + Jets,  $\gamma$  + Jets, Multi Jets (MJ)



- What if we use data directly to describe those samples ?
  - analyses using this technique.



• The agreement between Data and Monte Carlo (MC) simulated samples for  $\gamma$  + Jets and MJ is not satisfying and the statistics is too low for the training of subsequent discriminants.

We would like to improve the data driven approach used in the previously published





# Overview

### **I.** A data driven estimation of the background

### I. Training a GAN

- a. Generative Adversarial Network (GAN)
- b. Evaluation procedure

### **III.** Generating a full object (misidentified photon)

- a. Optimization of training
- Results b.

## **IV. Conclusions and outlooks**

# I. A data driven estimation of the background

In an event each photon is given a score (photon ID) representing its likelihood to be a photon. Control region in data based on photon ID is used to replace MC  $\gamma$  + Jets / MJ samples (better agreement, more statistics).



#### Need one photon with very low photon ID : probably a misidentified photon $\gamma$ (as opposed to a prompt photon $\gamma$

Many analyses already use data driven background estimation. By reverting the cut on the min photon ID, one needs either to get rid of the photon ID variable or to generate a new min photon ID !

This procedure was used in published analysis from CMS experiment [1], new ID was generated by :

- 1. Deriving a 1D probability density function (PDF) from the misidentified photon ID distribution
- 2. Generating a random min photonID following this PDF, in the signal region but below the max photonID
- 3. However correlations are not preserved

We propose a new method to generate a suitable photon (not only ID) taking into account these correlations thanks to ML and more specifically GAN (Generative Adversarial Networks)

[1] Measurements of ttH production and the CP structure of the Yukawa interaction between the Higgs boson and the top quark in the diphoton decay channel, CMS collaboration

#### CHEP 2023 - V. Lohezic



## I. A data driven estimation of the background

In an event each photon is given a score (photon ID) representing its likelihood to be a photon. Control region in data based on photon ID is used to replace MC  $\gamma$  + Jets / MJ samples (better agreement, more statistics).

Need one photon with very low photon ID : probably a misidentified photon  $\gamma$  (as opposed to a prompt photon  $\gamma$ 



Many analyses already use data driven background estimation. By reverting the cut on the min photon ID, one needs either to get rid of the photon ID variable or to generate a new min photon ID !

This procedure was used in published analysis from CMS experiment [1], new ID was generated by :

- 1. Deriving a 1D probability density function (PDF) from the misidentified photon ID distribution
- 2. Generating a random min photonID following this PDF, in the signal region but below the max photonID
- 3. However correlations are not preserved

We propose a new method to generate a suitable photon (not only ID) taking into account these correlations thanks to ML and more specifically GAN (Generative Adversarial Networks)

[1] Measurements of ttH production and the CP structure of the Yukawa interaction between the Higgs boson and the top quark in the diphoton decay channel, CMS collaboration

CHEP 2023 - V. Lohezic







# Overview

#### **A data driven estimation of the background**

### I. Training a GAN

- a. Generative Adversarial Network (GAN)
- **Evaluation procedure** b.

## **III.** Generating a full object (misidentified photon)

- a. Optimization of training
- b. Results

**V. Conclusions and outlooks** 

## II. Training a GAN II.a - Generative Adversarial Networks (GANs)

Would it be possible to create an algorithm capable of learning underlying correlations and capable of generating a sample statistically independent from the training sample?

Goodfellow et al. suggested a model consisting of two neural networks competing against each other : 

- the "discriminator" sorts samples between real and generated ones *i.e.* discriminates fakes
- the "generator" tries to produce samples which will fool the discriminator lacksquare











 Usually, monitoring the loss of a neural network is enough to evaluate its performance. It is not the case for GAN where both networks need to perform well against the other so their loss stays flat.

We need to set up a more elaborate evaluation procedure







- Usually, monitoring the loss of a neural network is enough to evaluate its well against the other so their loss stays flat.
- We need to set up a more elaborate evaluation procedure



performance. It is not the case for GAN where both networks need to perform



#### **II.b - Evaluation procedure**

training epoch on the training sample and on a validation sample :







• To evaluate the performance of a given model, we rely on different metrics computed for each



 $n_k$ : sum of the weights of **GANed events** in bin k N<sub>k</sub> : sum of the weights of **original events** in bin k

For the NLL we histogram our events in 4D :

- $\mathchar`-$  transverse momentum of misidentified photon  $p_{T_{\nu}}$
- pseudorapidity of misidentified photon  $\eta_{\gamma}$
- $p_T$  of diphoton pair over its mass  $\frac{PT_{\gamma\gamma}}{M}$  $m_{\gamma\gamma}$
- ID of misidentified photon  $\mathrm{ID}_\gamma$

Takes into account correlations by construction











II. Training a GAN

Metrics are fluctuating a lot !

fluctuations can be reduced by increasing the number of generation per event :



- a closer look at its performance



# $p_k$ (see slide 9) estimation is statistically limited creating fluctuations in the NLL. These

Seeing how the fluctuations decrease, we decide to go to 100 generation per event

Then we can find epochs where the model is reaching minima for these metrics and take

# Overview

### A data driven estimation of the background

### I. Training a GAN

a. Generative Adversarial Network (GAN)b. Evaluation procedure

## III. Generating a full object (misidentified photon)

- a. Optimization of training
- b. Results

V. Conclusions and outlooks

## III. Generating a full object (misidentified photon) III.a - Optimization of training

 $\bullet$ learn correlations





CHEP 2023 - V. Lohezic

#### III. Generating a full object (misidentified photon) III.a - Optimization of training -0.070

 $\bullet$ learn correlations





CHEP 2023 - V. Lohezic



- could help
- quantile transformation



 $\Rightarrow$  Transformation helps the GAN recover the gaps in  $\eta$  and the core of the ID and  $p_T$  distributions



#### **Example from scikit-learn's documentation**

#### **III.** Generating a full object

### II.b - Results

- GAN is able to generate a full misidentified object that would bass the selection criteria (see 1D distributions on diagonal)
- GAN learns correlations bet observables of the objects (see contours on off-diagonal plots also correlations with the rest event (see distance correlation coefficients matrix)
- This method could be used as a general tool to generate other objects for other use cases



-0.2

[GeV]

 $\mu_{\mathcal{A}}$ 



#### **III.** Generating a full object





#### II.b - Results

GAN is able to generate a full misidentified object that would pass the selection criteria (see 1D distributions on diagonal)

• GAN learns correlations between observables of the objects (see contours on off-diagonal plots) but also correlations with the rest of the event (see distance correlation coefficients matrix)

This method could be used as a general tool to generate other objects for other use cases



**III.** Generating a full object

#### **Original correlations**



**GANed correlations** 



# **IV. Conclusions and outlooks**

- We developed an evaluation procedure to test the GAN's generator performance and pick the best performing one
- Thanks to GAN we can generate a misidentified photon mimicking the behaviour of an object passing the photon selection criteria
- The sample produced for this showcase could be used for any H  $\rightarrow \gamma\gamma$  analysis
- This method can be used as a general tool to generate other objects for analyses dealing with background coming from misidentified objects







