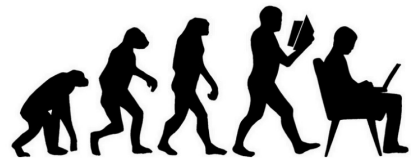# The A(i)DAPT program
## AI for Data Analysis and Preservation

Tommaso Vittorini

*on behalf of A(i)DAPT Working Group*



**A(i)DAPT**

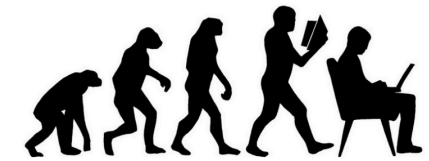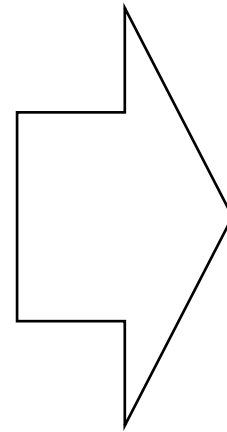**AI for Data Analysis and PreservaTion**

- Motivation and advantages of the deployed techniques

- Generative Adversarial Network overview

- Our approach towards reproducing experimental data

- Work in progress

- Conclusions

- Data collected by NP/HEP experiments are (always) affected by the detector's effects
- Before starting physics analysis the detector's effect unfolding is required
- Traditional observables may not be adequate to extract physics in multidimensional space (multi-particles in the final state)
- At High-Intensity frontiers, data sets are large and difficult to manipulate/preserve

**Should AI support NP/HEP experiments to extract physics from data in more efficient way?**



**A(i)DAPT**
**AI for Data Analysis and PreservaTion**

**Develop AI – supported procedures to:**
- Accurately fit data in multiD space
- Unfold detector effects
- Compare synthetic (AI-generated) to experimental data
- Quantify the uncertainty (UQ)

**Collaborative effort (regular meeting)**
- ML experts (ODU, Jlab)
- Experimentalists (Jlab Hall-B)
- Theorists (JPAC, JAM)

# Detector unfolding

- Detector effects make measured observables (detector-level) different from the 'true' observables (vertex level)

    **Acceptance:** Any measurement can access only a limited portion of the phase space. What can we say about these unmeasured regions?

    ➢ Interpolation: deal with the holes in the phase space
    ➢ Extrapolation: extend our coverage from the borders of measured regions

    **Resolution:** Any measurement has an experimental resolution that may modify cover up effects that we're looking for

    ➢ Spikes may be concealed behind the detector resolution
    ➢ Measurements could be extended to unphysical regions

- Mitigation strategy:
    ➢ Acceptance: 'Fiducial volumes' to exclude unmeasured regions and extend the covered measured of the phase space
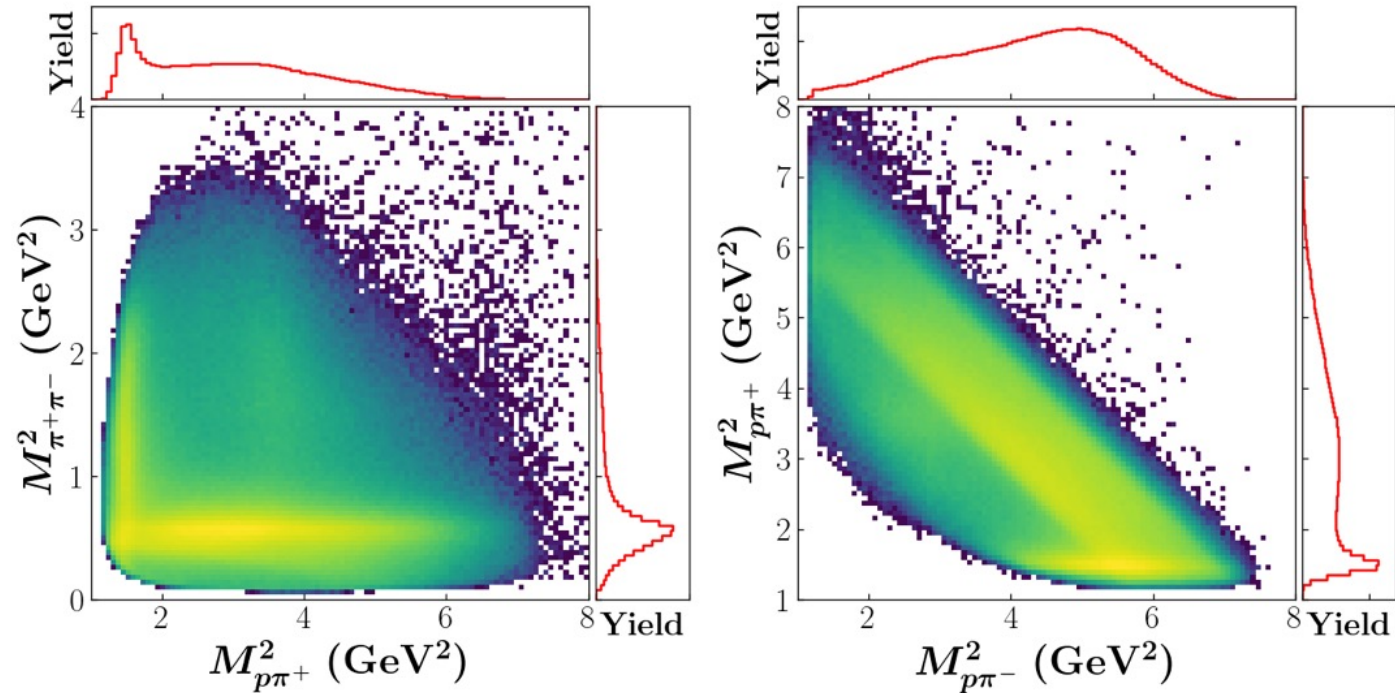    ➢ Resolution: build and validate ML-models to unfold resolution effects

$\gamma p \rightarrow \pi^+\pi^- p$ **(unpolarized)**

- Initial state: Fully known
- Final state: 3x3 indipendent variables
- Indipendent variables: (3x3) − 4 = 5 ($E_\gamma$ fixed)
- Many possible choices, such as $M_{\pi\pi}^2$, $M_{p\pi}^2$ $\theta_\pi, \alpha, \phi$

CLAS g11 $2\pi$ photoproduction

- $E_\gamma = (3 - 3.8) \, GeV$
- Dataset analyses on $\gamma p \rightarrow p\pi^+(\pi^-)$ with small contamination from $\gamma p \rightarrow p\pi^+$ (more than a single missing $\pi^-$)
- Complicated dynamics due to the overlap of $(p\pi)$ to form $\Delta$ baryon resnoances and $(\pi\pi)$ to form meson resonances

$$\frac{d\sigma \, (\gamma \, p \rightarrow p \, \pi^+\pi^-)}{dM_{\pi\pi} \, dM_{p\pi} \, d\cos(\theta_\pi) \, d\alpha \, d\phi}$$



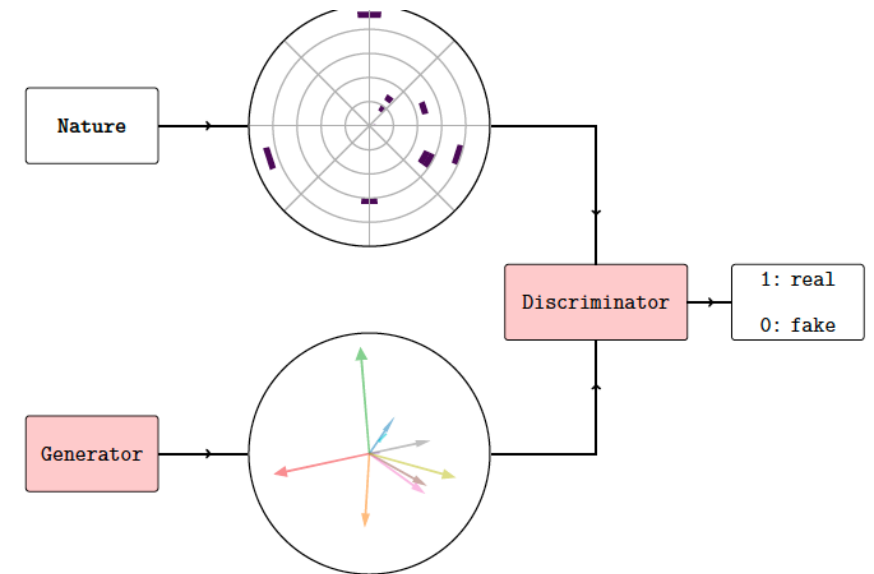AI could provide a new way to look at data and to extract observables and physics interpretation

Credit: Y.Alanazi Awadh, , P..Ambrozewicz, G. Costantini A.Hiller Blin, E. Isupov, T. Jeske, Y.Li, L.Marsicano W. Menlnitchouk, V.Mokeev, N.Sato, A.Szczepaniak, T.Viducic

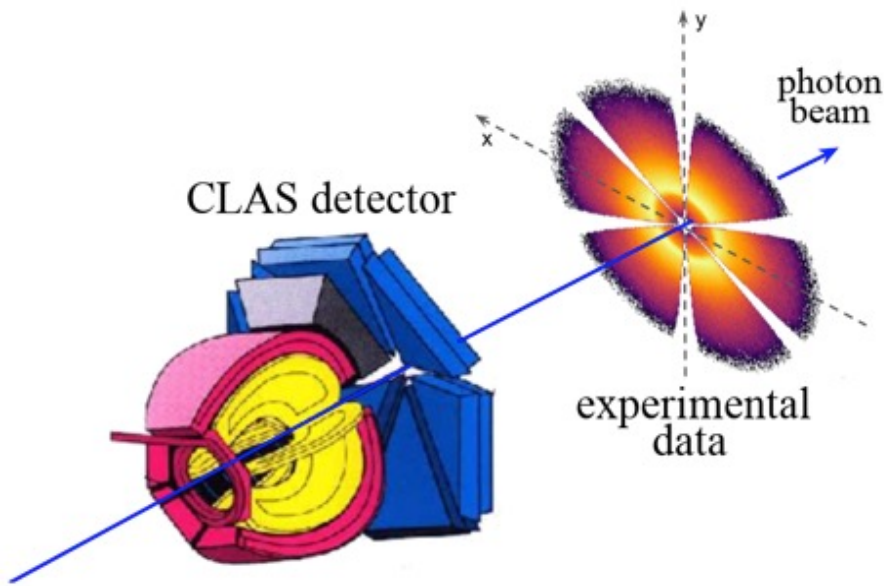# Generative Adversarial Networks (GANs)



- Generative model based on the competition between two Neural Networks: Generator vs Discriminator
  - **Generator** produces synthetic data which progressively reproduce realistic data and the **Discriminator** has to distinguish between synthetic and realistic data
  - **Generator** can be used to retain high dimensional correlations (detector proxies)
  - **Generator** can be used to provide highly realistic pseudo-data in an extremely fast way

CLAS g11 kinematics
- Dataset used by CLAS Collaboration for many publications
- Fiducial cuts $(p, \theta, \phi)$ as used in published analyses
- Focus on $\gamma p \to p\pi^+(\pi^-)$
- Final exclusive $2\pi$ state identified by missing mass technique (variables are reconstructed by energy/momentum conservation)
- Multi-pion background comes from $\gamma p \to p\omega^0 \to p\pi^+\pi^-\pi^0$
- At $E_\gamma = (3-4)$GeV reaction dynamics are dominated by $\rho^0$ photproduction through $\gamma p \to p\rho^0$ and $\Delta^{++}$ resonance excitation through $\gamma p \to \Delta^{++}\pi^-$
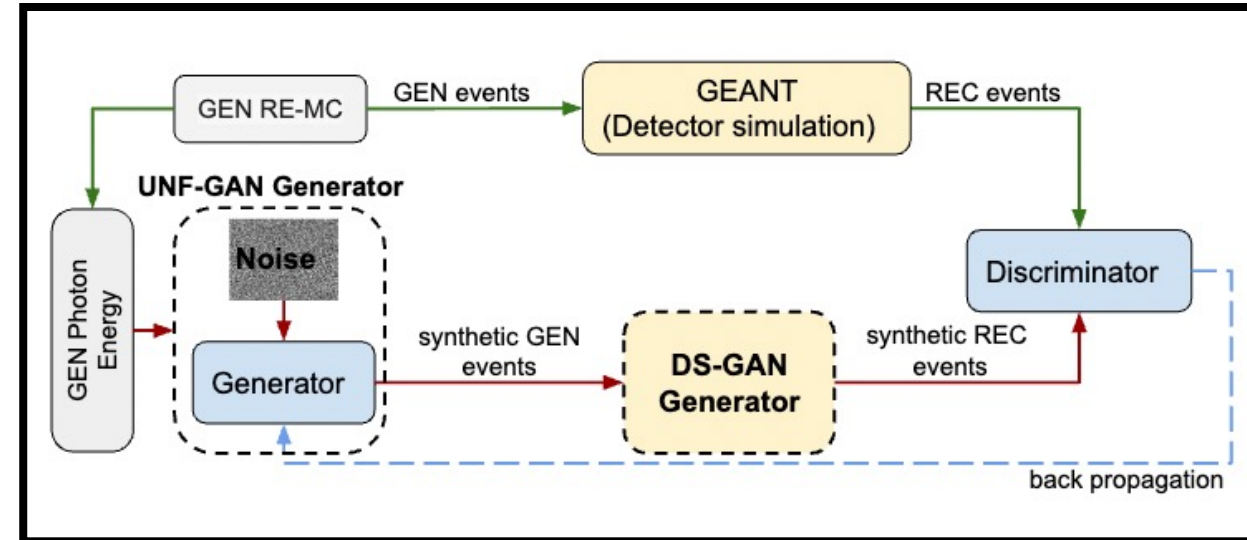
## $2\pi$ photoproduction closure test

- CLOSURE TEST:

    Demonstrate that GANs reproduce 'true' multi-d correlations, unfolding CLAS detector effects, comparing vertex-level (GEN) events with GAN GEN SYNT events, trained at detector-level and unfolded with a (GAN-based) detector proxy

1. Generate events with a (realistic) Monte Carlo $2\pi$ photoproduction model (RE-MC GEN pseudodata)

2. Apply detector effects (acceptance and resolution) via GSIM-GEANT (RE-MC REC pseudodata)

3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an indipendent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)

4. Deploy the unfolding GAN (UNF-GAN) that includes the DS-GAN, and train it with RE-MC REC pseudodata

5. Compare UNF-GAN GEN SYNT data to RE-MC GEN pseudodata

6. Replace RE-MC REC pseudo data with CLAS data in the training to unfold the vertex-level experimental distributions
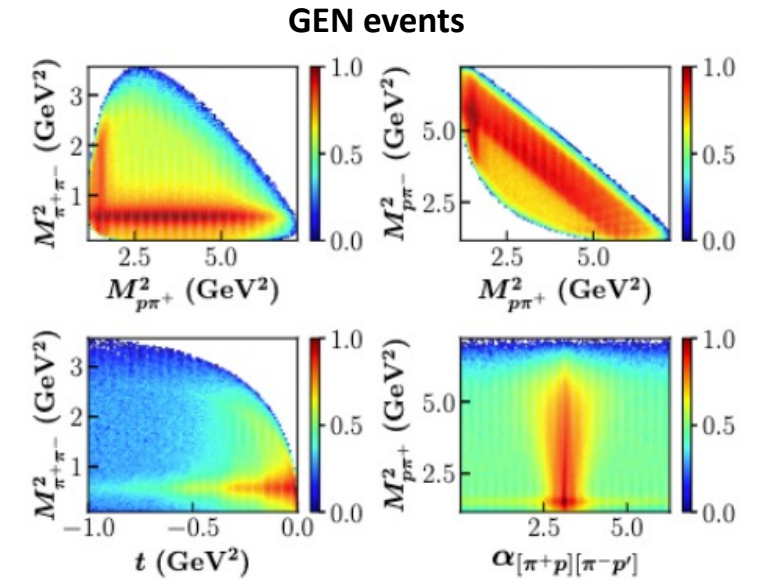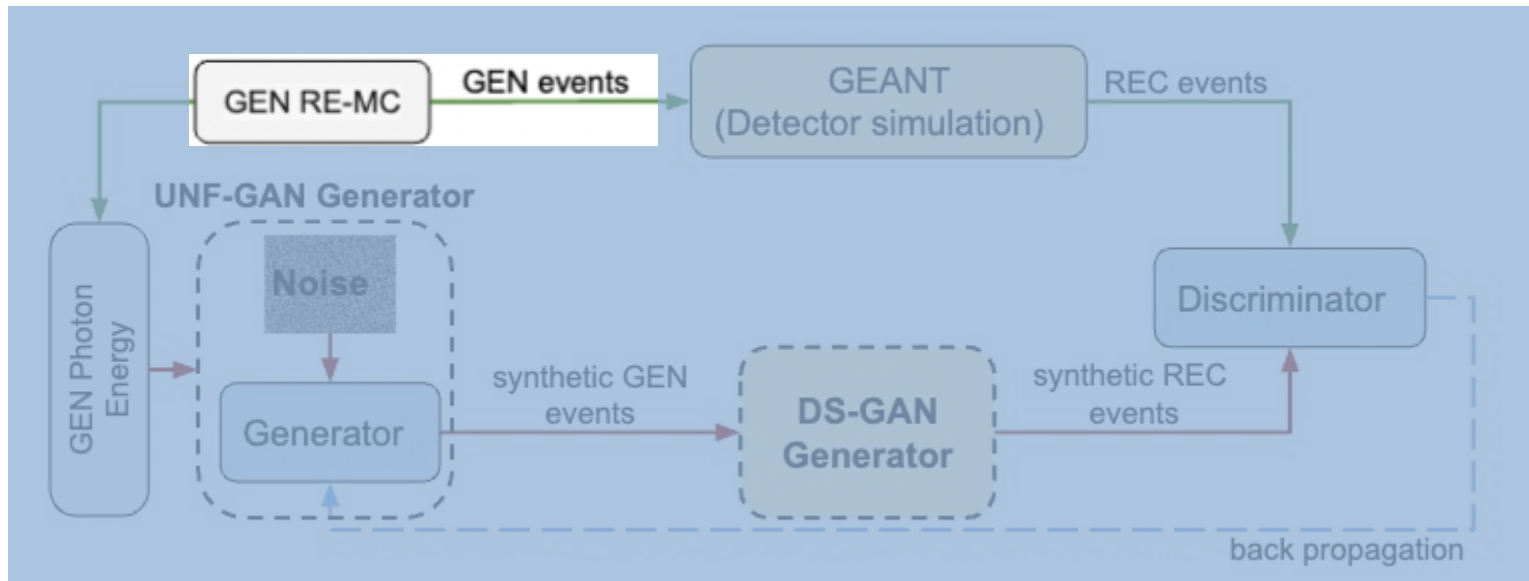


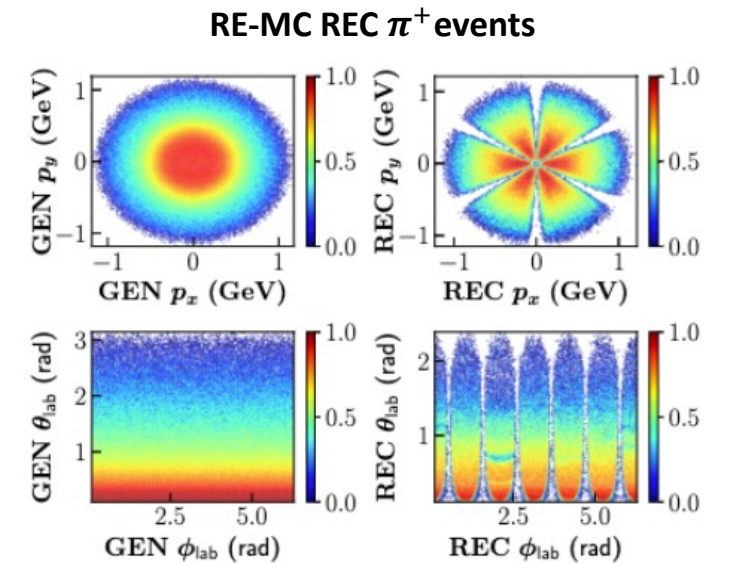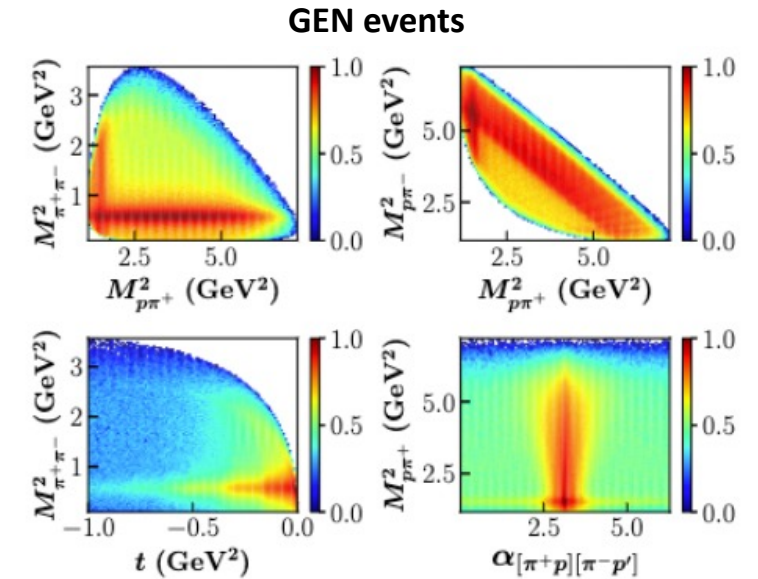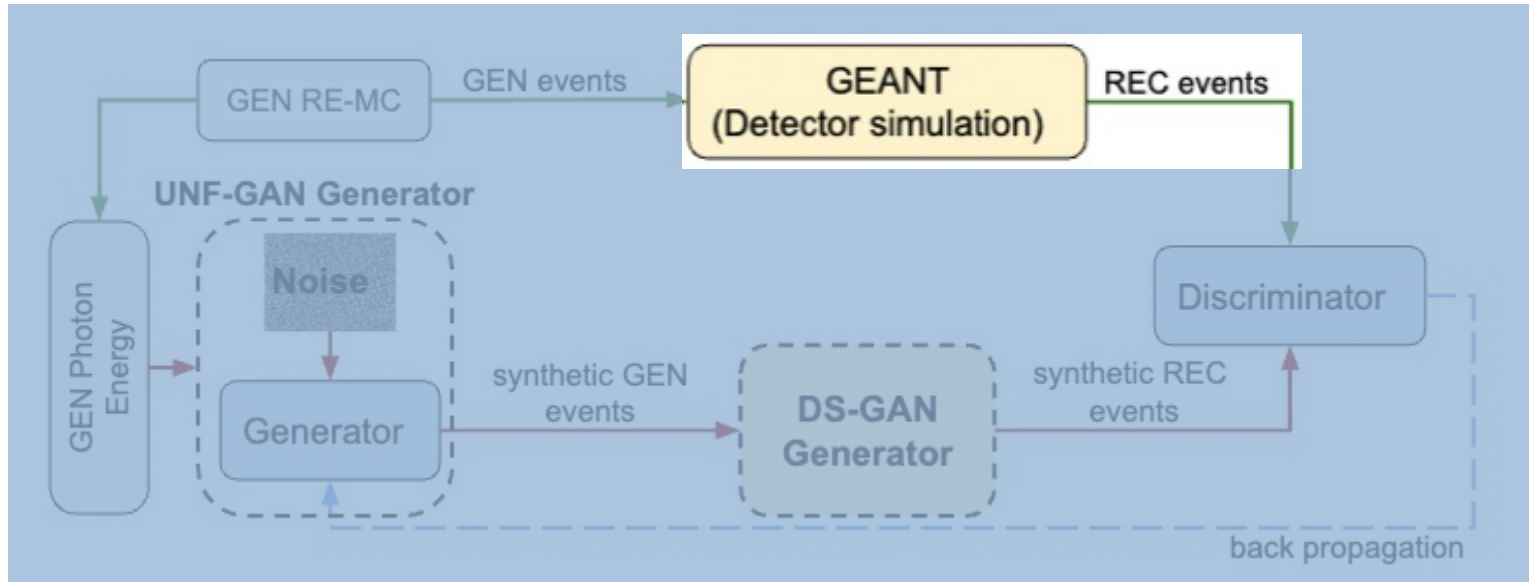Credit: T. Alghamdi et al. Phys. Rev. D **108**, 094030

1. Generate events with a (realistic) Monte Carlo $2\pi$ photoproduction model (RE-MC GEN pseudodata)

- RE-MC realistic Monte Carlo event generator to mimic real data. Includes measured cross-sections, angular distributions and decay of dominant mechanisms ($\rho^0, \Delta^{++}, \Delta^0$ + a contact term)
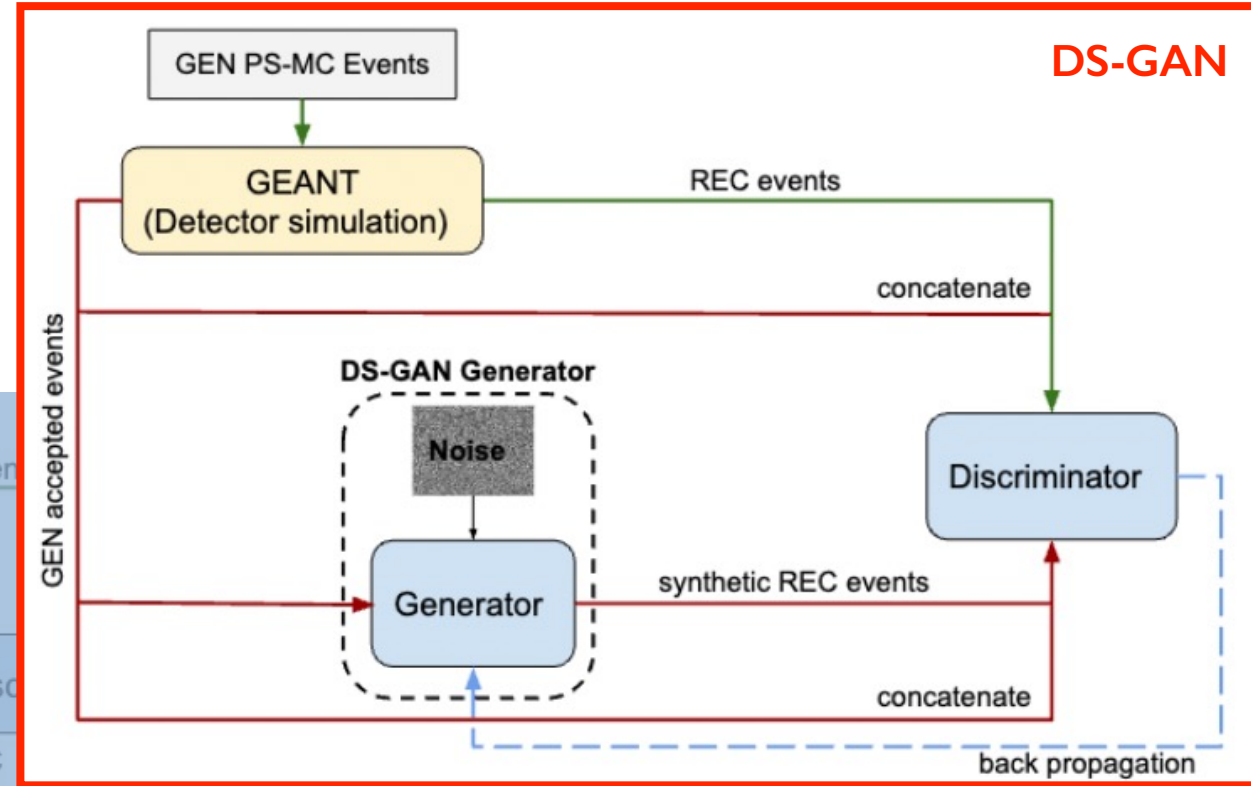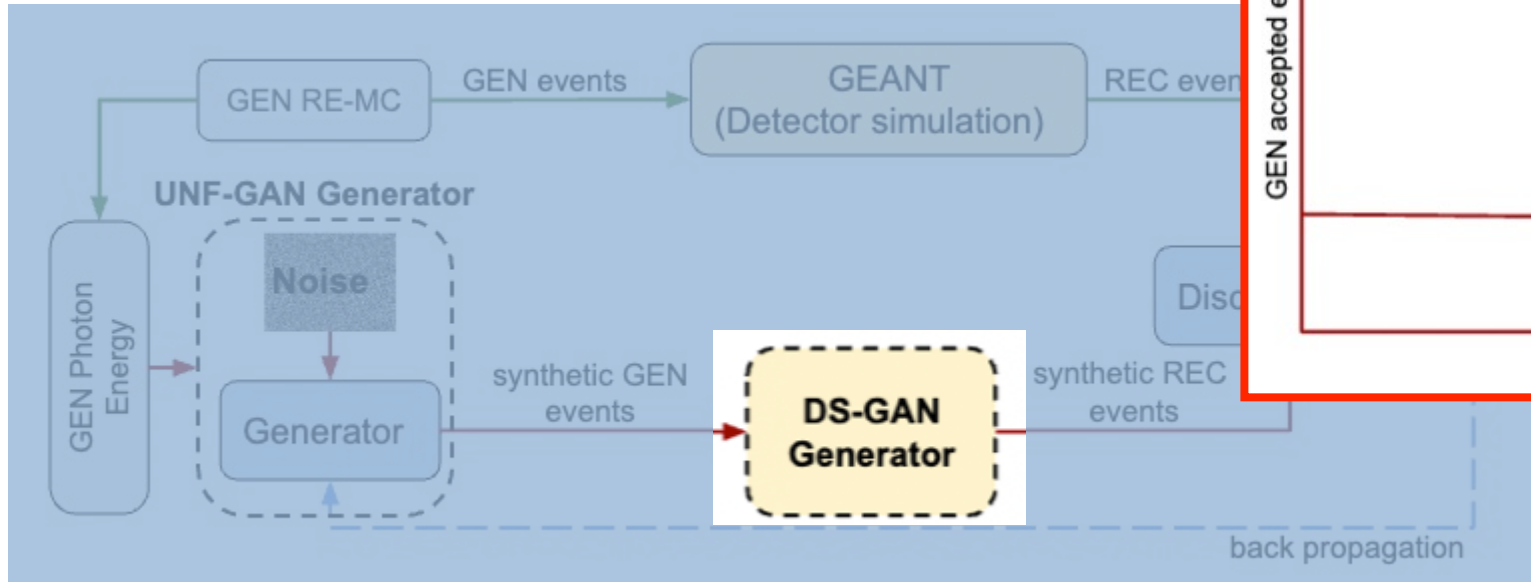
**GEN events**

2. Apply detector effects (acceptance and resolution) via GISM-GEANT (RE-MC REC pseudodata)

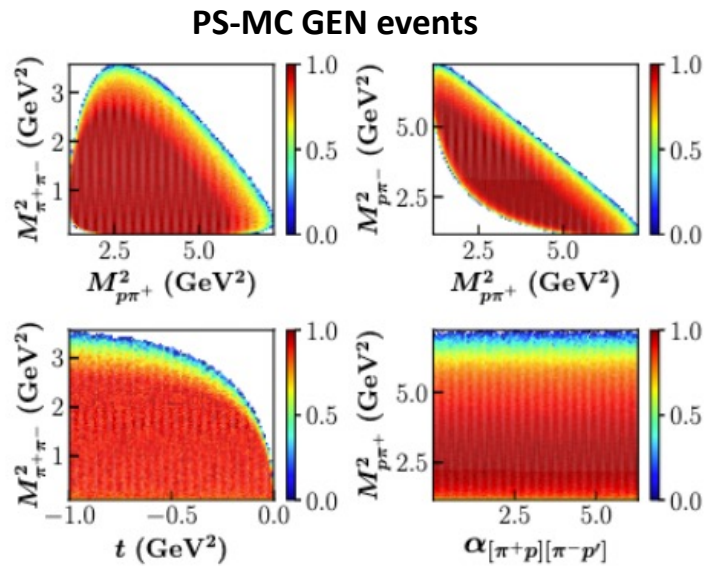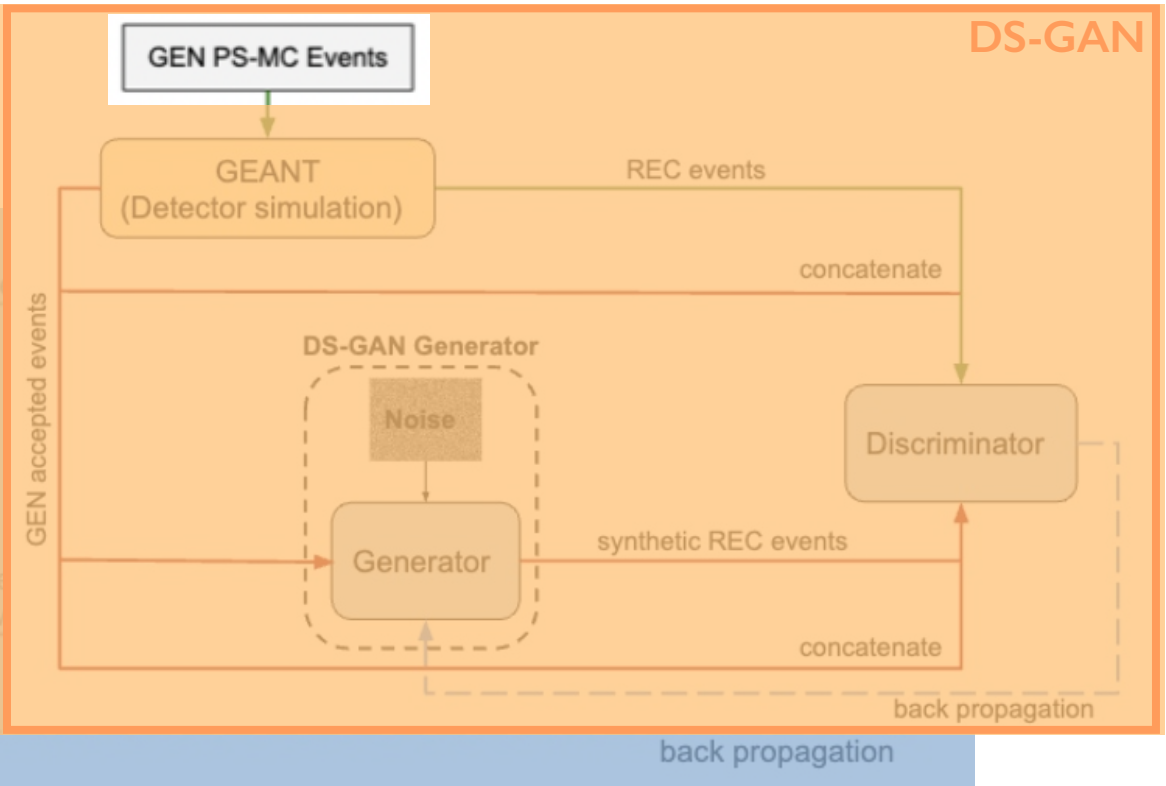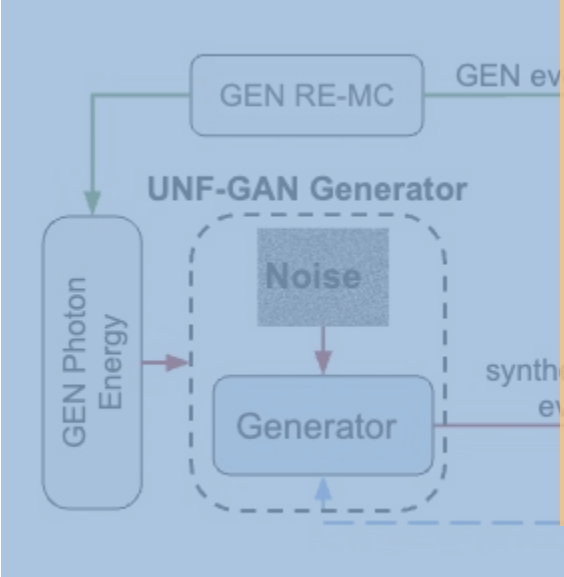- GSIM: detector simulation package to simulate CLAS detector effects based on GEANT3

3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an indipendent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)

3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an indipendent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)
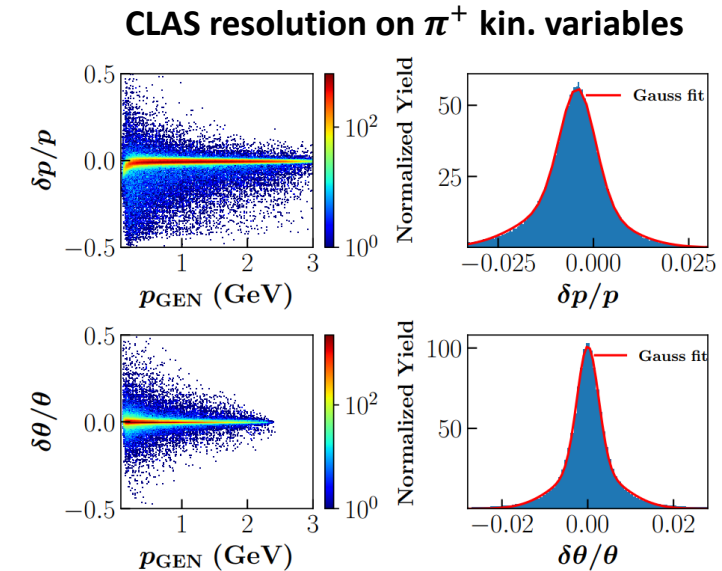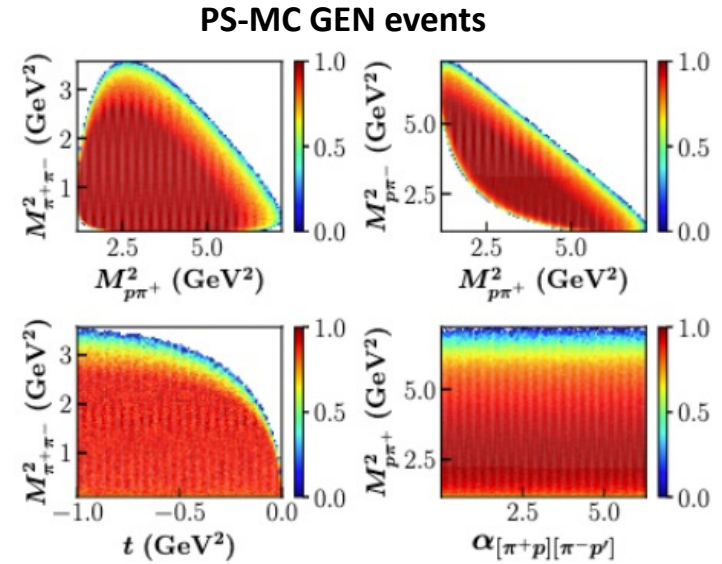
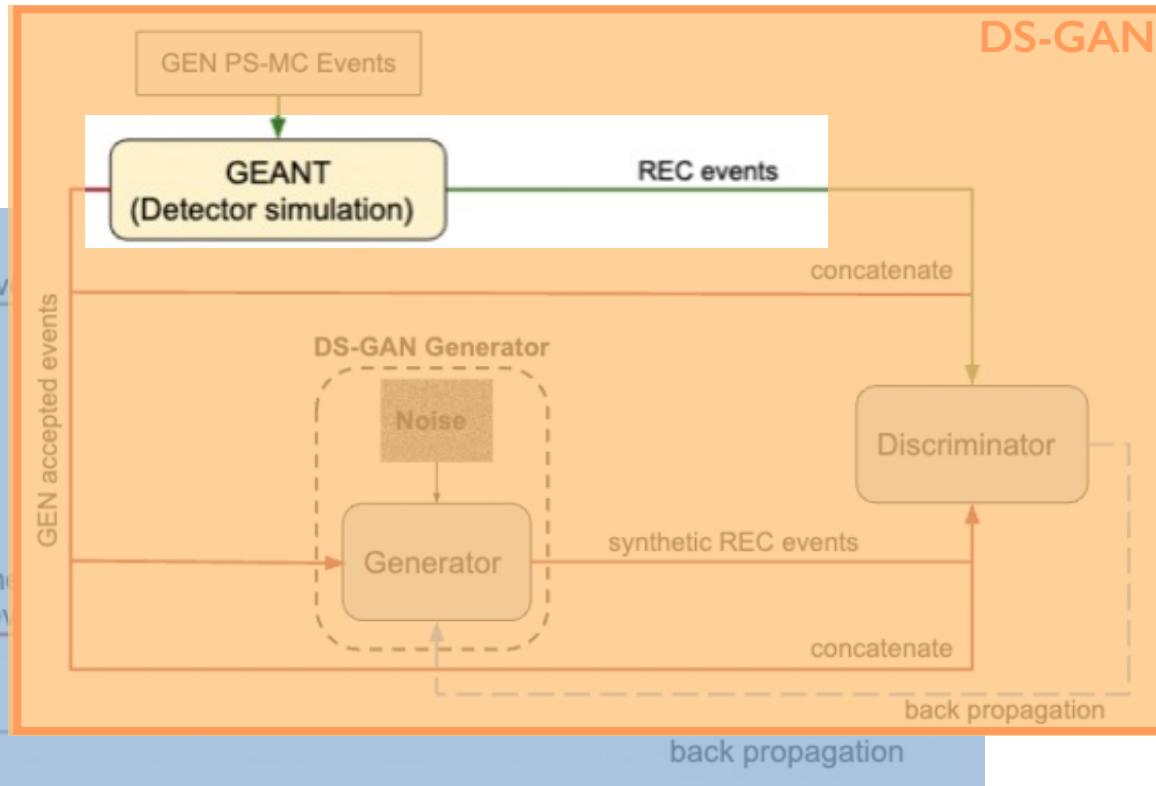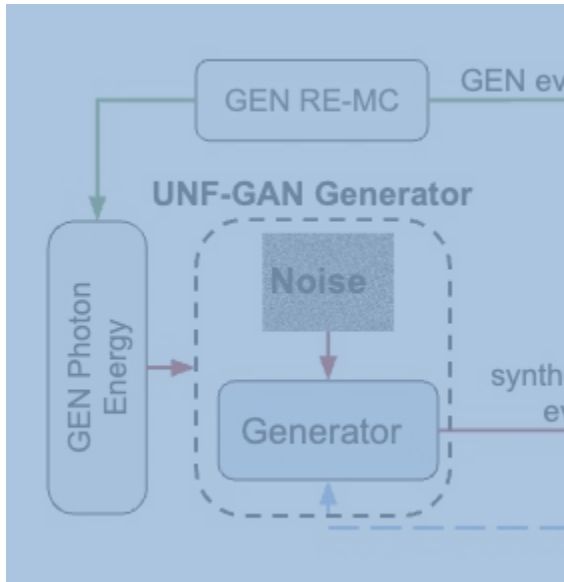- PS-MC: Phase space Monte Carlo event generator

**PS-MC GEN events**

3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an indipendent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)

- GSIM-GEANT to simulate CLAS acceptance and resolution



**PS-MC GEN events**

**CLAS resolution on $\pi^+$ kin. variables**

**MC REC pseudodata vs. DS-GAN synthetic data**

3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an indipendent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)

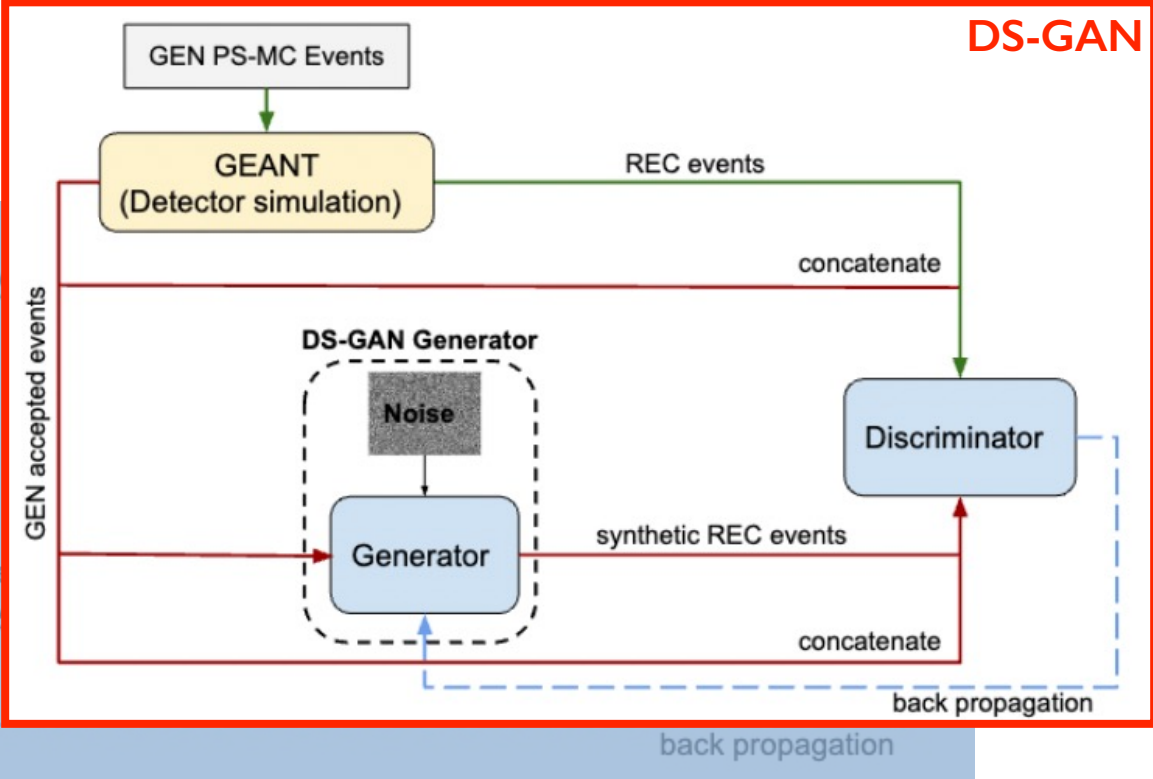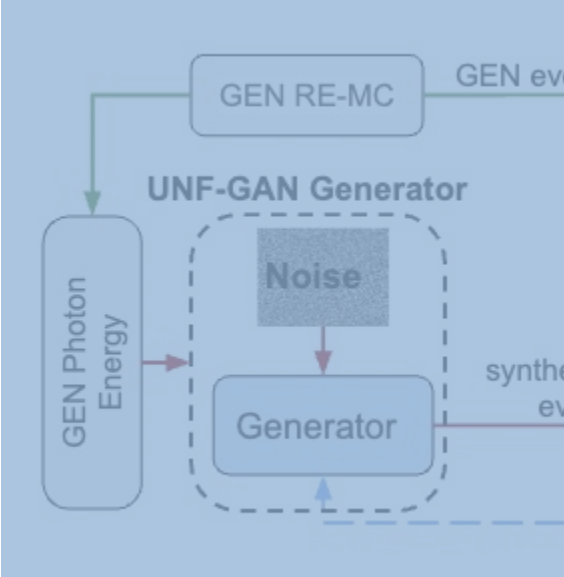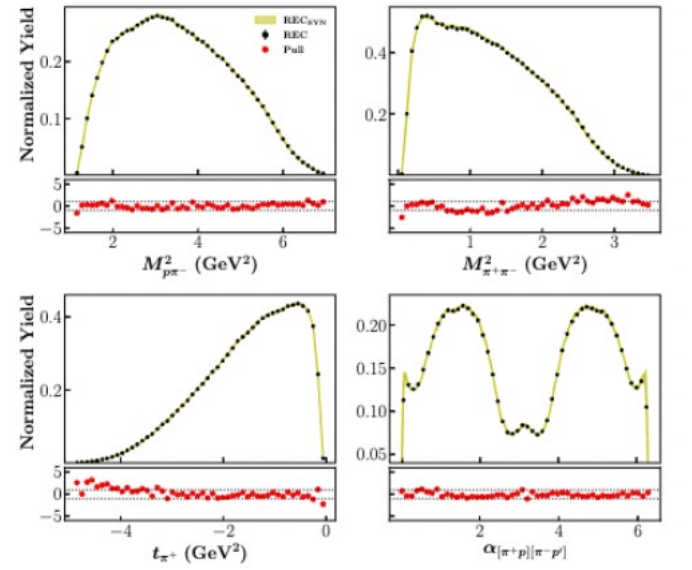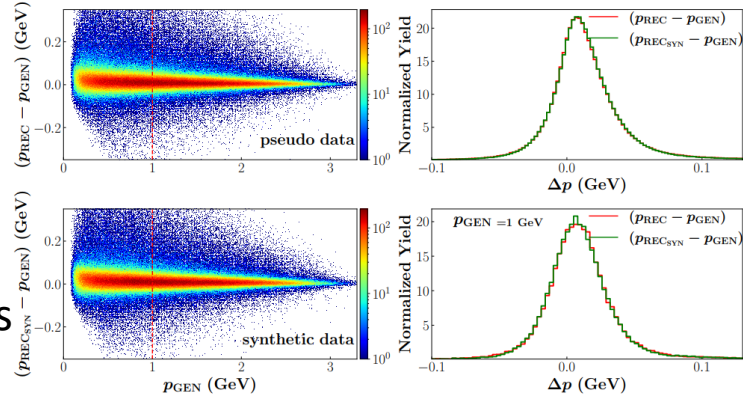- GSIM-GEANT to simulate CLAS acceptance and resolution



**DS-GAN**



**CLAS resolution**



Uncertainty quantification via **pull** calculation: Bootstrap with 20 indipendently trained GANs

**DS-GAN learned the CLAS detector effects!**

- UNF-GAN trained with REC-MC pseudodata (experimental data proxy)
- DS-GAN used to unfold CLAS detector effects (within acceptance)



**RE-MC GEN pseudodata vs. UNF-GAN SYN data**



5. Compare UNF-GAN GEN SYNT to RE-MC GEN pseudodata

Good agreement ($\pm 1\sigma$) for vertex-level training variables!

- Systematic of the full procedure (two-GANs) estimated by bootstrap with 20+20 independently trained GANs

4. Deploy the unfolding GAN (UNF-GAN) that includes the DS-GAN and train it with RE-MC REC pseudodata
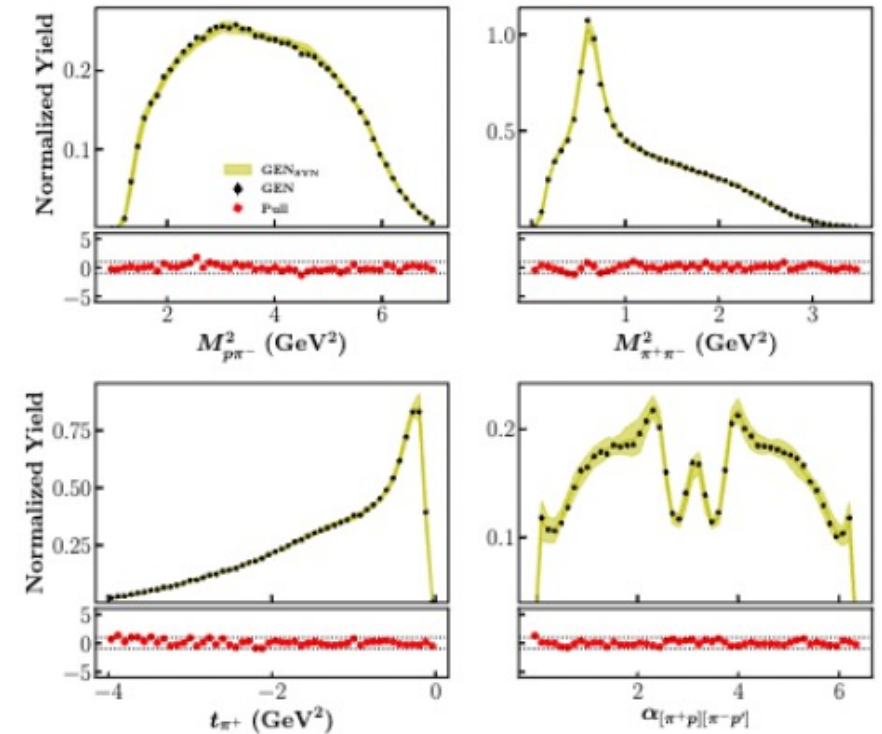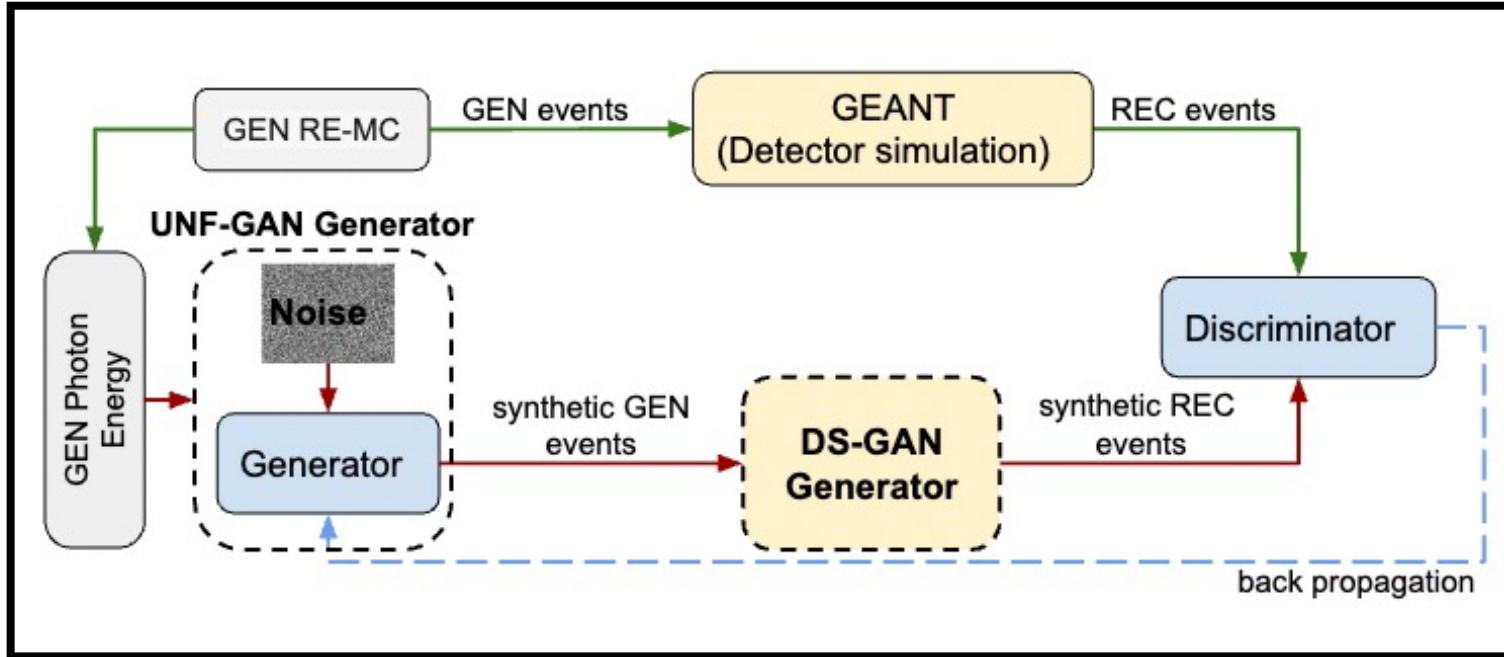
   - UNF-GAN trained with REC-MC pseudodata (experimental data proxy)
   - DS-GAN used to unfold CLAS detector effects (within acceptance)



**RE-MC GEN pseudodata vs. UNF-GAN SYN data**



**2D pulls**



5. Compare UNF-GAN GEN SYNT to RE-MC GEN pseudodata

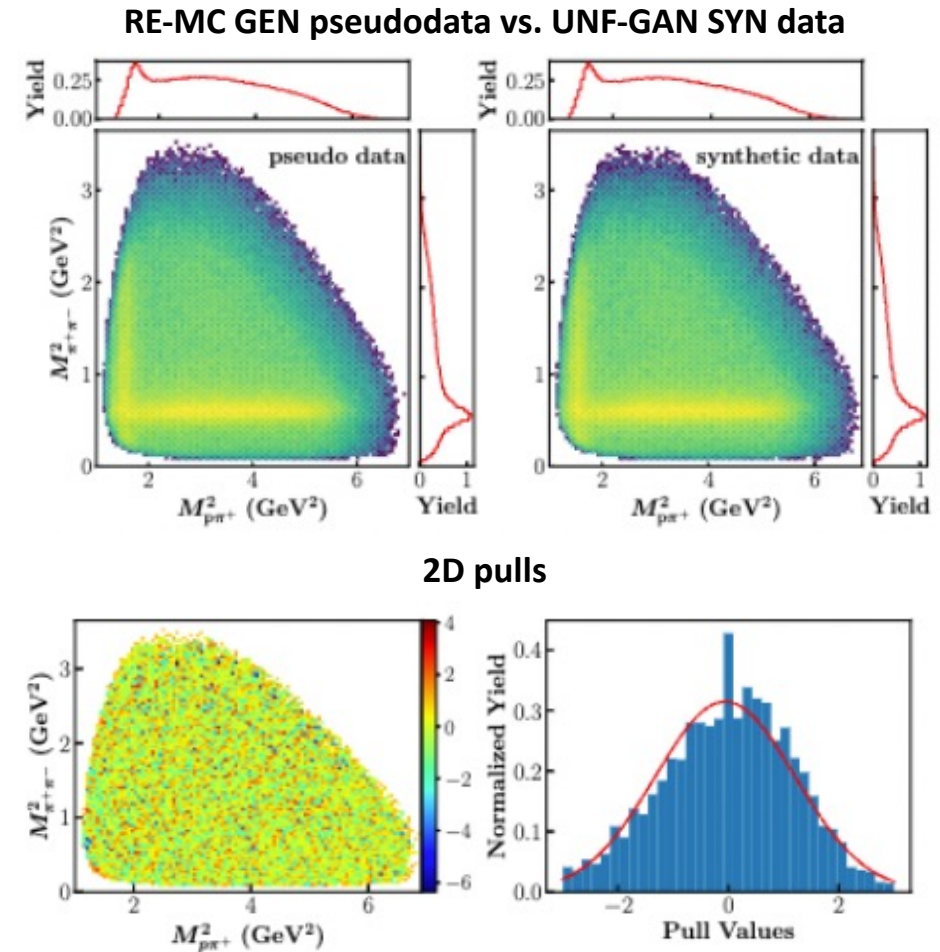Good agreement (±1σ) for 2D distributions (correlations)

4. Deploy the unfolding GAN (UNF-GAN) that includes the DS-GAN and train it with RE-MC REC pseudodata

   - UNF-GAN trained with REC-MC pseudodata (experimental data proxy)
   - DS-GAN used to unfold CLAS detector effects (within acceptance)



**Distribution in 4D bins**

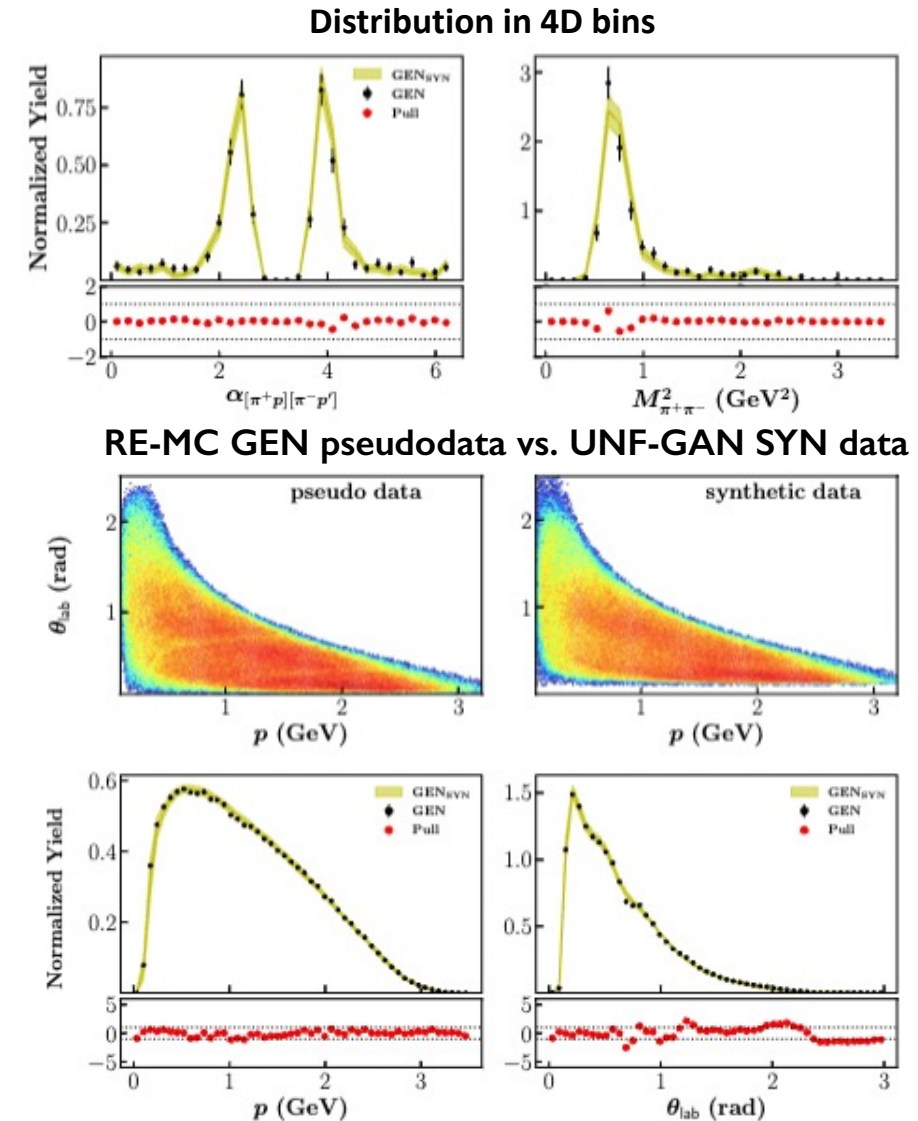**RE-MC GEN pseudodata vs. UNF-GAN SYN data**

5. Compare UNF-GAN GEN SYNT to RE-MC GEN pseudodata

Good agreement ($\pm 1\sigma$) for lab variables and in 4D bins

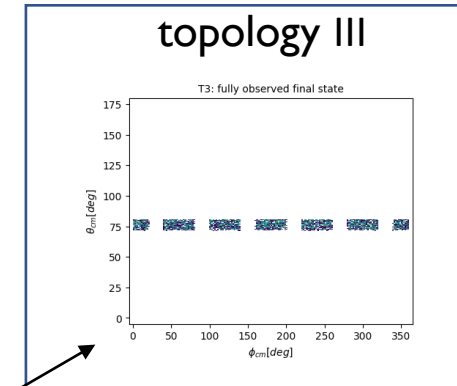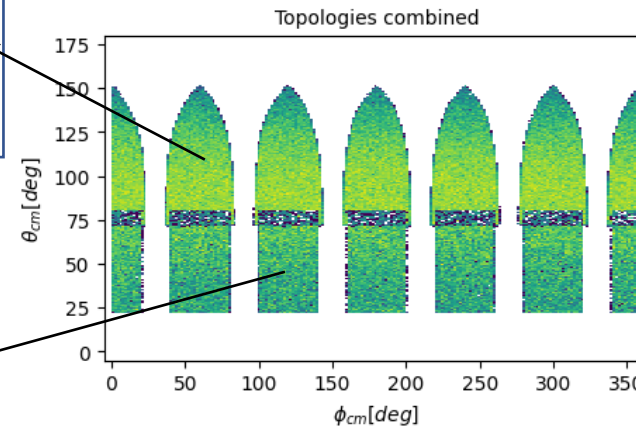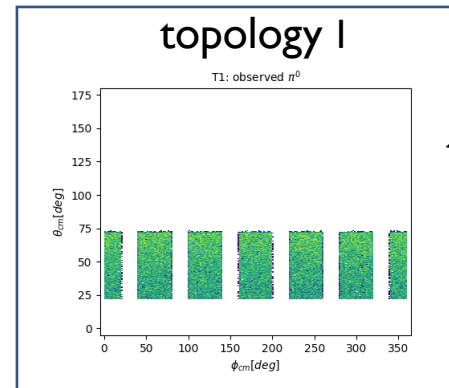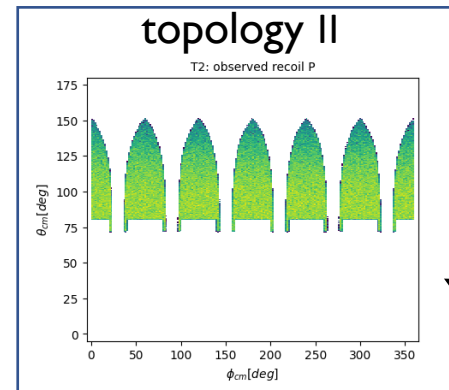- Simple 2-body process: $\gamma p \to \Delta^+(1232) \to \pi^0 p$
- Two independent variables (at fixed energy): $\theta_{cm}$ and $\phi_{cm}$
- Monte Carlo eventgenerator
- Simple model: Breit-Wigner with two parameters: $m_\Delta$ and $\Gamma_\Delta$

$$\frac{d\sigma}{d\Omega} \propto \frac{p_f}{p_i\,s} \sum_{\lambda_\gamma \lambda_p \lambda'_p} \left| (-)^{\lambda_\gamma} H_{|\lambda_\gamma - \lambda_p|} \frac{d^{3/2}_{\lambda_\gamma - \lambda_p, -\lambda'_p}(\theta)}{m_\Delta^2 - s - i\Gamma_\Delta m_\Delta} \right|^2$$

$$\propto \frac{p_f}{p_i\,s} \frac{3\,|H_{3/2}|^2 + 5\,|H_{1/2}|^2 - 3\cos 2\theta \left( |H_{3/2}|^2 - |H_{1/2}|^2 \right)}{(m_\Delta^2 - s)^2 + \Gamma_\Delta^2 m_\Delta^2}$$

- Detector acceptance (CLAS) implemented via fiducial cuts (coils, minimum proton momentum and angle in the lab frame)
  - topology 1: $\gamma\, p \to (p)\, \pi^0$ (proton missing)
  - topology 2: $\gamma\, p \to p\, (\pi^0)$ ($\pi^0$ missing)
  - topology 3: $\gamma\, p \to p\, \pi^0$ (all detected)
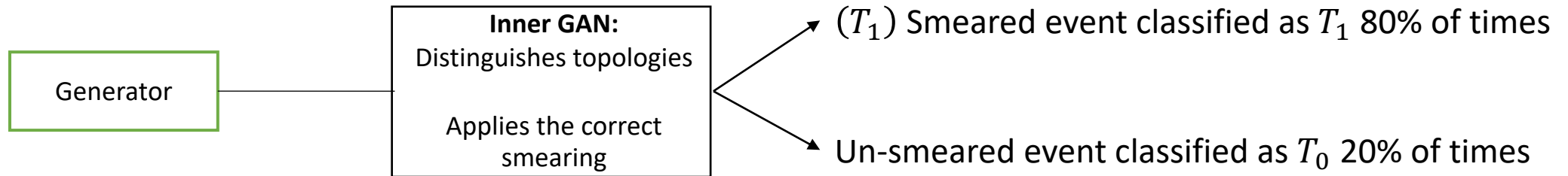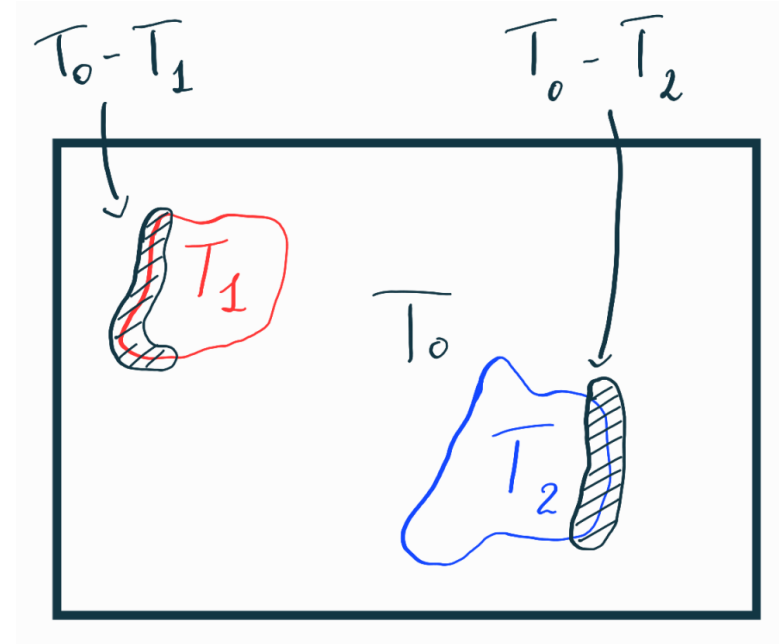  - [topology 0: unmeasured]

**Build a single Network able to generate in the full phase space according to the correct distributions**



Credit: T. Vittorini , T. Alghamdi, Y. Li

- Some evets may sometimes be classified in the wrong category:

  - Defective paddles can lead to missing events ($T_1$ which become $T_0$)

  - Reconstructed events can get missclassified ($T_1$ which become $T_2$ and viceversa)

- Existence of regions where a given generated event can go into different topologies with a given probability: Mixed events
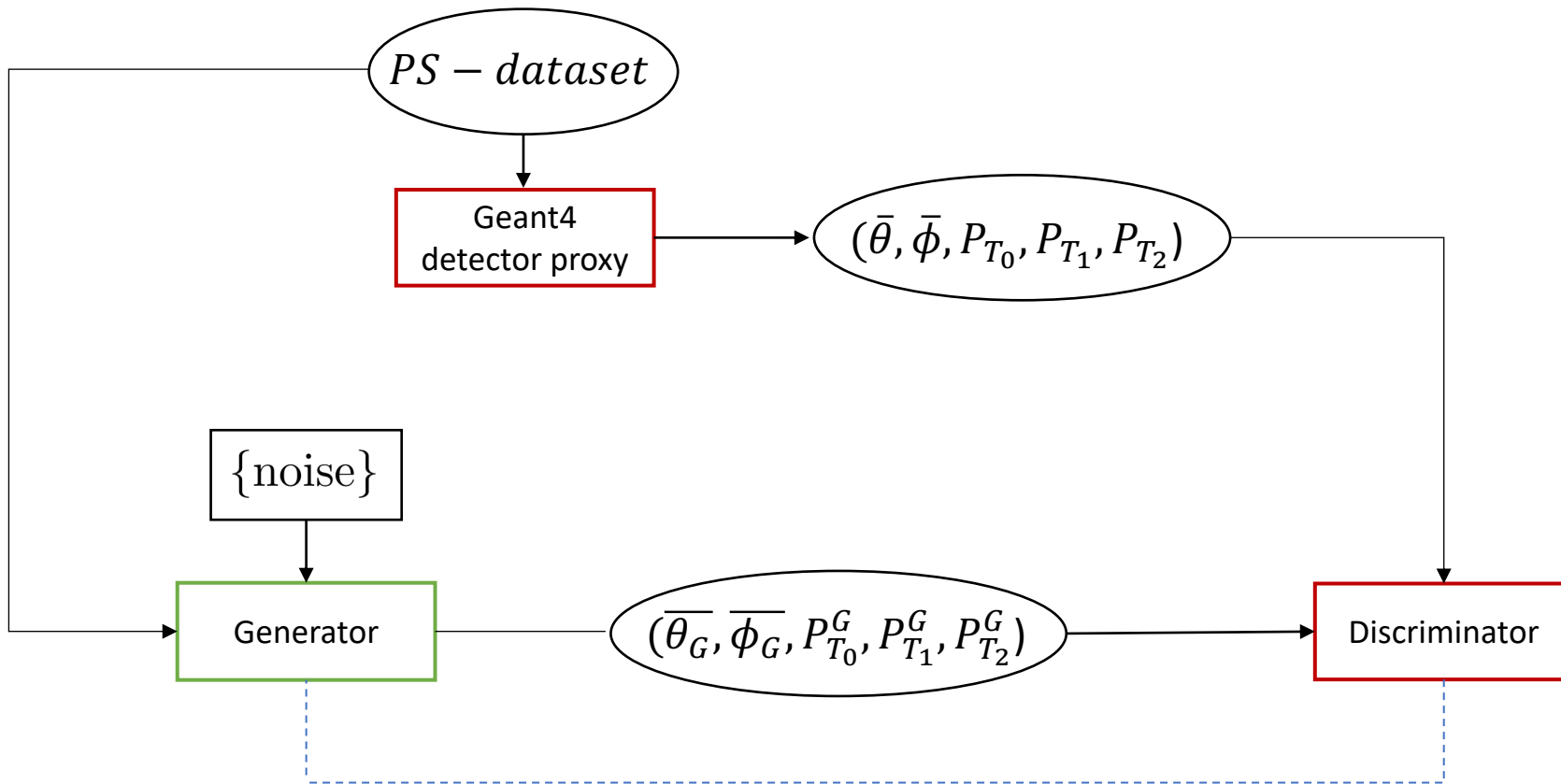


| Generator |
|---|

| **Inner GAN:** Distinguishes topologies <br><br> Applies the correct smearing |
|---|

$(T_1)$ Smeared event classified as $T_1$ 80% of times

Un-smeared event classified as $T_0$ 20% of times

Credit: T. Vittorini , T. Alghamdi, Y. Li

- Here we are assuming that the smearing applied to mixed events is similar: We will apply a common smearing to all the events being classified with multiple topologies

- **Inner GAN**

- The process of generating training $P_i$ will be carried out outside of the dataset training

- Each PS-event will be associated with different probabilities

- The inner GAN training dataset will consist of $\theta$ and $\phi$ values to identify the event + the probabilities of being classified into a given topology

- The output of the $P_i^G$ could be distribution peaked around the expected value
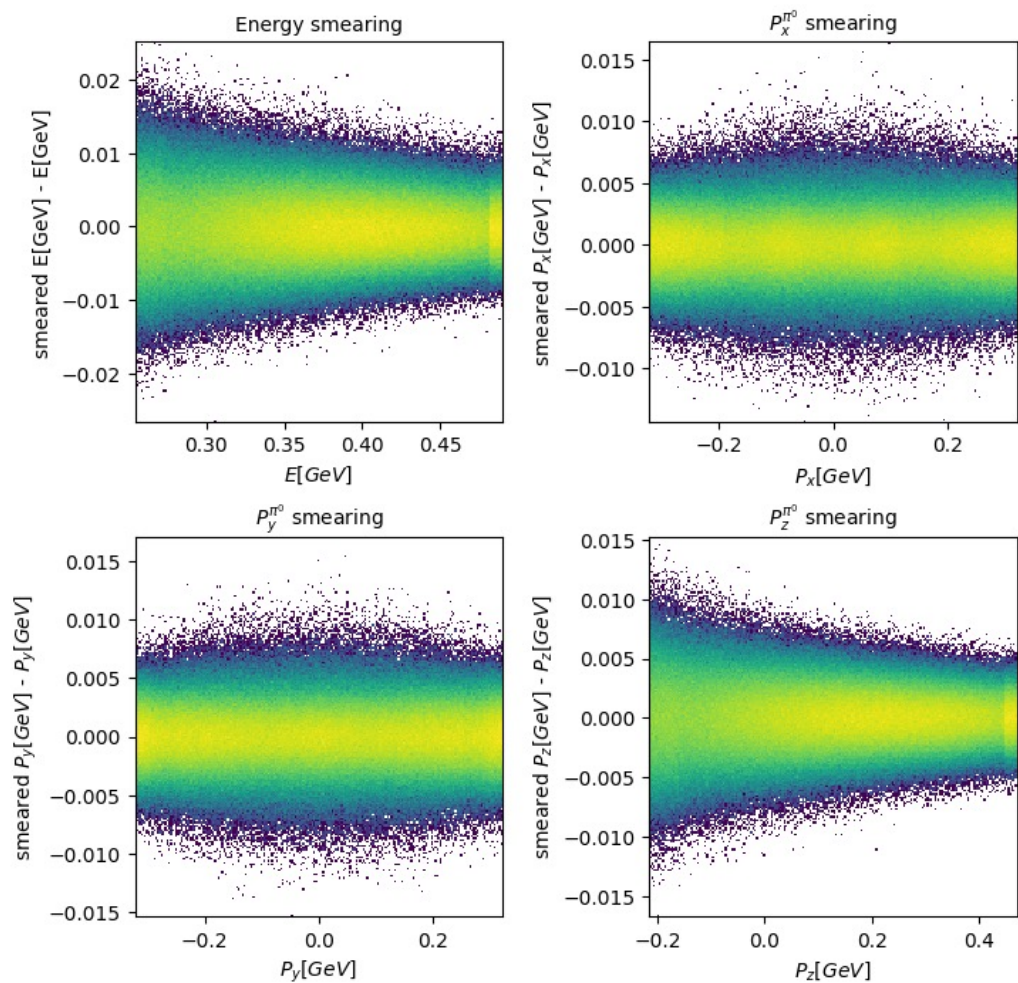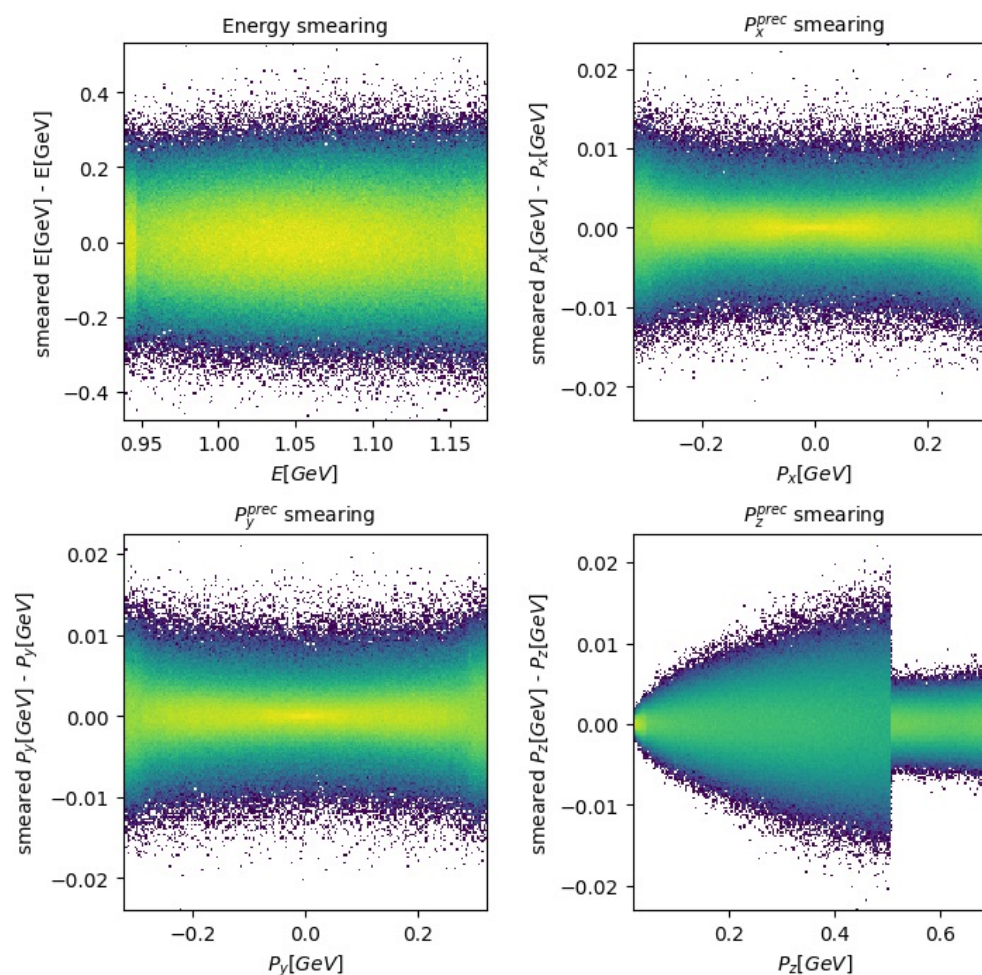
Credit: T.Vittorini , T.Alghamdi, Y. Li

- Reasonable smearing applied to the training variables:

$\pi^0$ 4-momenta

Recoil proton 4-momenta



Credit: T. Vittorini , T. Alghamdi, Y. Li

The A(i)DAPT program

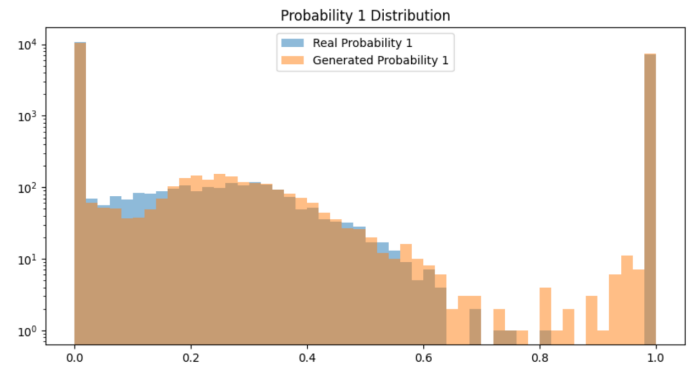- Here we are assuming that the smearing applied to mixed events is similar: We will apply a common smearing to all the events being classified with multiple topologies

- **Outer GAN**



- The outer GAN trained generator will be able to generate vertex level events in the full phase space, with the relative distributions derived according to experimental data (in the measured regions) and some realistic model (in the unmeasured region)

Credit: T. Vittorini , T. Alghamdi, Y. Li

- Working towards the application of the developed machinery to CLAS12 pseudodata

- Pseudo-data for the reaction $ep \rightarrow e'p\,\pi^+\pi^-$
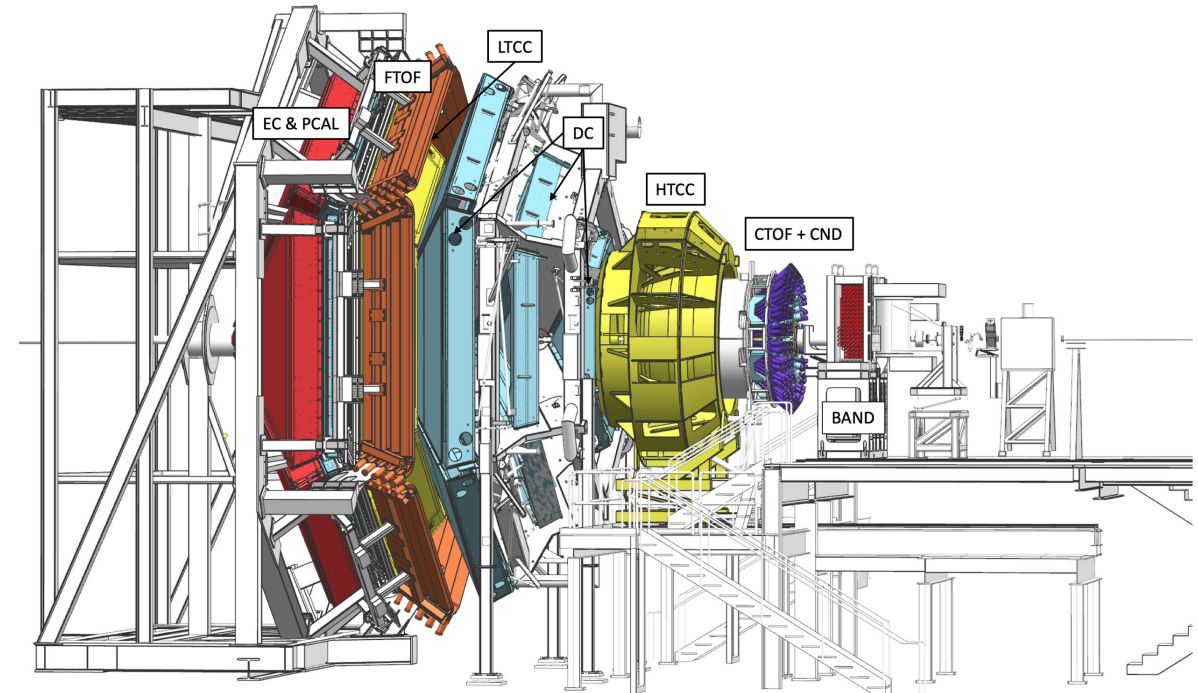
- Focus on the topology with final state $e'p\,\pi^+(\pi^-)$

**Test the inner – GAN architecture on CLAS12 detector**

- If this procedure works well on CLAS and CLAS12 data the architecture robustness is guaranteed

- We can put together in a coherent way information from different kinematic regions
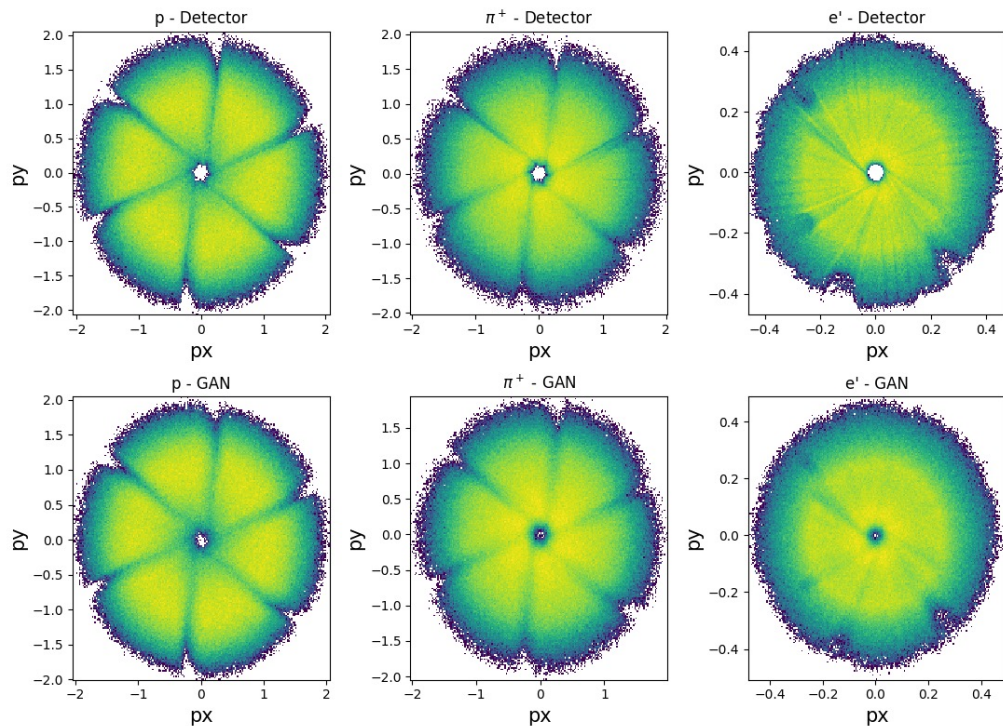


EC & PCAL
FTOF
LTCC
DC
HTCC
CTOF + CND
BAND

Credit: D. Glazier, T. Alghamdi, M. Spreafico

## $P_x$ vs $P_y$ comparison



## Resolution comparison



## Smeared training variables comparison



Credit: D. Glazier, T. Alghamdi, M. Spreafico

# Summary

**A(I)DAPT program aims to demonstrate a novel way to extract and interpret physics observables**
- Multi-step program
- We performed a positive closure test on $2\pi$ photoproduction
- We demonstrated that GANs are a viable tool to unfold detector effects (smearing) to generate a synthetic copy of data
- We demonstrated that the original correlations are preserved
- Preserve data in alternative compact and efficient form

**We are working on:**
- Quantifying the systematic error introduced by the detector acceptance
- Implementing this architecture into jlab software in order to make it easily available to everyone
- Further verify that this procedure is well defined confronting the results obtained analysing CLAS data with traditional analysis
- Make this procedure an efficient way to analyse CLAS12 $2\pi$ data

**There is still a long way to go to be able to use AI to extract physics from data in an efficient way, but we are moving towards the right direction!**
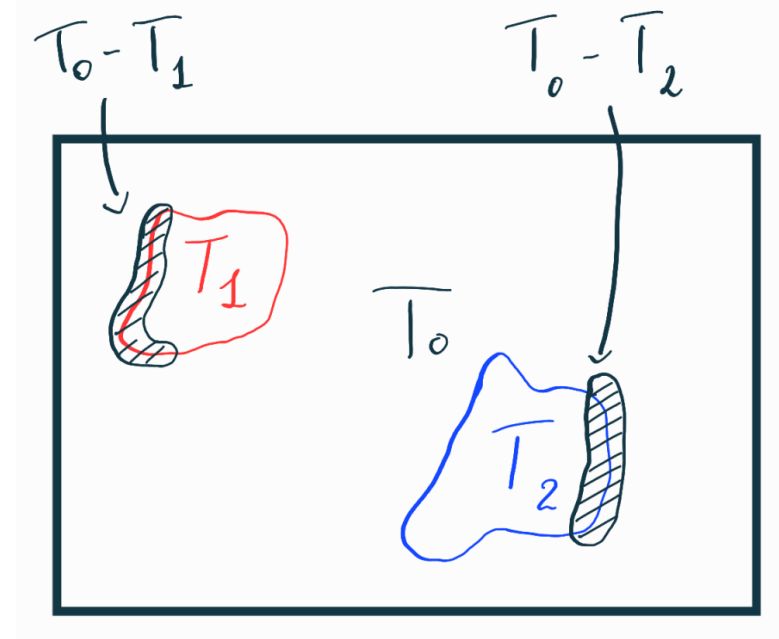
# Thank you!

- How do we define these Mixed events?

  - We generate a dataset of vertex level events in all the phase space

  - We pass **each** event in our dataset through a realistic geant4 proxy of our detector multiple times (for example 100)



| Vertex | $\#T_0$ | $\#T_1$ | $\#T_2$ |
|---|---|---|---|
| event 0 | 20 | 80 | 0 |
| ... | ... | ... | ... |
| event N | 30 | 0 | 70 |

$P_i$ is the probability of a given event to end up in a specific topology

- Each event can then be defined as $(\theta, \phi, P_{T_0}, P_{T_1}, P_{T_2})$

- Event 0 would look like $(\theta, \phi, 0.2, 0.8, 0)$

- An event which always belongs to $T_1$ will look like this: $(\theta, \phi, 0, 1, 0)$

Credit: T. Vittorini , Y. Alanazi, T. Alghamdi, Y. Li