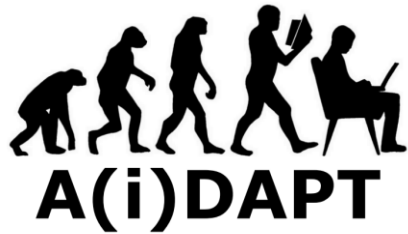


The A(i)DAPT program

AI for Data Analysis and Preservation

Alessandro Pilloni

on behalf of A(i)DAPT Working Group



AI for Data Analysis and Preservation

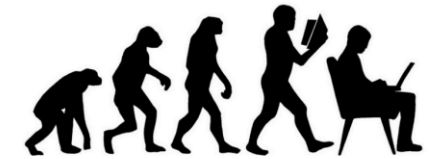
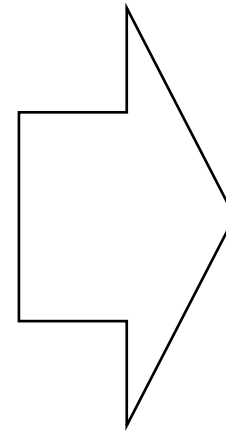


Università
degli Studi di
Messina



- Data collected by NP/HEP experiments are (always) affected by the detector's effects
- Before starting physics analysis, unfolding the detector effect 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 uncertainties (UQ)
- Move from cross sections to amplitudes

Collaborative effort (regular meetings)

- ML experts (ODU, JLab)
- Experimentalists (JLab Hall-B)
- Theorists (JPAC, JAM)

Y. Alanazi, T. Alghamdi, M. Battaglieri, M. Bondi, Ł. Bibrzycki, A.V. Golda, A.N. Hiller Blin, E.L. Isupov, Y. Li, L. Marsicano, W. Melnitchouk, V.I. Mokeev, G. Montaña, A. Pilloni, N. Sato, M. Spreafico, A.P. Szczepaniak, T. Vittorini



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

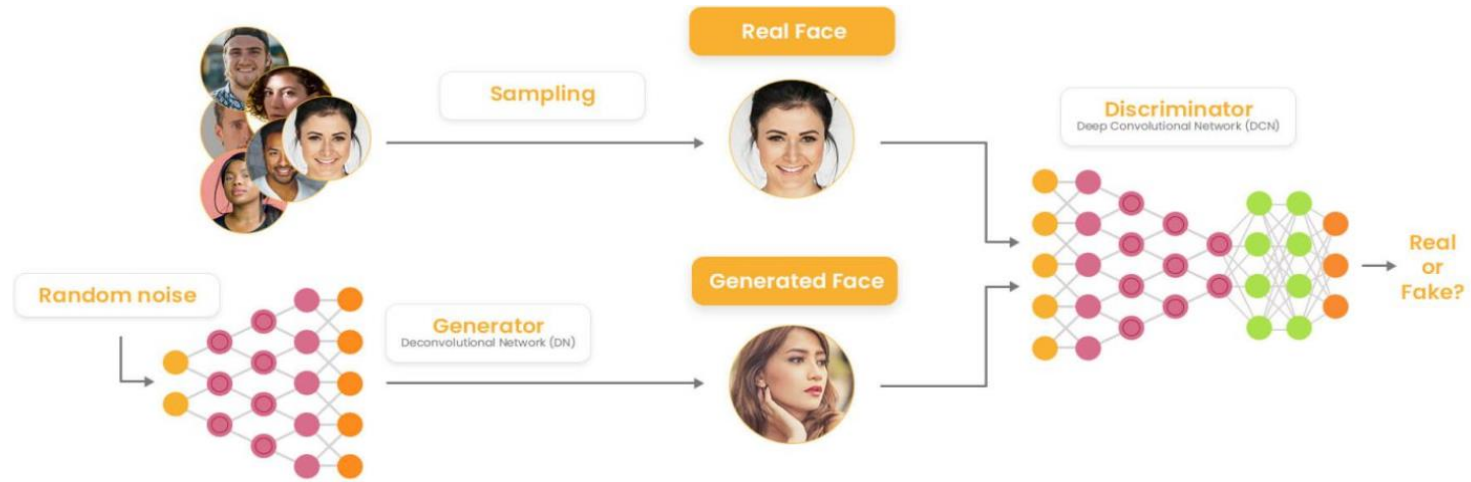
- Spikes may be concealed behind the detector resolution
- Measurements could be extended to unphysical regions

- Mitigation strategy:

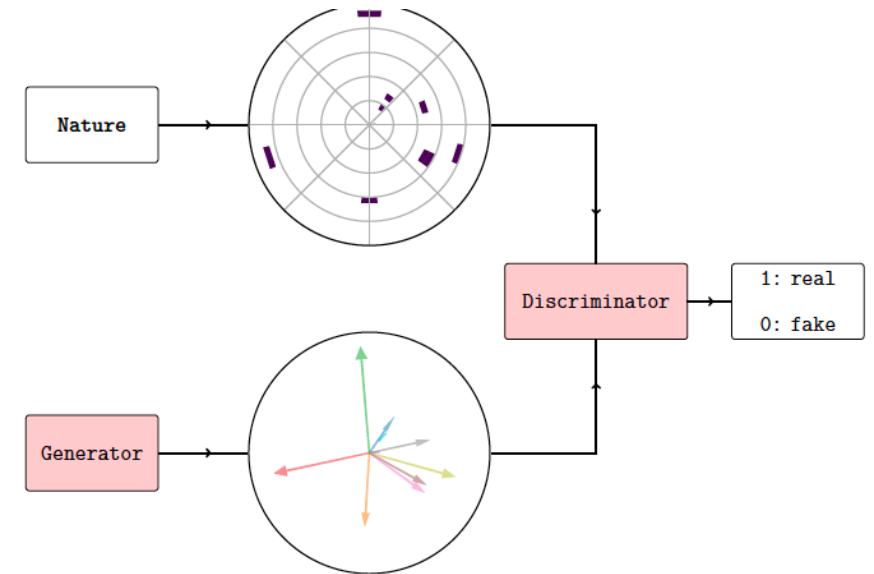
- Acceptance: ‘Fiducial volumes’ to exclude unmeasured regions and extend the covered measured of the phase space
- Resolution: build and validate ML-models to unfold resolution effects



Generative Adversarial Networks (GANs)



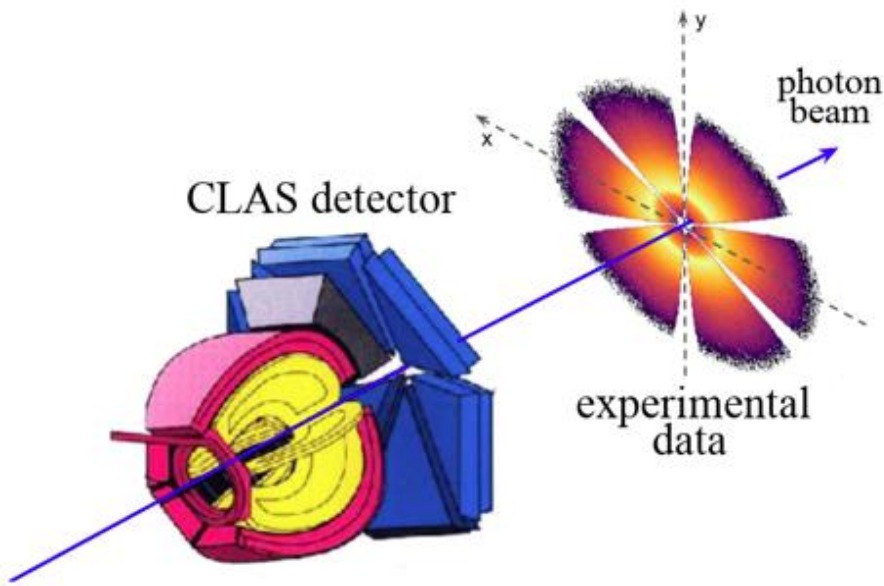
- Generative model based on the competition between two Neural Networks: Generator vs Discriminator
 - **Generator** produces synthetic data which progressively reproduce realistic data and the **Discriminator** has to distinguish between synthetic and realistic data
 - **Generator** be used to retain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



Multi-d cross-section: exclusive 2π photoproduction

M. Battaglieri *et al.* (CLAS Collaboration)
Phys. Rev. Lett. 102, 102001

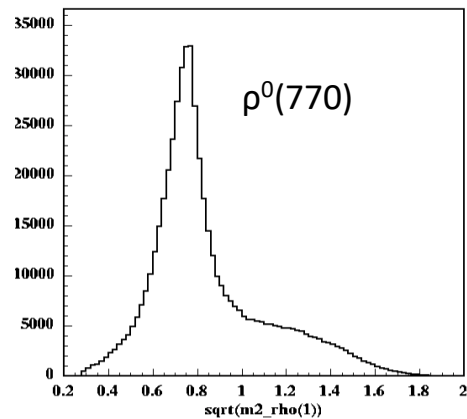
M. Battaglieri *et al.* (CLAS Collaboration)
Phys. Rev. D 80, 072005



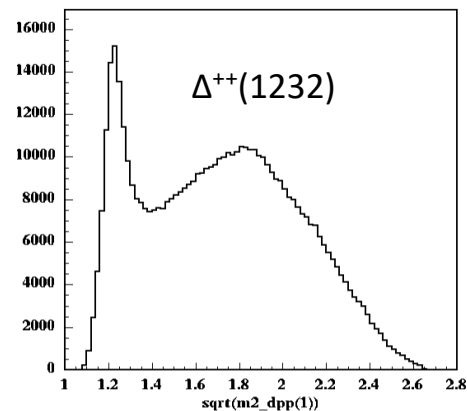
CLAS g11 kinematics

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

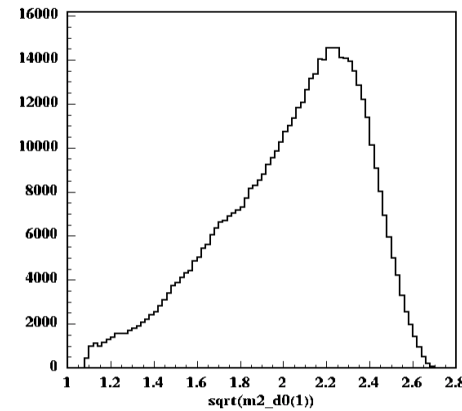
$M(\pi\pi)$



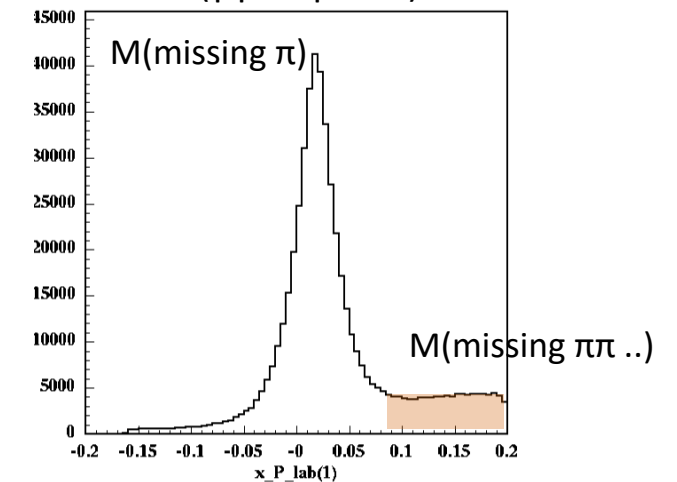
$M(p\pi^+)$



$M(p\pi^-)$

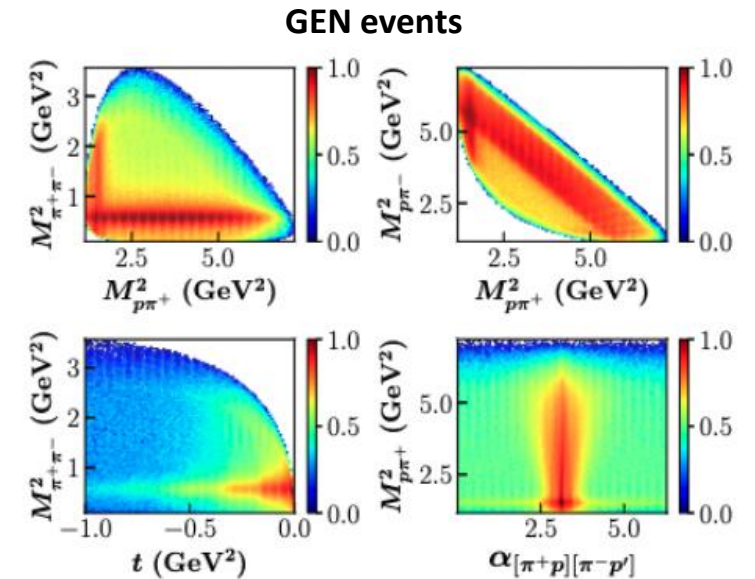
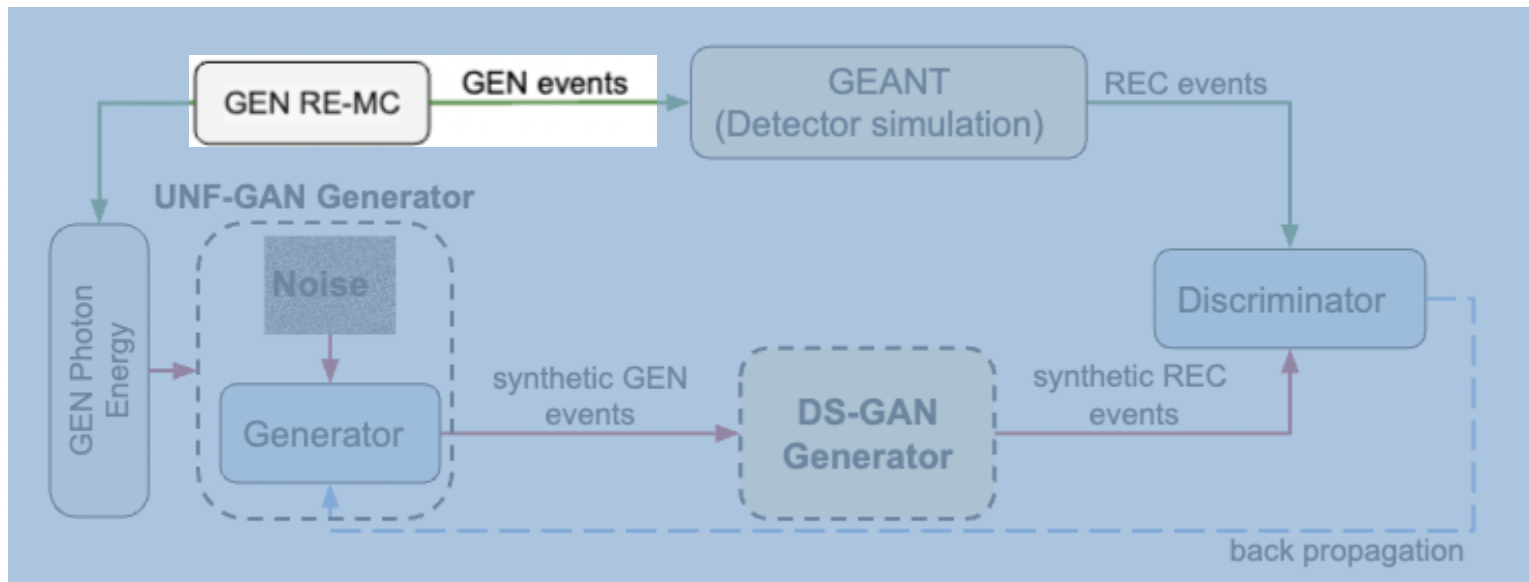


$M(\gamma p \rightarrow p\pi^+ X)$



2π photoproduction closure test

1. Generate events with a (realistic) Monte Carlo 2π 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 (ρ^0 , Δ^{++} , Δ^0 + a contact term)

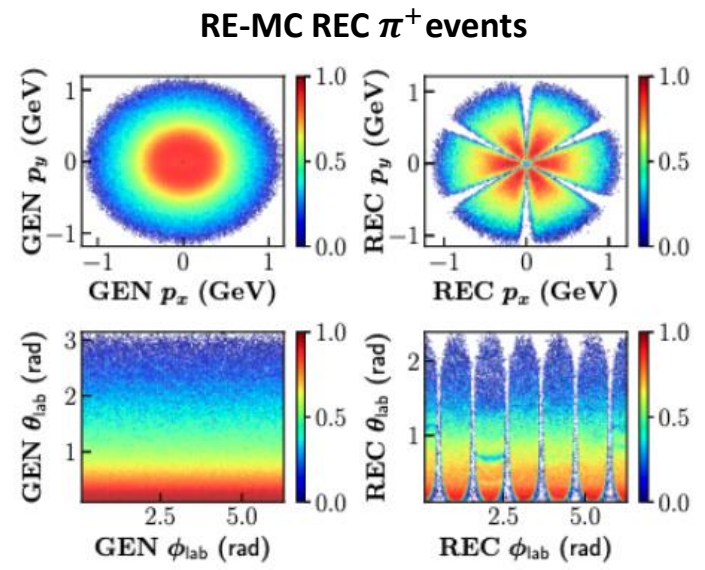
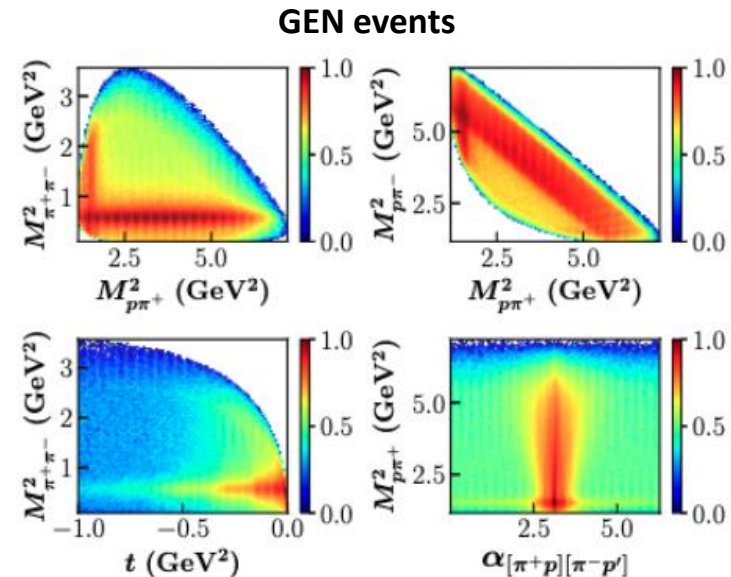
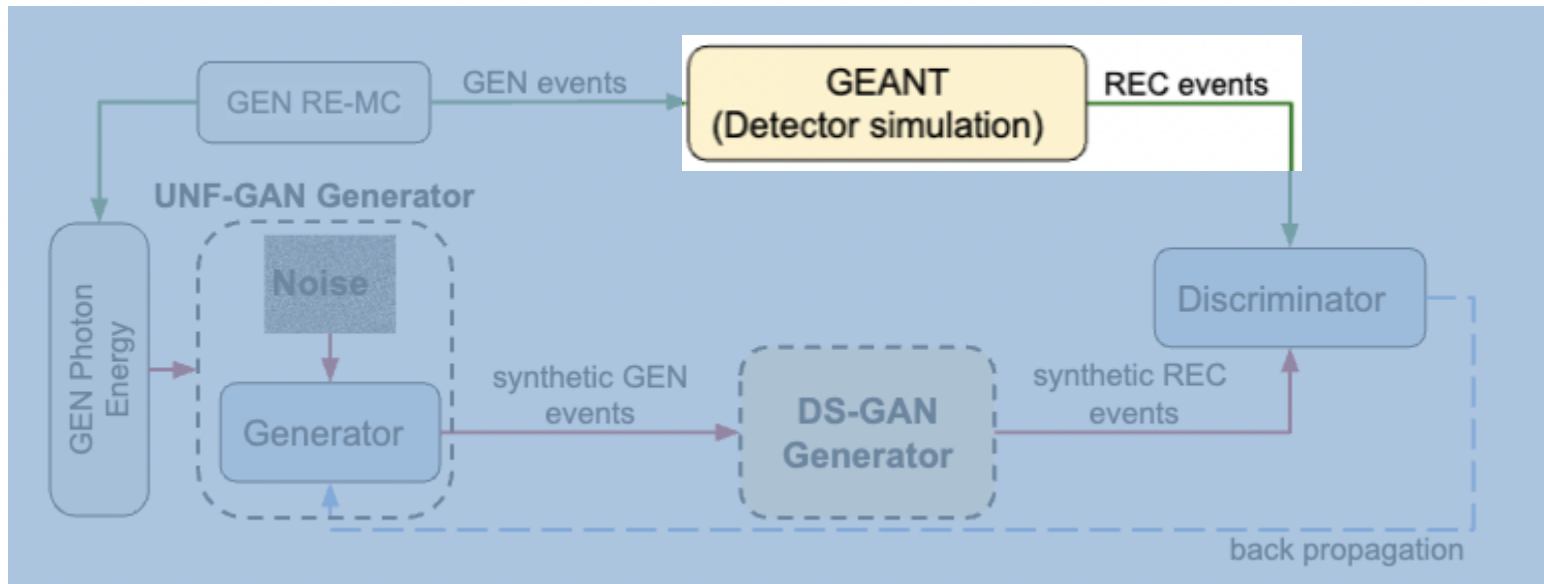


Credit: T. Alghamdi, Y. Alanazi, M. Battaglieri, L. Bibrzycki, A.V. Golda, A.N. Hiller Blin, E.L. Isupov, Y. Li, L. Marsicano, W. Melnitchouk, V.I. Mokeev, G. Montana, A. Pilloni, N. Sato, A.P. Szczepaniak, T. Vittorini, *PRD108, 094030 (2023)*



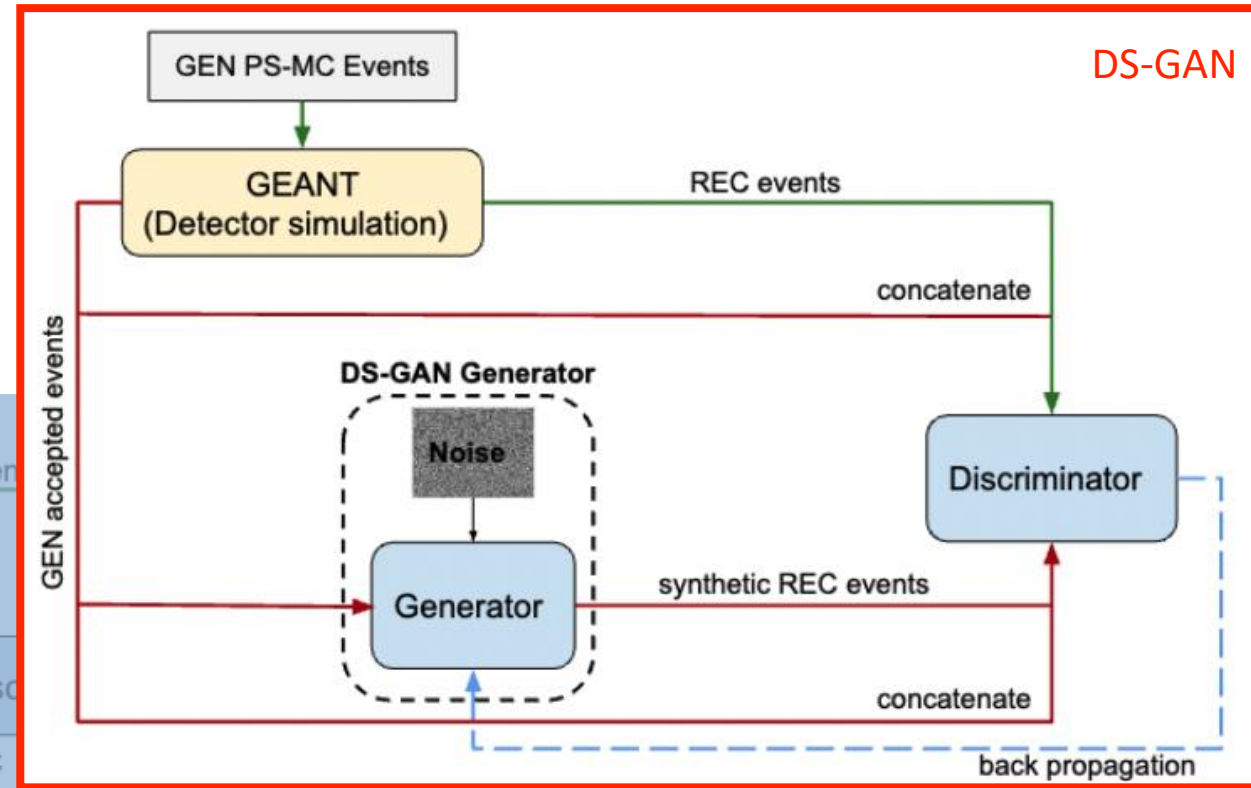
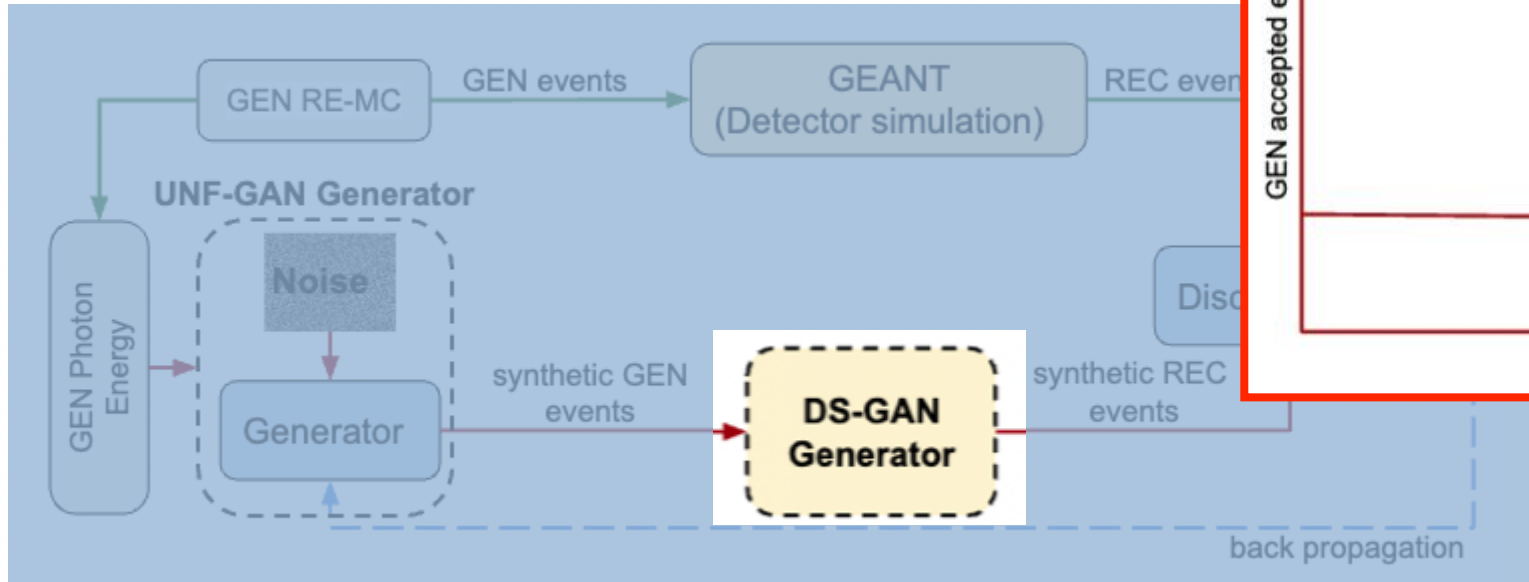
2π photoproduction closure test

- 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



2π photoproduction closure test

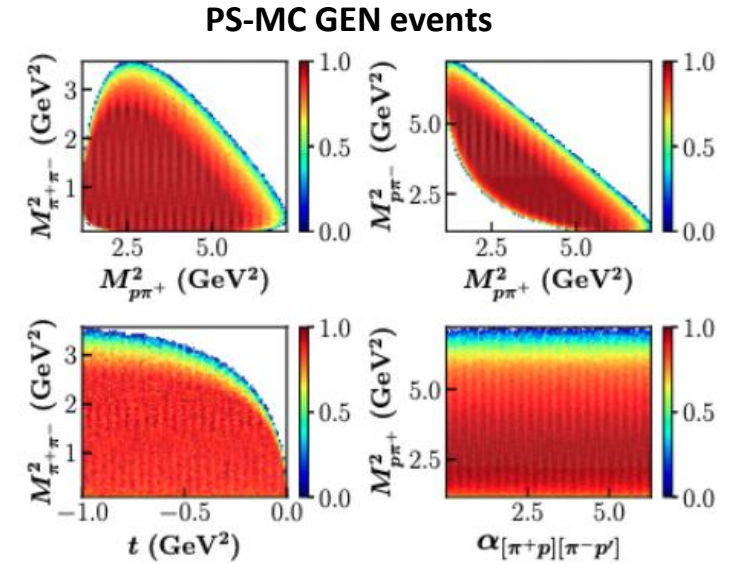
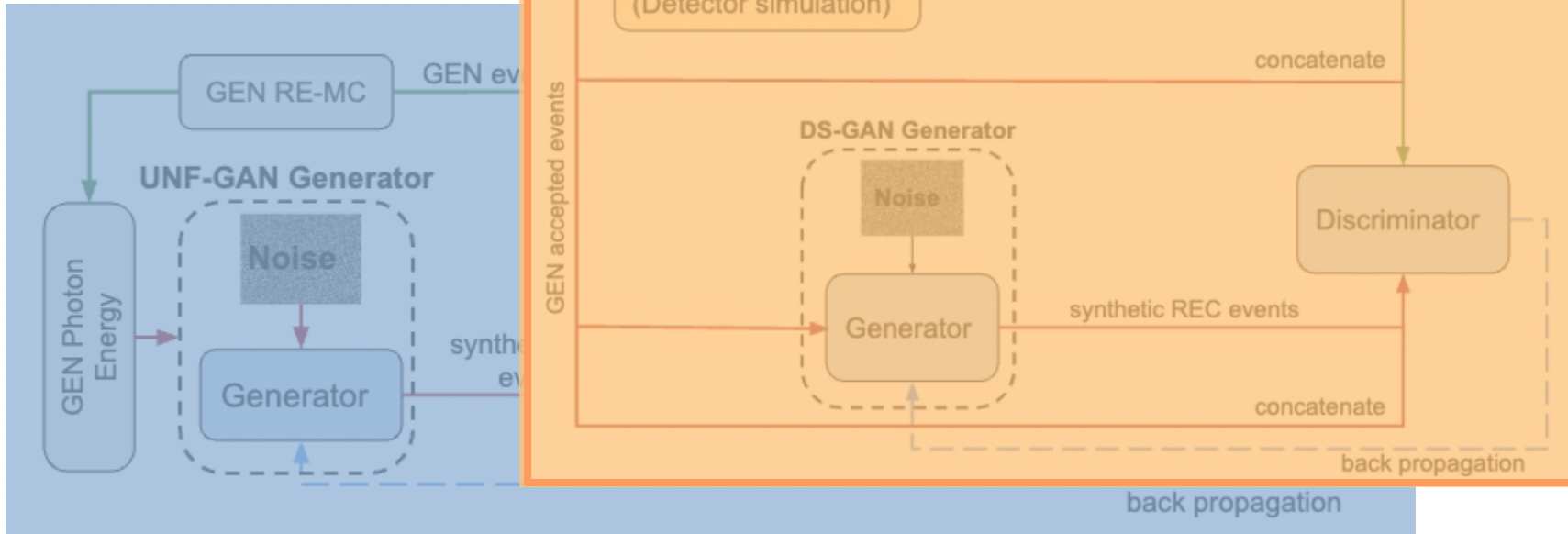
- 3. Deploy a secondary GAN (DS-GAN) to learn detector effects using an independent MC event generator (PS-MC) + GSIM-GEANT (GEN and REC pseudodata)



2π photoproduction closure test

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

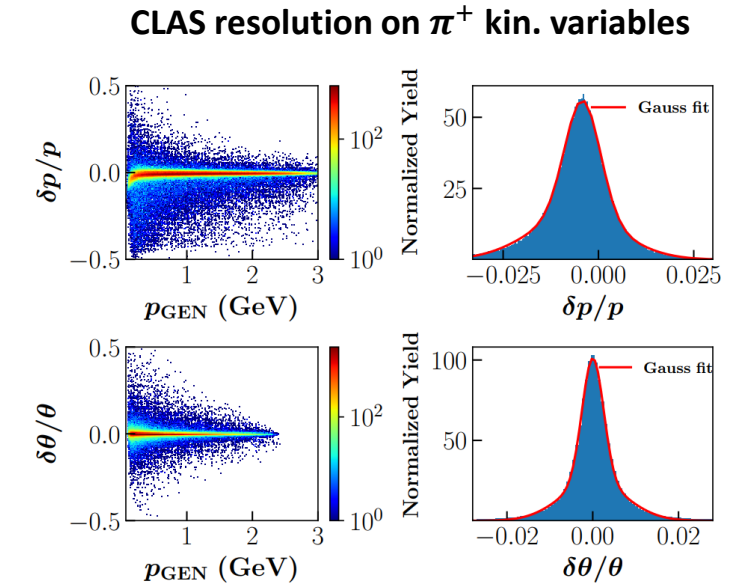
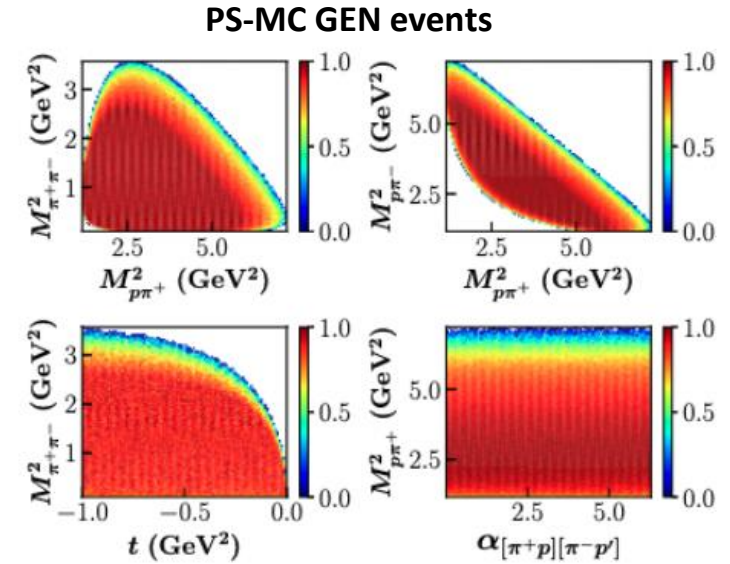
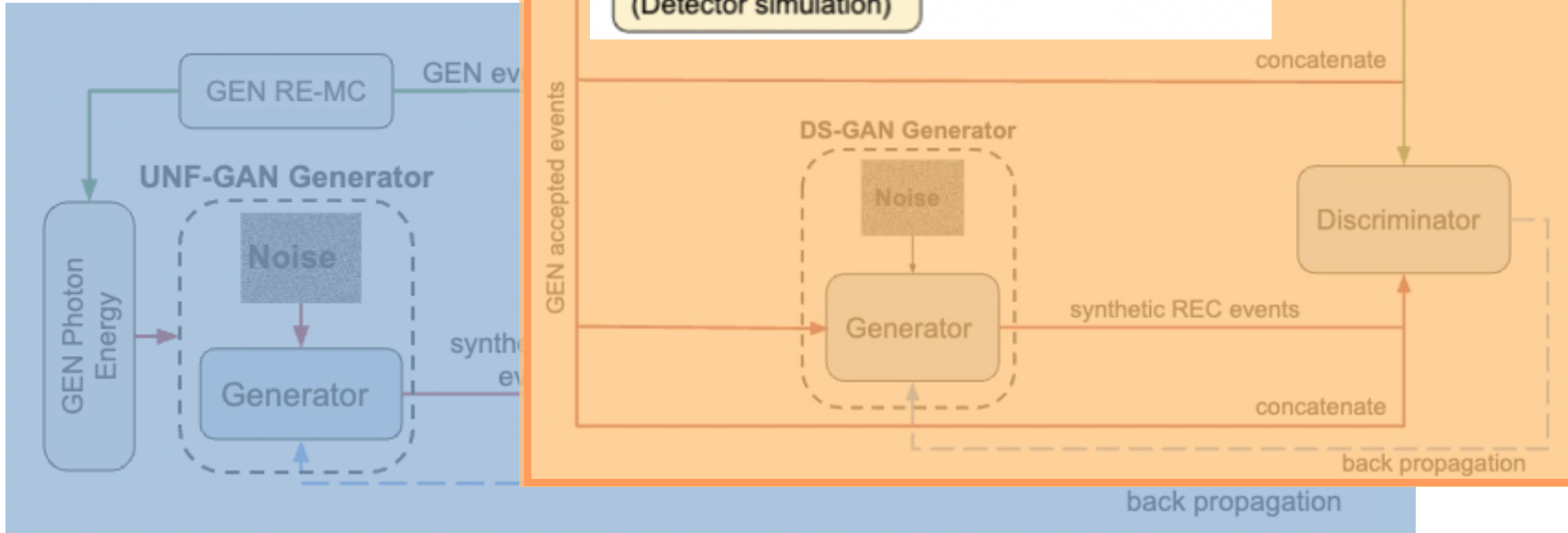
- PS-MC: Phase space Monte Carlo event generator



2π photoproduction closure test

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

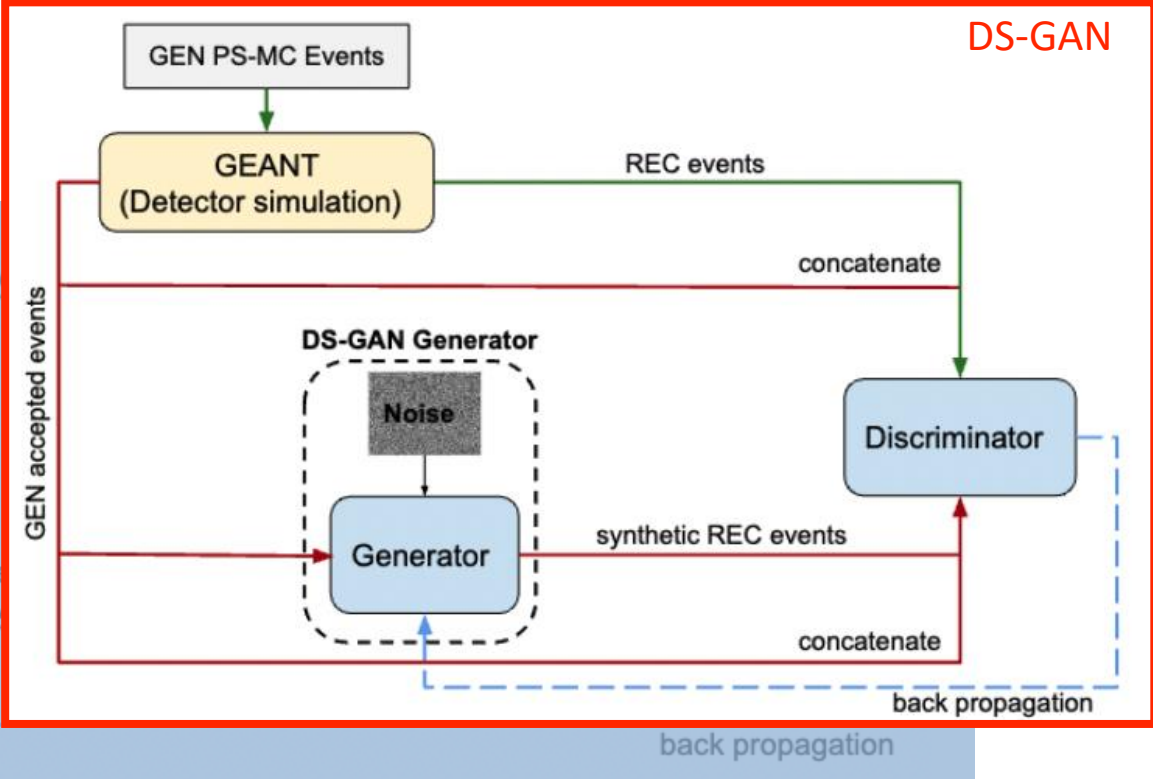
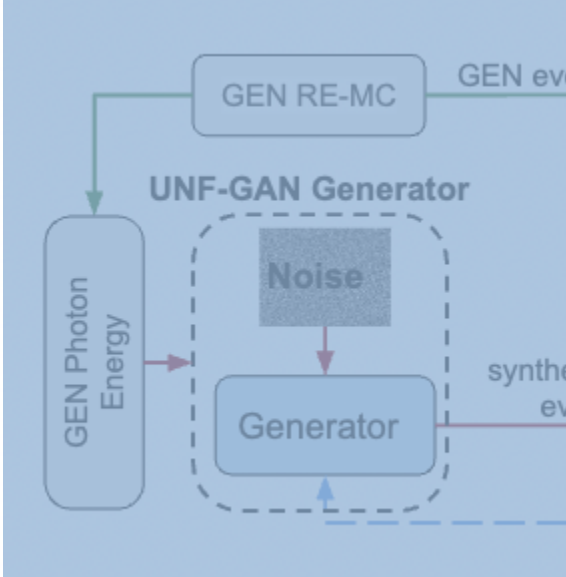
- GSIM-GEANT to simulate CLAS acceptance and resolution



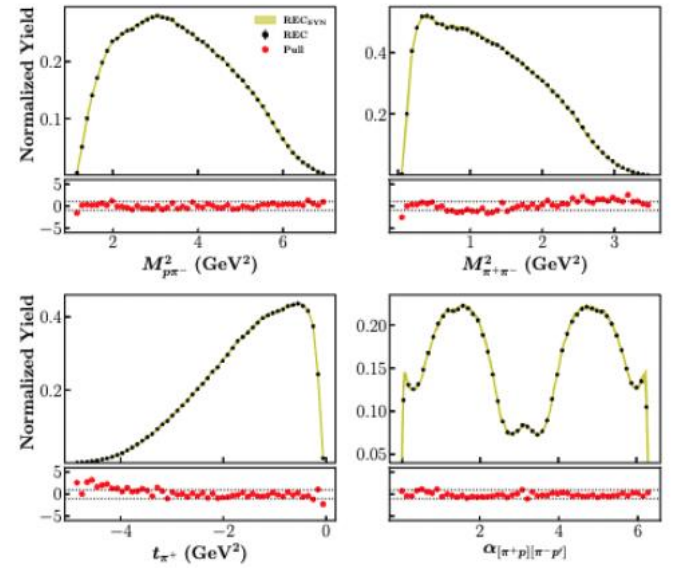
2π photoproduction closure test

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

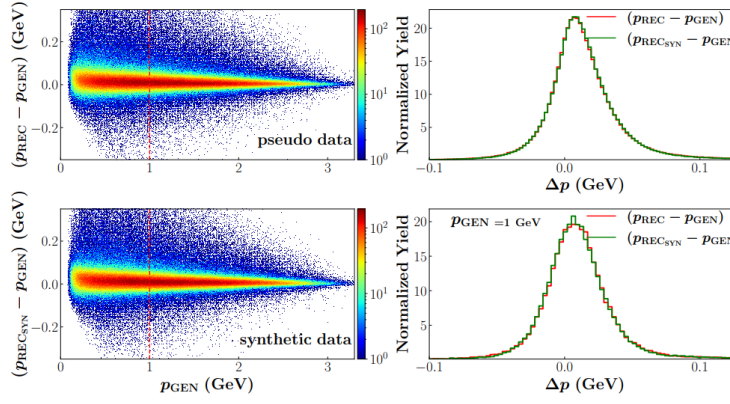
- GSIM-GEANT to simulate CLAS acceptance and resolution



MC REC pseudodata vs. DS-GAN synthetic data



CLAS resolution



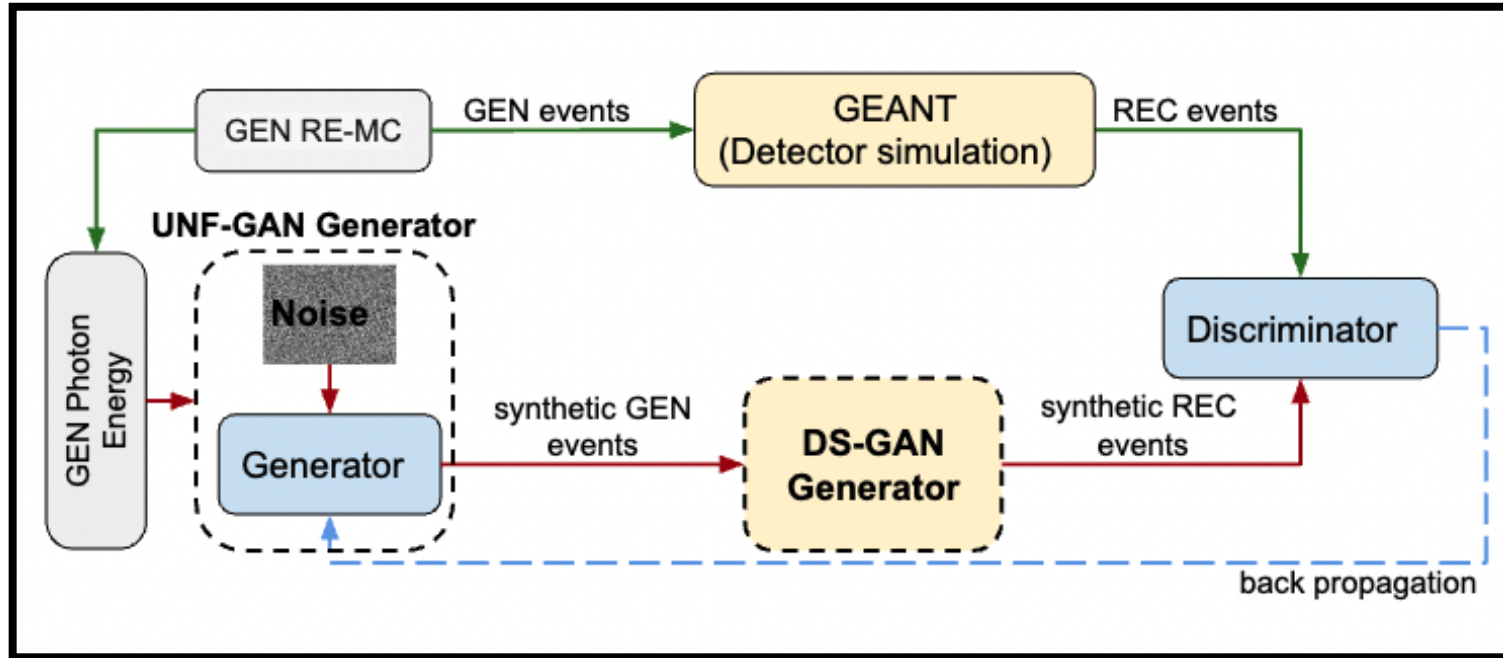
DS-GAN learned the CLAS detector effects!



2π photoproduction closure test

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

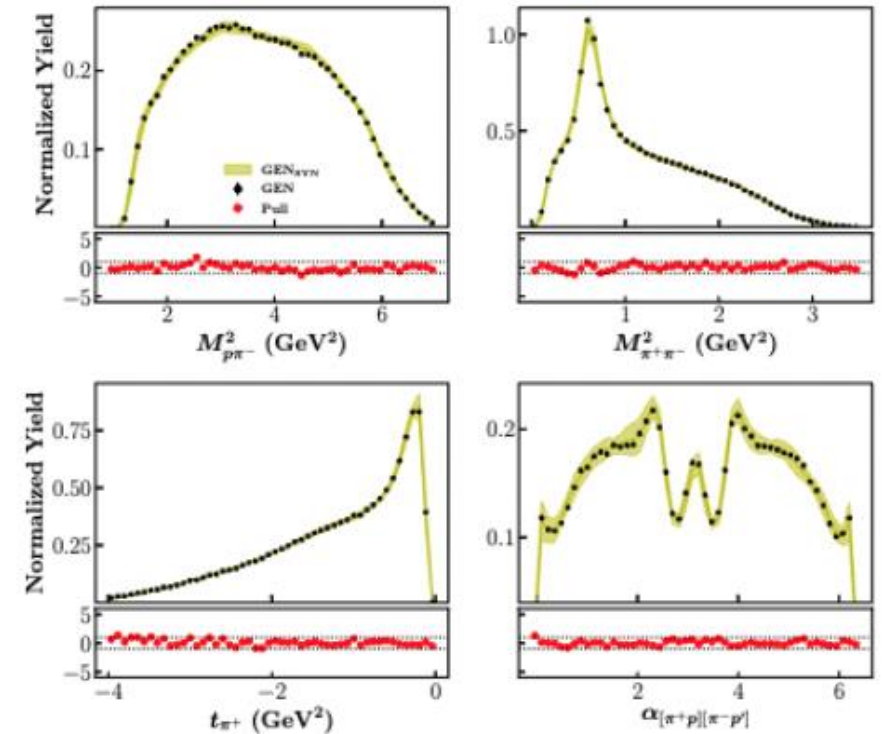
- UNF-GAN trained with REC-MC pseudodata (experimental data proxy)
- DS-GAN used to unfold CLAS detector effects (within acceptance)



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

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

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



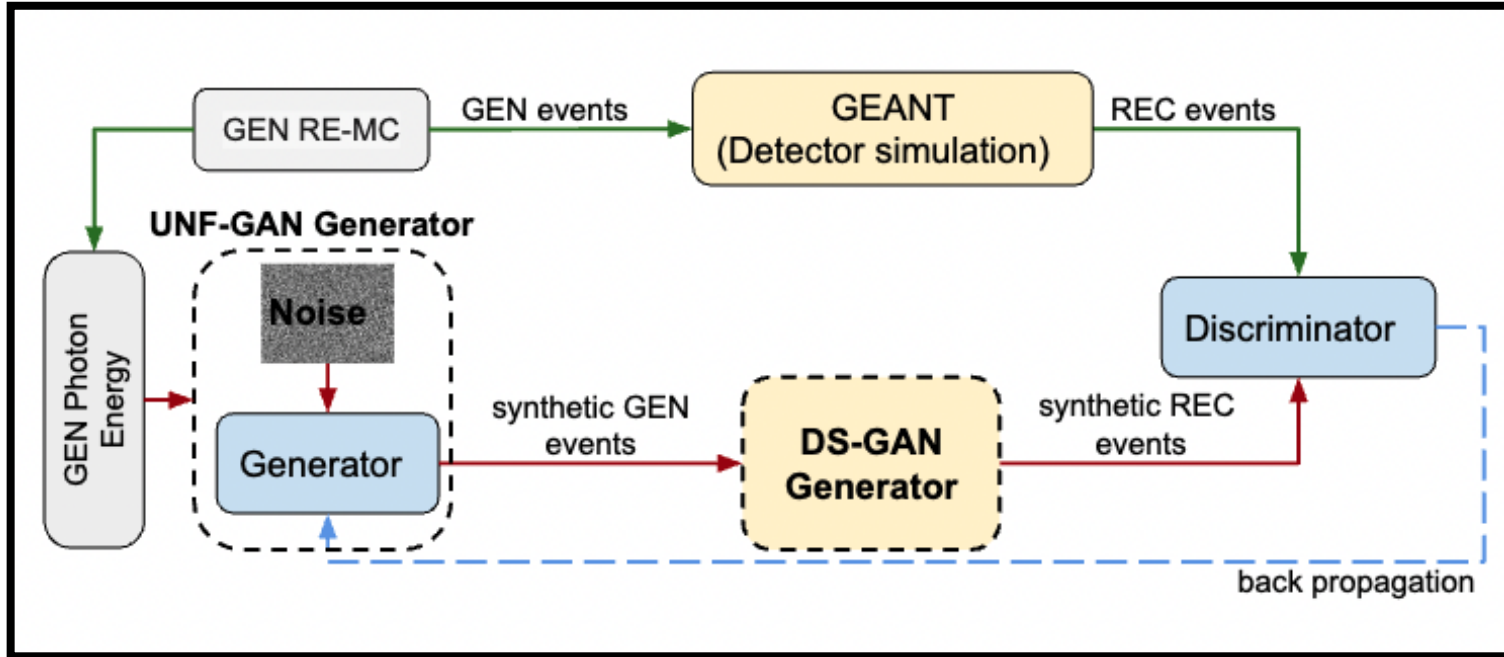
- Systematic of the full procedure (two-GANs) estimated by bootstrap with 20+20 independently trained GANs



2π photoproduction closure test

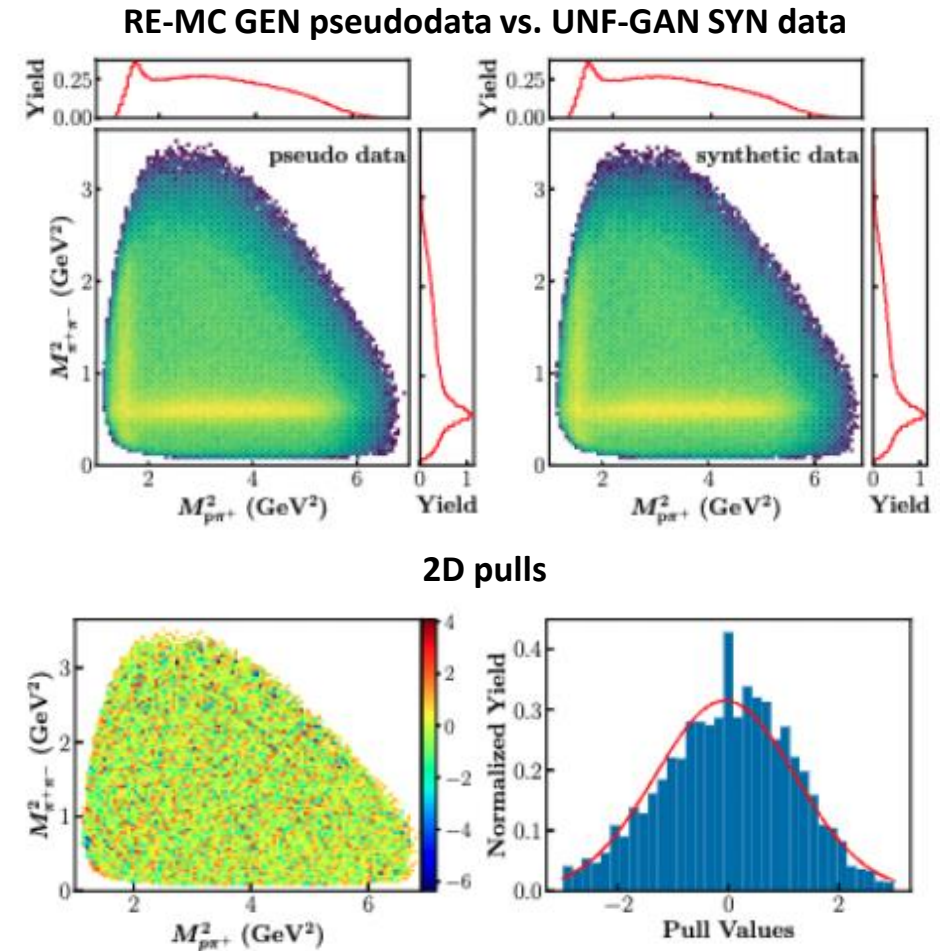
4. Deploy the unfolding GAN (UNF-GAN) that includes the DS-GAN and train it with RE-MC REC pseudodata

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



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

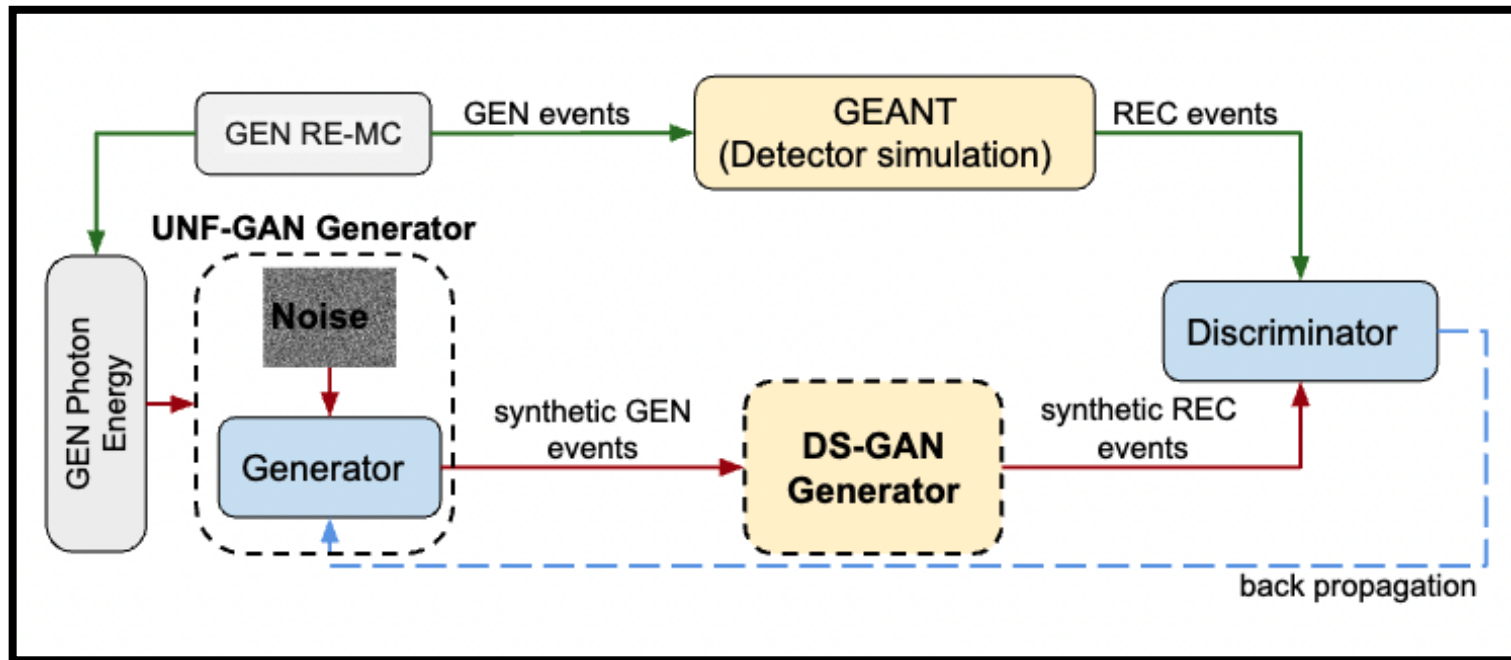
Good agreement ($\pm 1\sigma$) for 2D distributions (correlations)



2π photoproduction closure test

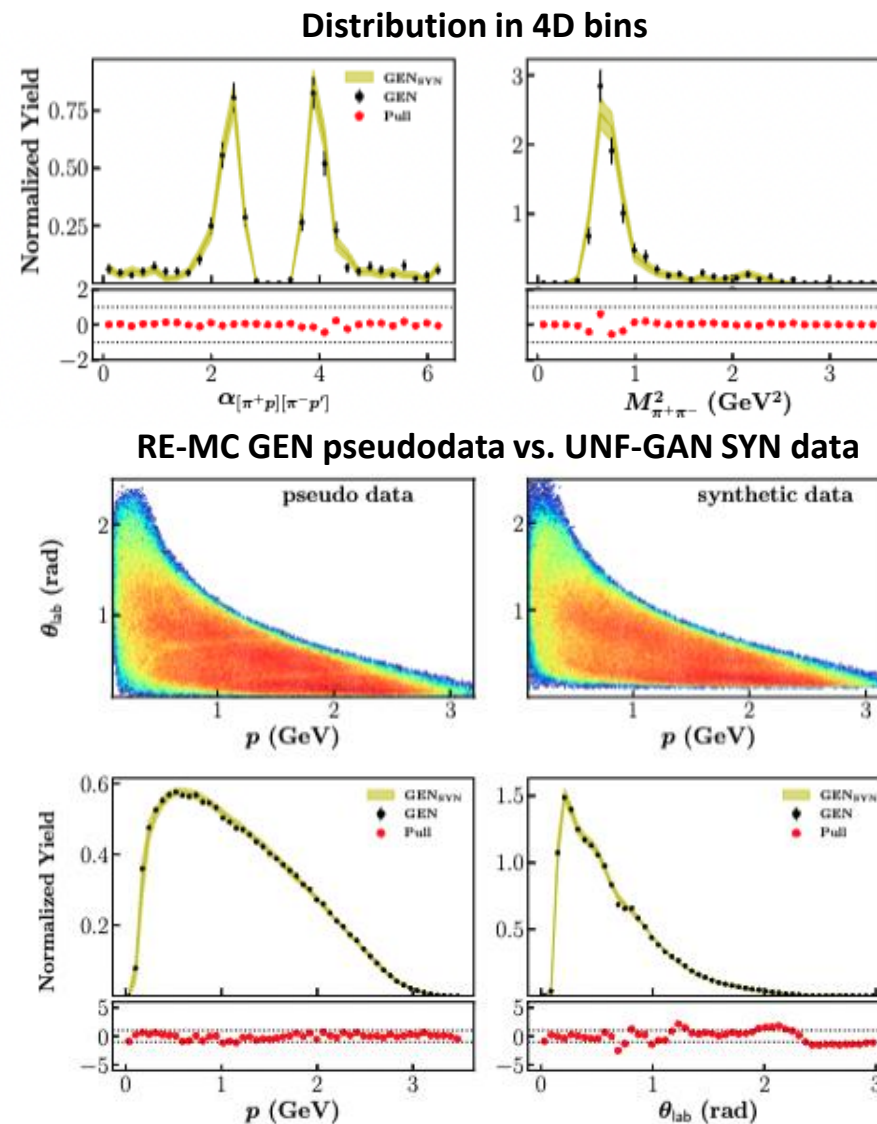
4. Deploy the unfolding GAN (UNF-GAN) that includes the DS-GAN and train it with RE-MC REC pseudodata

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



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

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



Introducing Acceptance

Credit: T. Vittorini, *in progress*

- Simple 2-body process: $\gamma p \rightarrow \Delta^+(1232) \rightarrow \pi^0 p$
- Simple model: Breit-Wigner amplitude with parameters m_Δ and Γ_Δ

$$\begin{aligned} \frac{d\sigma}{d\Omega} &\propto \frac{p_f}{p_i s} \sum_{\lambda_\gamma \lambda_p \lambda'_p} \left| (-)^{\lambda_\gamma} H_{|\lambda_\gamma - \lambda_p|} \frac{d_{\lambda_\gamma - \lambda_p, -\lambda'_p}^{3/2}(\theta)}{m_\Delta^2 - s - i\Gamma_\Delta m_\Delta} \right|^2 \\ &\propto \frac{p_f}{p_i s} \frac{3 |H_{3/2}|^2 + 5 |H_{1/2}|^2 - 3 \cos 2\theta \left(|H_{3/2}|^2 - |H_{1/2}|^2 \right)}{(m_\Delta^2 - s)^2 + \Gamma_\Delta^2 m_\Delta^2} \end{aligned}$$

- Two independent variables for the process: We chose them to be the scattering $\theta_{lab}^{\pi^0}$ angle and the azimuthal $\phi_{lab}^{\pi^0}$ angle in the lab frame

- **Goals:**

- Build a single GAN model which includes all the available phase space of the same reaction in order to understand if extending the range of the measured phase space would improve our knowledge of relevant observables
- Quantify the model dependence on the unmeasured regions of the phase space



Introducing Acceptance

- Implementing CLAS acceptance cuts we define three different measured topologies for the reaction $\gamma p \rightarrow \pi^0 p$

- Topology 0: $\gamma p \rightarrow (\pi^0 p)$ (Unmeasured)
- Topology 1: $\gamma p \rightarrow \pi^0(p)$
- Topology 2: $\gamma p \rightarrow (\pi^0)p$
- Topology 3: $\gamma p \rightarrow \pi^0 p$

- Build a single GAN model which includes all the measured phase space regions + different models in the unmeasured region

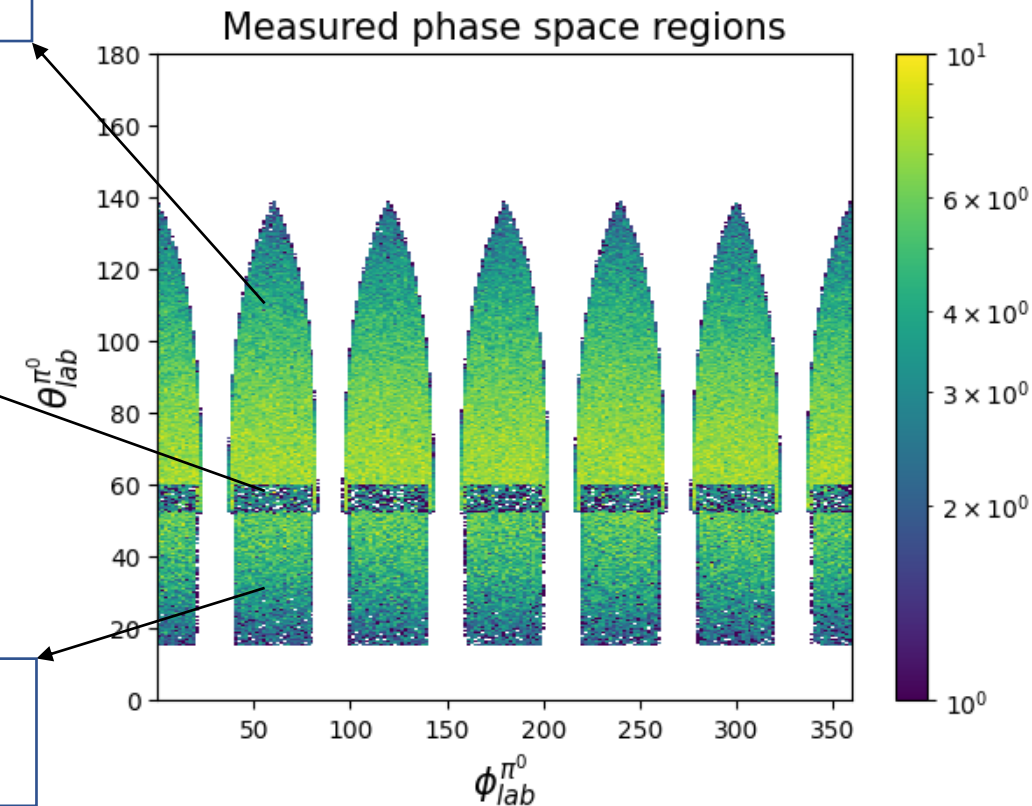
Understand the advantages of a multi-topology approach

1° topology:
Observed π^0

3° topology:
Both particles in
the final state are
observed

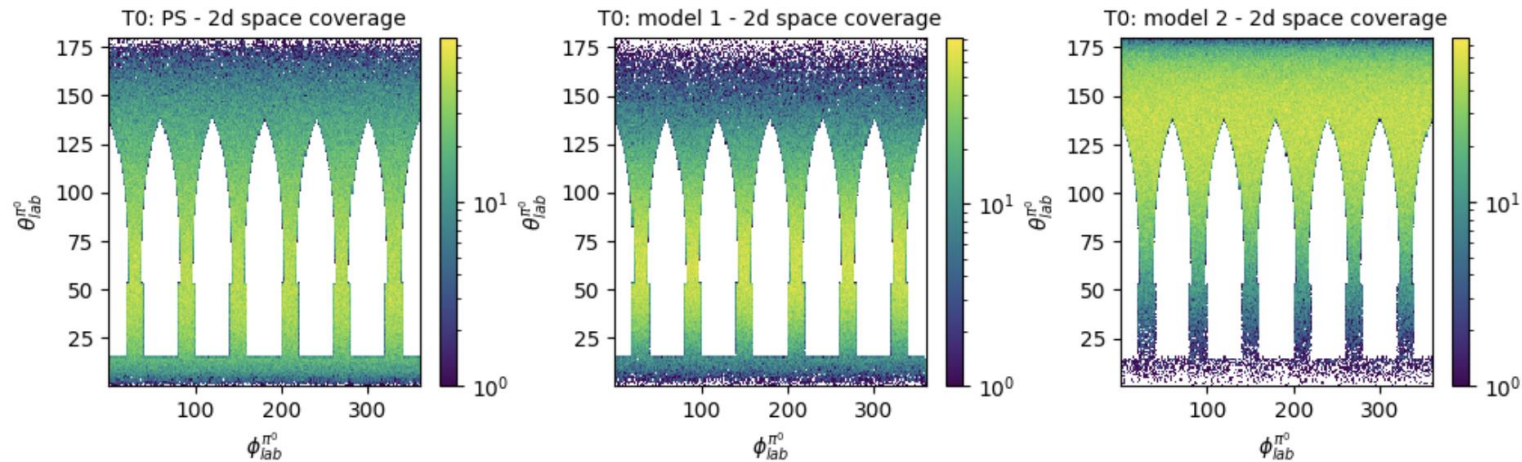
2° topology:
Observed recoil
proton

Combining the topologies in the lab frame
looking at π^0 variables ($\theta_{lab}^{\pi^0}, \phi_{lab}^{\pi^0}$):



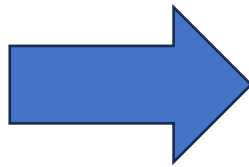
Introducing Acceptance

- Distinguish different models, built modifying the amplitudes, which will cover the unmeasured region of the phase



- Build two different datasets (PS and Model2) to train the same GAN architecture with them to check the model dependence

- $Model1 = (M1_{meas})$ will be considered the equivalent of observed data
- $PS = (PS_{T0}, PS_{meas})$
- $Model 2 = (M2_{T0}, M2_{meas})$



- Histograms considering just the measured quantities (and check that for both dataset the 'true' observed variables are recovered)
- Histogram considering **all the phase space** (measured and unmeasured) to quantify the model dependence
- Nonsharp cuts are currently under study



Towards amplitude reconstruction

Credit: G. Montaña, *in progress*

- Elastic scattering $\pi^+\pi^- \rightarrow \pi^+\pi^-$

$$A(s, \cos \theta) = \sum_{\ell=0}^n (2\ell + 1) f_{\ell}(s) P_{\ell}(\cos \theta)$$

- Breit-Wigner type partial waves $f_{\ell}(s)$, $\ell = 0, 1$

$$f_0(s) = \frac{m_{\sigma} \Gamma_{\sigma}}{m_{\sigma}^2 - s - i \Gamma_{\sigma} m_{\sigma}} \quad m_{\sigma} \approx 0.475 \text{ GeV}, \Gamma_{\sigma} \approx 0.55 \text{ GeV}$$

$$f_1(s) = \frac{m_{\rho} \Gamma_{\rho}}{m_{\rho}^2 - s - i \Gamma_{\rho} m_{\rho}} \quad m_{\rho} = 0.775 \text{ GeV}, \Gamma_{\rho} = 0.147 \text{ GeV}$$

$$\longrightarrow A(s, \cos \theta) = f_0(s) + 3f_1(s) \cos \theta$$

- Differential cross section

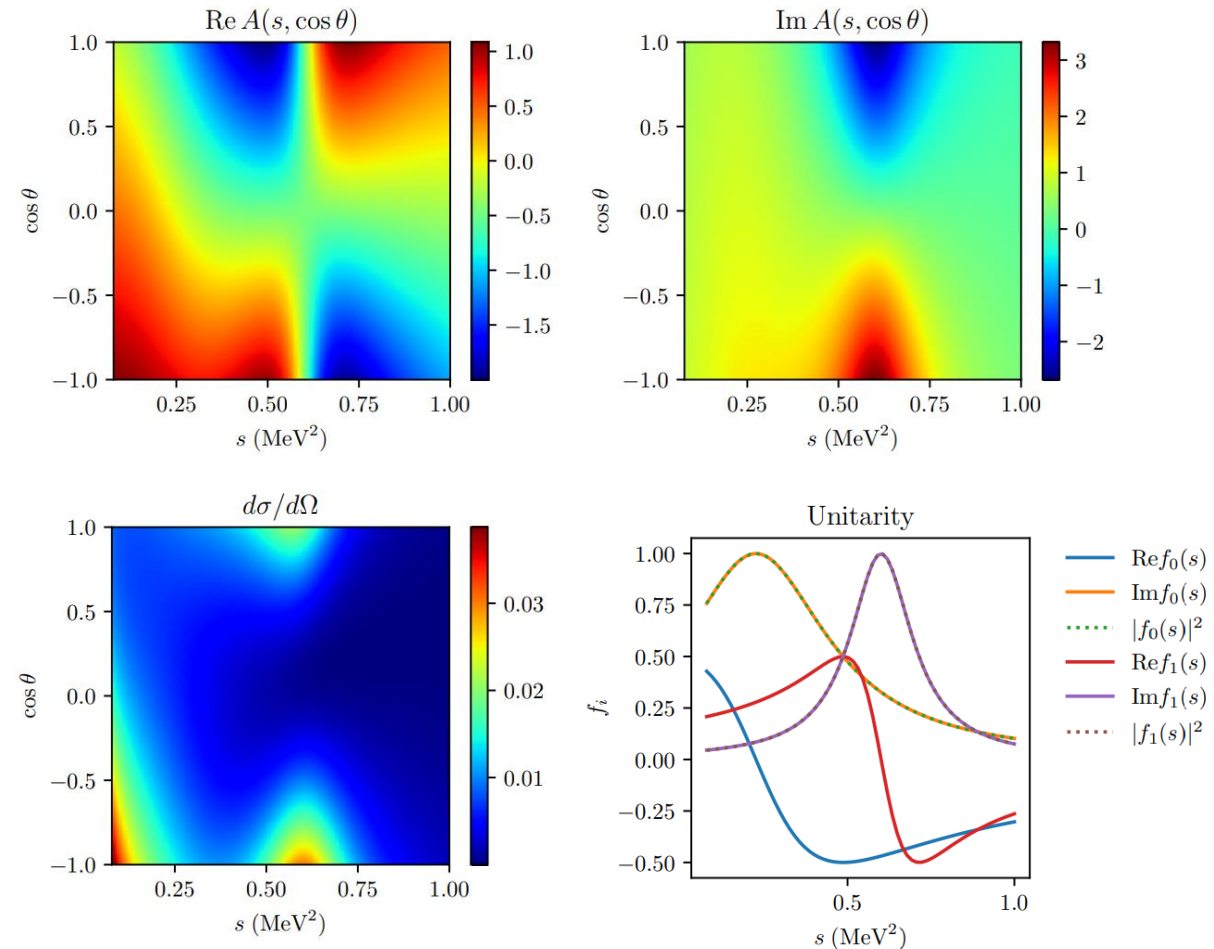
$$\frac{d\sigma}{d\Omega} = \frac{1}{64\pi^2} \frac{1}{s} |A(s, \theta)|^2$$

- Physics constraint: Unitarity of the partial waves

$$\text{Im} f_0(s) = |f_0(s)|^2$$

$$\text{Im} f_1(s) = |f_1(s)|^2$$

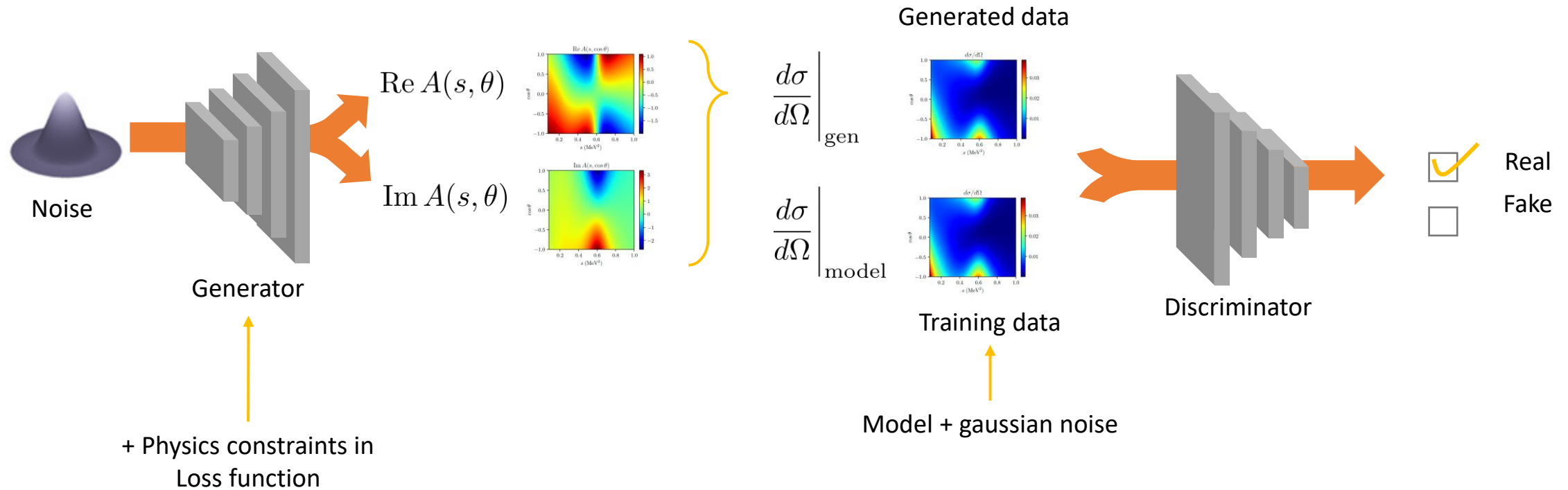
$$f_{\ell}(s) = \frac{1}{2} \int_{-1}^{+1} d(\cos \theta) P_{\ell}(\cos \theta) A(s, \theta)$$



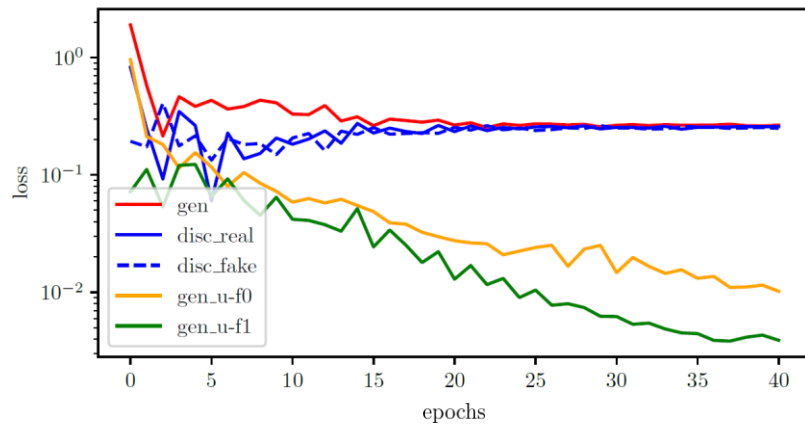
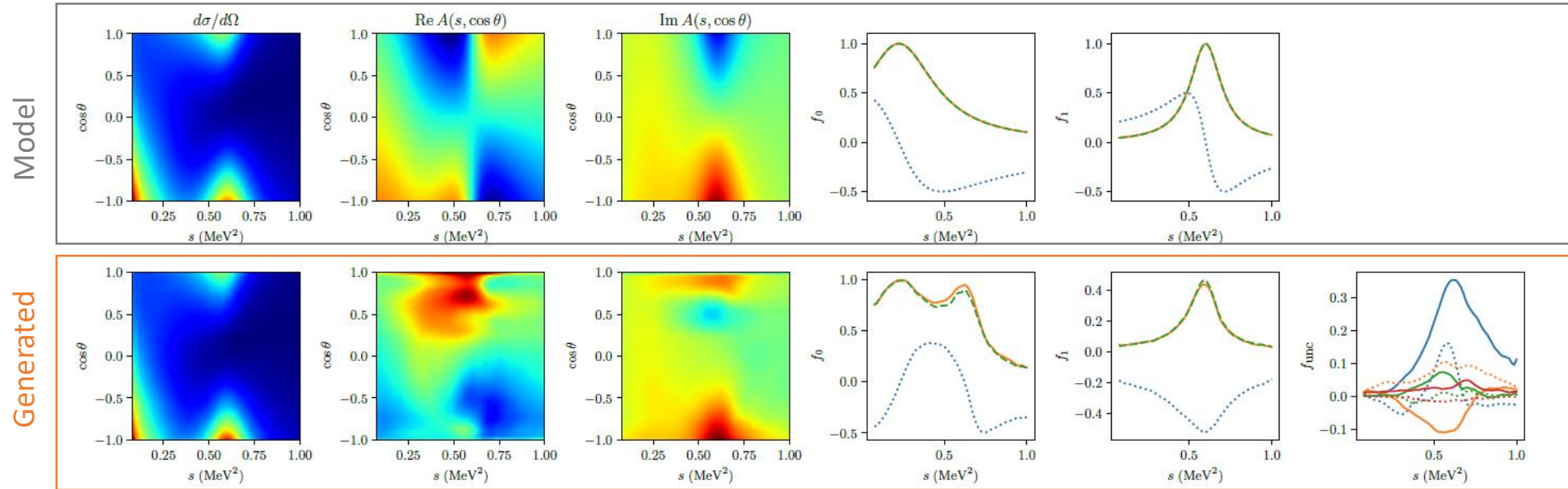
Generative Adversarial Network (GAN) with constraints

Two neural networks, the **generator** and the **discriminator**:

- The **generator** needs to capture the data distribution
- The **discriminator** estimates the probability that a sample comes from the training data rather than from the generator



Preliminary results (i)



- Cross section is reproduced qualitatively
- Unitarity constraint is satisfied
- Partial waves $\ell \geq 2$ are large



More physics constraints

- Unitarity of the partial waves $f_\ell(s)$, $\ell = 0, 1$

$$\text{Im}f_0(s) = |f_0(s)|^2$$

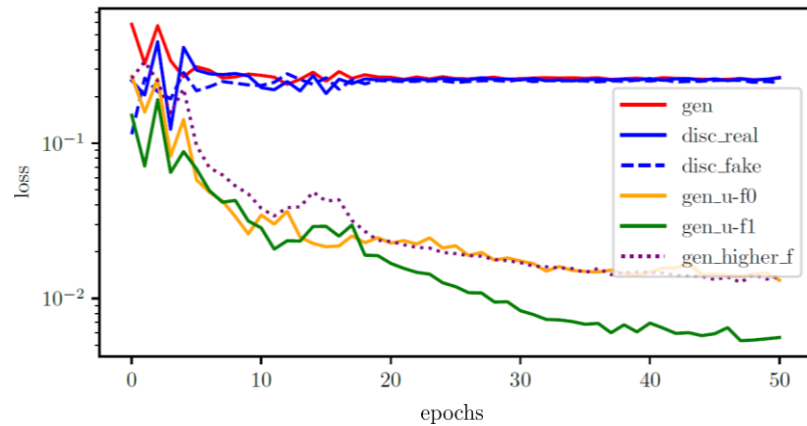
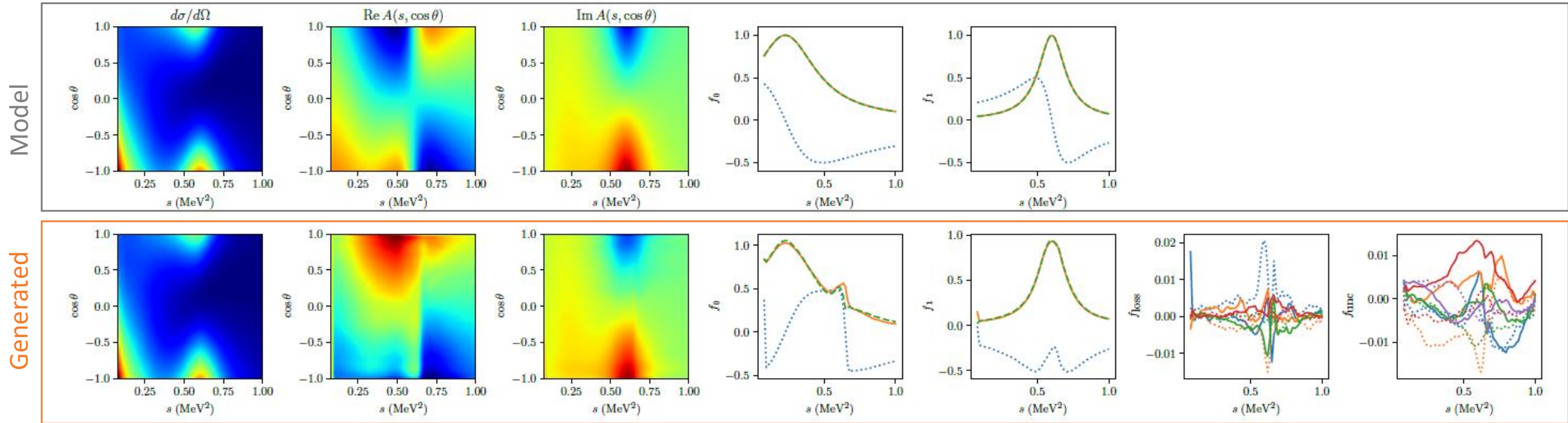
$$\text{Im}f_1(s) = |f_1(s)|^2$$

- Suppression of higher partial waves

$$f_\ell(s) = 0, \ell \geq 2$$



Preliminary results (ii)



- Cross section is reproduced qualitatively
- Unitarity constraint is satisfied
- Partial waves $\ell \geq 2$ are suppressed
- Ambiguity in the sign of the real part



More physics constraints

- Unitarity of the partial waves $f_\ell(s)$, $\ell = 0, 1$

$$\text{Im}f_0(s) = |f_0(s)|^2$$

$$\text{Im}f_1(s) = |f_1(s)|^2$$

- Suppression of higher partial waves

$$f_\ell(s) = 0, \ell \geq 2$$

- Positive derivative of the phase shift $\delta_\ell(s) = \text{atan} \left(\frac{\text{Im}f_\ell(s)}{\text{Re}f_\ell(s)} \right)$

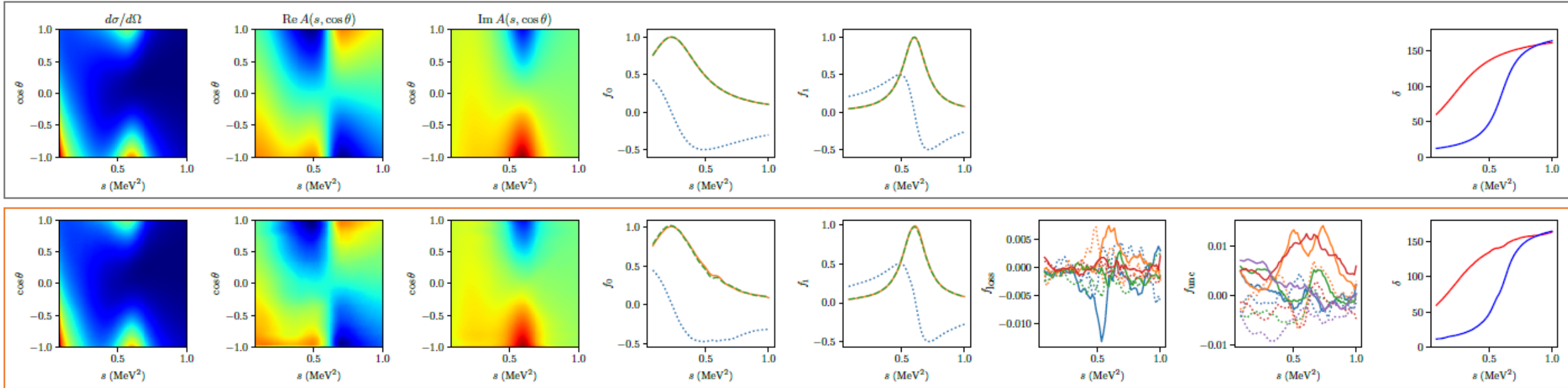
$$\frac{d}{ds} \delta_0(s) \geq 0$$

$$\frac{d}{ds} \delta_1(s) \geq 0$$

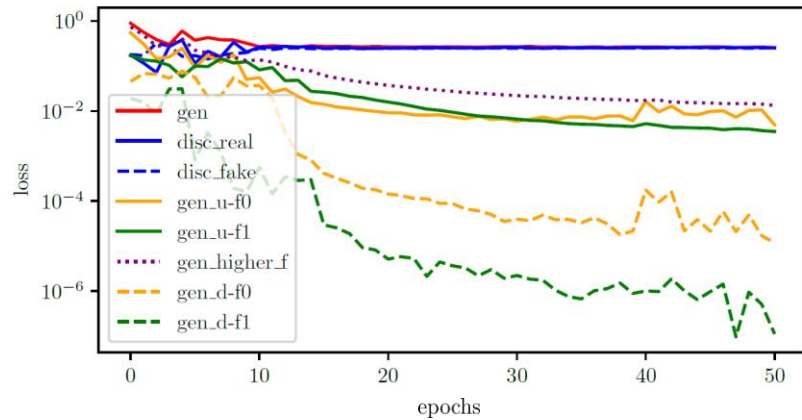


Preliminary results (iii)

Model



Generated



- Cross section is reproduced qualitatively
- Unitarity constraint is satisfied
- Partial waves $\ell \geq 2$ are suppressed
- The real part takes the right sign



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 2π 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 in order to extract a 4D cross-section
- Reach the amplitudes in a model-independent way
- Make this procedure an efficient way to analyse CLAS12 2π 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!



BACKUP

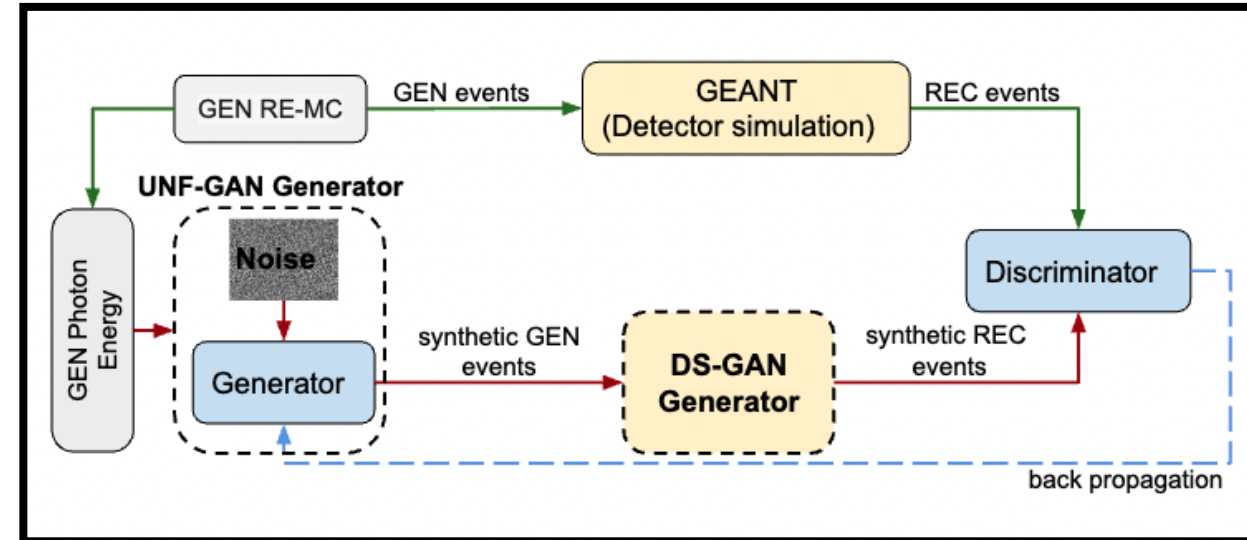


2π photoproduction closure test

- 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 2π 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 independent 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



Exclusive reactions: 2 → 3

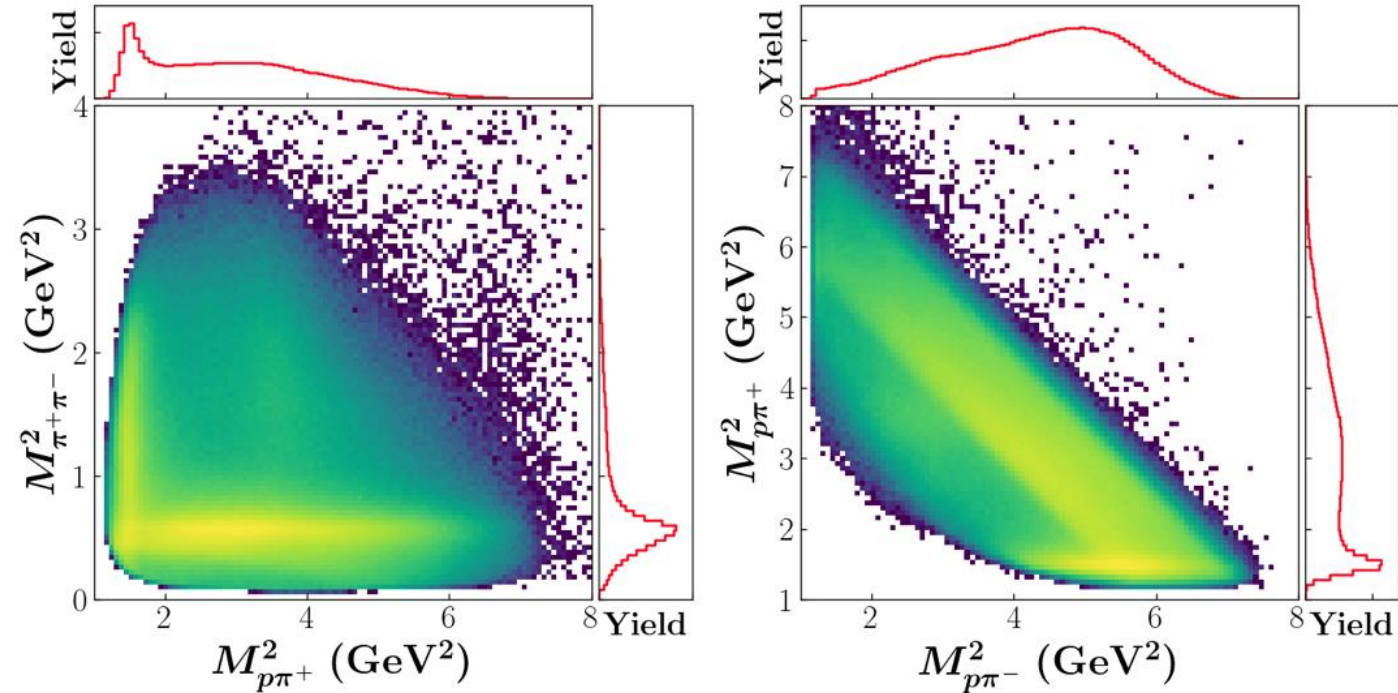
$\gamma p \rightarrow \pi^+ \pi^- p$ (unpolarized)

- Initial state: Fully known
- Final state: 3x3 independent variables
- Independent variables: $(3 \times 3) - 4 = 5$ (E_γ fixed)
- Many possible choices, such as $M_{\pi\pi}^2$, $M_{p\pi}^2$, θ_π , α , ϕ

CLAS g11 2π photoproduction

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

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

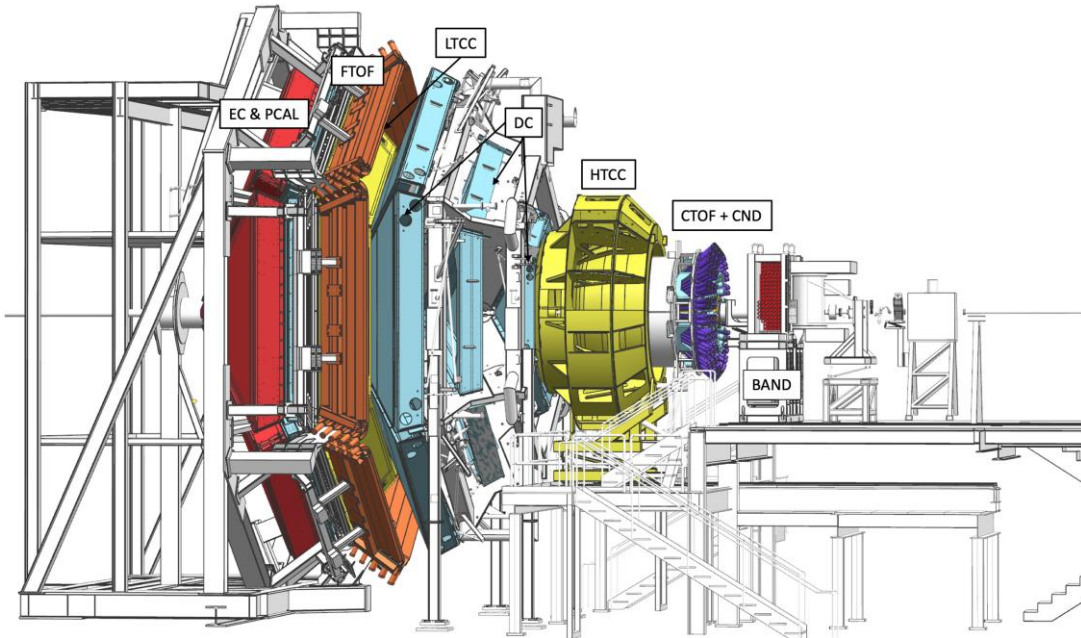


AI could provide a new way to look at data
and to extract observables and physics
interpretation

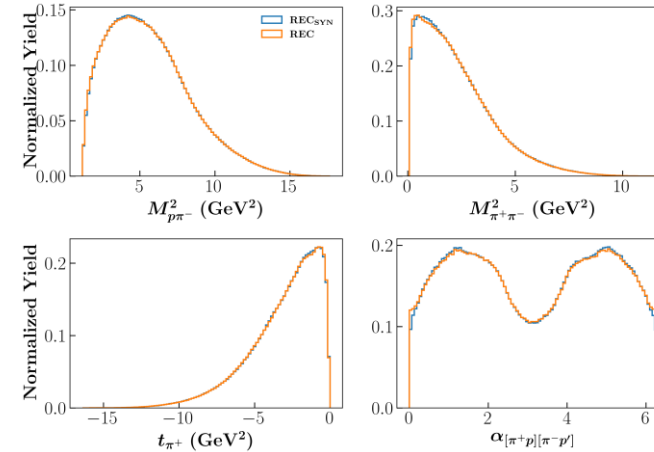


CLAS12 application

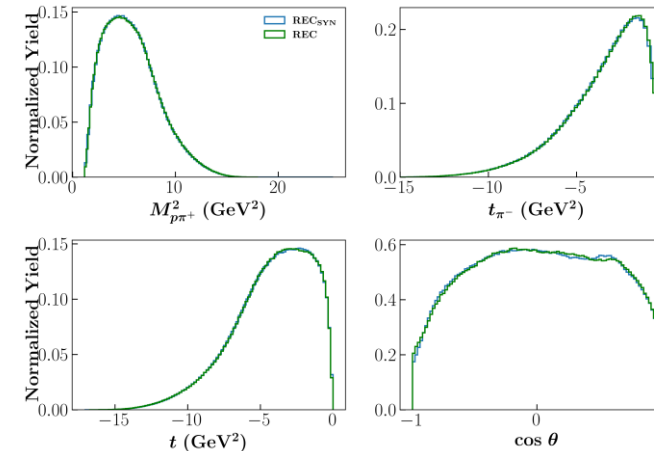
- Working towards the application of the developed machinery to CLAS12 pseudodata
- If this procedure works well on CLAS and CLAS12 data the architecture robustness is guaranteed
- We can put together in a coherent way information from different kinematic regions



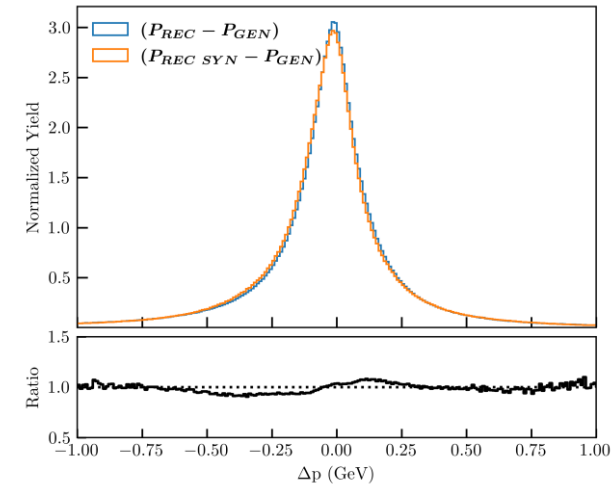
REC SYN vs REC pseudodata training variables



REC SYN vs REC pseudodata derived variables



CLAS12 resolution



Good agreement for training and derived variables

Credit: Derek Glazier, Tareq Alghamdi

