

e- Cab 12

INFN

QNP2022 - The 9th International Conference on Quarks and Nuclear Physics

5-9 September 2022

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

M.Battaglieri (INFN)

on behalf of A(i)DAPT Working Group





$$rac{\mathrm{d}\sigma}{\mathrm{d}\Omega} = (2\pi)^4 m_i m_f rac{p_f}{p_i} ig| T_{fi} ig|^2$$

Differential solid angle d  $\Omega$ 

**S**Lab12

- The cross section is related to the transition probability between an initial to a final state
- In case of scattering, cross sections provides information about the elementary interaction
- Cross section is expressed as squared sum of scattering amplitudes (complex functions) interaction properties
- It is derived by measuring the momentum distributions of reaction particle (at different CM energy)
- Correlations between particles in the final state reflects the underlying dynamics
- Cross sections fully replaces the 4-mom data sample in a compact and efficient way
- Cross section is the starting point for any higher level physics analysis



depending on the kinematic Lorentz-invariant of the problem and embedding the

• Traditional approach: particles (4-momenta) measured into the detector, extract the relevant observables, extract physics mechanisms

• Cross section **preserves** this information as replacement for the original particle-by-particle scattering information

Exclusive reactions:  $2 \rightarrow 2$ 



A(i)DAPT AI for Data Analysis and Preservation

e- 🚱 Lab 12

## $2 \rightarrow 2$ scattering (no polarisation)

- Initial state: known
- Final state: 2 x 3
- Parameters:  $(2 \times 3) 4 = 2$
- Possible choice: -t and  $\phi$
- the physics depends only on one variable (-t)
- It worked (and still works!) well if limited to channels with a single variable
- Xsec, Polarization observables, angular distribution, decay matrix, ...

- $2 \rightarrow 3$  scattering (no polarization)
- Initial state: known
- Final state:  $3 \times 3$
- Parameters:  $(3 \times 3) 4 = 5$
- Possible choice:  $M^2_{\pi\pi}$ ,  $M^2_{\rho\pi}$ ,  $\theta_{\pi}$ ,  $\alpha$ ,  $\phi$

## **CLAS gII** $2\pi$ photo production

- $E_{V} = (3.0 3.8) \text{ GeV}$ 
  - $\gamma p \rightarrow p \pi^+ \pi^-$  exclusive reaction
  - data set analyses so far  $\gamma p \rightarrow p \pi^+$ ( $\pi$ ) + small contamination of  $\gamma p \rightarrow$ p π<sup>+</sup> (more than a missing π<sup>-</sup>)
  - complicated dynamic for the overlap of  $(p\pi)$  to form  $\Delta$  baryon resonances and  $(\pi\pi)$  to form meson resonances



- It does not work (in practice) when you have several independent variables: multi-particle final states (spectroscopy) or multi-variable correlations (SIDIS)
- In the integration to reduce to I-dim all correlations are lost



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



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

# **Detector unfolding**

- Detector effects make measured observables (detector-level) DIFFERENT from 'true' observables (vertex-level)
  - Acceptance: any measurements only access a limited region of the phase space. How to recover the unmeasured region?
    - Interpolation: holes in the phase space
    - Extrapolation: border of the accessible phase space
  - **Resolution**: any measurements introduce an experimental resolution that may hide or washout the effect searched for
    - A spike could be not resolved, the measurement may extend in an unphysical region (e.g. negative squared missing mass)
- For both effects, one needs to quantify the systematic errors introduced to the vertex-level observables
- Mitigation strategy:
  - Acceptance: *fiducial volumes*' to exclude unmeasured regions verifying the training convergence
  - Resolution: closure test with a reasonable model of the detector using a detector proxy (parametric or GEANTbased)



## Generative Adversarial Network (GANs)

• The colored boxes are built using NNs

e- Cab12

- Discriminator is trained to output "real" for Nature samples
- Generator is trained to fool the discriminator
- The Generator can be used as data compression tool
- Typical size for the Generator: O(IMB) to be compared to NP/HEP experiments data set O(IGB/TB)
- Simple to distribute instead of stored events on tapes

6



### https://doi.org/10.24963/ijcai.2021/588



## **ML Event Generator GAN scheme**



- I00-d white noise entered at 0, unit standard dev.
- Generator: 5 hidden layers / 512 neurone per layer, ReLU activation function. Last layer connected to 2 neurons output to generate  $V_1$  and  $V_2$  variables
- Discriminator: same NN architecture as for the generator
- Detector proxy: similar architecture
- Least Squares GAN (LSGAN)
- Trained adversarially for 100000 epochs (pass through the training data set)
- Adam's optimizer

A(i)DAPT AI for Data Analysis and Preservation



• eic-smear: parametric smearing routine for the Electron Ion Collider detectors (no GEANT-based simulations) Parameters tuned to reproduce ZEUS/H1 detectors

• Full  $4\pi$  acceptance

•

## I) GAN training w/o detector effects

Pseudo-data sample (JAM)

e- Cab 12

- Inclusive electron DIS generated at  $E_{CM}$ =318.2 GeV (HERA kinematics)
- 2-dim differential cross section  $d\sigma/dxdQ^2$
- Lorentz boosted from CM to Lab (+ uniform azimuthal angle)
- To reduce violation of momentum conservation on the edge of the phase space due to smearing effects, electron momentum is replaced by new variables:

$$u_1 = \ln \left( (k'_0 - k'_z)/1 \,\text{GeV} \right),$$
 $\nu_2 = \ln \left( (2E_e - k'_0 - k'_z)/1 \,\text{GeV} \right),$ 

## Uncertainty Quantification via pull calculation

- Metric: pull  $pull = \frac{E[\mathcal{P}(\mathcal{O}|bin)]_{GAN} E[\mathcal{P}(\mathcal{O}|bin)]_{JAM}}{\sqrt{V[\mathcal{P}(\mathcal{O}|bin)]_{GAN} + V[\mathcal{P}(\mathcal{O}|bin)]_{JAM}}}$
- Bootstrap with 10 independently trained GANs





## I) GAN training w/o detector effects

Pseudo-data sample (JAM)

e- 🕄 Lab12

INFN

- Inclusive electron DIS generated at  $E_{CM}$ =318.2 GeV (HERA kinematics)
- 2-dim differential cross section  $d\sigma/dxdQ2$
- Lorentz boosted from CM to Lab (+ uniform azimuthal angle)
- To reduce violation of momentum conservation on the edge of the phase space due to smearing effect, electron momentum is replaced by new variables:

$$u_1 = \ln \left( (k'_0 - k'_z)/1 \,\text{GeV} \right),$$
 $\nu_2 = \ln \left( (2E_e - k'_0 - k'_z)/1 \,\text{GeV} \right),$ 

## Uncertainty Quantification via *pull* calculation

• Metric: 
$$pull$$
  $pull = \frac{E[\mathcal{P}(\mathcal{O}|bin)]_{GAN} - E[\mathcal{P}(\mathcal{O}|bin)]_{JAM}}{\sqrt{V[\mathcal{P}(\mathcal{O}|bin)]_{GAN} + V[\mathcal{P}(\mathcal{O}|bin)]_{JAM}}}$ 

• Bootstrap with 10 independently trained GAN





### **No Detector Effects**

# 0.8**II) GAN training WITH detector effects** • eic-smear introduces significant distortions to the 0.6 detector level sample in particular on $V_2$ 0.40.20.0 0.810.6 0.40.2 0.0-2

e Cab12

INFN



## **II)** GAN training with detector effects



A(i)DAPT AI for Data Analysis and Preservation



12

e-Cab12

INFN

ERA M data	• The
	• 'gen gene
	• 'reco on r
	• 'ge recc bars
1e-2 x = 2e-2	the • The reco
x = 3e-2 $x = 5e-2$	
104	

## Conclusions

- The 'closure test' was successful
- 'generated' events (trained on generated)
- 'reconstructed' events (trained on reconstructed)
- 'generated' (trained on reconstructed) with larger error bars, in particular on the edge of the phase space
- The PDFs are correctly recovered

## CLAS gll data set

e- 💰 Lab 12

- Same data set used by CLAS Collaboration for many publications
- Fiducial cuts ( $p, \Theta, \phi$ ) as used in published analysis
- All four topologies available but only focused on  $\gamma p \rightarrow p \pi^+ (\pi)$
- Final exclusive  $2\pi$  state identified by missing mass technique (variables reconstructed by energy/momentum conservation)
- Multipion background comes from  $\gamma p \rightarrow p \omega^0 \rightarrow p \pi^+ \pi^- \pi^0$
- At  $E_{\gamma}=3-4$  GeV reaction dynamics dominated by  $\rho^0$  photo production ( $\gamma p \rightarrow p \rho^0$ ) and  $\Delta^{++}$  resonance excitation ( $\gamma p \rightarrow \Delta^{++} \pi^-$ )



13

## Pseudo-data event generator

### • Two event generators

I) REALISTIC used to mimic real data. Includes measured xsec, angular distributions, dominant mechanism ( $\rho^0$  and  $\Delta^{++}$  channels) II) PHASE SPACE used to train the detector proxy GAN. Based on phase space with uniform distributions

• Reconstructed events are passed through the full analysis chain (GEANT detector simulation + CLAS data reconstruction code)



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



A(i)DAPT AI for Data Analysis and Preservation



## GAN training on data (no detector unfolding)

- Train a NN to generate events (synthetic) with the SAME correlations of experimental data
- Replace the xsec with a NN (synthetic data are equivalent to data)
- Light GAN-based event generator can generate any statistics

**GAN** architecture



- INVARIANTS
- detector acceptance/resolution
- whole data set

- observables will be extracted from  $Y_{LM}$ )



**S**Lab12

• Tested different options: 4-moms in CM or LAB or using kinematics

• Implemented a folding/unfolding procedure to take into account

• Training performed on 10% of the data set to use the others for systematic studies. When optimised the training will be done on the

• Results are validated by comparing data/synt I-dim projected and selected bins in 5-d space in LAB, CM and INVARIANTS space

• Error quantified to include statistical (via bootstrap) and systematic (different data samples, different reference systems)

• The final check is performed comparing  $Y_{LM}$  moments extracted from data and synt-data (in the assumption that any further physics

# Synthetic vs data (no detector unfolding)







16

A(i)DAPT AI for Data Analysis and Preservation

 $(\mathbf{b})$ 

0.5

100

 $(\mathbf{c})$ 



## Moments of the angular distribution (no detector unfolding)



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

17

e-Slab12



• Moments extracted from synthetic data copy moments extract from data

## tSNE analysis (no detector unfolding)







e-Cab12

MC-Realistic (gen)

### • t-SNE analysis strategy

- explore correlations with multi-pion final states
- Select the prominent  $\Delta$ ++(1232) peak and check correlations
- Select areas of t-SNE1.vs.t-SNE2 space to demonstrate meson/baryon systems separation



18

### orrelations rate meson/baryon systems separation

M.Battaglieri - INFN



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



A(i)DAPT AI for Data Analysis and Preservation

## **Detector proxy**



## **Closure test (using simulations)**







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



A(i)DAPT AI for Data Analysis and Preservation

# **Closure test (using simulations)**



A(i)DAPT AI for Data Analysis and Preservation

21

e- Cab12

INFN

## ML to access multi-d cross section



### Goals

- Cross section: embed multi-d cross section information (correlation) in a data-trained event generator
- <u>Preserve</u> data in an alternative compact and efficient form (to be applied to current JLab physics program)
- <u>Statistics</u>: use the NN to determine the necessary statistics for a given analysis
- <u>Statistics</u>: overcome statistics limitation exploiting ML super-resolution
- <u>Detector Efficiency</u>: folding/unfolding detector effects to extract physics at vertex level (via sim or data)
- Physics analysis: incorporate Universality (of scattering amplitudes) training the NN with different kinematics of the same final state or different final states (coupled channels)
- <u>Physics analysis</u>: extract from the NN features related to the underlying physics
- <u>Physics analysis</u>: structure the NN to reflect amplitudes properties (poles, cuts, dynamics, ...)
- Collaborative effort ML Data manipulation Validation Unfolding
  - Theory
- Regular weekly meeting

- develop a procedure to best fit data (ML Group)
- develop a procedure to compare synt-data to data (Validation Group)
- develop a procedure to quantify the error associated to sent-data (Data manipulation Group)
- develop a procedure to take into account the detector effect (Folding/Unfolding Group)
- extract physics form data and sent-data and compare (Theory Group)
- Wiki page: https://clasweb.jlab.org/wiki/index.php/A%28I%29DAPT\_-\_AI\_for\_Data\_Analysis\_and\_PresevaTion

A(i)DAPT AI for Data Analysis and Preservation