



Unfolding Data with Machine Learning



vmikuni@lbl.gov

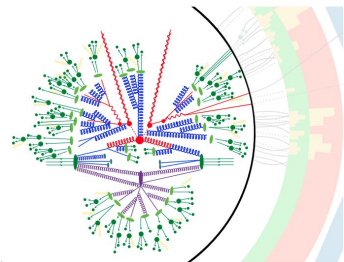


vinicius-mikuni

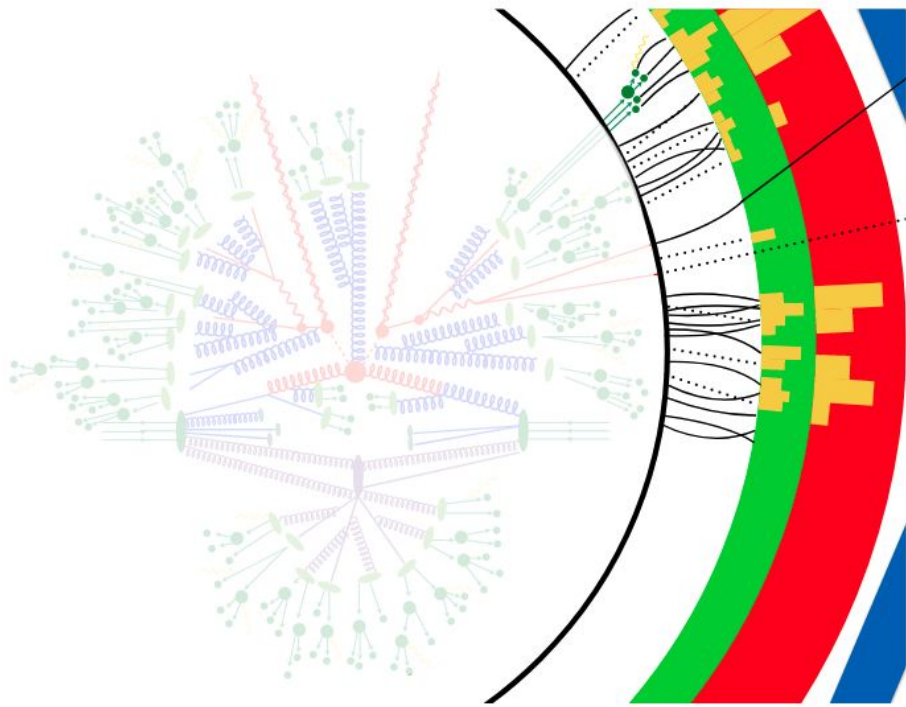
Vinicius Mikuni



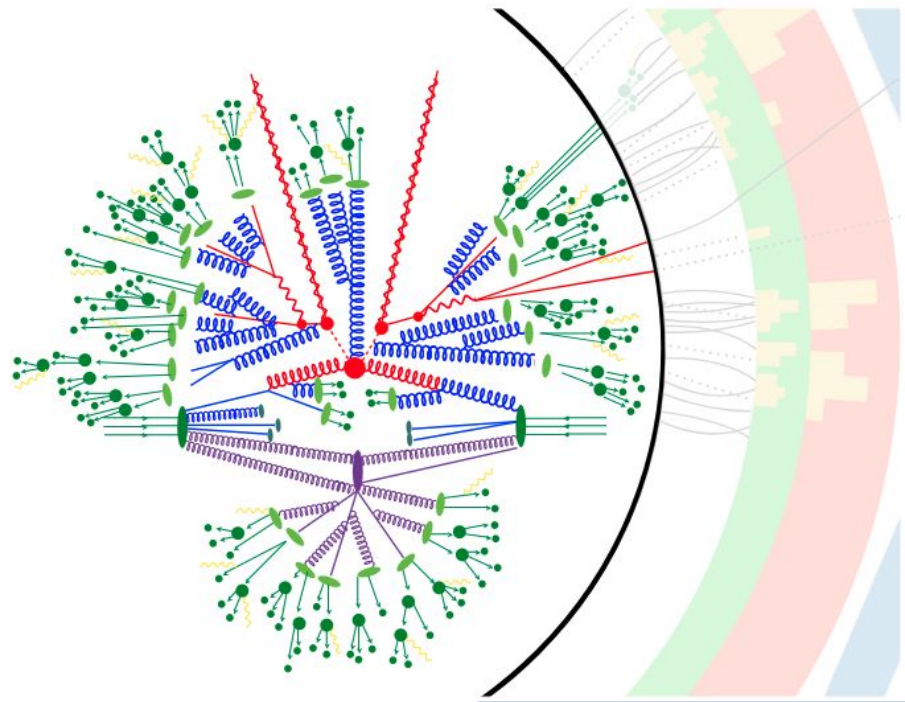
Unfolding



What we measure

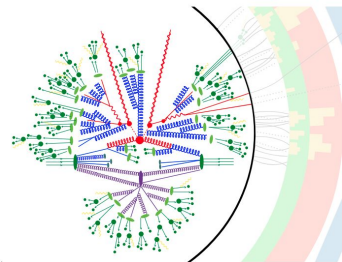


What we want

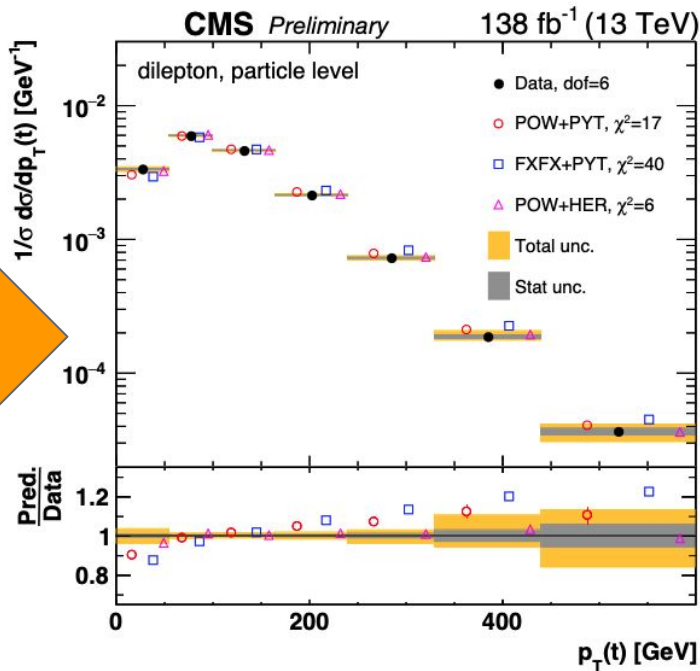
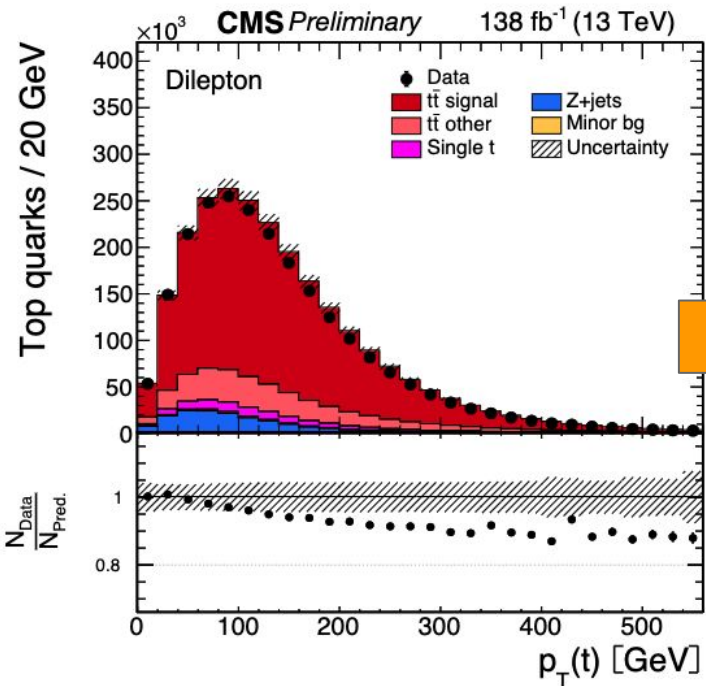




Unfolding



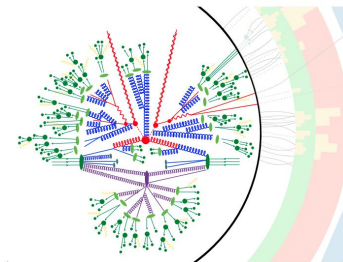
Source: CMS-PAS-TOP-20-006



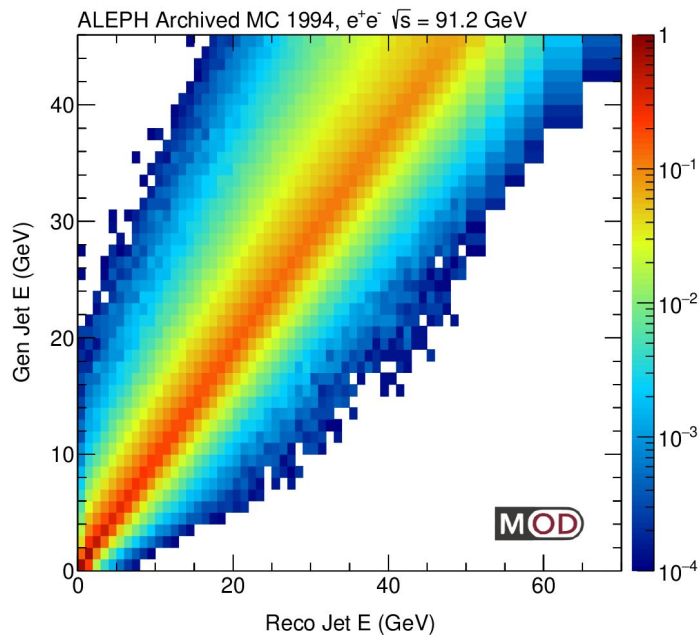
Traditional methods for **unfolding** use **histograms**



Unfolding



Source: MITHIG-MOD-21-001



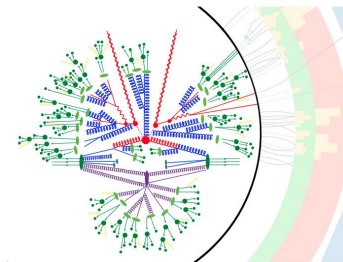
$$a_i = R_{ij} b_j$$

R_{ij} is the response matrix: **P(observed in bin i | true in bin j)**

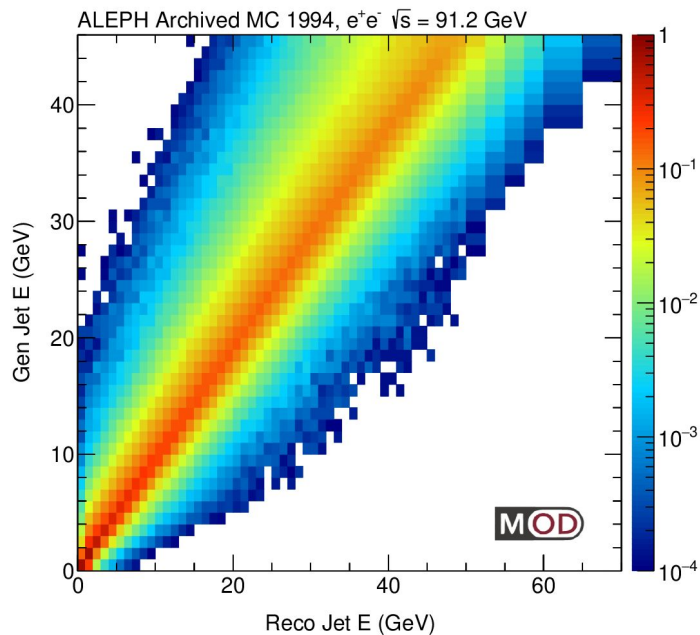
Traditional methods for **unfolding** use **histograms**



Unfolding



Source: MITHIG-MOD-21-001



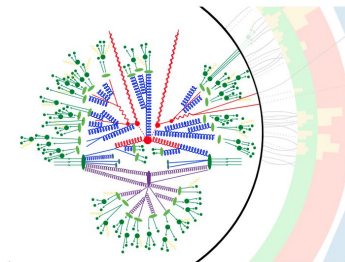
$$a_i = R_{ij} b_j$$

- R_{ij} is the response matrix: **P(observed in bin i | true in bin j)**
- Traditional unfolding is all about inverting the matrix R_{ij}

Traditional methods for **unfolding** use **histograms**



Unfolding



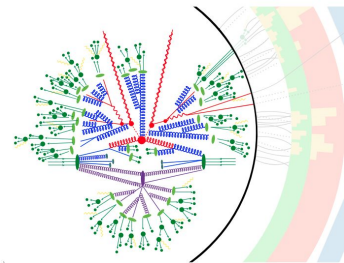
How to define the **optimal binning**?

- Choice depends on the **distribution** and **phase space**
- Need to compromise when **combining** results from **different experiments**





Unfolding



How to define the **optimal binning**?

- Choice depends on the **distribution** and **phase space**
- Need to compromise when **combining** results from **different experiments**

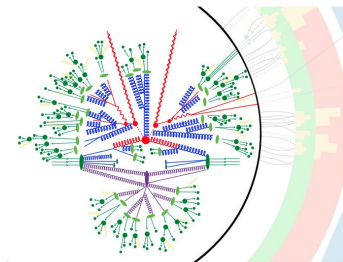
How to include **multiple distributions**?

- Histograms are hard to scale: **curse of dimensionality**
- Unfolding uncertainties can be reduced using **additional observables**





Unfolding



How to define the **optimal binning**?

- Choice depends on the **distribution** and **phase space**
- Need to compromise when **combining** results from **different experiments**

How to include **multiple distributions**?

- Histograms are hard to scale: **curse of dimensionality**
- Unfolding uncertainties can be reduced using **additional observables**

How to unfold distributions that are **not** defined for each event?

- Moments of distributions
- Energy Correlators

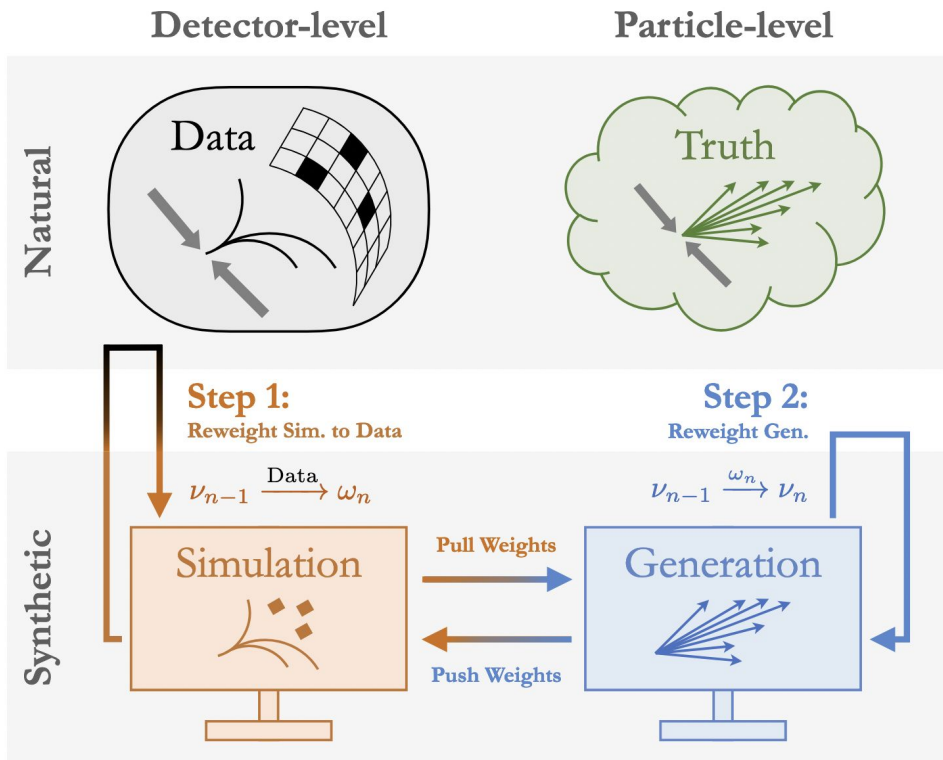




Going beyond histograms: Omnifold*

For unfolding using **invertible networks** see:

- SciPost Phys. 9 (2020) 074 e-Print: [2006.06685](https://arxiv.org/abs/2006.06685)



ML is used to define a method for unfolding that is unbinned and can use multiple distributions at a time

2 step iterative approach

- Simulated events after detector interaction are reweighted to match the data
- Create a “new simulation” by transforming weights to a proper function of the generated events

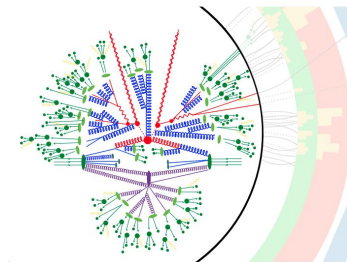
Machine learning is used to approximate **2** likelihood functions:

- **reco MC to Data** reweighting
- **Previous** and **new Gen** reweighting

* Andreassen et al. PRL 124, 182001 (2020)



Omnifold



EUROPEAN ORGANISATION FOR NUCLEAR RESEARCH (CERN)



Submitted to: Phys. Rev. Lett.



CERN-EP-2024-132
May 31, 2024

A simultaneous unbinned differential cross section measurement of twenty-four Z+jets kinematic observables with the ATLAS detector

The ATLAS Collaboration

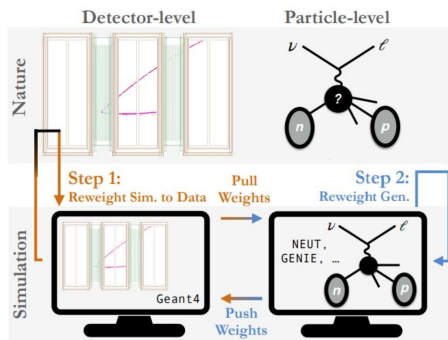
CMS Physics Analysis Summary

Contact: cms-pag-conveners-smp@cern.ch

2024/06/03

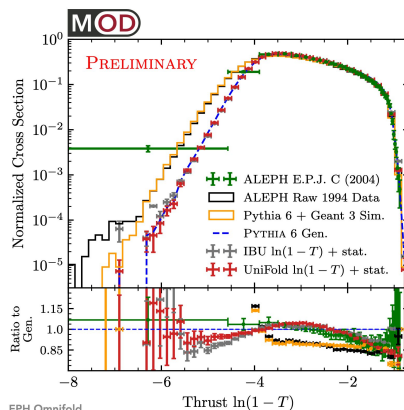
Measurement of event shapes in minimum bias events from pp collisions at 13 TeV

The CMS Collaboration



Measurement of lepton-jet correlation in deep-inelastic scattering with the H1 detector using machine learning for unfolding

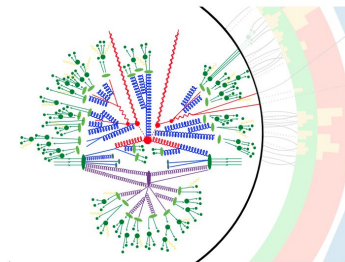
V. Andreev,²³ M. Arratia,³⁵ A. Baghadasaryan,⁴⁶ A. Baty,¹⁶ K. Begzsuren,³⁹ A. Belousov,^{23,*} A. Bolz,¹⁴ V. Boudry,⁵¹ G. Brandt,¹² D. Britzger,²⁶ A. Buniatyan,⁶ L. Bystritskaya,²² A.J. Campbell,¹⁴ K.B. Cantun Avila,⁴⁷ K. Cerny,²⁸ V. Chekelian,²⁸ Z. Chen,³⁷ J.G. Contreras,⁴⁷ L. Cunqueiro Mendez,²⁷ J. Cvach,³³ J.B. Dainton,¹⁹ K. Daim,⁴⁵ A. Deshpande,³⁸ C. Diaconu,³¹ G. Eckerlin,¹⁴ S. Egli,⁴³ E. Elen,¹⁴ L. Favart,⁴ A. Fedotov,⁵² J. Fellesse,¹² M. Fleischer,¹⁴ A. Fomenko,²³ C. Gal,³⁹ J. Gayler,¹⁴ L. Goerlich,¹⁷ N. Gogitidze,²³ M. Gouzevitch,⁴² C. Grab,⁴⁹ T. Greenshaw,¹⁹ G. Grindhammer,²⁶ D. Haidt,¹⁴ R.C.W. Henderson,¹⁹ J. Hessler,²⁹ J. Hladky,³³ D. Hoffmann,²¹ R. Horstberger,⁴³ T. Hrousof,⁵ F. Huber,¹⁵ P.M. Jacobs,⁵ M. Jaquet,²⁹ T. Janssen,⁴ A.W. Jung,⁴⁴ H. Jung,¹⁴ M. Kapichine,¹⁰ J. Katzy,¹⁴ C. Kiesling,²⁶ M. Klein,¹⁹ C. Kleinwort,¹⁴ H.T. Klest,³⁸ R. Kogler,¹⁴ P. Kostka,¹⁹ J. Kretschmar,¹⁹ D. Krücker,¹⁴ K. Krüger,¹⁴ M.P.J. Landon,²⁰ W. Lange,⁴⁸ P. Laycock,⁴¹ S.H. Lee,³ S. Levonian,¹⁴ W. Li,¹⁶ J. Lin,¹⁶ K. Lipka,¹⁴ B. List,¹⁴ J. List,¹⁴ B. Lobodzinski,²⁶ E. Malinovsky,²³ H.-U. Martyn,¹ S.J. Maxfield,¹⁹ A. Mehta,¹⁹ A.B. Meyer,¹⁴ J. Meyer,¹⁴ S. Mikocki,¹⁷ M.M. Mondal,³⁸ A. Morozov,¹⁰ K. Müller,⁵⁰ B. Nachman,⁵ Th. Naumann,⁴⁸ P.R. Newman,⁶ C. Niebuhr,¹⁴ G. Nowak,¹⁷ J.E. Olsson,¹⁴ D. Ozerov,⁴³ S. Park,³⁸ C. Pascaud,²⁹ G.D. Patel,¹⁹ E. Perez,¹¹ A. Petrukhin,⁴² I. Picuric,³² D. Pitzl,¹⁴ R. Polifka,³⁴ S. Preins,³⁵ V. Radescu,³⁰ N. Raicevic,³² T. Ravandani,³⁹ P. Reimer,²³ E. Rizvi,²⁰ P. Robmann,⁵⁰ R. Roosen,⁴ A. Rostovtsev,²⁵ M. Rotaru,⁷ D.P.C. Sankey,⁸ M. Sauter,¹⁵ E. Sauvan,^{21,2} S. Schmitt,¹⁴ B.A. Schmookler,³⁸ L. Schoeffel,¹² A. Schöning,¹⁵ F. Seifow,¹⁴ S. Shushkevich,²⁴ Y. Soloviev,²³ P. Sopicki,¹⁷ D. South,¹⁴ V. Spaskov,¹⁰ A. Specka,³¹ M. Steder,¹⁴ B. Stella,³⁶ U. Straumann,⁵⁰ C. Sun,³⁷ T. Sykora,³⁴ P.D. Thompson,⁹ D. Traynor,²⁰ B. Tsepeldorj,^{39,40} Z. Tu,⁴¹ A. Valkárová,³⁴ C. Vallée,²¹ P. Van Mechelen,⁴ D. Wegener,⁹ E. Wünsch,¹⁴ J. Žáček,³⁴ J. Zhang,³⁷ Z. Zhang,²⁹ R. Zlebčík,³⁴ H. Zohrabyan,⁴⁶ and F. Zomer²⁹ (The H1 Collaboration)



Increasing adoption by experimental collaborations: **ATLAS, CMS, H1, T2K, Aleph**



Omnifold



Reco level

● Data ○ MC

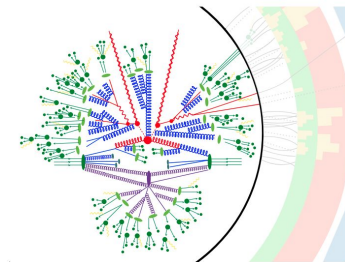


Generator level

● Data (○) MC



Omnifold



Reco level

● Data ○ MC

Iteration 1



Step 1:

- Train a classifier to separate **data** from **MC** events
- Reweight **reco level MC** with weights:

$W(\text{reco}) =$

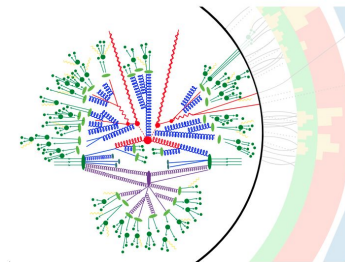
$$p_{\text{Data}}(\text{reco}) / p_{\text{MC}}(\text{reco})$$

Generator level

● Data (○) MC



Omnifold



Reco level

● Data ○ MC

Iteration 1



Step 2:

- Pull weights from **step 1** to generator level events
- Train a classifier to separate **initial MC at gen level** from **reweighted MC** events
- Define a **new simulation** with weights that are a **proper function of gen level kinematics**

$$W(\text{gen}) = \frac{p_{\text{weighted}}}{p_{\text{MC}}(\text{gen})}$$

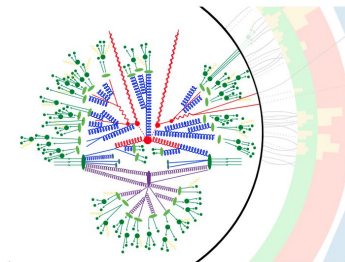


Generator level

● Data (○) MC (○) MC reweighted



Omnifold



Reco level

● Data ○ MC

Iteration 1



Start again from **step 1** using the **new simulation** after **pushing** the weights from **step 2**

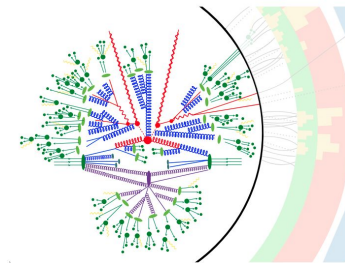
- Guaranteed convergence to the maximum likelihood estimate of the generator-level distribution when number of iterations go to infinite
- In practice, less than 10 iterations are enough to achieve convergence

Generator level

● Data () MC



Omnifold



Reco level

● Data ○ MC

Iteration N



Start again from **step 1** using the **new simulation** after **pushing** the weights from **step 2**

- **Guaranteed convergence** to the maximum likelihood estimate of the generator-level distribution when number of iterations goes to infinite
- In practice, **less than 10 iterations** are enough to achieve convergence

Generator level

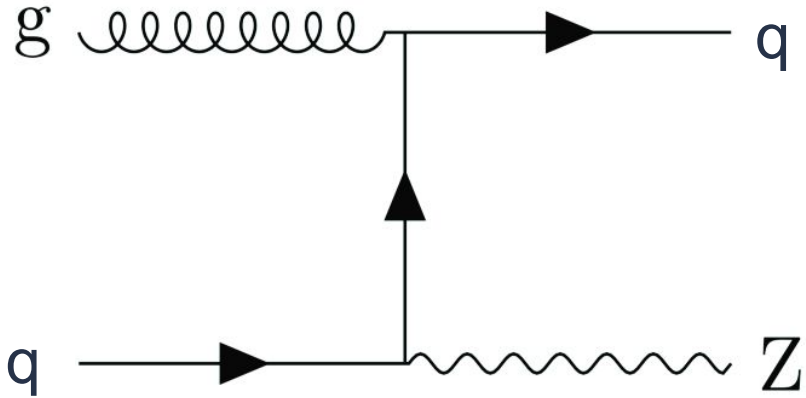
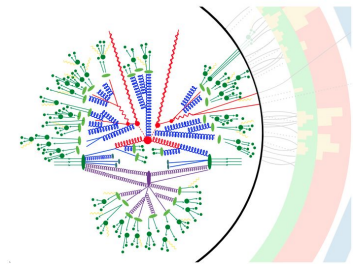
● Data (○) MC

Part 2

Applications



OmniFold dataset

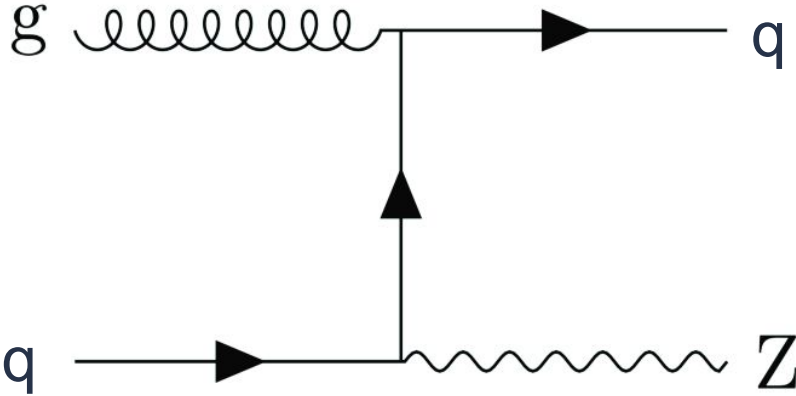
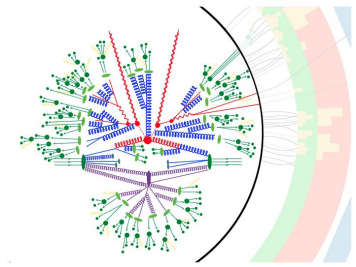


We are going to unfold **6 jet substructure observables** simultaneously using **OmniFold**

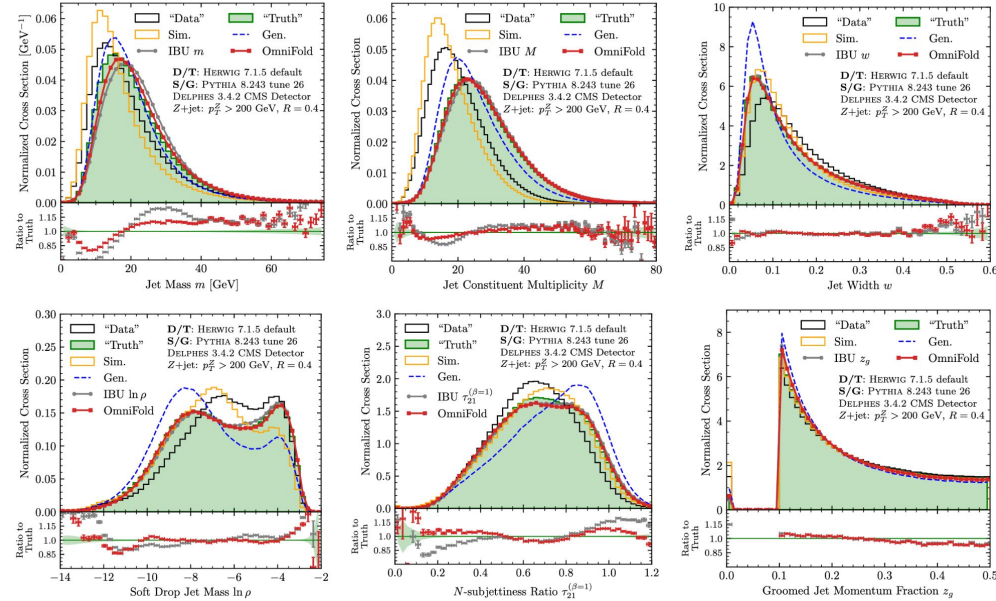
Only consider Z decaying to neutrinos:
mostly a single jet per event.



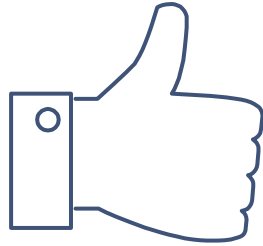
OmniFold dataset



Only consider Z decaying to neutrinos:
mostly a single jet per event.



Phys. Rev. Lett. 124, 182001 (2020)



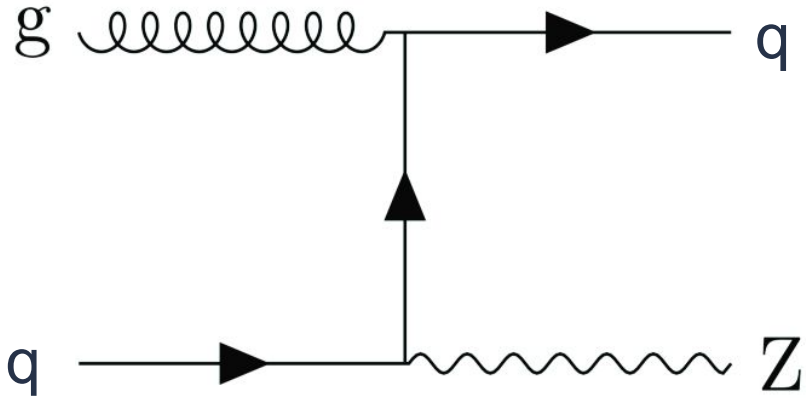
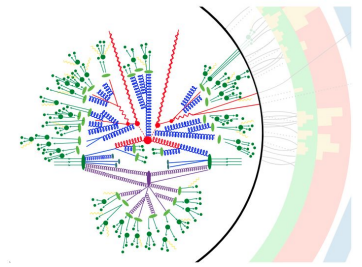
THANKS!

Any questions?

Backup



OmniFold dataset



Only consider Z decaying to neutrinos:
mostly a single jet per event.

Observables:

- ▷ Jet mass
- ▷ Particle Multiplicity
- ▷ $\tau_{21} = \tau_2 / \tau_1$ see Energ. Phys. 2012, 93 (2012).
- ▷ Jet width (τ_1)
- ▷ $\log \rho = 2 \log M_{SD} / p_T$
- ▷ Momentum fraction z_g after using Soft Drop