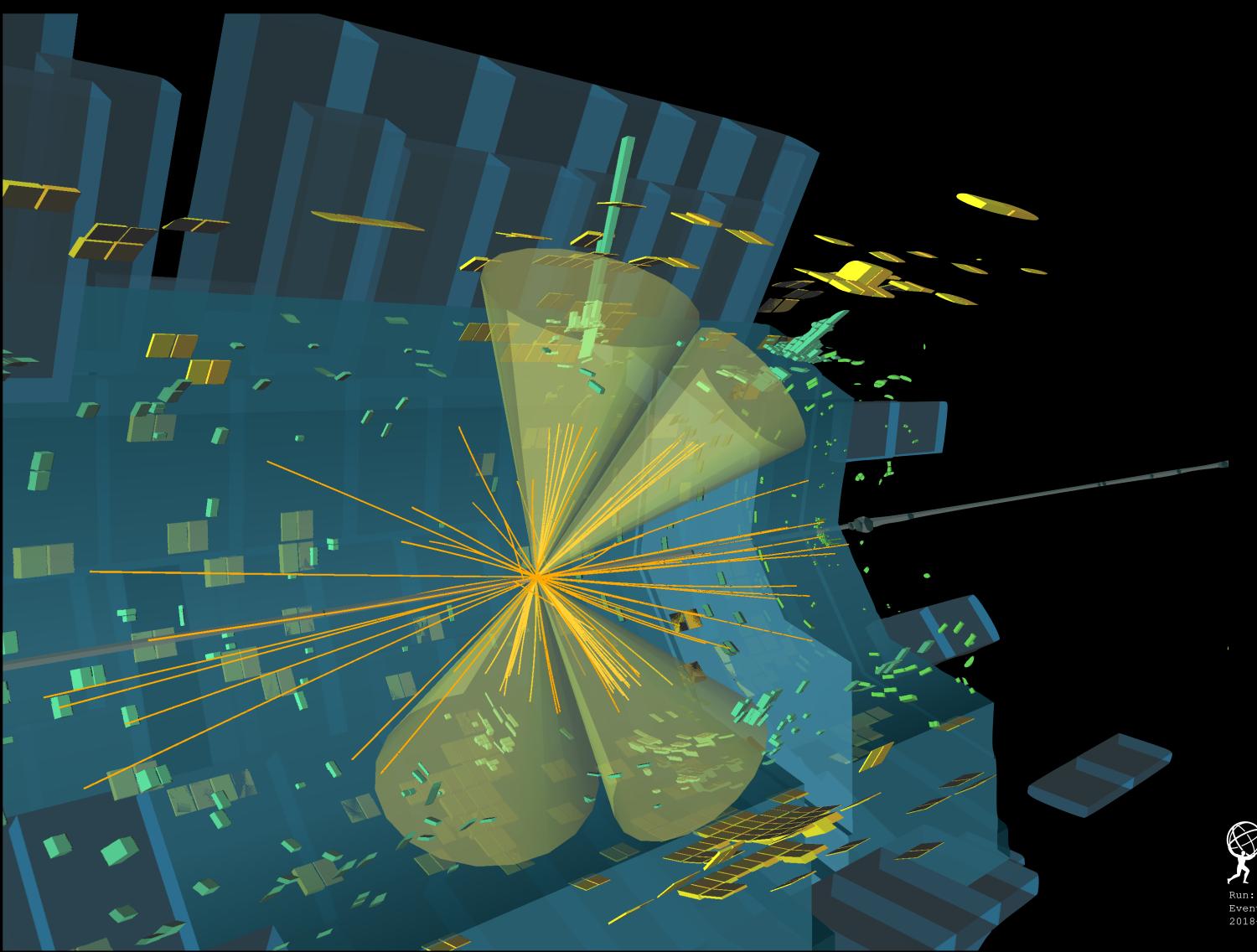
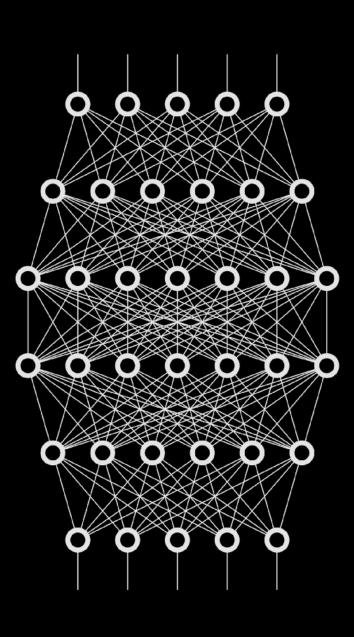
Machine Learning for Particle Physics



$$J = -\frac{1}{4} F_{nv} F^{nv} + i F_{D} V + h.c$$

$$+ V_{ij} V_{j} \phi + h.c$$

$$+ |D_{n} \phi|^{2} - V(\phi)$$



Nicole Hartman

nicole.hartman@tum.de

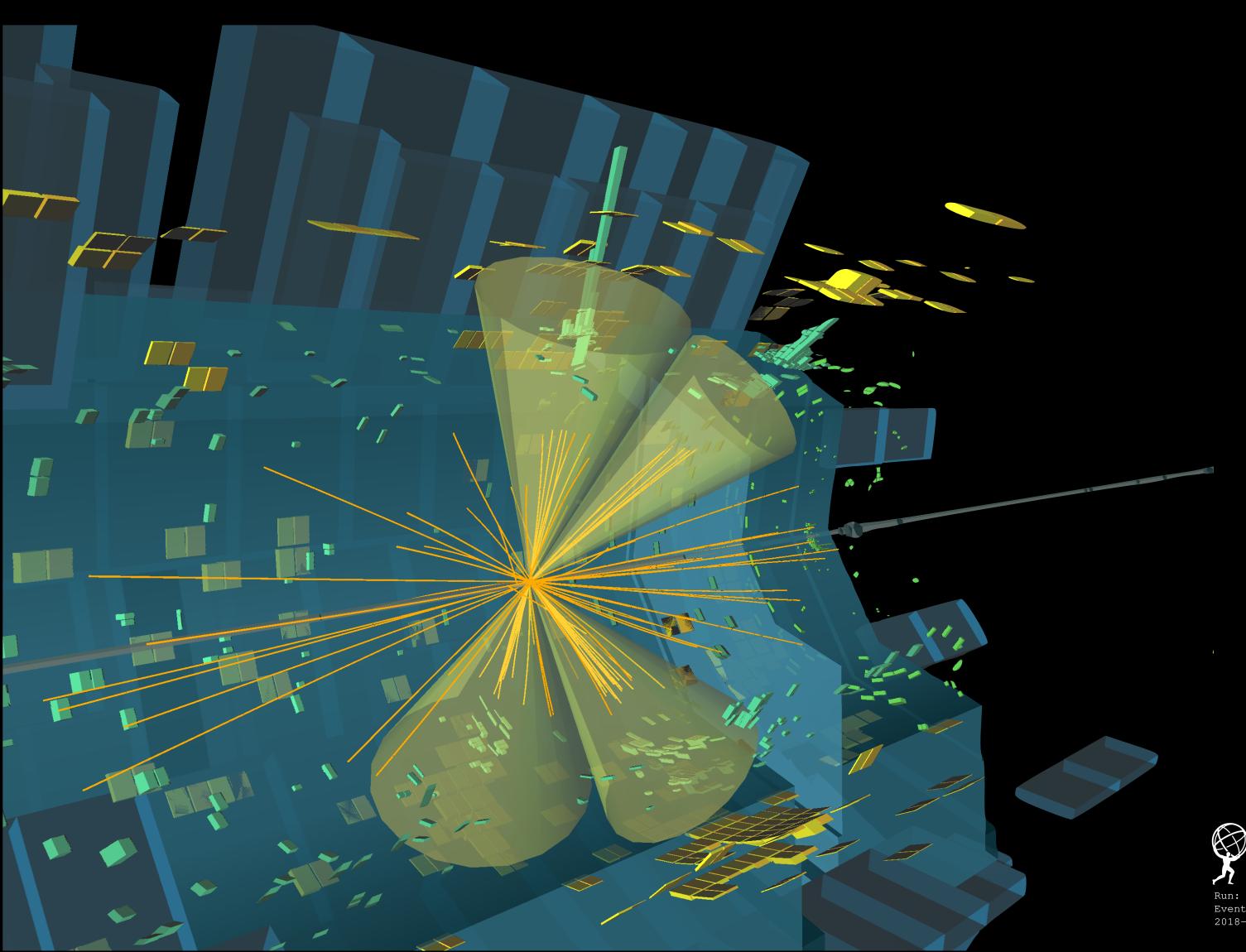
Belle II Germany meeting

University of Bonn, 8.9.2025

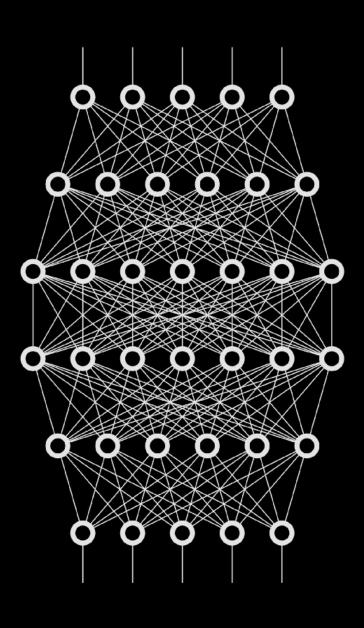




High Energy Physics Machine Learning for Particle Physics



$$J = -\frac{1}{4}F_{n}F^{n}$$
+ $iFDY + h.c$
+ $Y_{i}Y_{j}Y_{j} + h.c$
+ $|D_{n}P|^{2} - V(p)$



Nicole Hartman

nicole.hartman@tum.de

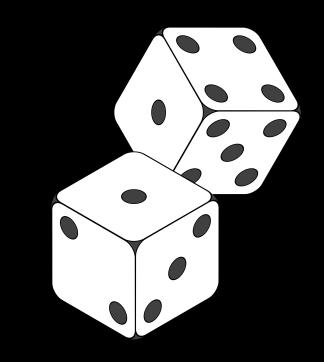
Belle II Germany meeting

University of Bonn, 8.9.2025





Particle physics: Probabilistic model to describe the world we live in.



$$J = -\frac{1}{4}F_{n}F^{n}$$
+ iFDY+h.c
+ $\frac{1}{4}Y_{ij}Y_{j} + h.c$
+ $\frac{1}{4}Y_{ij}Y_{j} + h.c$
+ $\frac{1}{4}Q_{n}P_{i}^{2} - V(\emptyset)$

$$p(x \mid \theta) = \frac{1}{\sigma} \int dx \mid \mathcal{M}(x) \mid^2$$
data SM

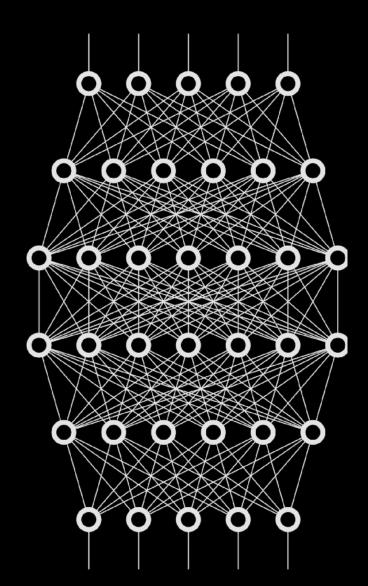
Machine Learning: Probabilistic model to describe

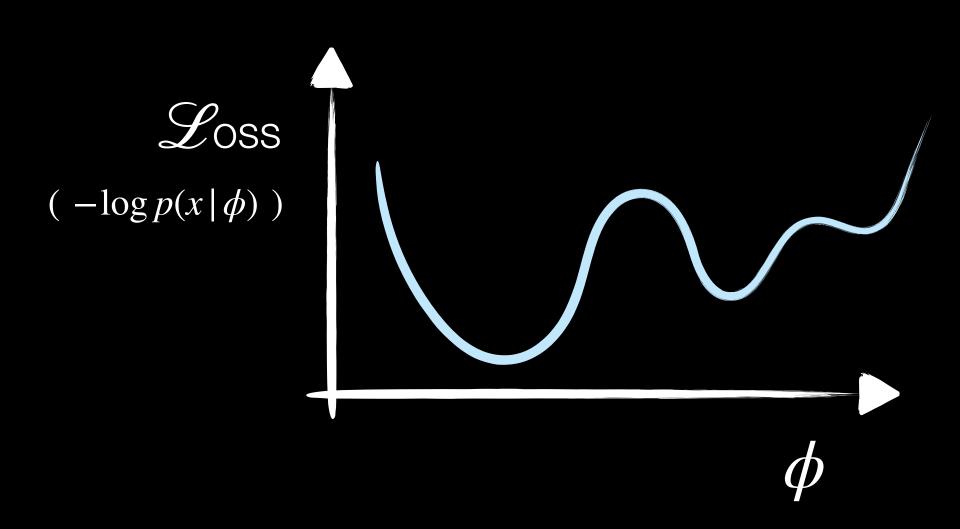
the world we live in.

$$p(x | \phi)$$
data

Parameters

$$\phi \leftarrow \phi - \lambda \nabla_{\phi} \mathcal{L}$$



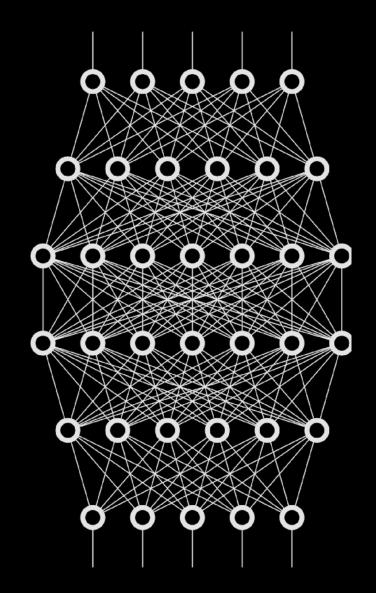


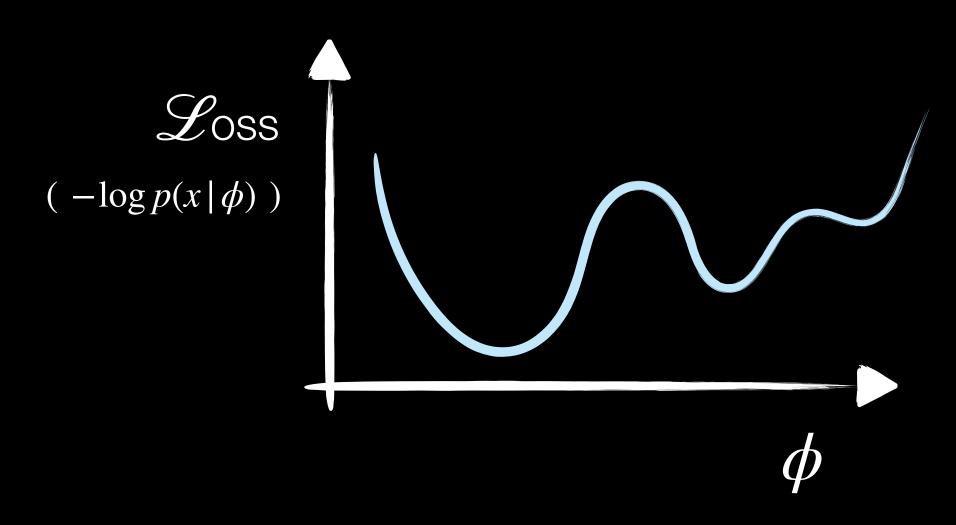
Machine Learning: Probabilistic model to describe

the world we live in.

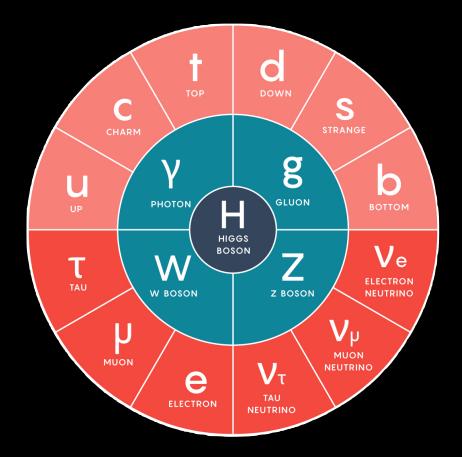
Data drives discoveries

$$\phi \leftarrow \phi - \lambda \nabla_{\phi} \mathcal{L}$$





Machine Learning for particle physics



W/Z 1983

High fidelity
MC simulations

Cern Courrier;
Discovery of the weak
neutral currents

top 1995

Probabilistic reconstruction

Phys. Rev. D 94, 032004 (2016)

A. Matrix Element Technique

This measurement uses the matrix element technique [13]. This method provides the most precise m_t measurement at the Tevatron in the ℓ +jets final state [1], and was applied in previous measurement of m_t in the dilepton final state using 5.3 fb⁻¹ of integrated luminosity [39]. The ME method used in this analysis is described below.

B. Event probability calculation

The ME technique assigns a probability to each event, which is calculated as

$$P(x, f_{t\bar{t}}, m_t) = f_{t\bar{t}} \cdot P_{t\bar{t}}(x, m_t) + (1 - f_{t\bar{t}}) \cdot P_{\text{bkg}}(x), (2)$$

where $f_{t\bar{t}}$ is the fraction of $t\bar{t}$ events in the data, and $P_{t\bar{t}}$ and P_{bkg} are the respective per-event probabilities calculated under the hypothesis that the selected event is either a $t\bar{t}$ event, characterized by a top quark mass m_t , or background. Here, x represents the set of measured observables, i.e., p_T , η , and ϕ for jets and leptons. We assume that the masses of top quarks and anti-top quarks are the same. The probability $P_{t\bar{t}}(x, m_t)$ is calculated as

$$P_{t\bar{t}}(x,m_t) = \frac{1}{\sigma_{\text{obs}}(m_t)} \int f_{\text{PDF}}(q_1) f_{\text{PDF}}(q_2) \times \frac{(2\pi)^4 |\mathcal{M}(y,m_t)|^2}{q_1 q_2 s} W(x,y) d\Phi_6 dq_1 dq_2,$$
(3)

B-factories 2000s

Measurements of Branching Fraction, Polarization, and Direct-CP-Violating Charge Asymmetry in $B^+ \to K^{*0} \rho^+$ Decays

The BABAR Collaboration

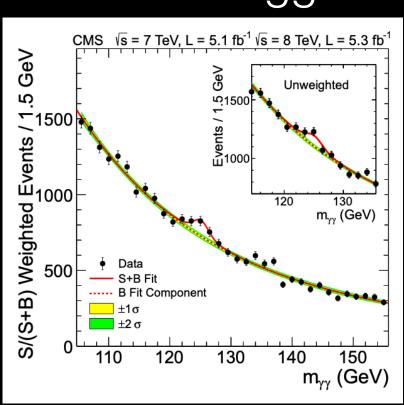
August 19 900

To discriminate signal from continuum background, we also use a neural network (NN) combining six variables: a Fisher discriminant made from two event-shape variables (see [18]); the cosine of the angle between the direction of the B and the collision axis (z) in the CM frame; the cosine of the angle between the B-thrust axis and the z axis; the cosine of the angle between the B-thrust axis and the thrust of the particles of the rest of the event; the angle between the direction of the π^0 and that of one of its daughter photons in the π^0 rest frame (π^0 decay angle); and the sum of transverse momenta relative to the z-axis of the particles in the rest of the event.

hep-ex/0408093

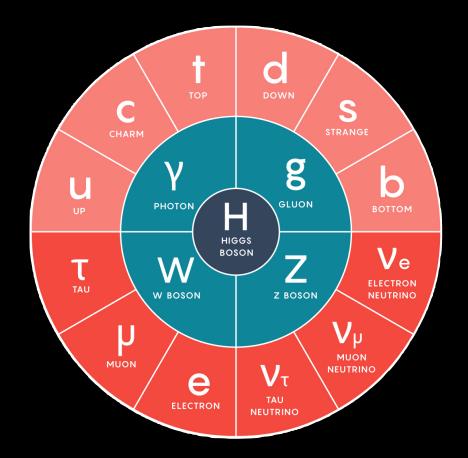


BDTs for Higgs discovery



Nature 560 (2018) Phys Lett B (2012) 30

Machine Learning for particle physics



W/Z 1983

High fidelity
MC simulations

Cern Courrier;
Discovery of the weak
neutral currents

top 1995

Probabilistic reconstruction

Phys. Rev. D 94, 032004 (2016)

A. Matrix Element Technique

This measurement uses the matrix element technique [13]. This method provides the most precise m_t measurement at the Tevatron in the ℓ +jets final state [1], [2], and was applied in previous measurement of m_t in the dilepton final state using 5.3 fb⁻¹ of integrated luminosity [39]. The ME method used in this analysis is described below.

B. Event probability calculation

The ME technique assigns a probability to each event, which is calculated as

$$P(x, f_{t\bar{t}}, m_t) = f_{t\bar{t}} \cdot P_{t\bar{t}}(x, m_t) + (1 - f_{t\bar{t}}) \cdot P_{\text{bkg}}(x), \quad (2)$$

where $f_{t\bar{t}}$ is the fraction of $t\bar{t}$ events in the data, and $P_{t\bar{t}}$ and $P_{\rm bkg}$ are the respective per-event probabilities calculated under the hypothesis that the selected event is either a $t\bar{t}$ event, characterized by a top quark mass m_t , or background. Here, x represents the set of measured observables, i.e., p_T , η , and ϕ for jets and leptons. We assume that the masses of top quarks and anti-top quarks are the same. The probability $P_{t\bar{t}}(x,m_t)$ is calculated as

$$P_{t\bar{t}}(x,m_t) = \frac{1}{\sigma_{\text{obs}}(m_t)} \int f_{\text{PDF}}(q_1) f_{\text{PDF}}(q_2) \times \frac{(2\pi)^4 |\mathcal{M}(y,m_t)|^2}{q_1 q_2 s} W(x,y) d\Phi_6 dq_1 dq_2,$$
(3)

B-factories 2000s

Measurements of Branching Fraction, Polarization, and Direct-CP-Violating Charge Asymmetry in $B^+ \to K^{*0} \rho^+$ Decays

The BABAR Collaboration

August 19 200/

To discriminate signal from continuum background, we also use a neural network (NN) combining six variables: a Fisher discriminant made from two event-shape variables (see [18]); the cosine of the angle between the direction of the B and the collision axis (z) in the CM frame; the cosine of the angle between the B-thrust axis and the z axis; the cosine of the angle between the B-thrust axis and the thrust of the particles of the rest of the event; the angle between the direction of the π^0 and that of one of its daughter photons in the π^0 rest frame (π^0 decay angle); and the sum of transverse momenta relative to the z-axis of the particles in the rest of the event.

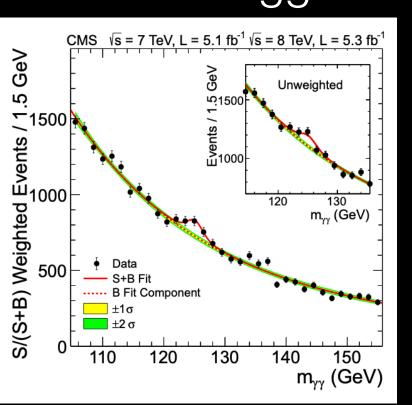
hep-ex/0408093







BDTs for Higgs discovery



Nature 560 (2018) Phys Lett B (2012) 30

The New Hork Times

Physicists Find Elusive Particle Seen as Key to Universe

Share full article









Scientists in Geneva on Wednesday applauded the discovery of a subatomic particle that looks like the Higgs boson. Pool photo by Denis Balibouse

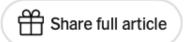
By Dennis Overbye

July 4, 2012

article

The New York Times

Scientists See Promise in Deep-Learning Programs







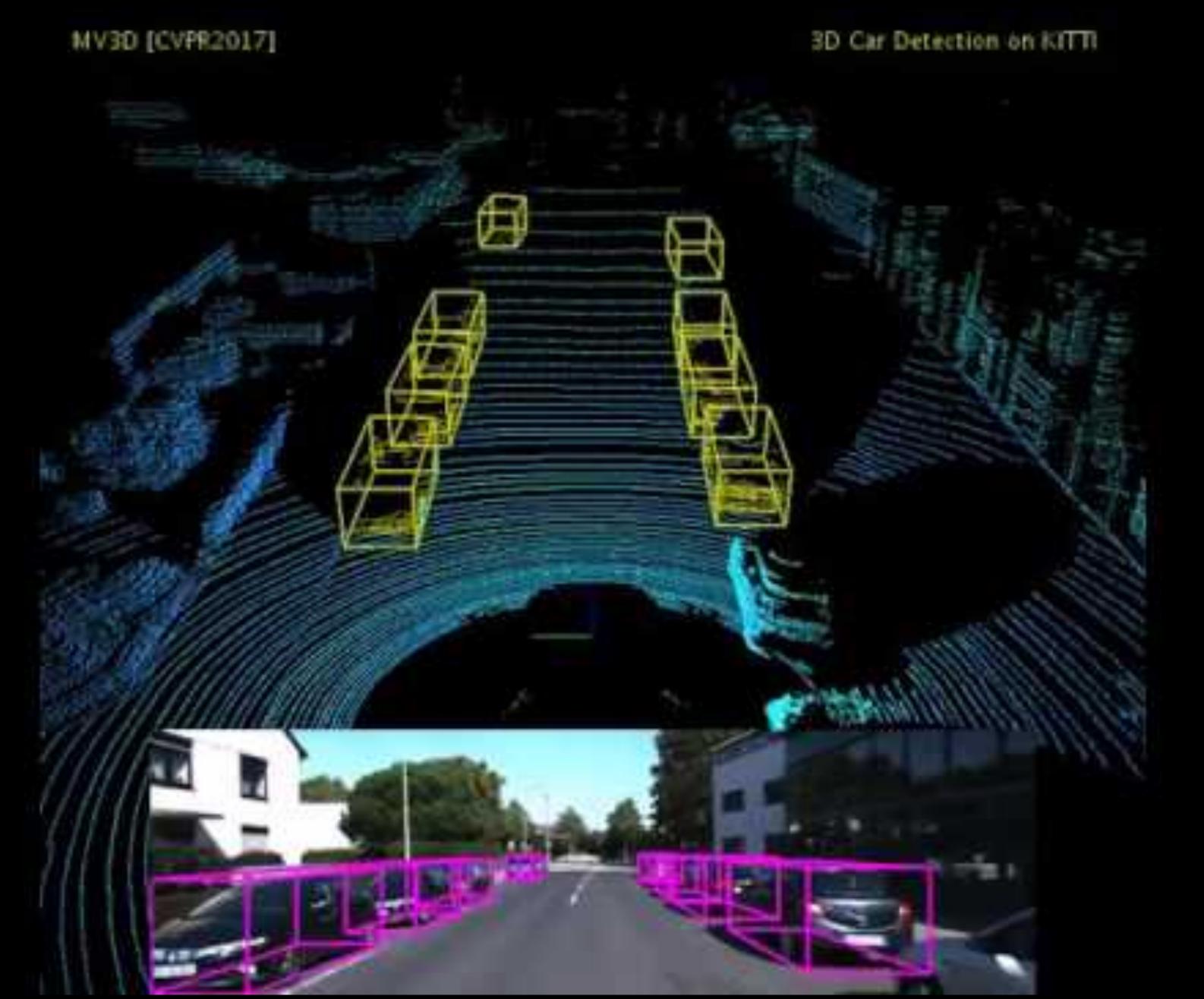


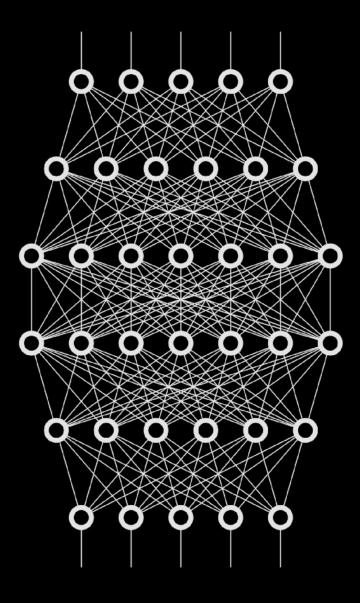
A voice recognition program translated a speech given by Richard F. Rashid, Microsoft's top scientist, into Mandarin Chinese. Hao Zhang/The New York Times

By John Markoff

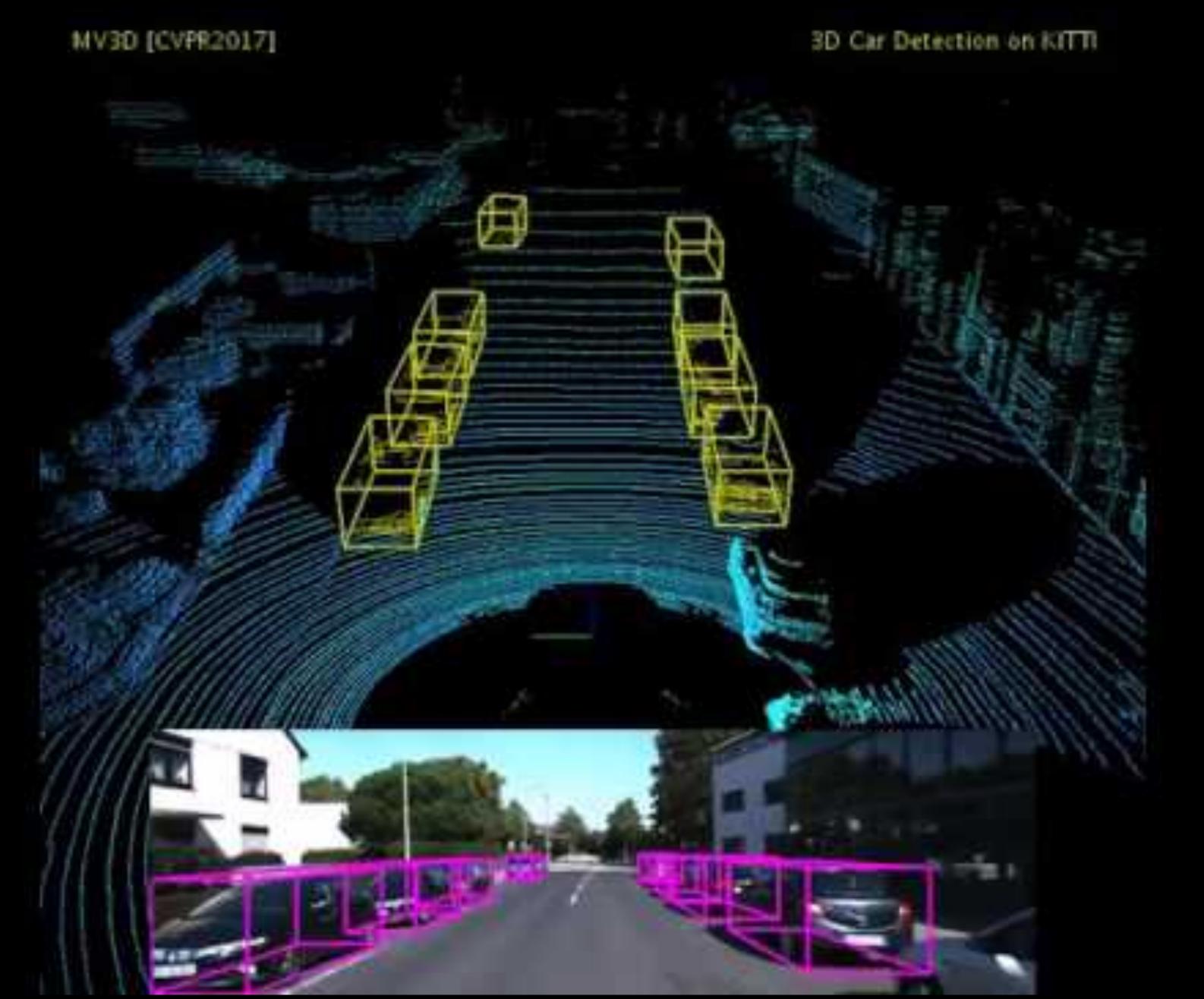
Nov. 23, 2012

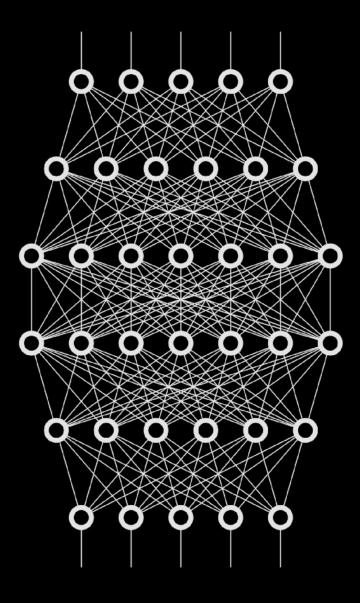






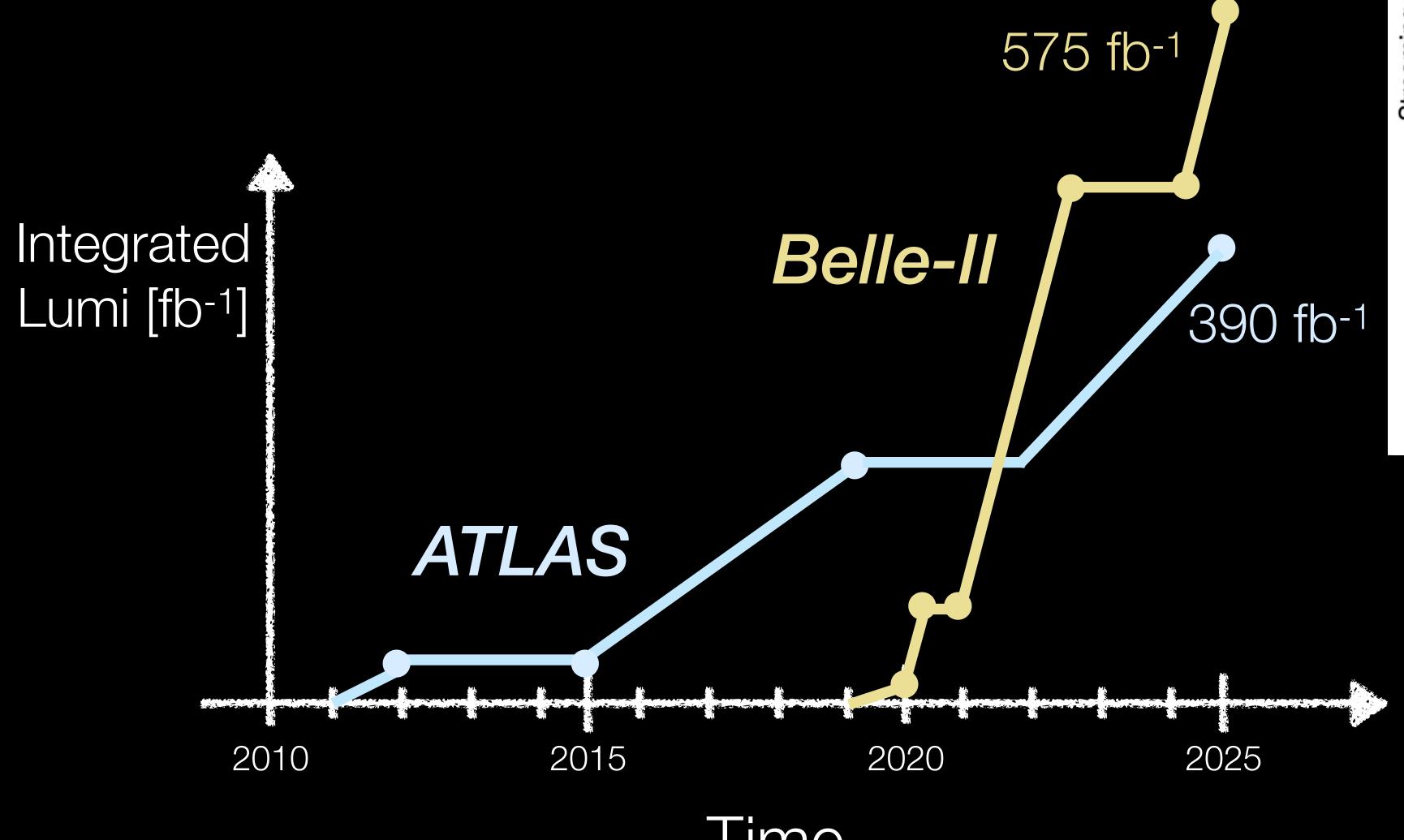
Networks get *deep* with CNNs; RNNs; graph networks; transformers

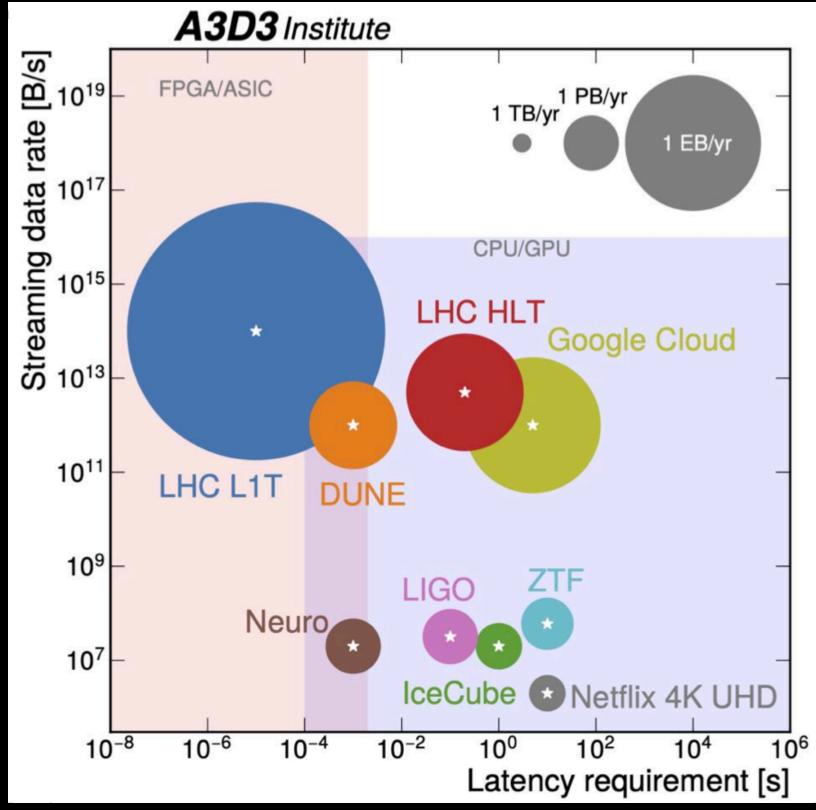




Networks get *deep* with CNNs; RNNs; graph networks; transformers

Big science, big data





https://a3d3.ai/about/

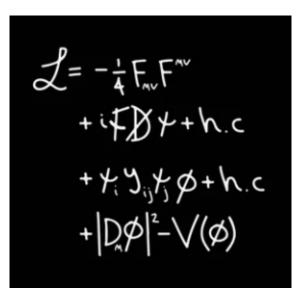
ATLAS: Run 3, HL-LHC <u>Belle-II lumi</u>

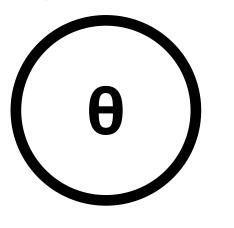
Time

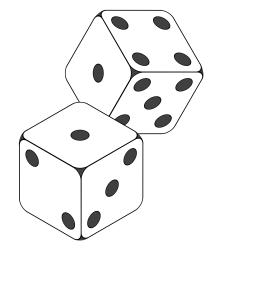
$p(x \mid \theta)$ data theory

Simulation

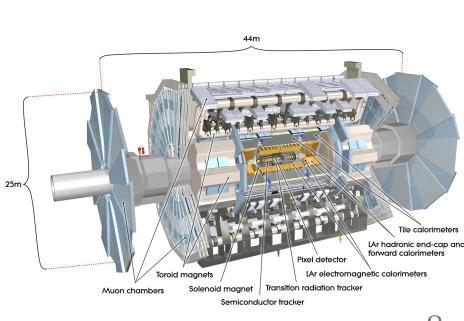
Generation / Forward modelling

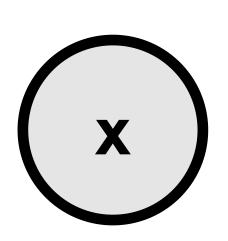


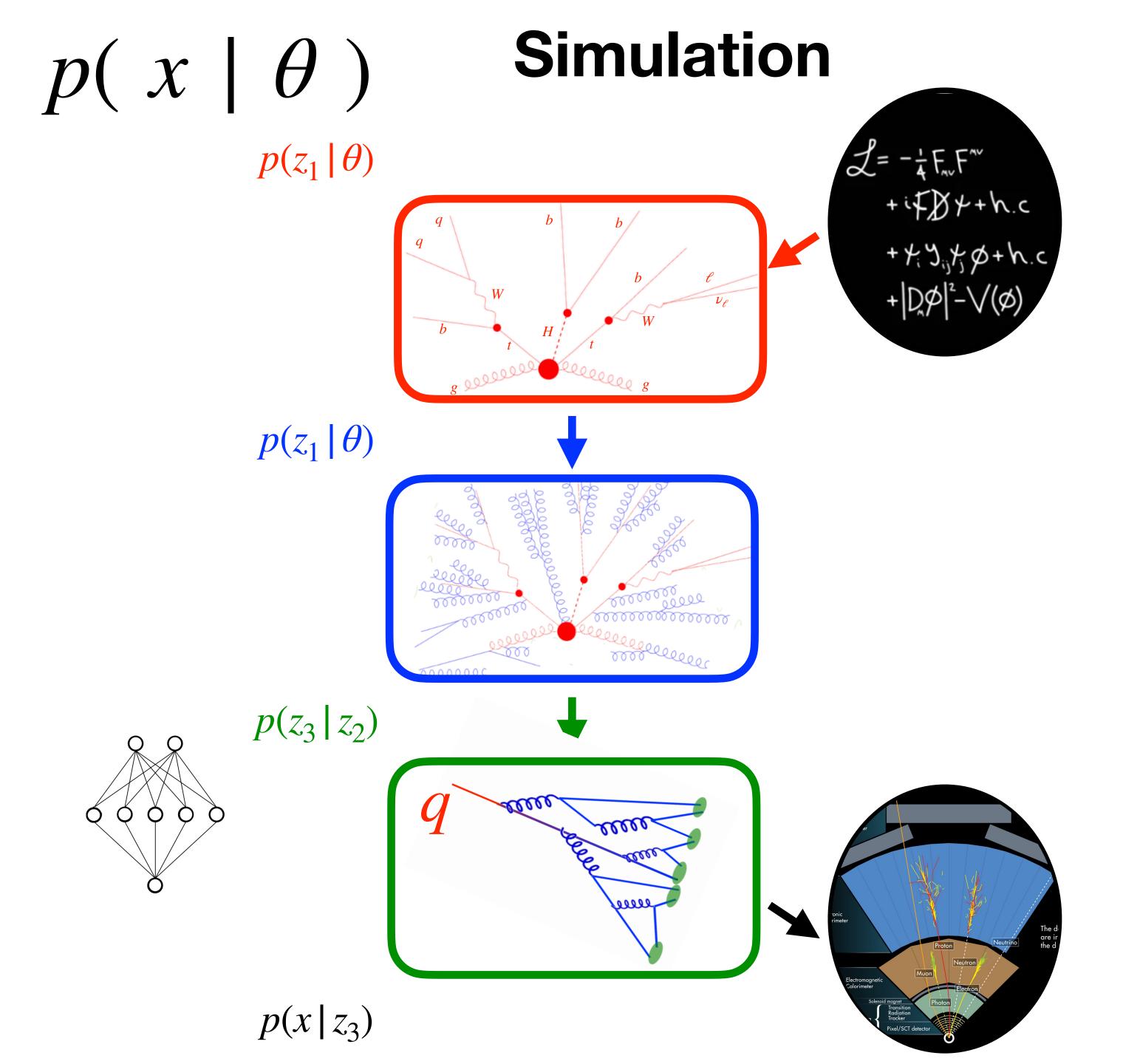


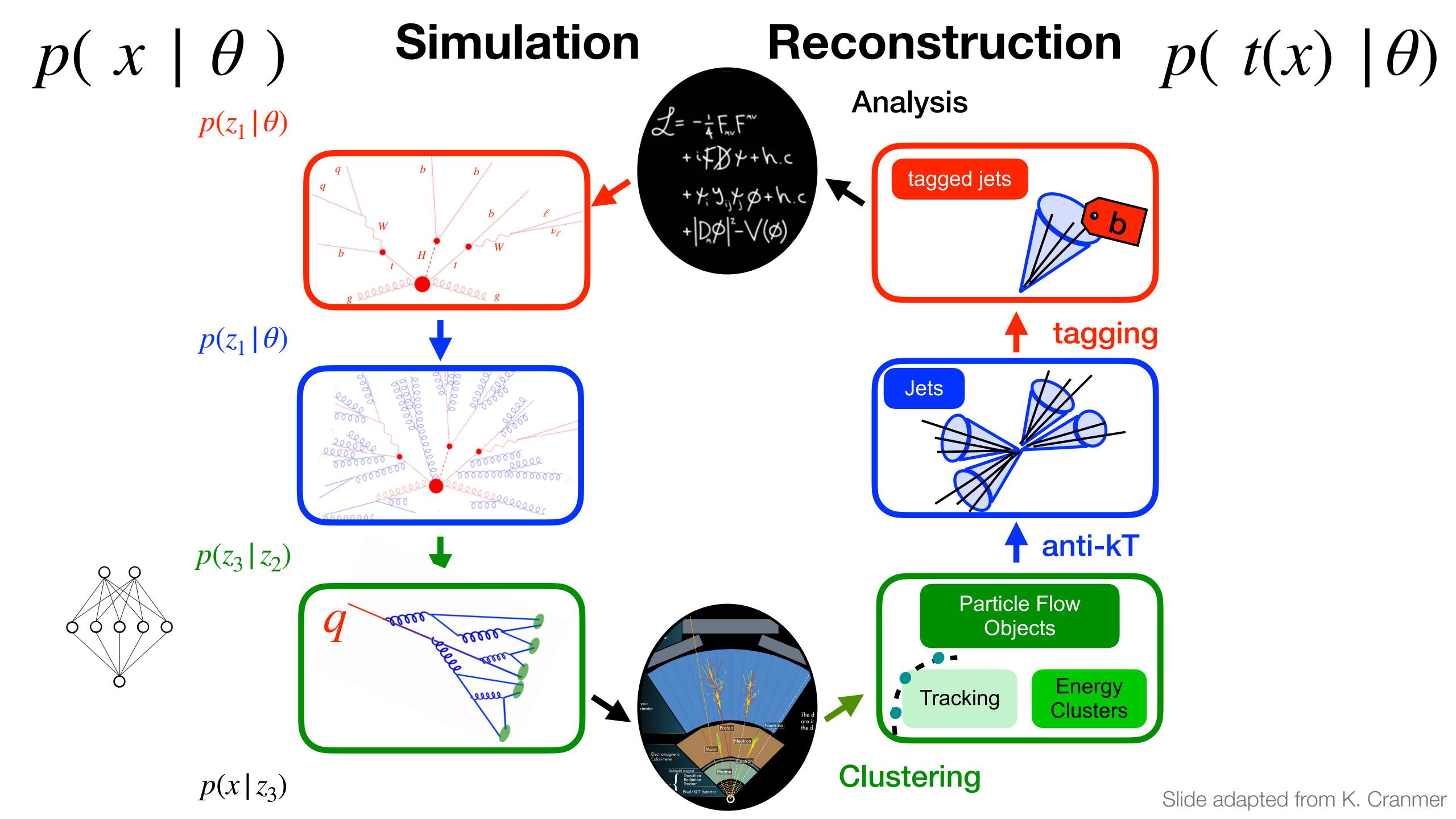


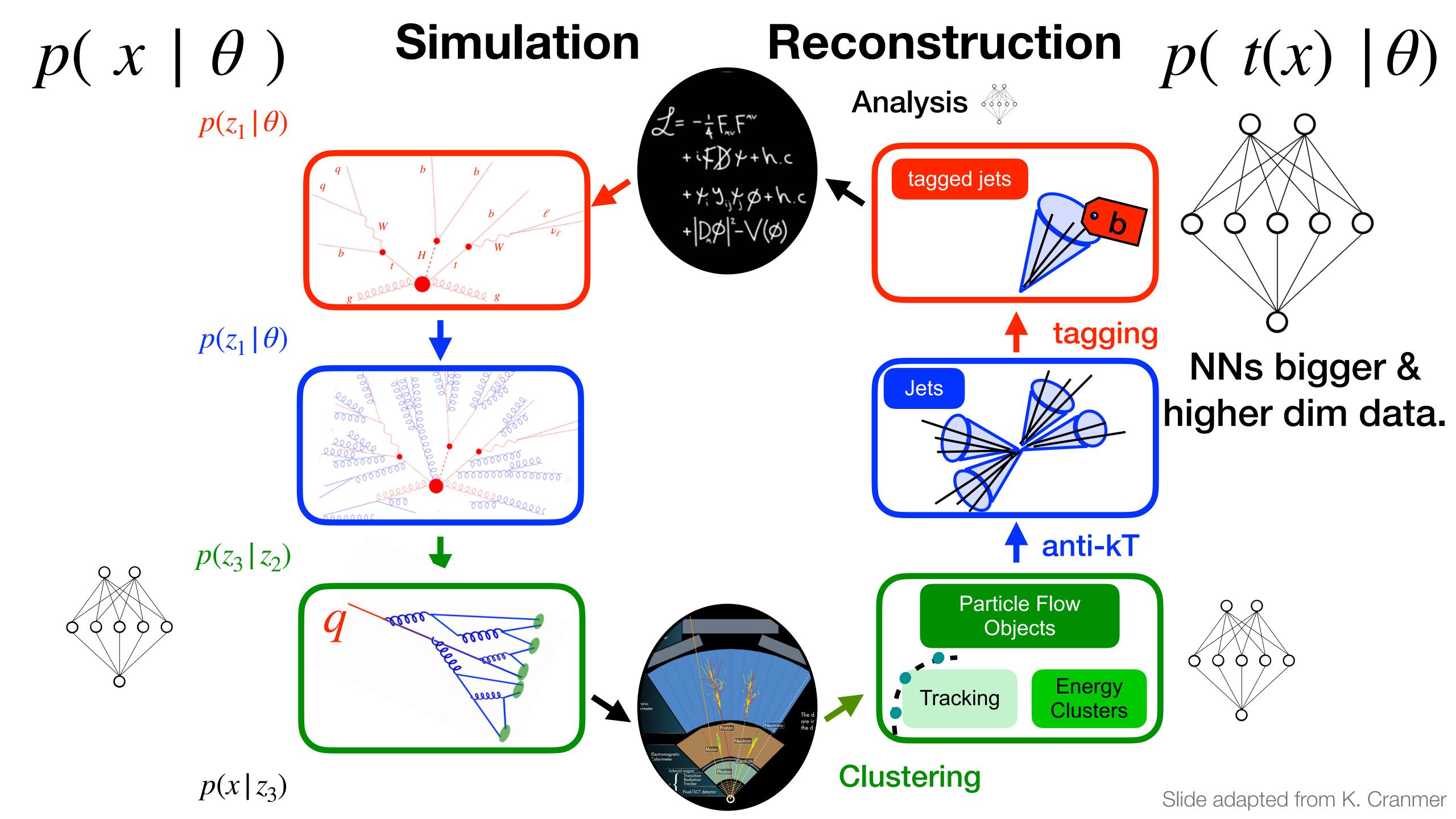








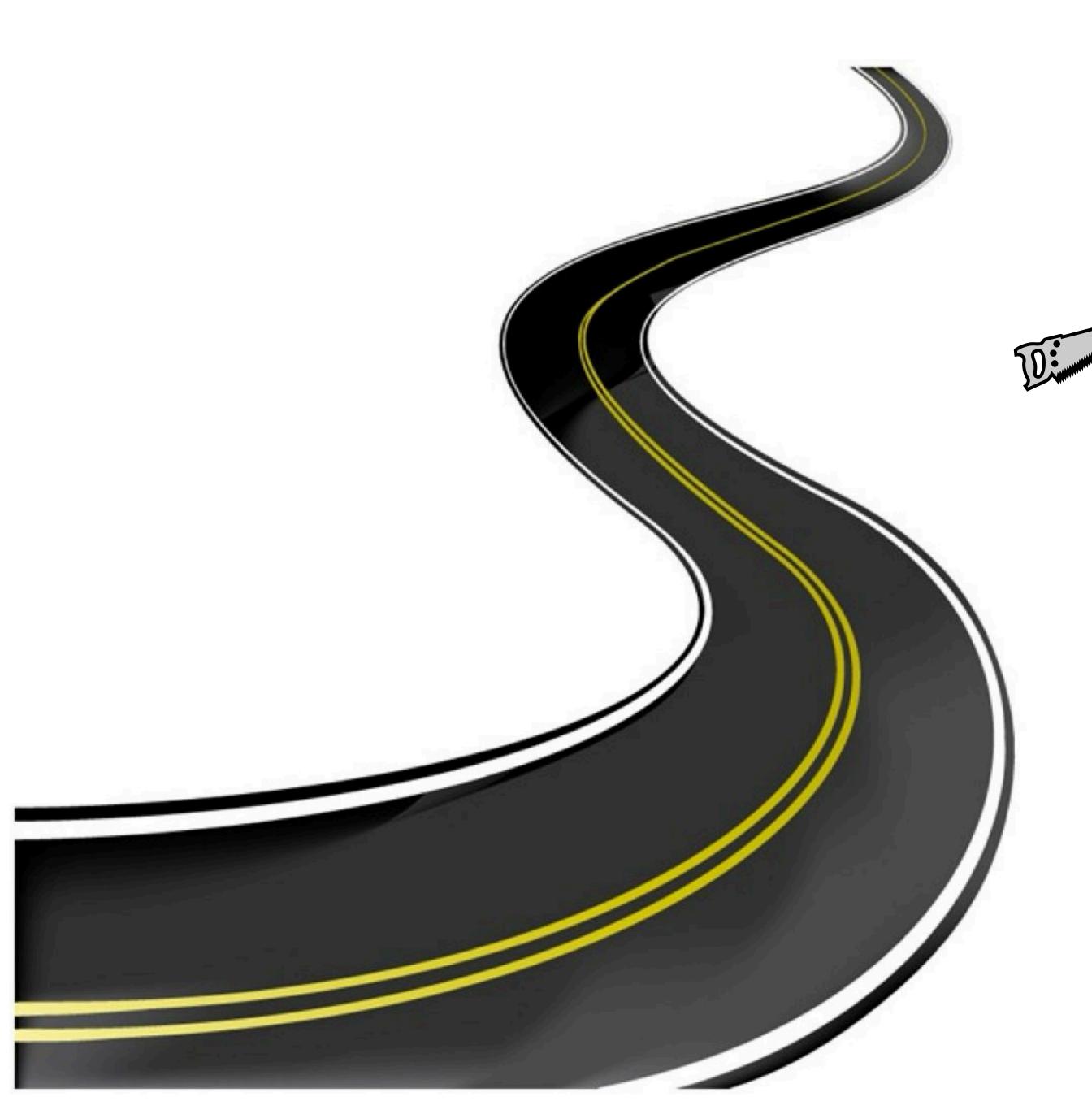






1) ML for our tools

2) ML for our future



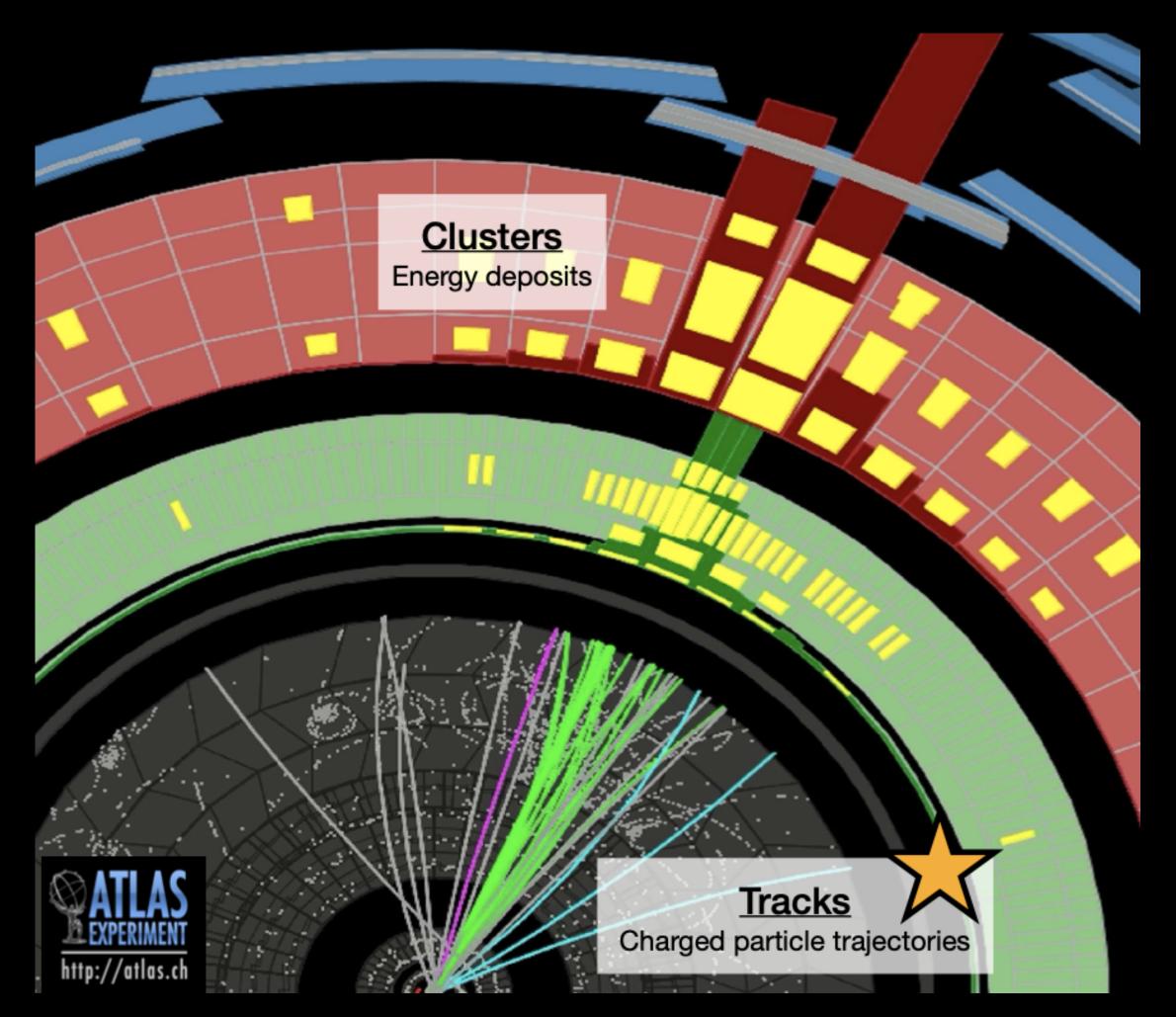
1) ML for our tools

Jet tagging

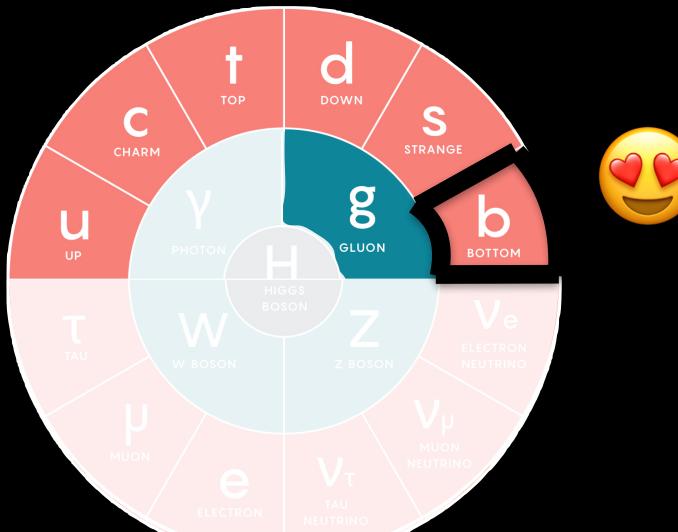


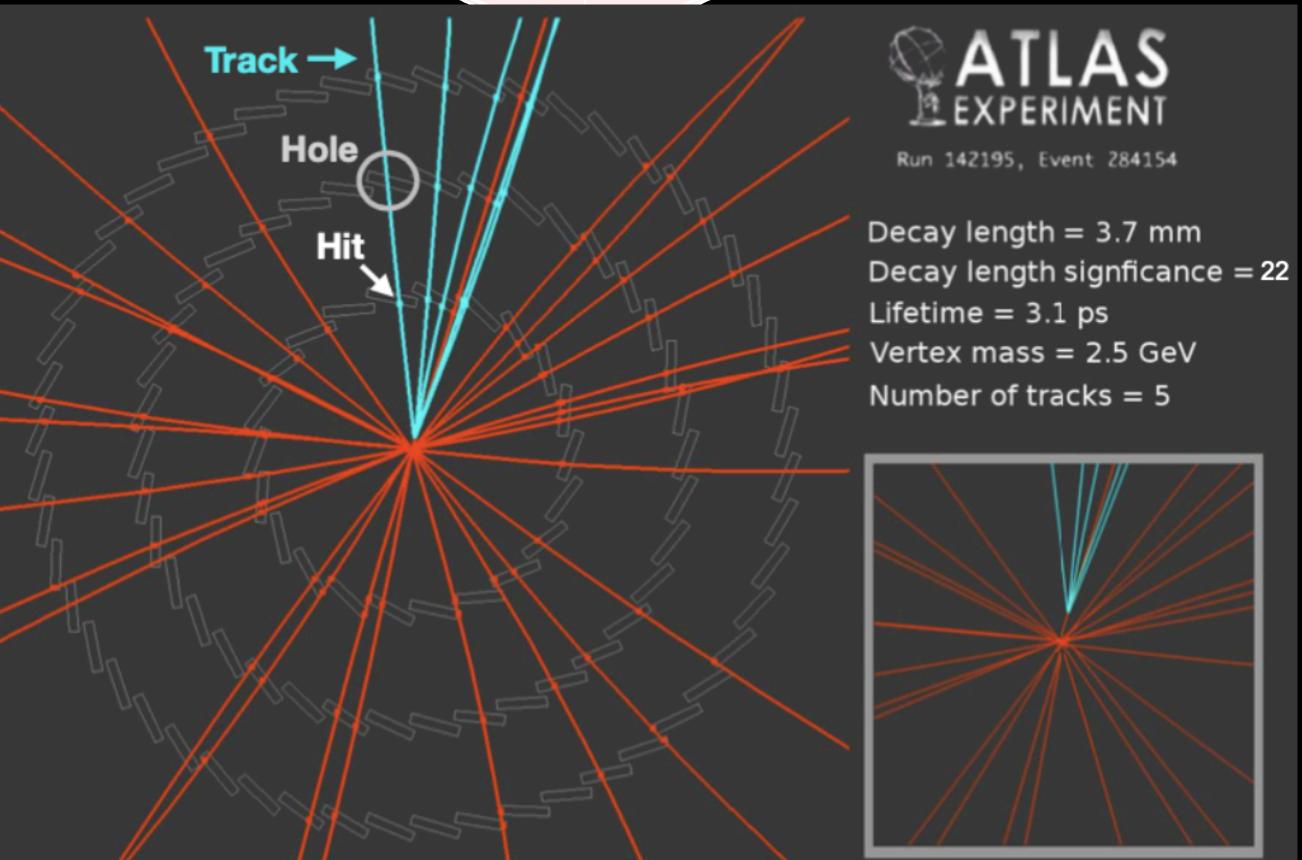
Generative networks

2) ML for our future

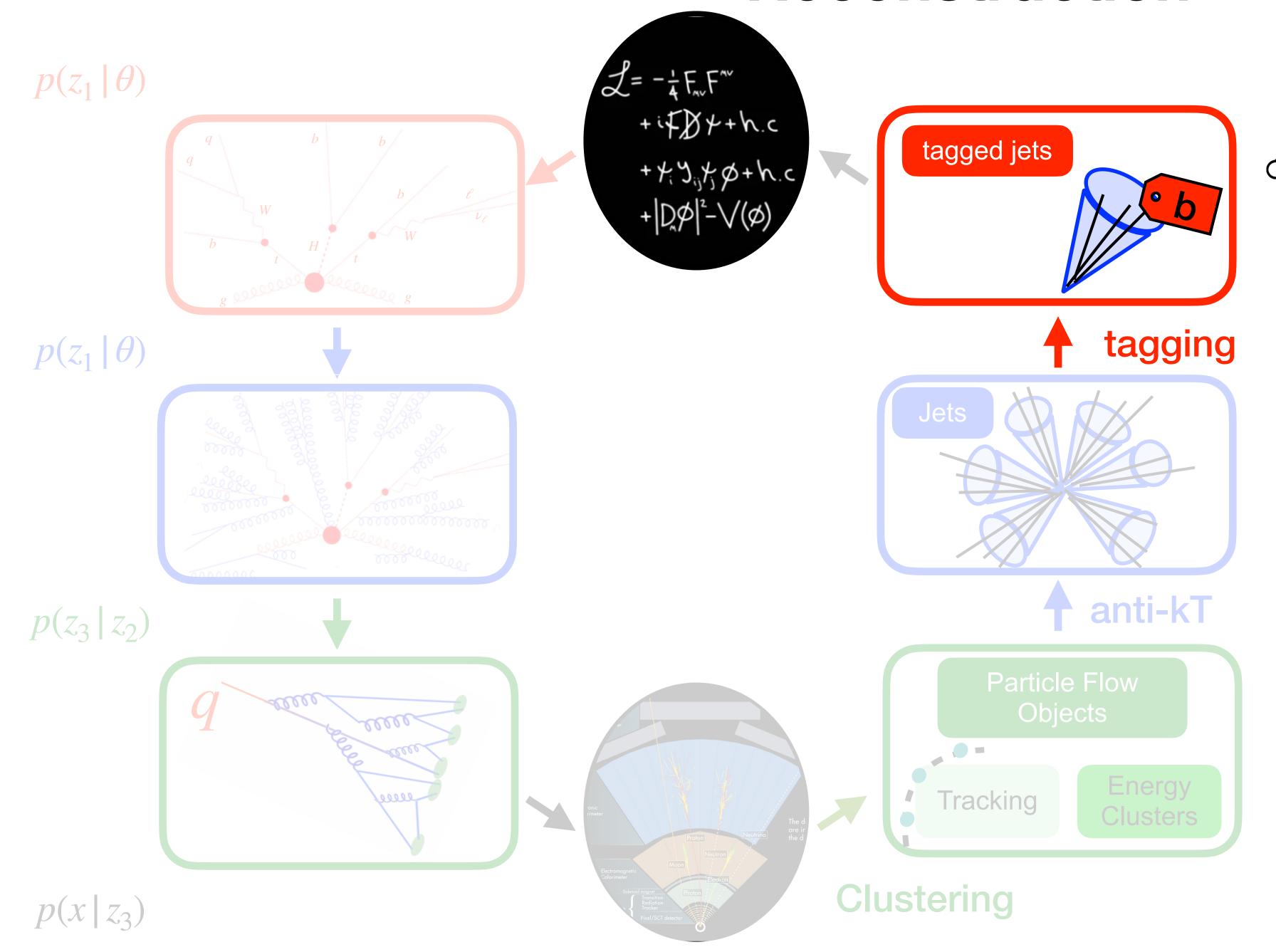


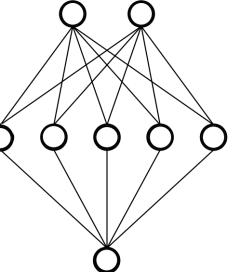
Eur. Phys Journal C Vol 73 3 (2013) 2304





Reconstruction





Deep learning for HEP

CNNs (2015)

RNNs (2017)

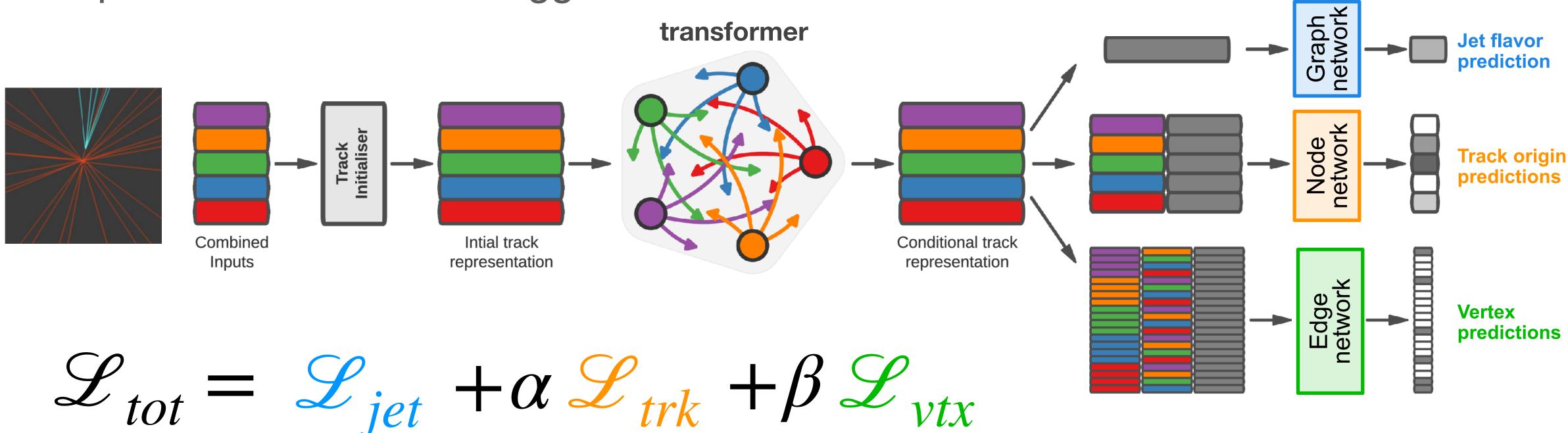
Deep Sets / Graphs (2018)

Transformers (2020)



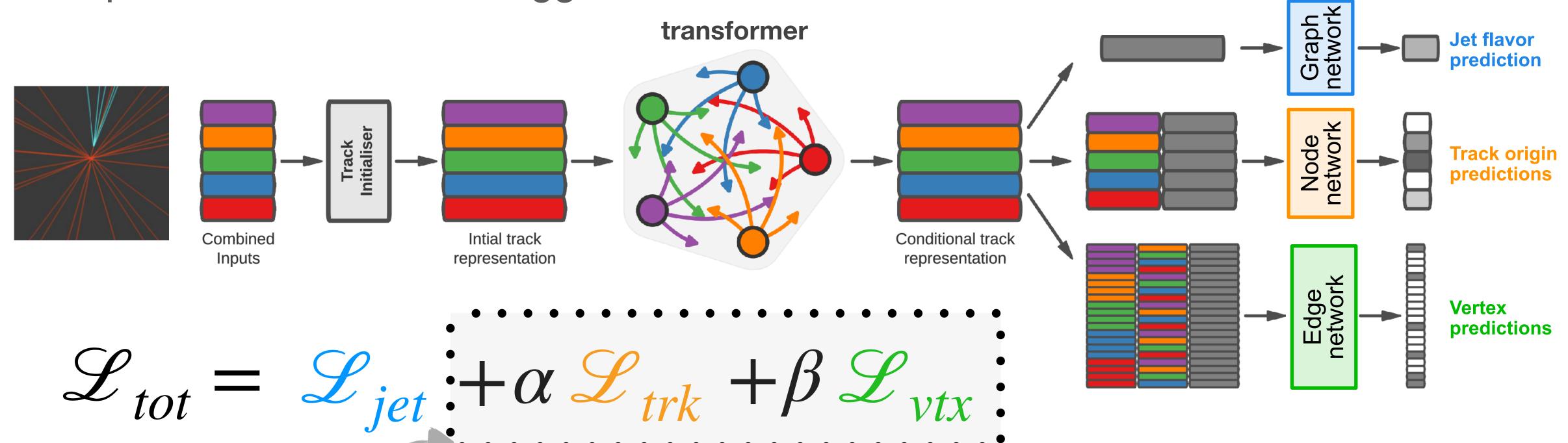
Jet tagging: multi-task learning!

Graph-net based flavour tagger



Jet tagging: multi-task learning!

Graph-net based flavour tagger



GN1 (2021): helped a lot

- 30m training jets
- 100% improvement ATL-PHYS-PUB-2022-027

FTAG-2023-01

Vertex finding

Group tracks with pair-wise compatibility > 0.5

ATLAS Simulation Preliminary

 $\sqrt{s} = 13 \text{ TeV}$ $t\bar{t}$ jets

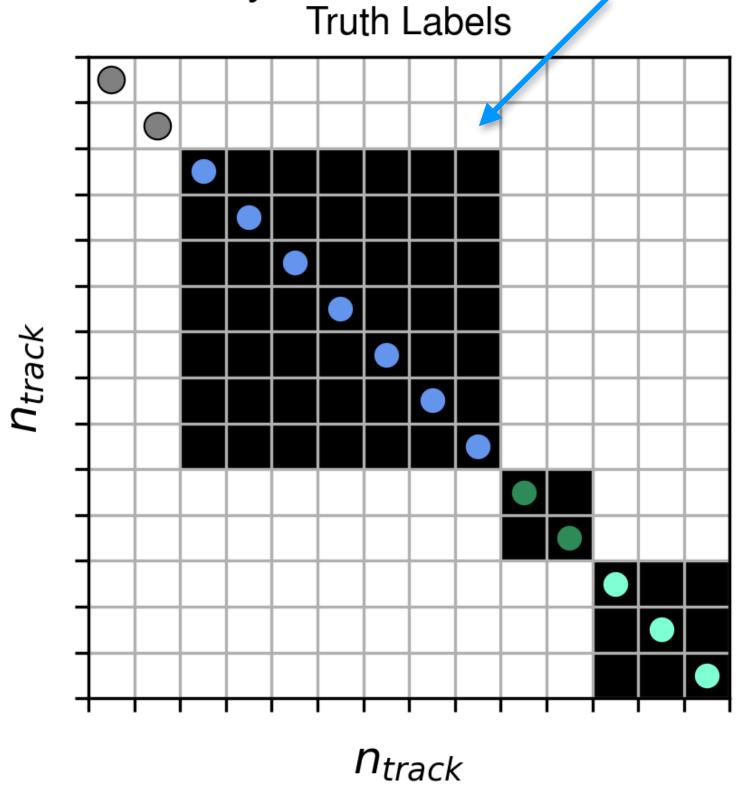
Truth b-jet

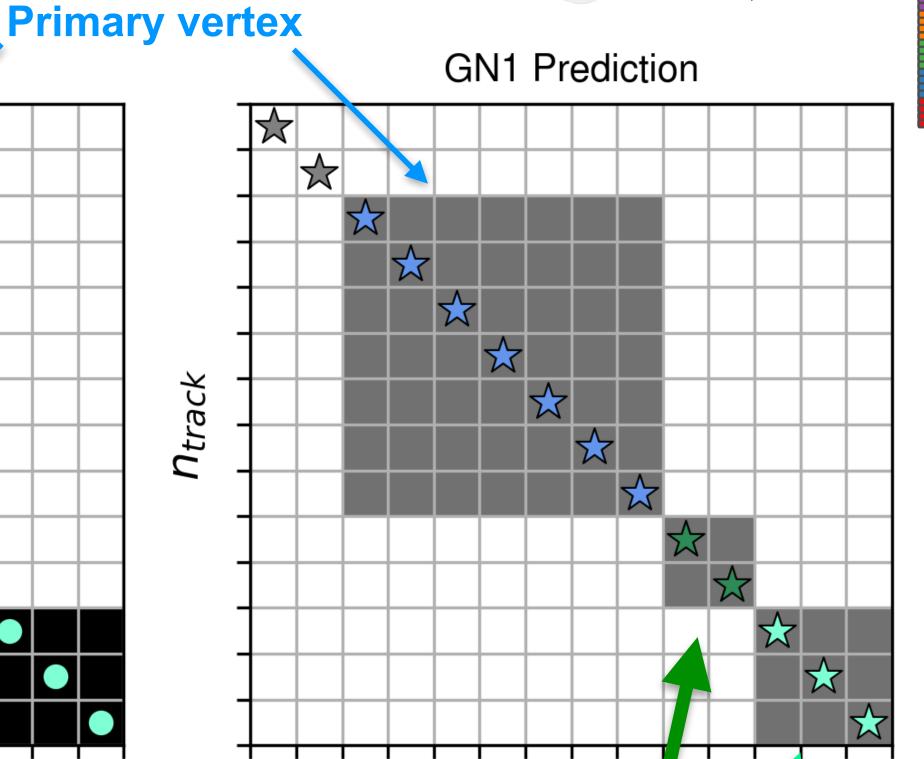
 $p_T = 134.1 \text{ GeV}$

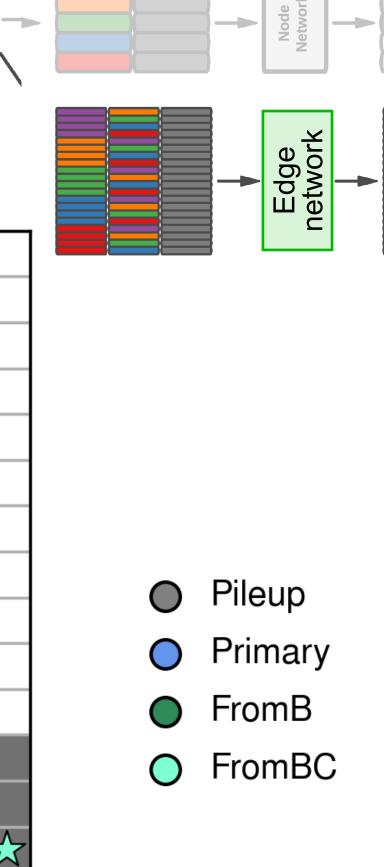
 $p_b = 0.995$

 $p_c = 0.005$

 $p_u = 0.000$







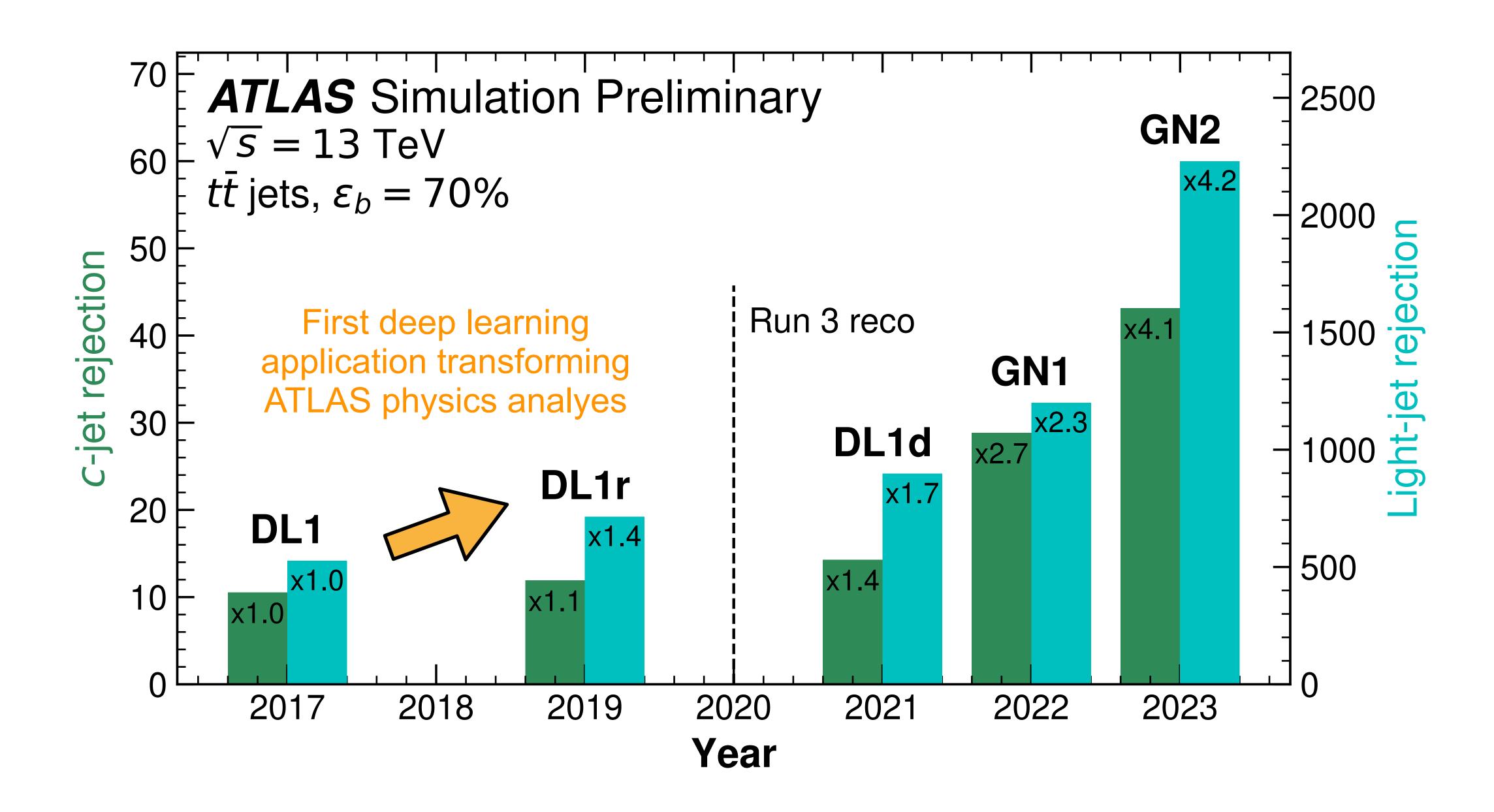
representation

From the B-decay vertex

Primary Vertex B

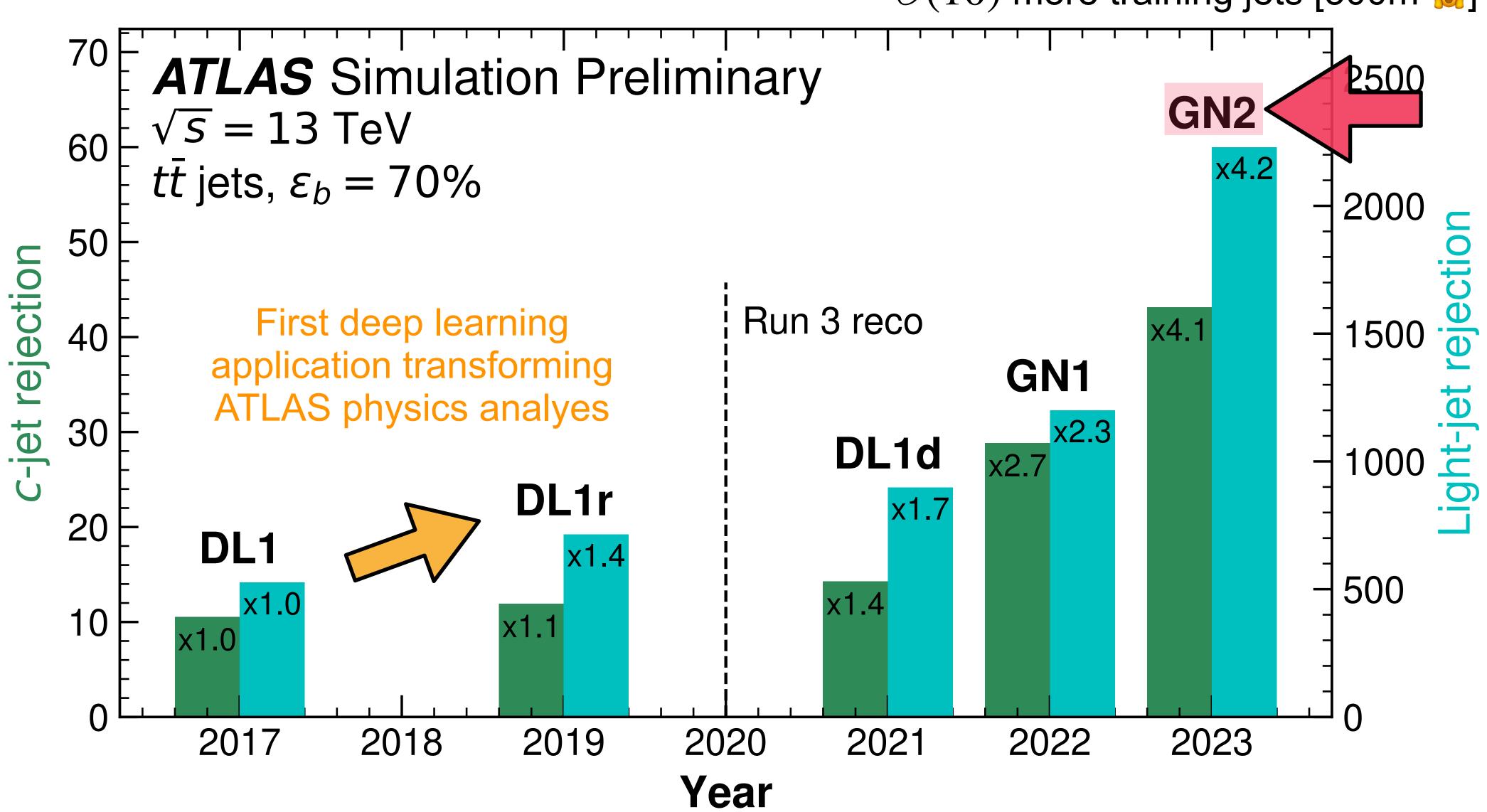
From tertiary Ddecay vertex

Т



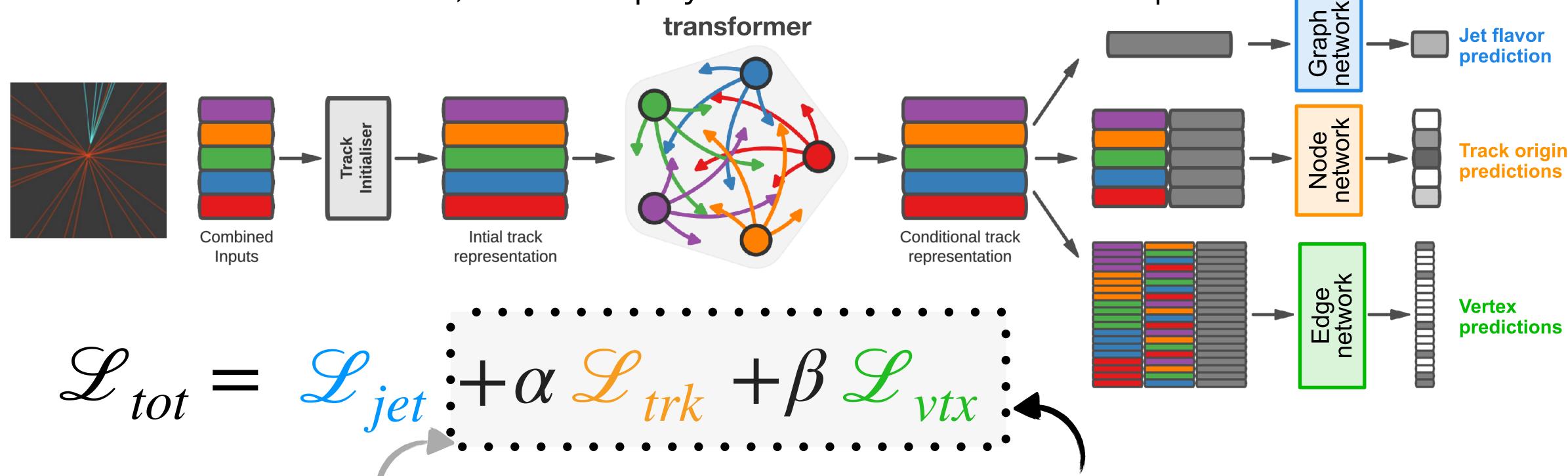
Transformer-based tagger

 $\mathcal{O}(10)$ more training jets [300m \mathbb{Q}]



More data, less physics

The more data we add, the less phyiscs inductive biases help.



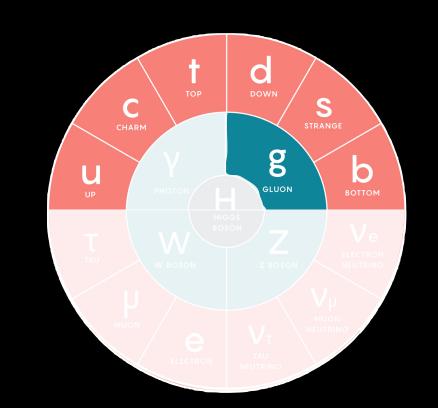
GN1 (2021): helped a lot

- 30m training jets
- 100% improvement ATL-PHYS-PUB-2022-027

GN2 (2023): help a little

- 300m training jets
- 15% improvement
 2505.19689

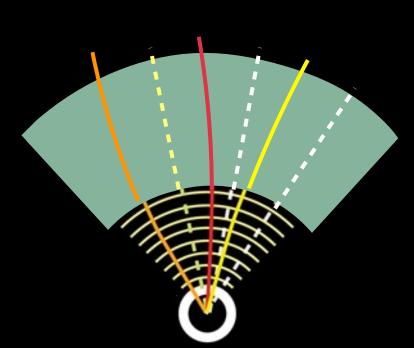
16

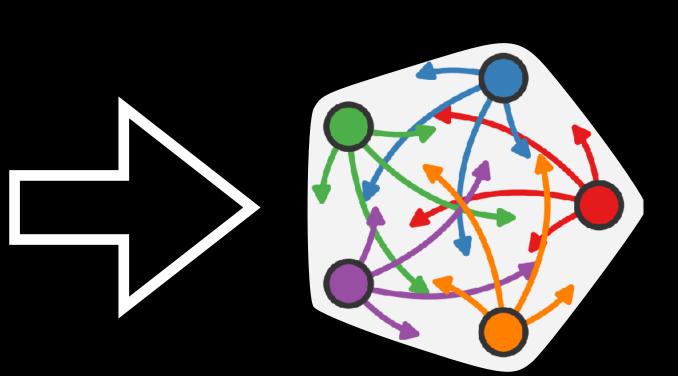


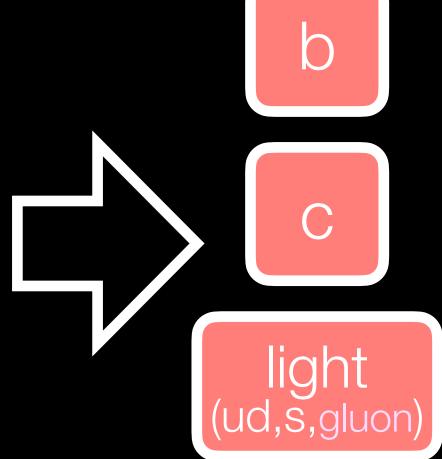




GN2

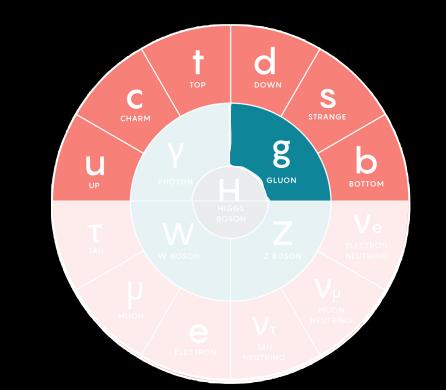


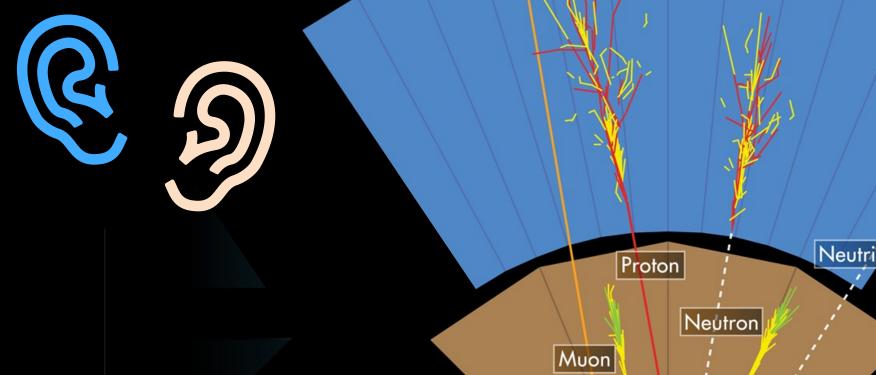


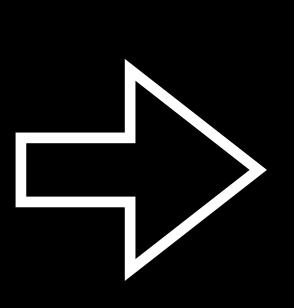


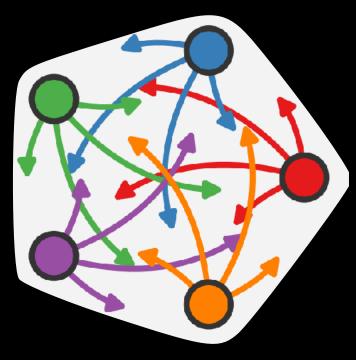
Classify other jets

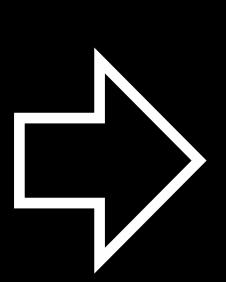










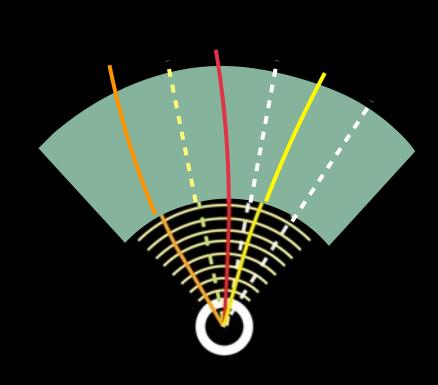


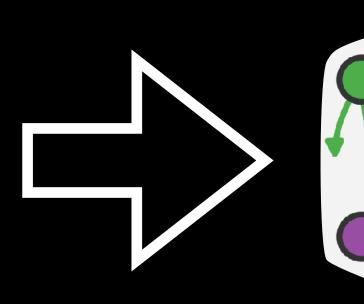


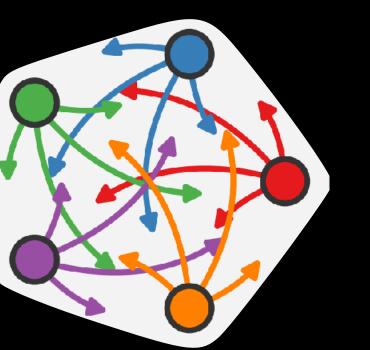
Classify b-jets

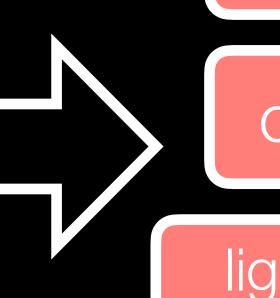


GN2



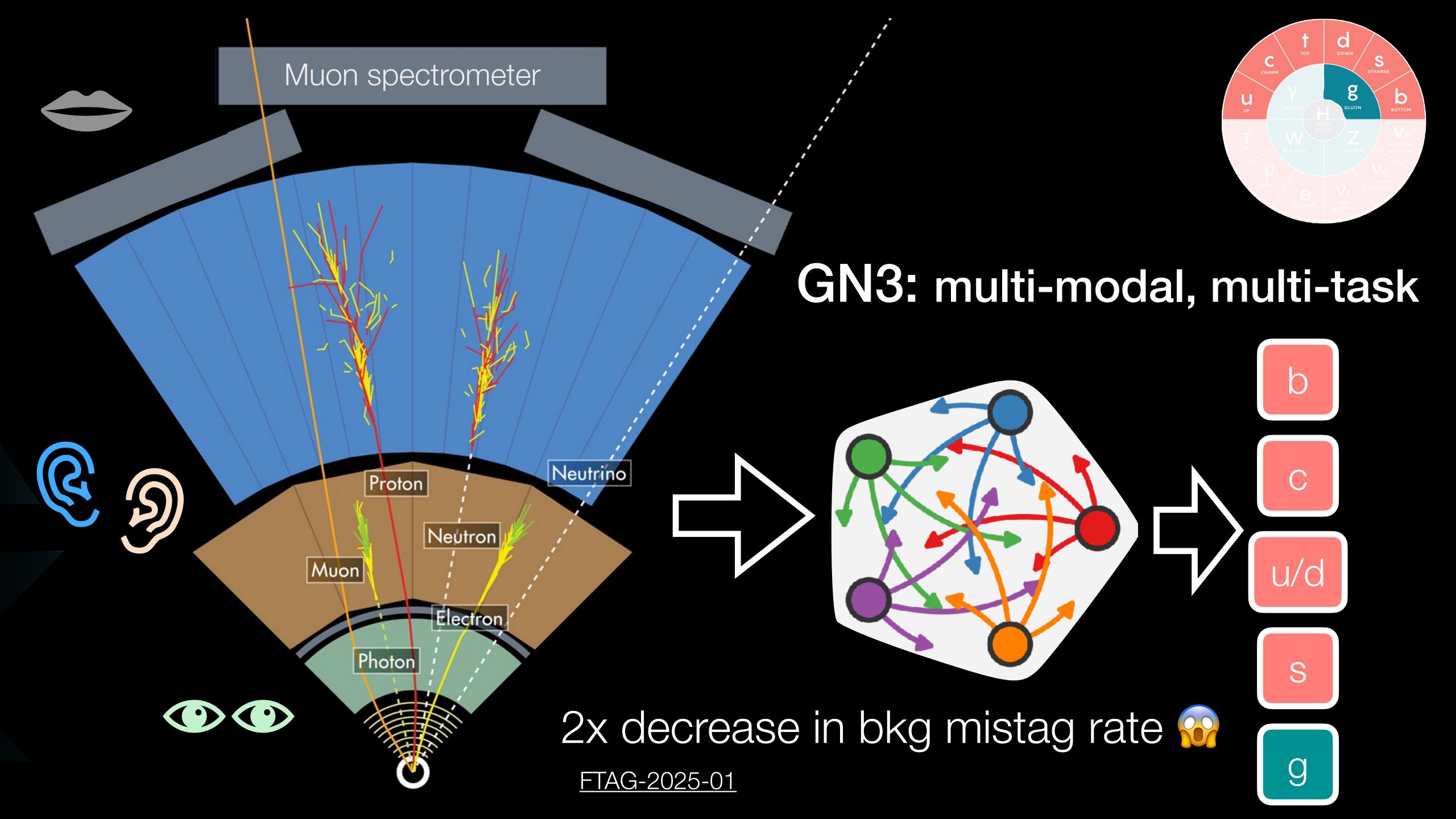






C

light (ud,s,gluon)



Consity ratio

Density ratio intro

Generalizes the ratio of two histograms to high dimensions

Histograms



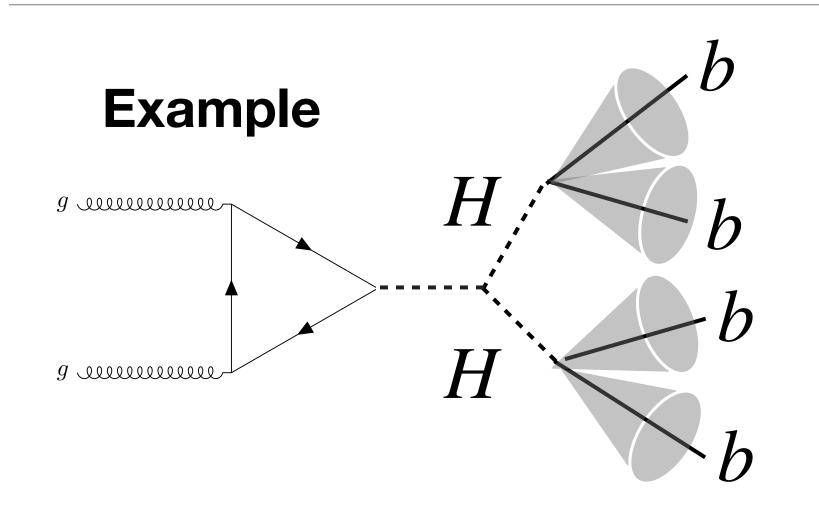
Neural density ratio estimation

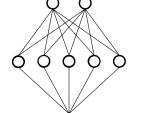
$$NN_{\Theta}(x) = \frac{p_{A}(x)}{p_{B}(x)}, \quad x \in \mathbb{R}^{d}$$

Loss: binary cross entropy

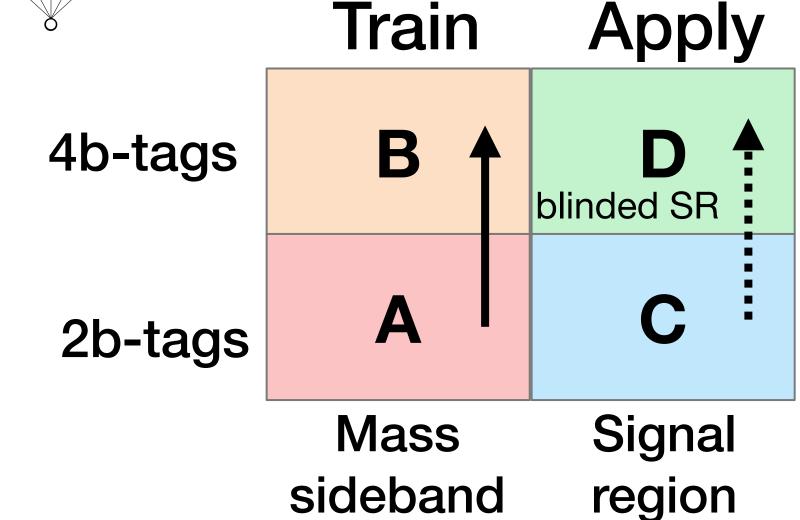
$$\mathcal{L} = -\sum_{x_i \in A} \log \mathsf{NN}_{\theta}(x_i) - \sum_{x_i \in B} \log \left(1 - \mathsf{NN}_{\theta}(x_i)\right)$$

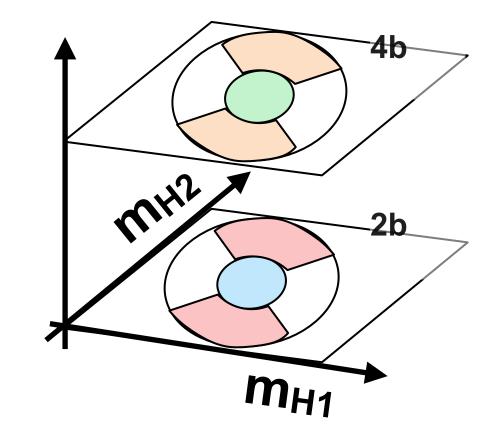
Background estimation

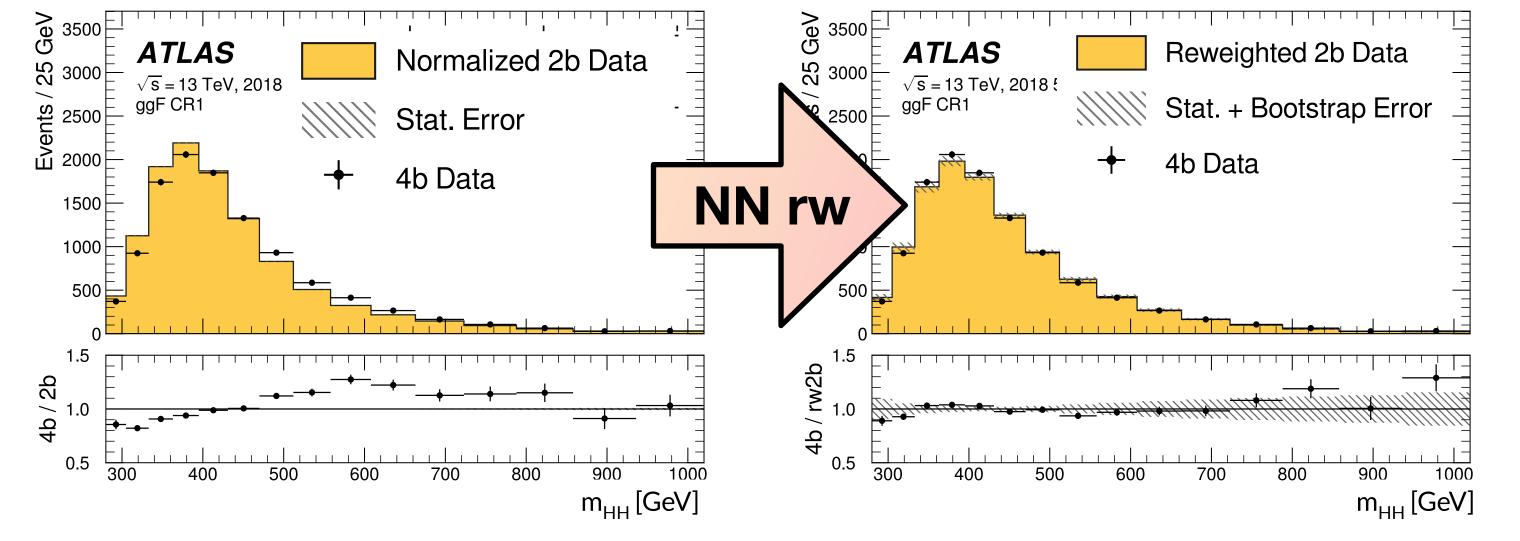




Generalized ABCD method







$$p_{4b} = w(x) \cdot p_{2b}(x)$$

Key assumption:

WCR(x) is a valid \approx of WSR(x).

1911.00405

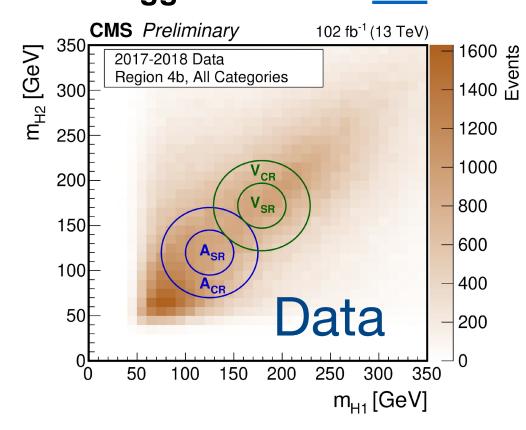
2202.07288

2301.03212

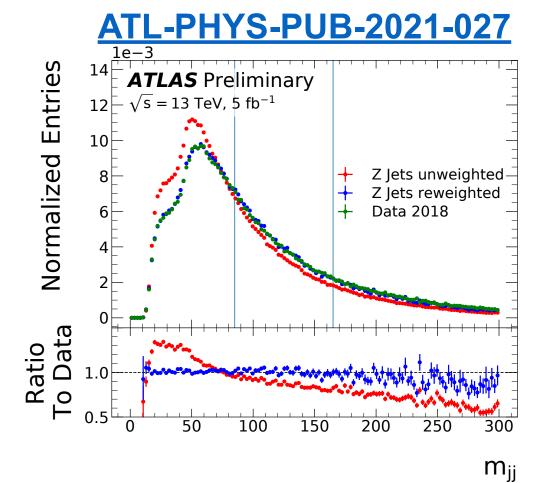
Wide spread adoption

CMS HH4b NR

BDT rw: <u>2202.09617</u> + Higgs Pairs 2022 <u>talk</u>

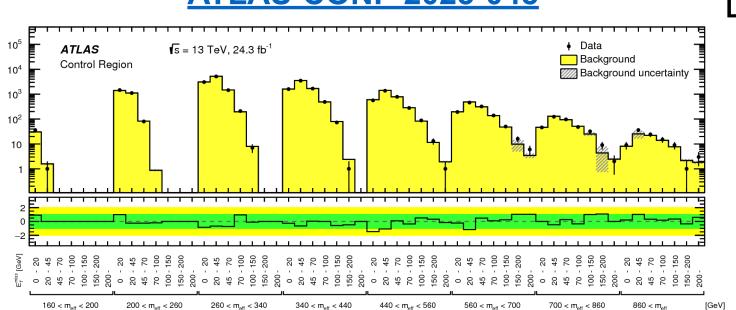


gg tagger Match MC -> data



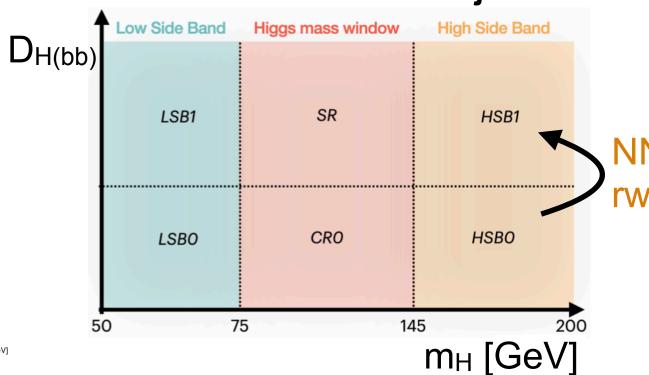
HH4b + MET

BDT rw <u>1806.04030</u> and <u>ATLAS-CONF-2023-048</u>



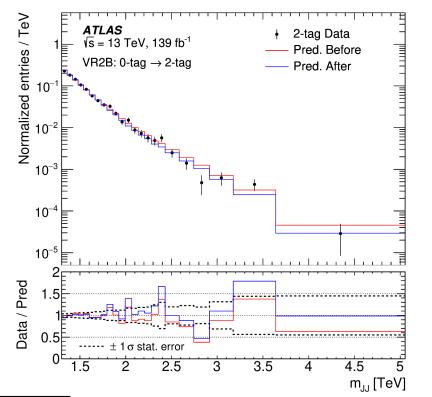
Y -> XH <u>2306.03637</u>

NN rw for anomalous jets



VH(qqbb)

BDT 2007.05293



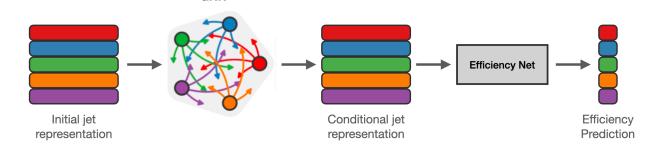
Search for generic resonances

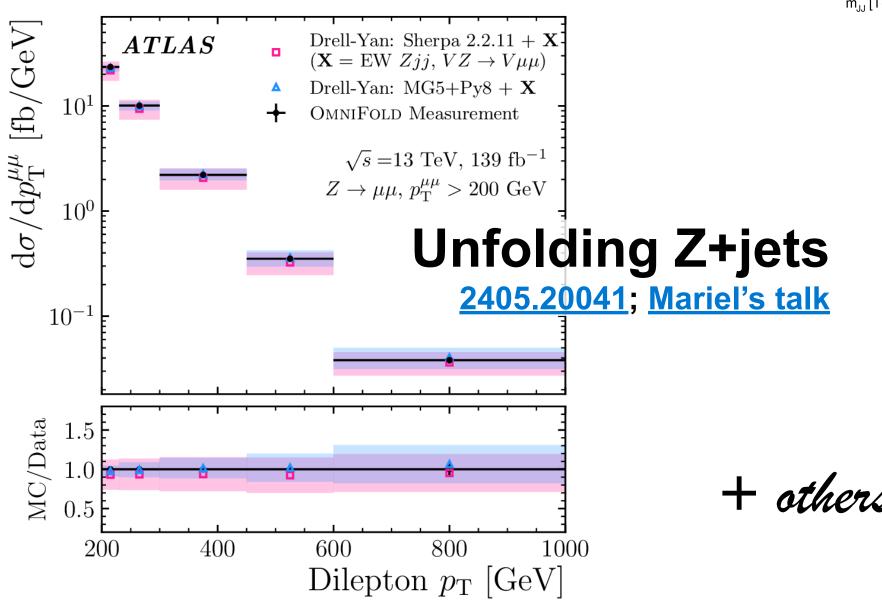
MC -> data reweighting in sideband 2502.09770



FTAG truth -> reco

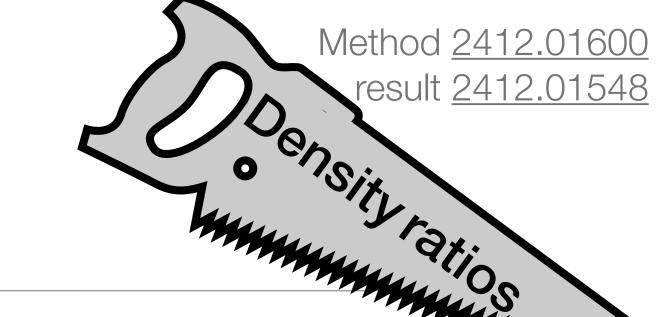
ATL-PHYS-PUB-2022-041

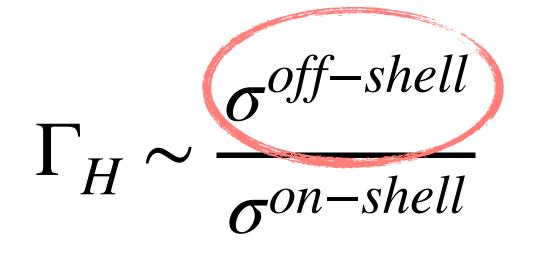


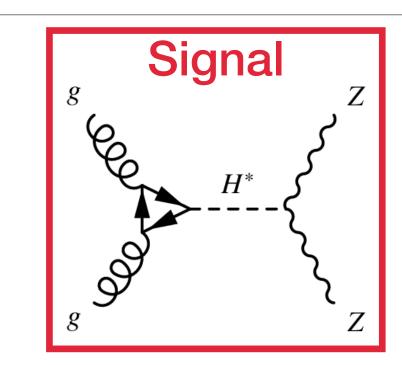


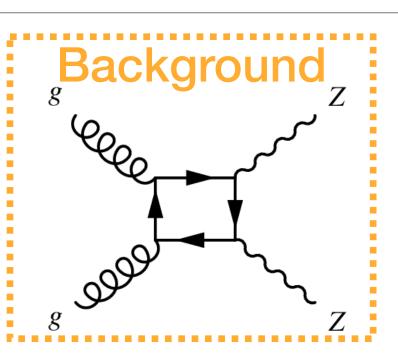


Simulation-based inference





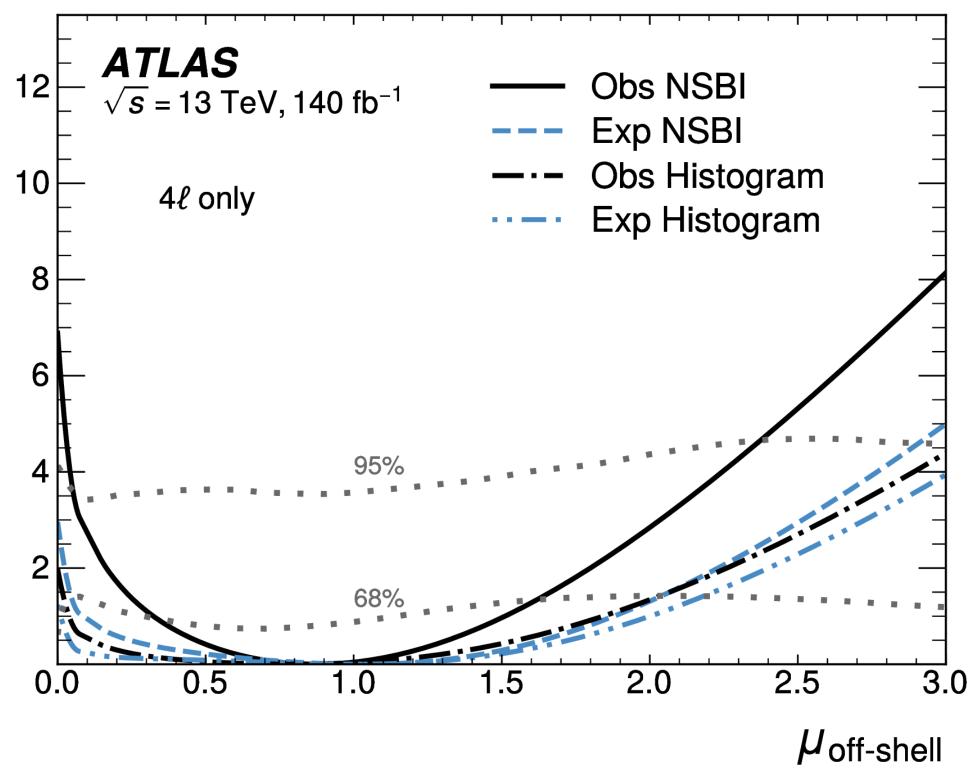




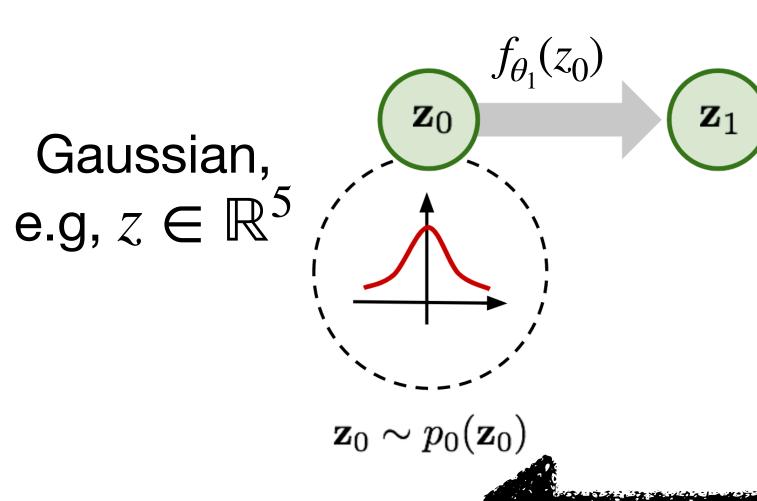
$$P_{off-shell}(\mu) = \mu P_S + \sqrt{\mu} P_I + P_B$$

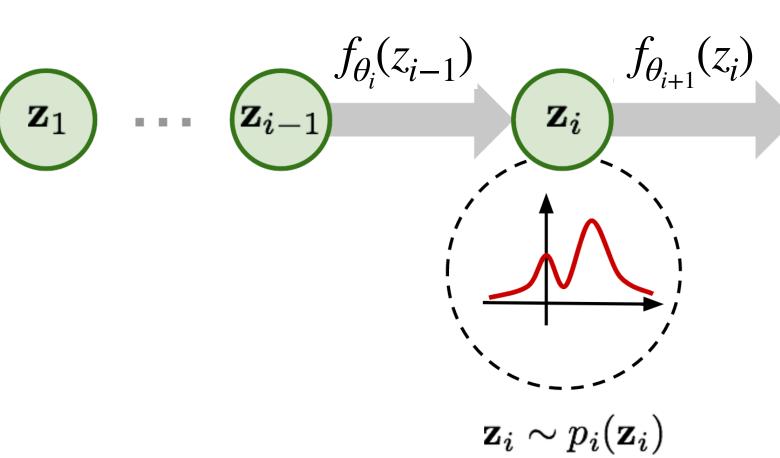
Kinemmatics depend on μ ... optimal classifier changes for each μ .

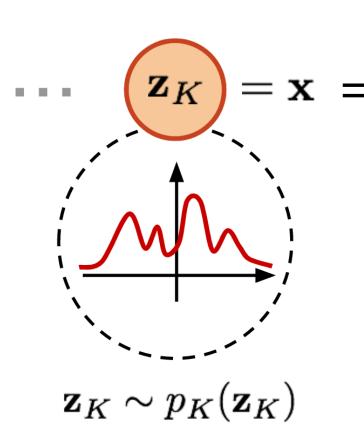
$$\frac{p(x|\mu)}{p'(x)} = \frac{1}{\sigma(\mu)} \left(\mu \sigma_s \frac{p_S(x)}{p'(x)} + \sqrt{\mu} \sigma_I \frac{p_I(x)}{p'(x)} + \sigma_B \frac{p_B(x)}{p'(x)} \right)$$
Eackup

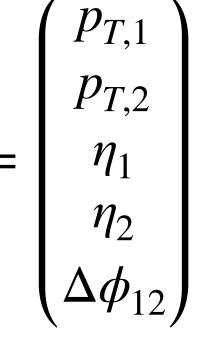


Learn $p_{\phi}(x)$









Predict momenta of particles 1,2

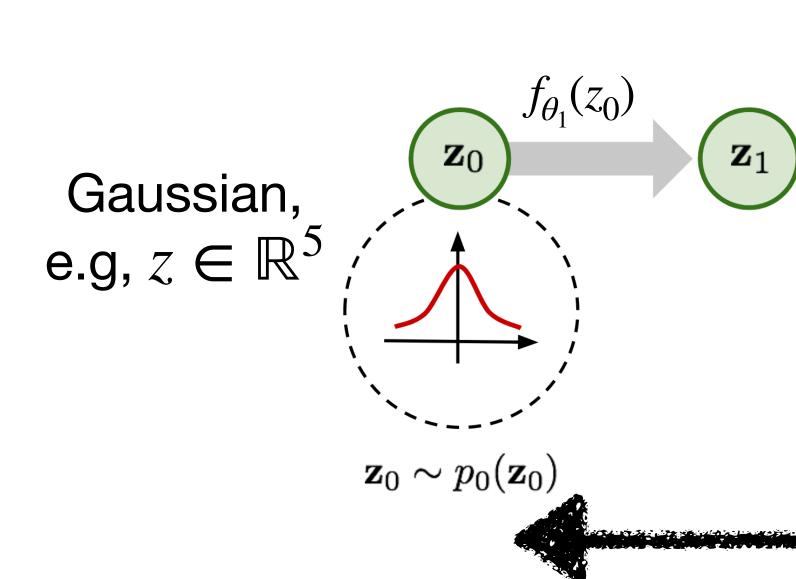
blog post

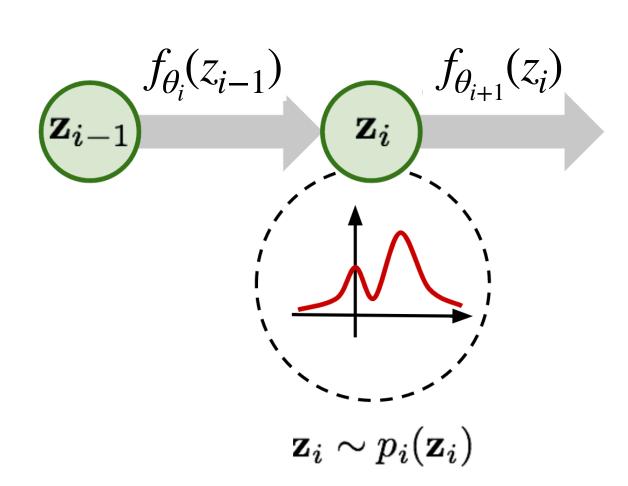
$$f_{\phi}^{-1} = f_{\phi_1}^{-1} \circ f_{\phi_2}^{-1} \circ \cdots \circ f_{\phi_L}^{-1}$$

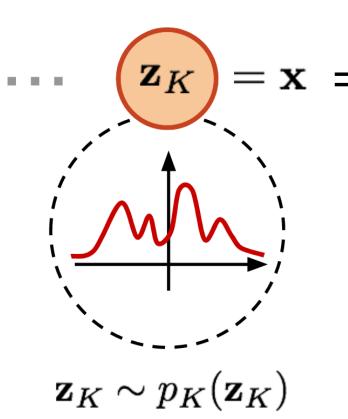
Invertible f_{ϕ_i} for density of training samples



Learn $p_{\phi}(x)$







$$= \begin{pmatrix} p_{T,1} \\ p_{T,2} \\ \eta_1 \\ \eta_2 \\ \Delta \phi_{12} \end{pmatrix}$$

Predict momenta of particles 1,2

blog post

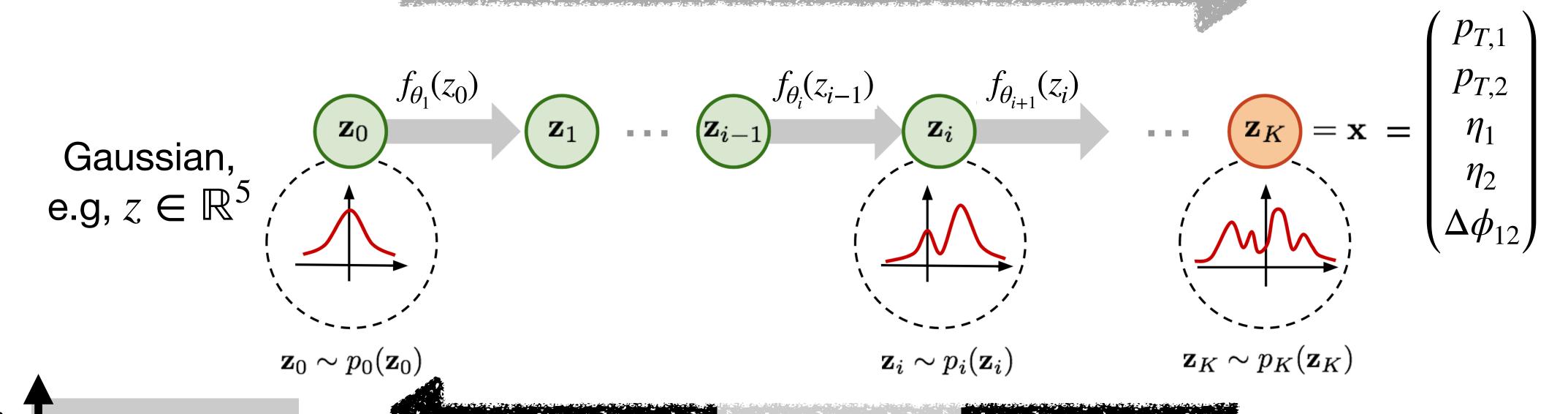
$$f_{\phi}^{-1} = f_{\phi_1}^{-1} \circ f_{\phi_2}^{-1} \circ \cdots \circ f_{\phi_L}^{-1}$$

Invertible f_{ϕ_i} for density of training samples



$$\mathscr{L}oss = -\log p_{\phi}(x) = -\log p_{z}(f_{\phi}^{-1}(x)) - \sum_{i=1}^{K} \left| \frac{\partial J_{\phi_{i}}}{\partial z_{i}^{T}} \right|$$

Learn $p_{\phi}(x|y)$



N₂ Sideband Signal region

 m_1

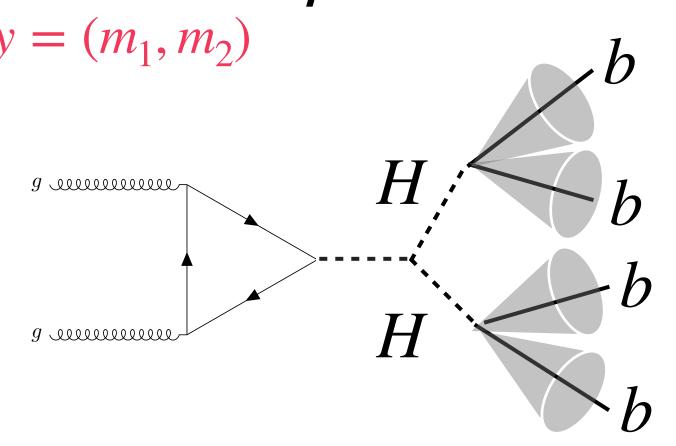
$$f_{\phi}^{-1} = f_{\phi_1}^{-1} \circ f_{\phi_2}^{-1} \circ \cdots \circ f_{\phi_L}^{-1}$$

Conditional density: $p_{\theta}(x \mid y)$, $y = (m_1, m_2)$

$$\mathcal{L}oss = -\log p_{\phi}(x|\mathbf{y}) = -\log p_{z}(f_{\phi}^{-1}(x|\mathbf{y})) - \sum_{i=1}^{|\mathcal{O}f_{\phi_{i}}|} \frac{\partial f_{\phi_{i}}}{\partial z_{i}^{T}}$$

blog post

Learn $p_{\phi}(x|y)$ $y = (m_1, m_2)$



Learns about preprocessing cut $(\Delta \eta_{HH} < 1.5)$

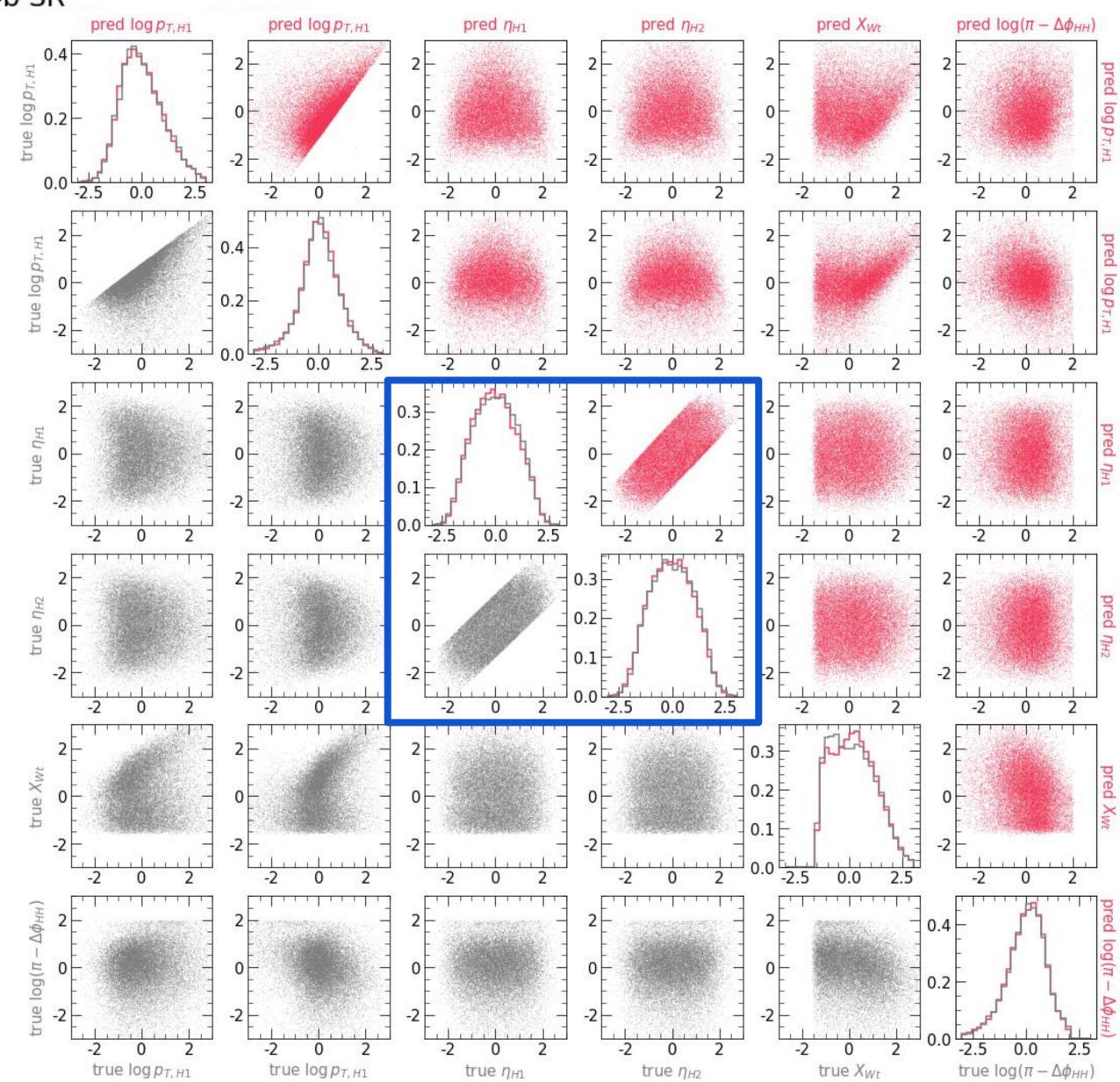
Sideband Signal region m₁

data

ATLAS Thesis

 $\sqrt{s} = 13 \text{ TeV}, 126 \text{ fb}^{-1}$ 4b SR

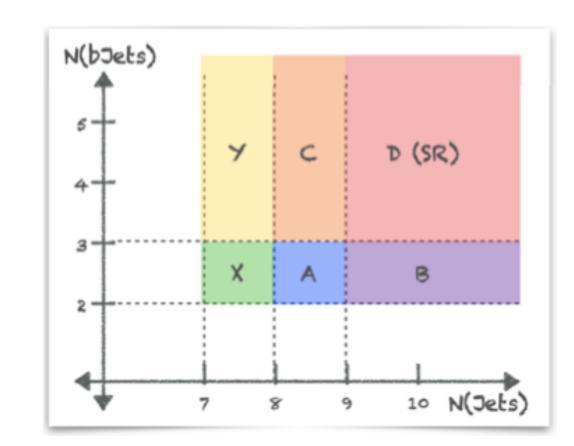
flow pred



Bkg estimation (when we don't trust QCD simulation)

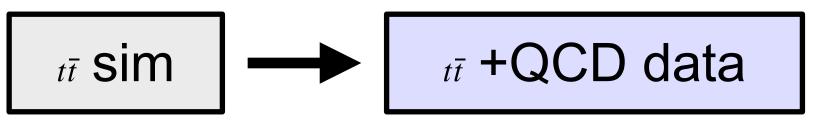
- · CMS 4 tops <u>evidence</u> (2023), M. Quinnan's ML4Jets <u>slides</u>
- Methods / pheno papers [1], [2]
- ATLAS genric search, <u>2502.09770</u>
- · CMS genric search, 2412.03747, seminar

Learn how to interpolate into the SR



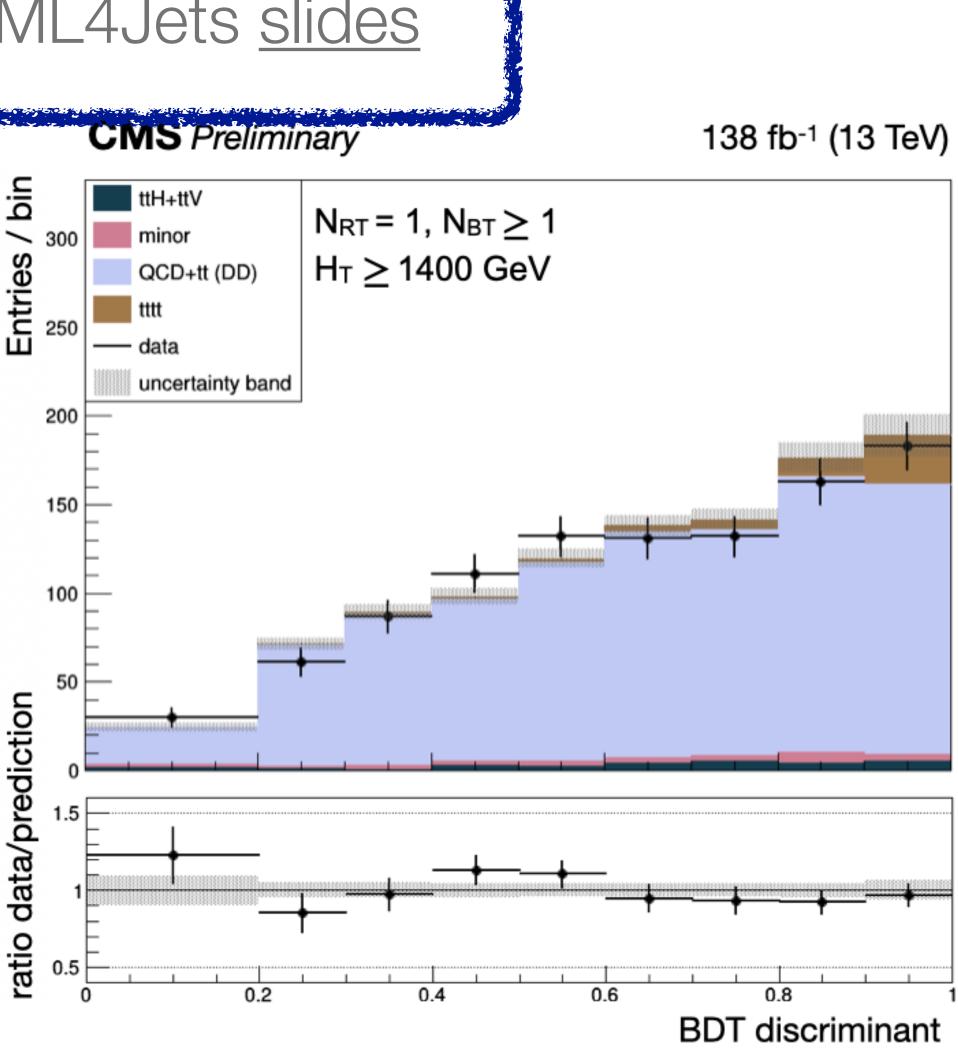
Shape

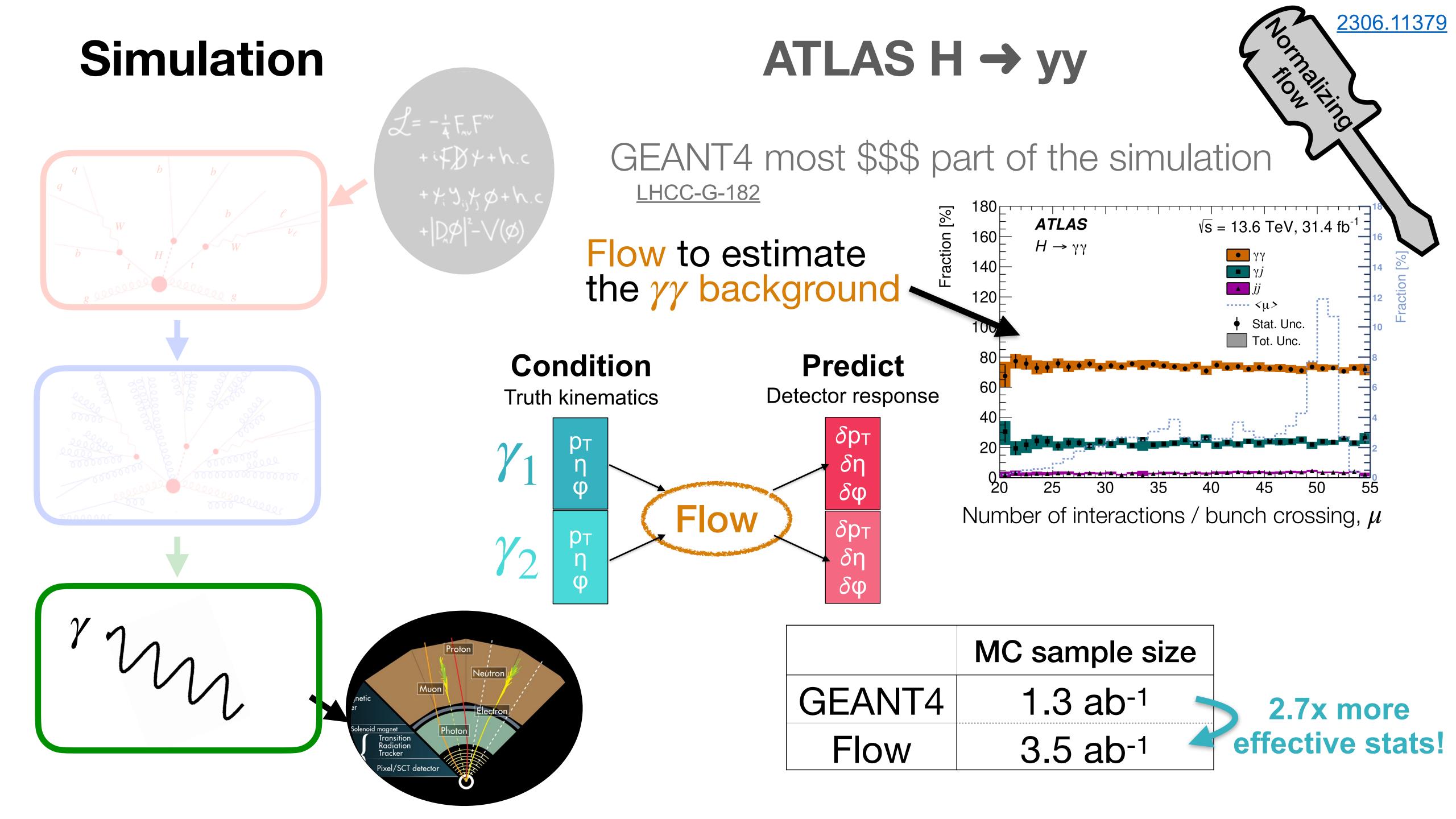
Based on a flow:





First public result with LHC data using flows!





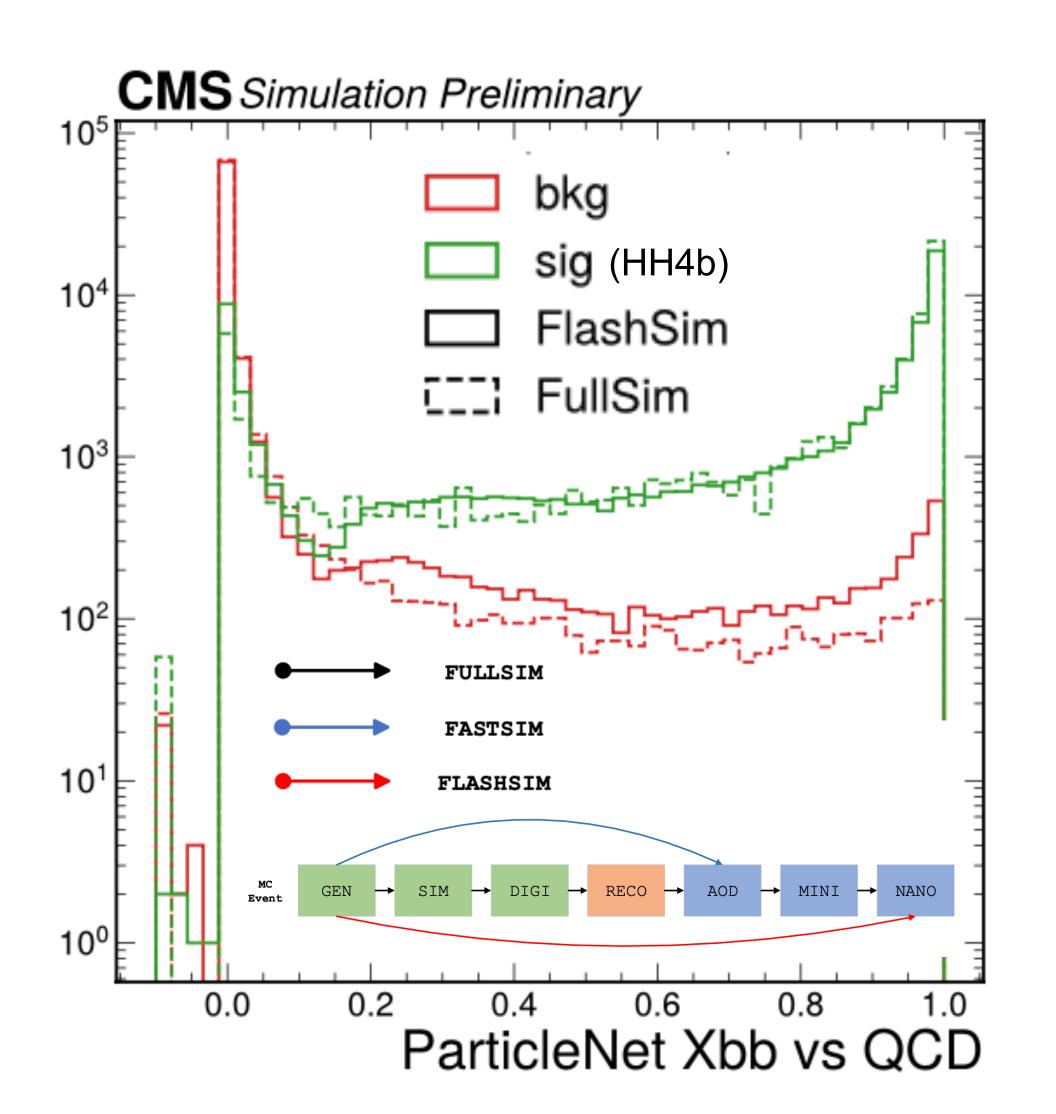
Simulation 2= - + Fr. F~ +iFD++h.c + + + y + p + h.c $+\left|\sum_{x}\phi\right|^{2}-\bigvee(\phi)$ Jeeses

CMS "flash sim"

GEANT4 most \$\$\$ part of the simulation

<u>flashSim</u>

slides





1) ML for our tools

2) ML for our future

- Foundation Models
- Reconstruction = HEP foundation model
- End-to-end analysis
- Self supervised learning
- Incorporating domain expertise

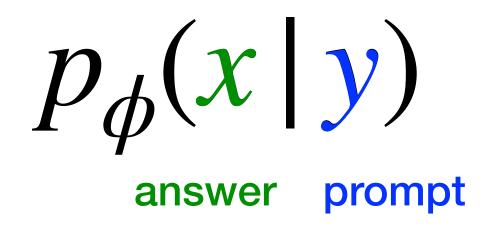
Computer vision

2204.06125

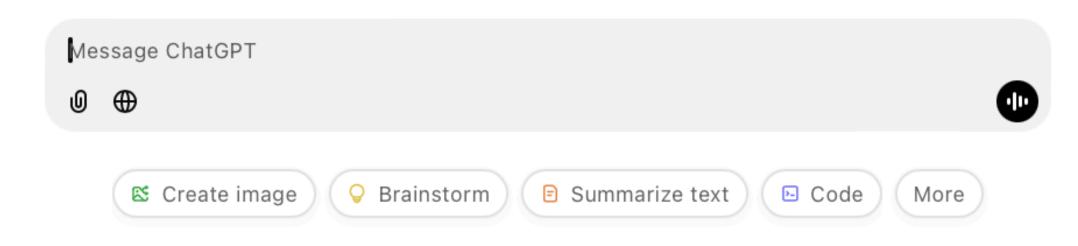


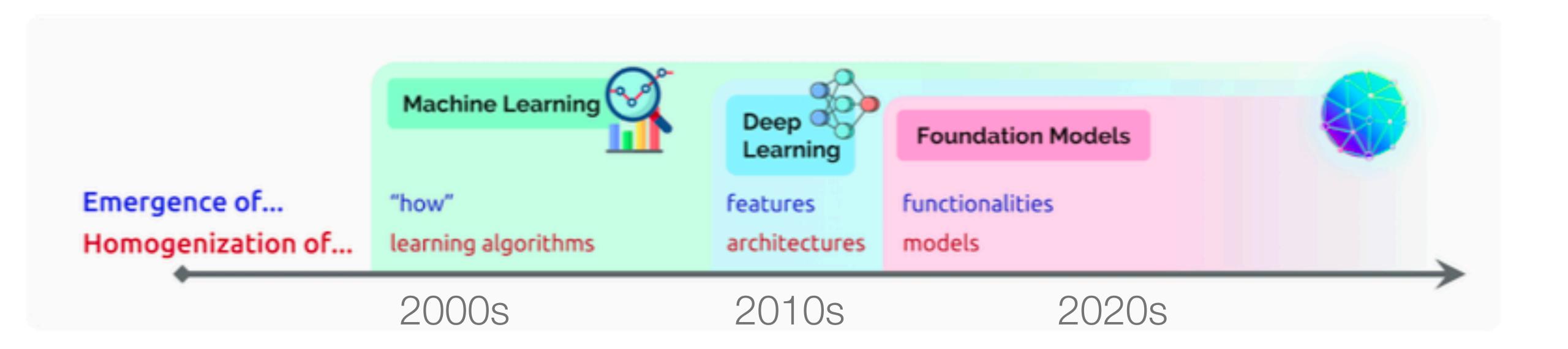
a teddy bear on a skateboard in times square

Natural language

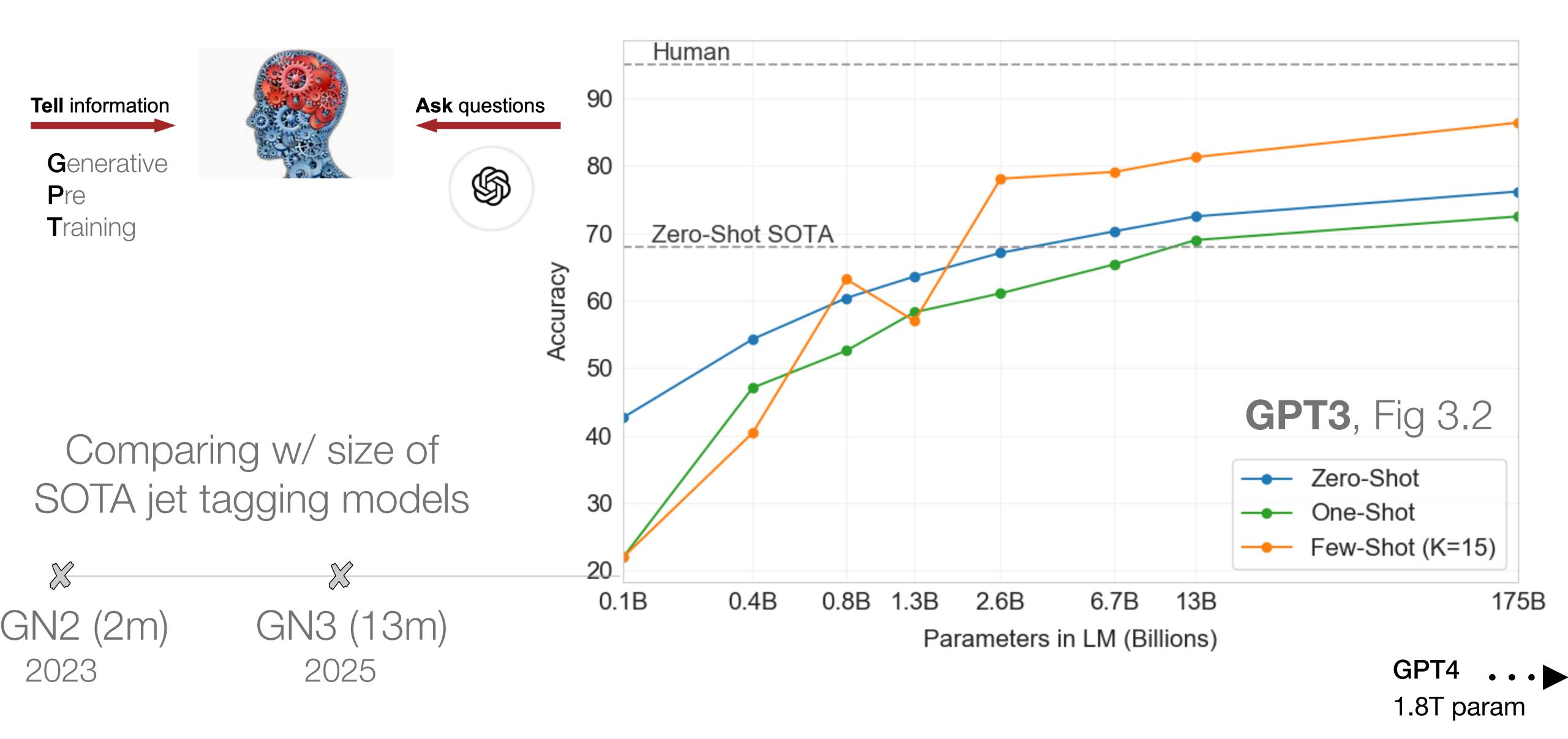


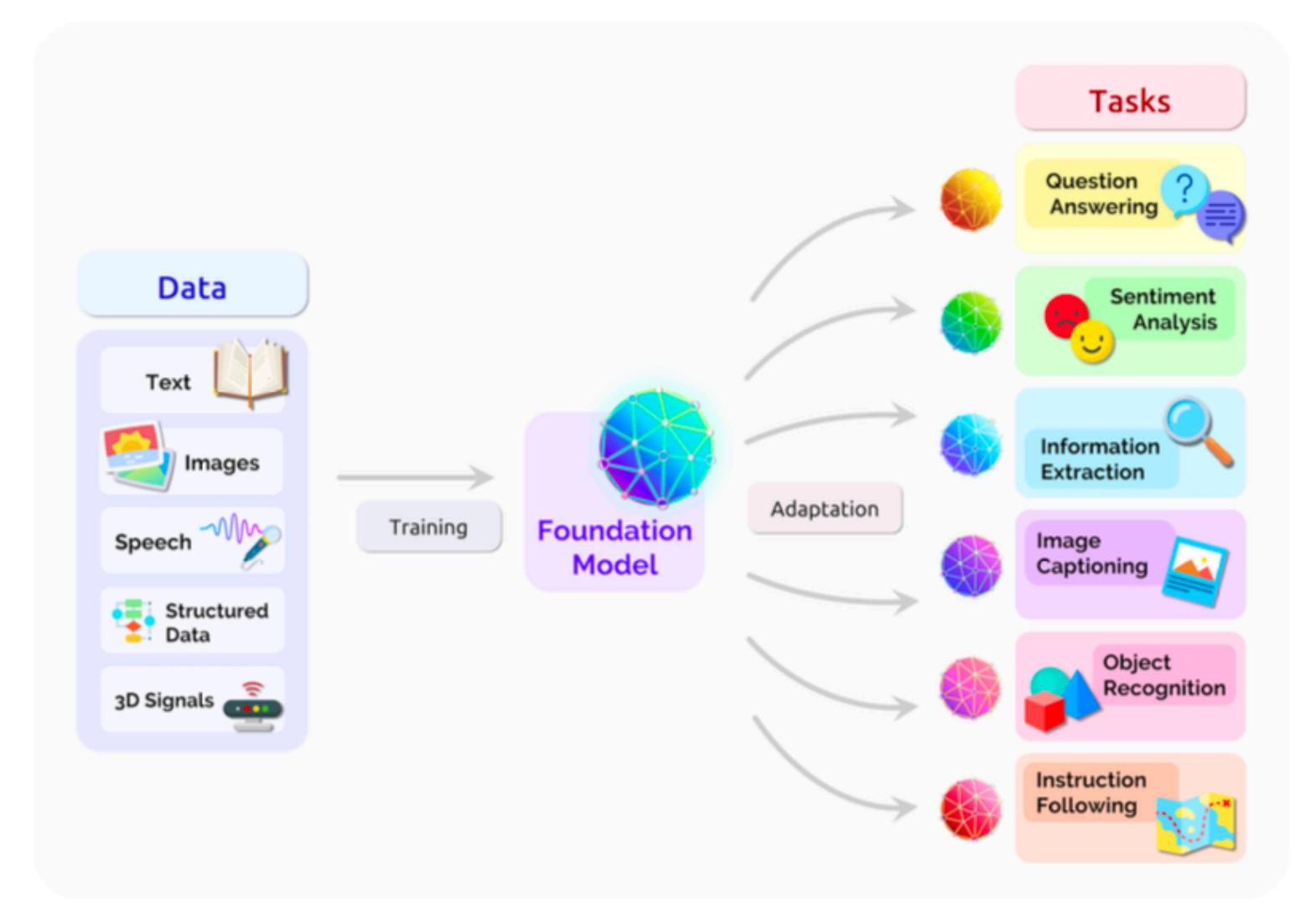
What can I help with?





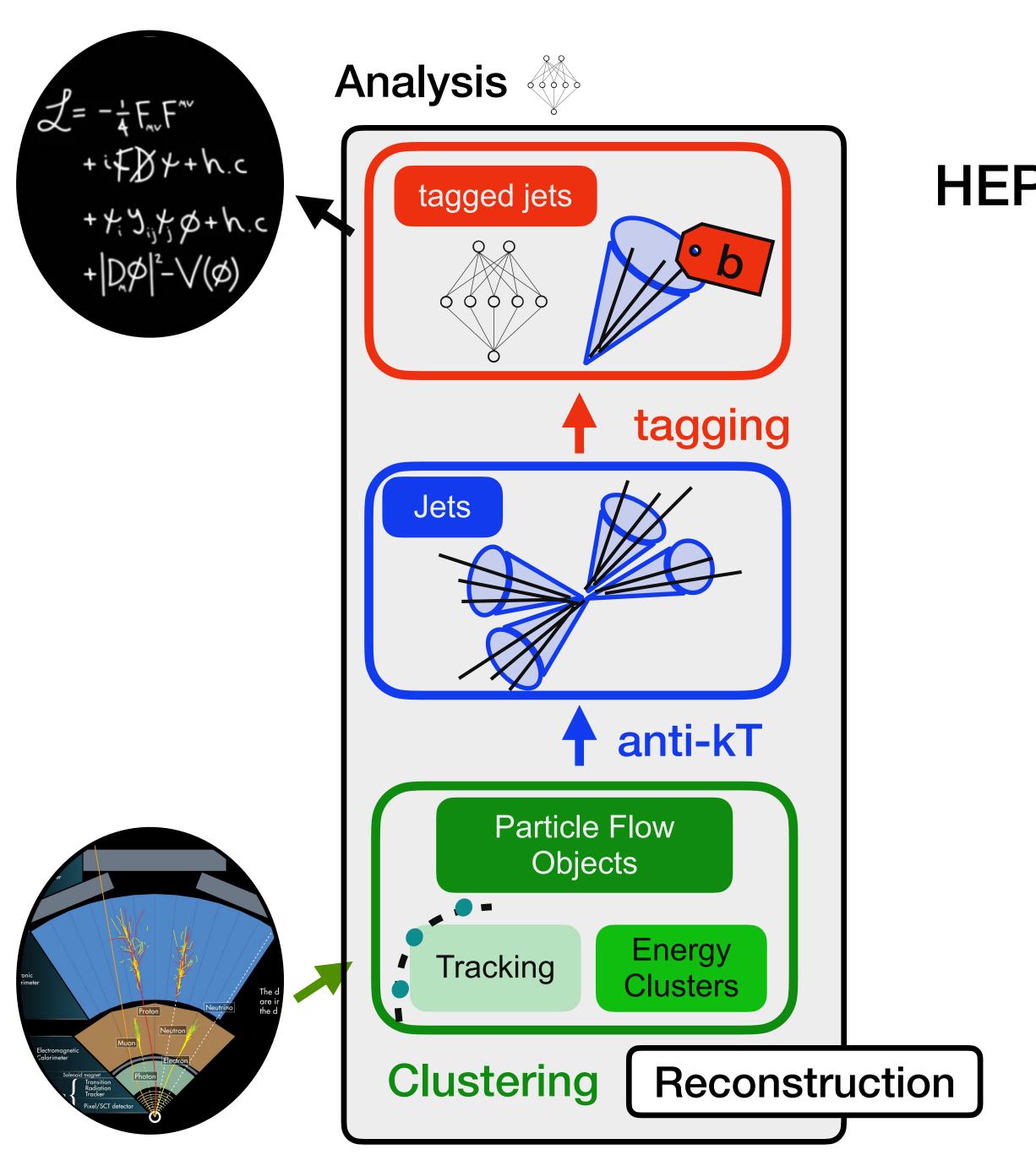
"From a technological point of view, **foundation models** are not new... however, the sheer scale and scope of foundation models from the last few years have stretched our imagination of what is possible."



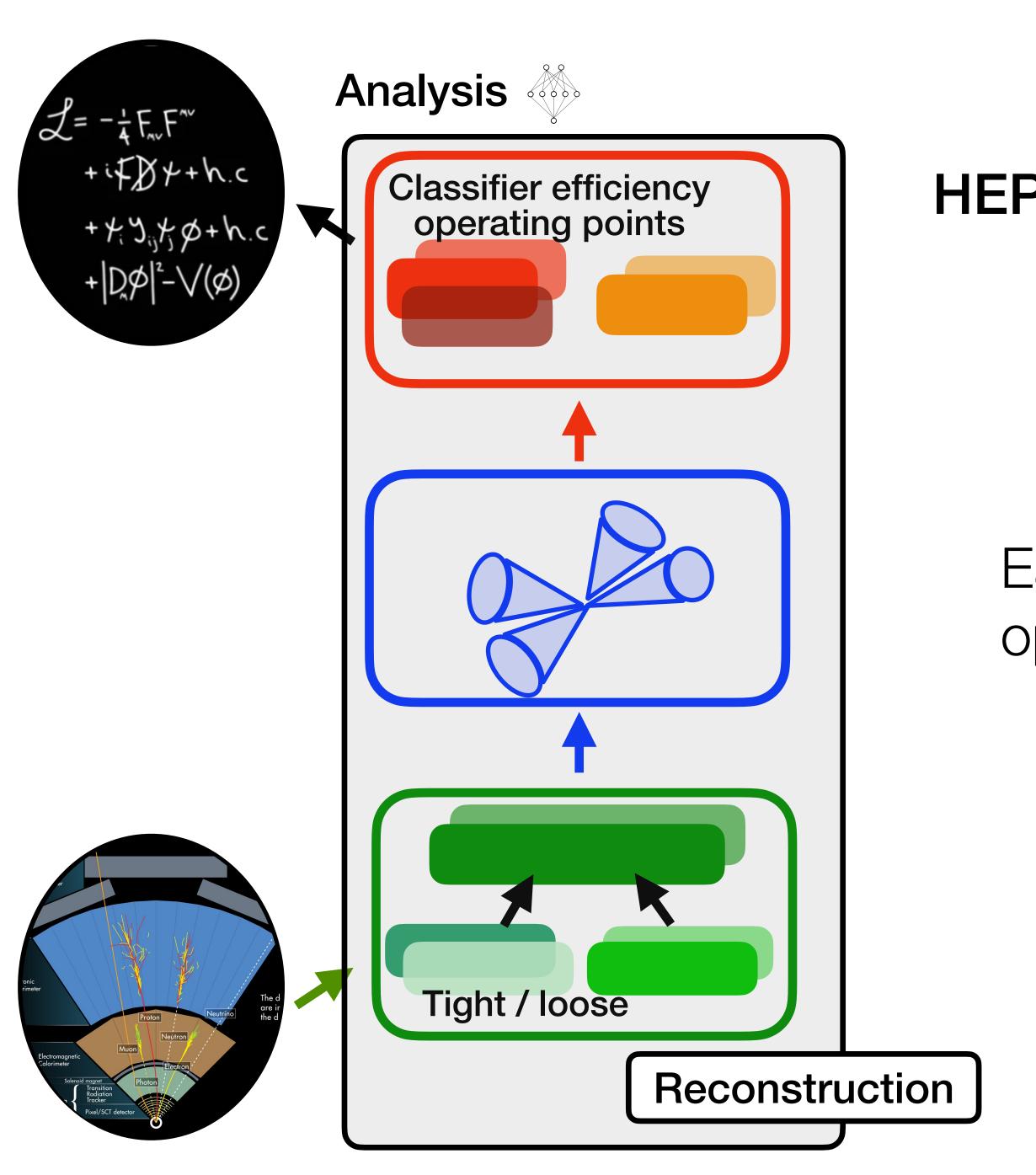




"A **foundation model** is any model that is trained on broad data (generally self-supervision at scale) that can be adapted (e.g, fine-tuned) to a wide range of downstream tasks."

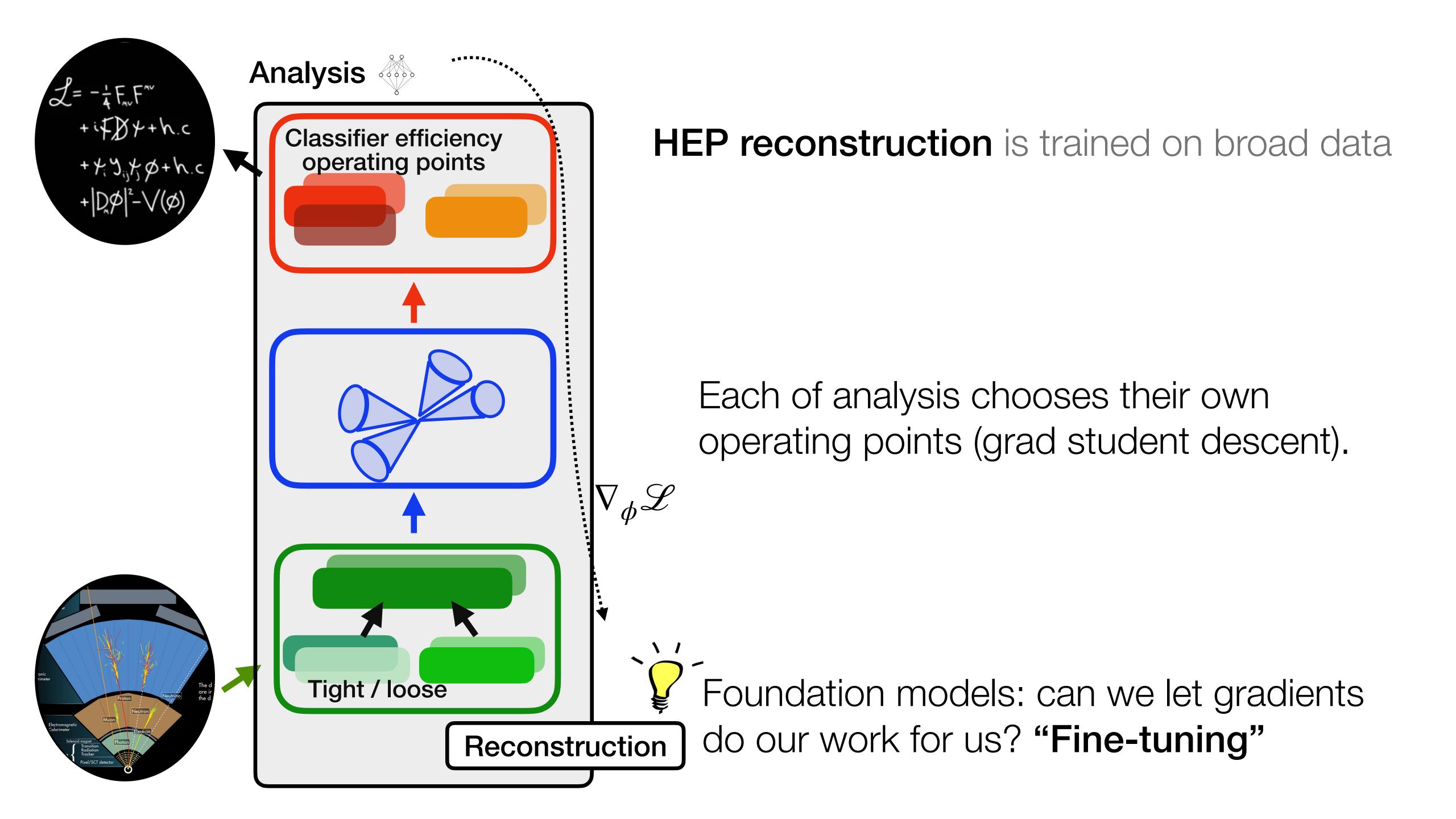


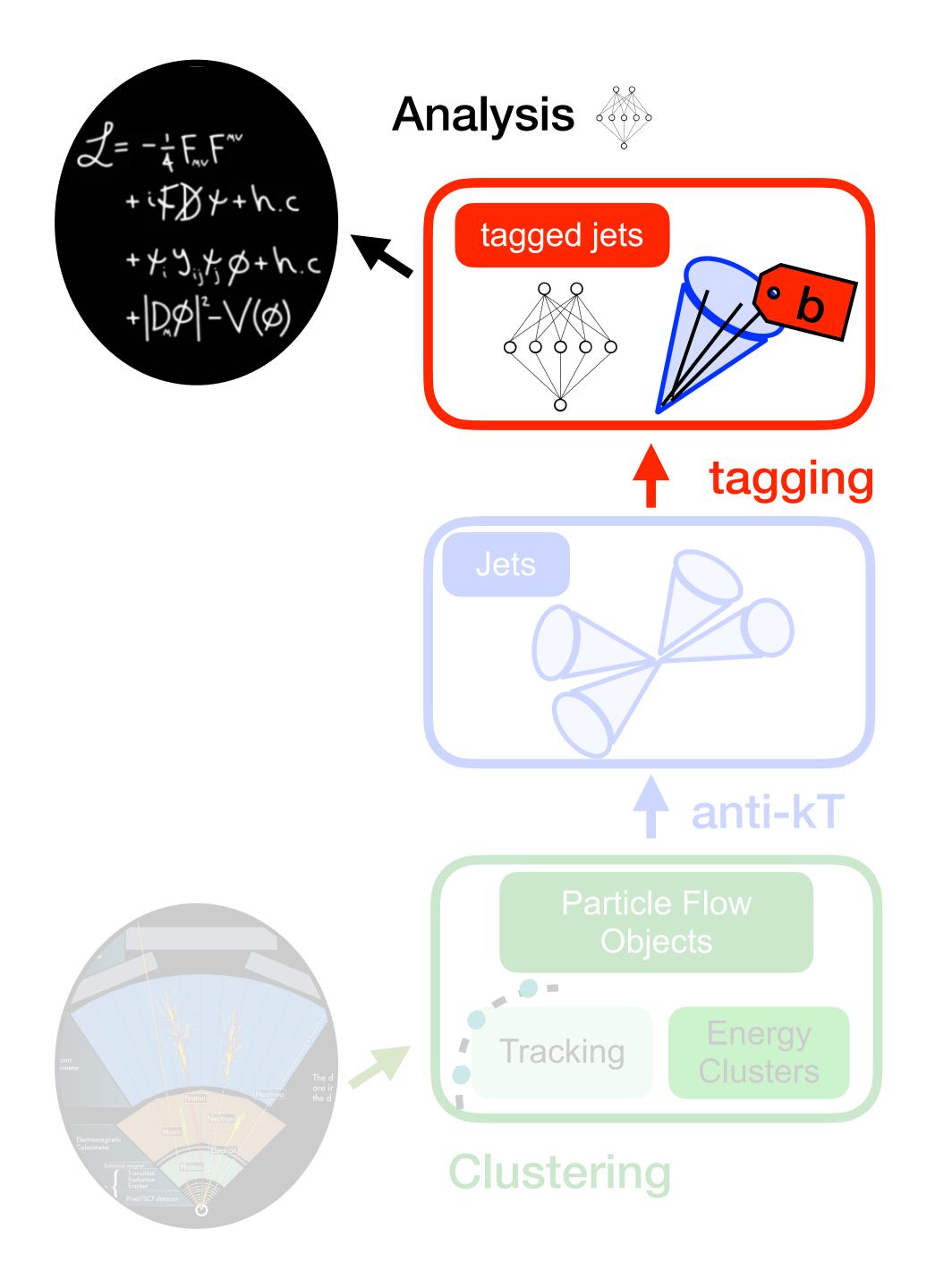
HEP reconstruction optimized on broad data



HEP reconstruction is trained on broad data

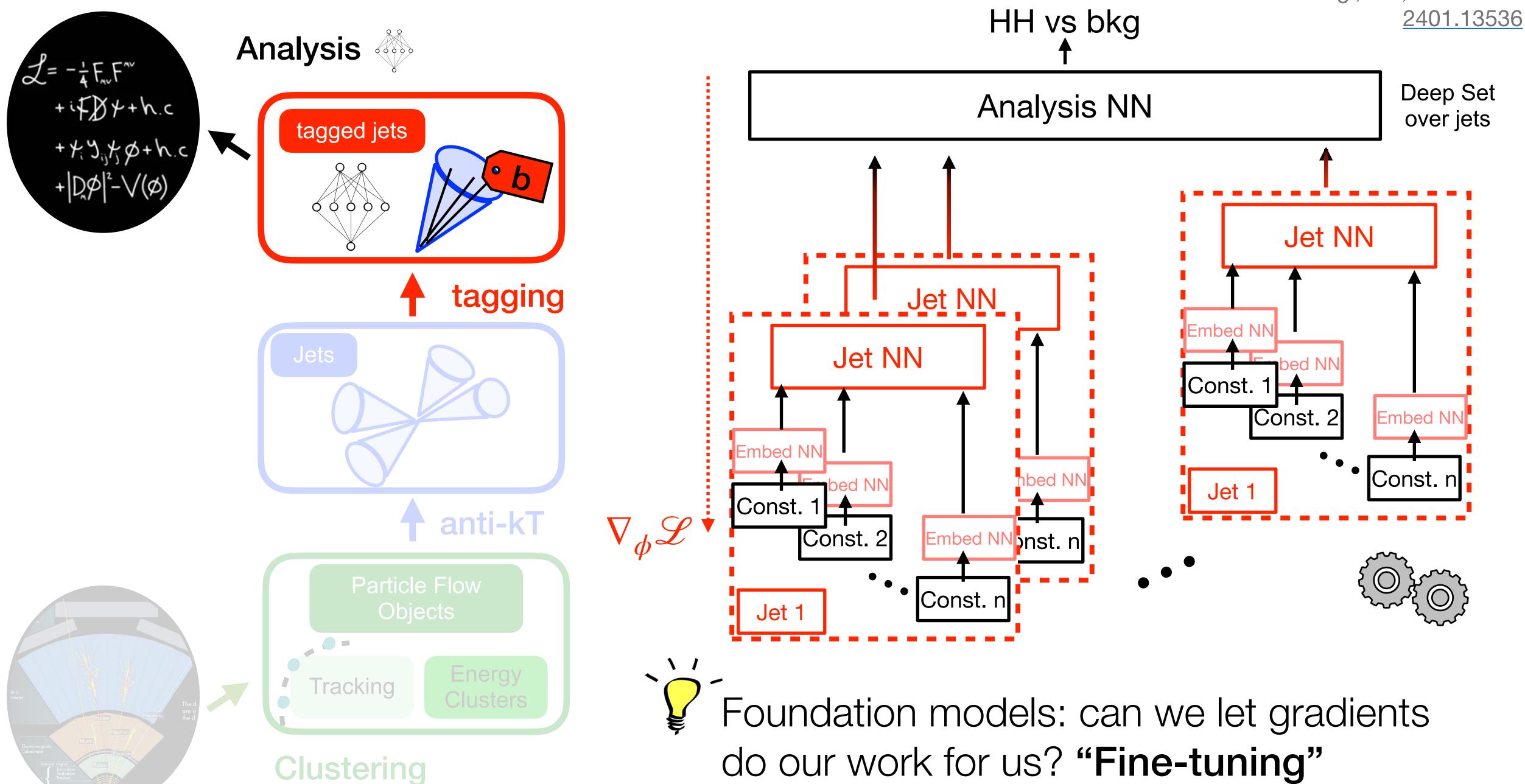
Each of analysis chooses their own operating points (grad student descent).

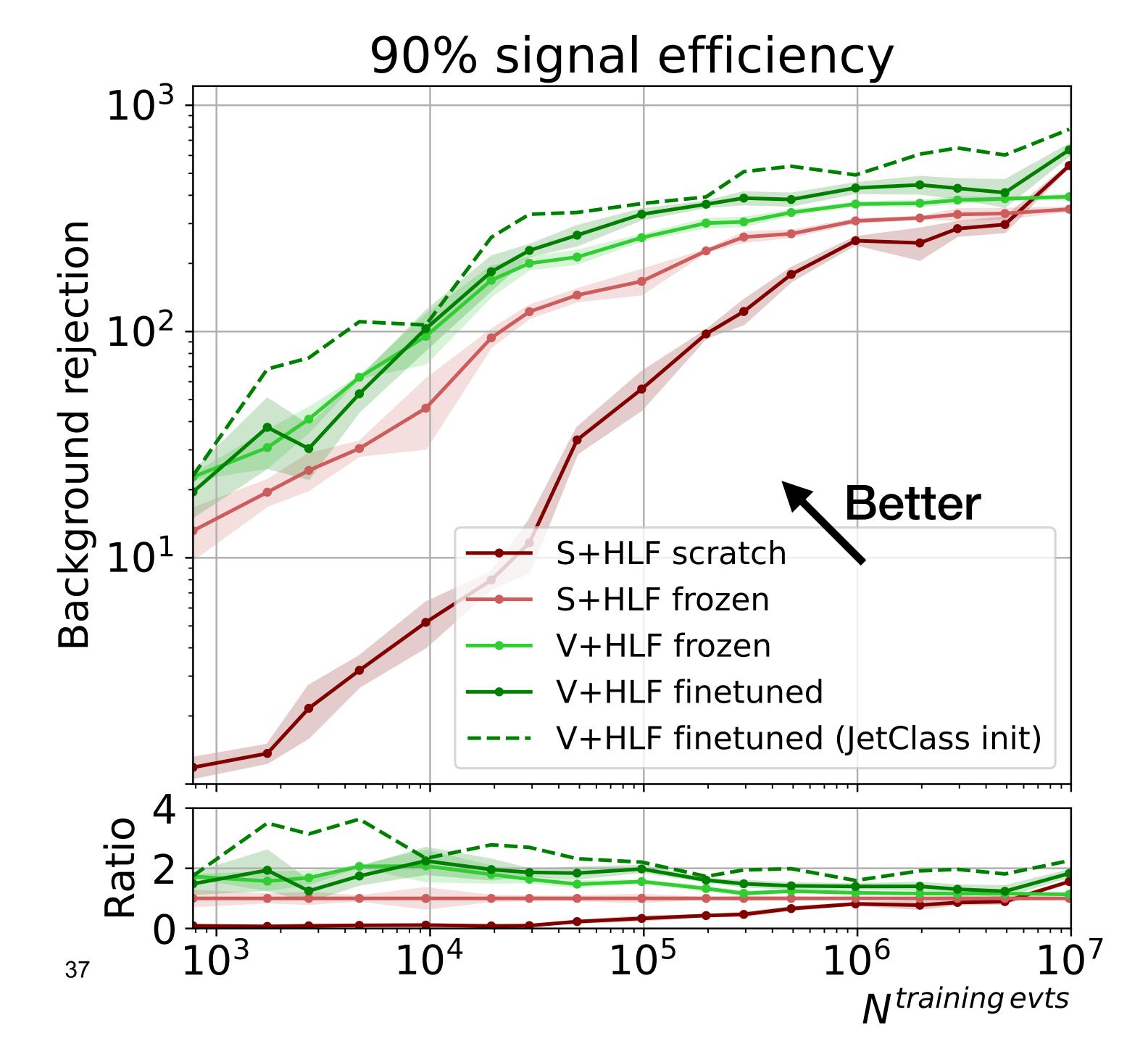




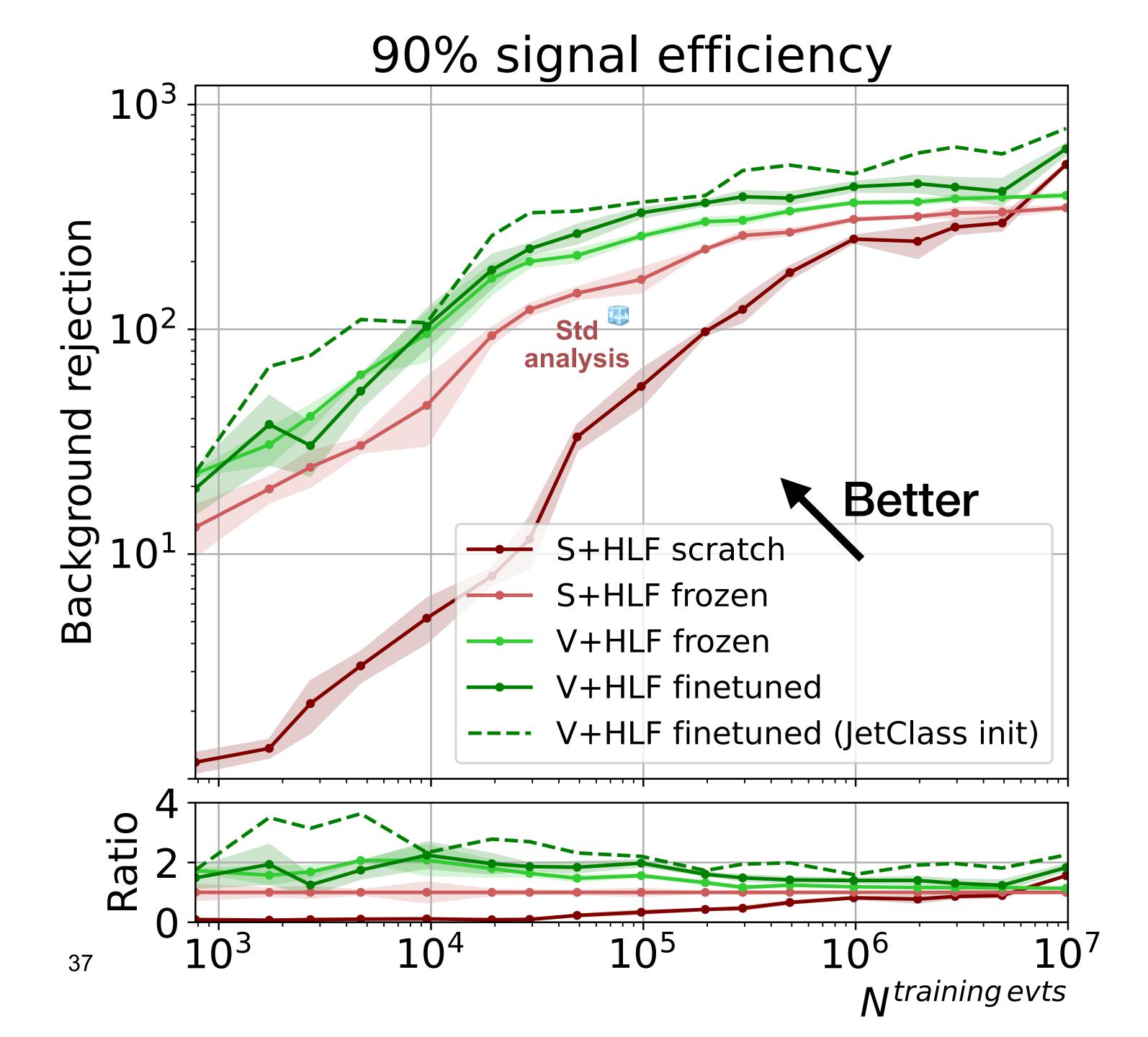


Foundation models: can we let gradients do our work for us? "Fine-tuning"

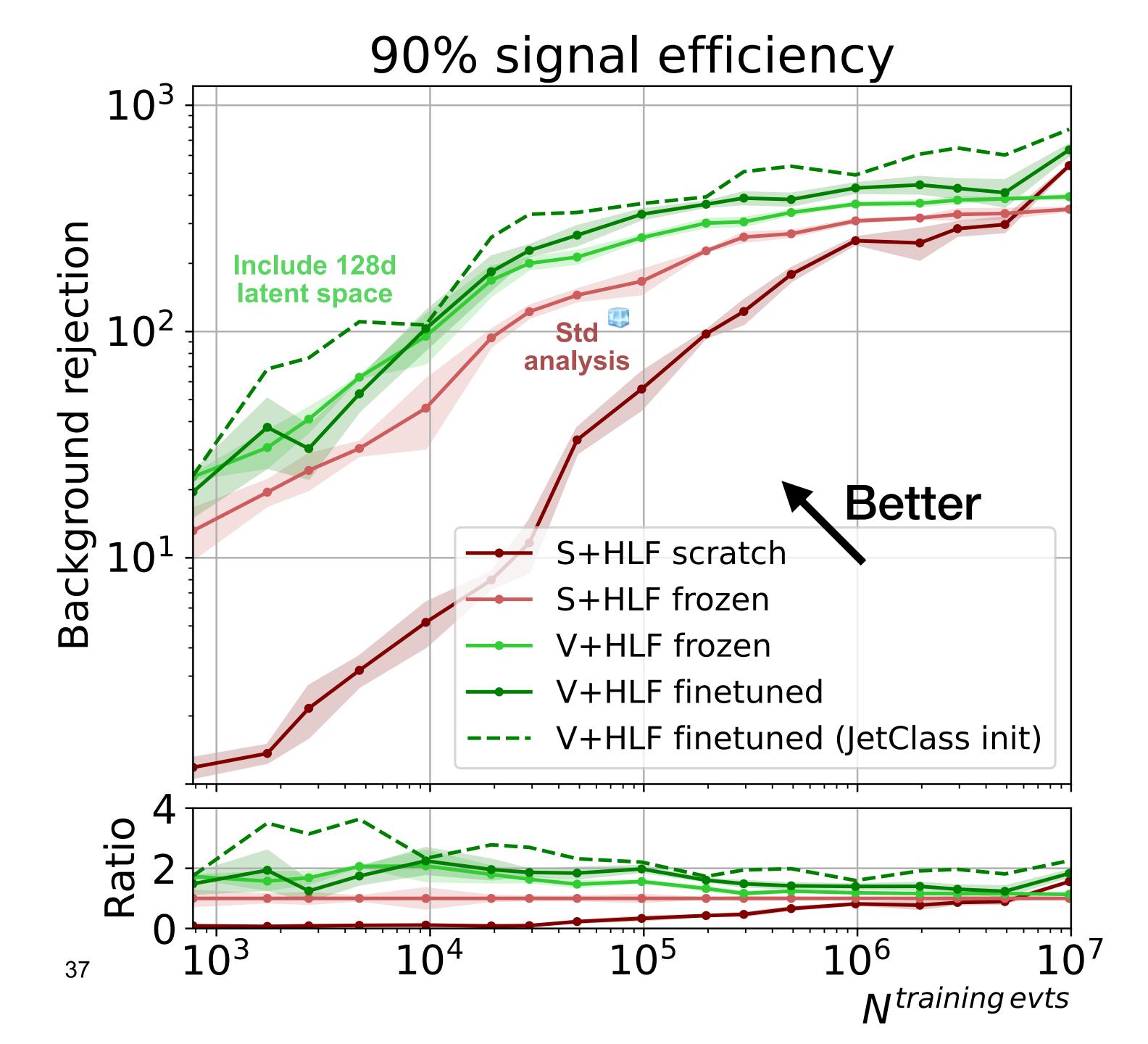






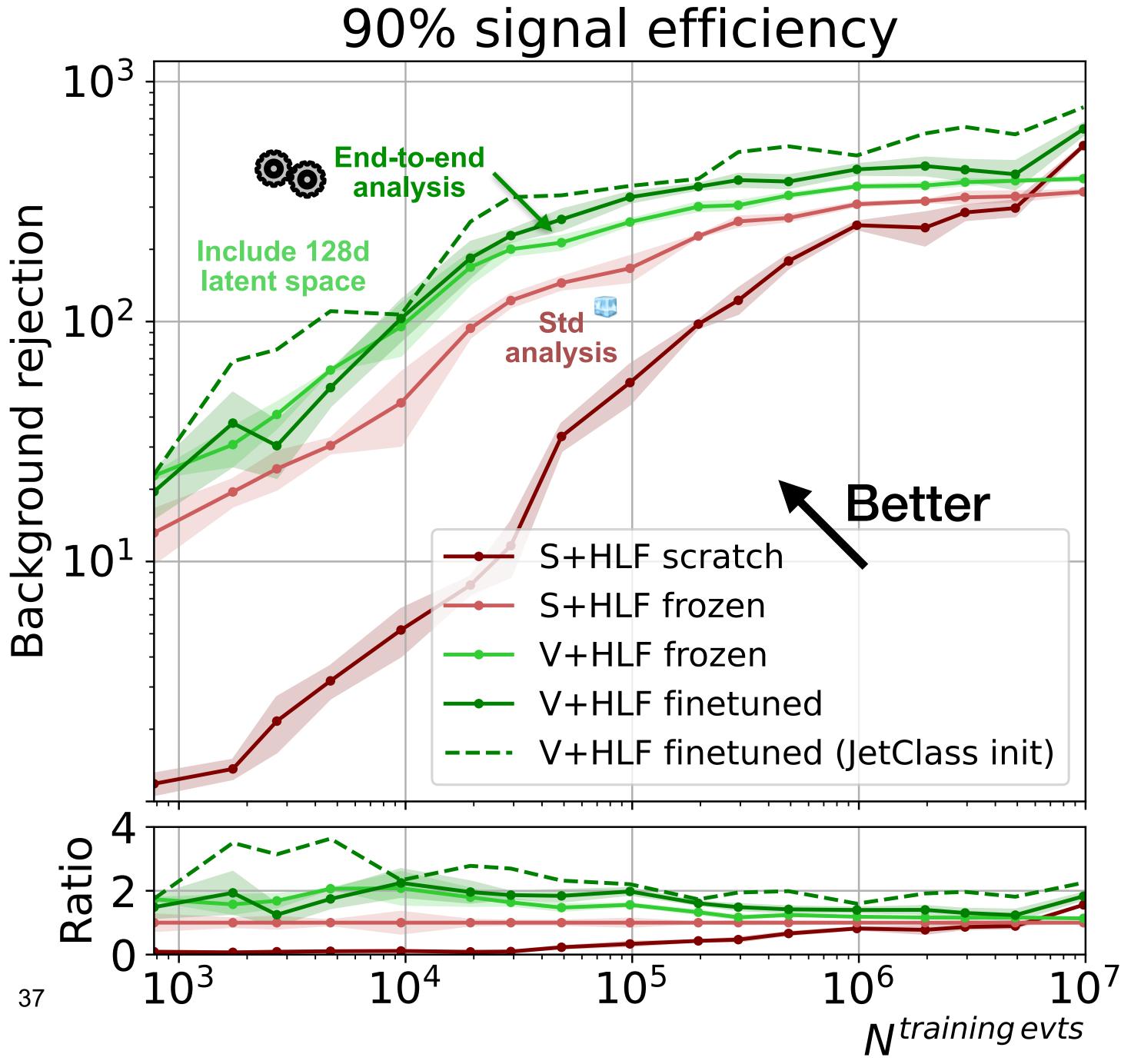








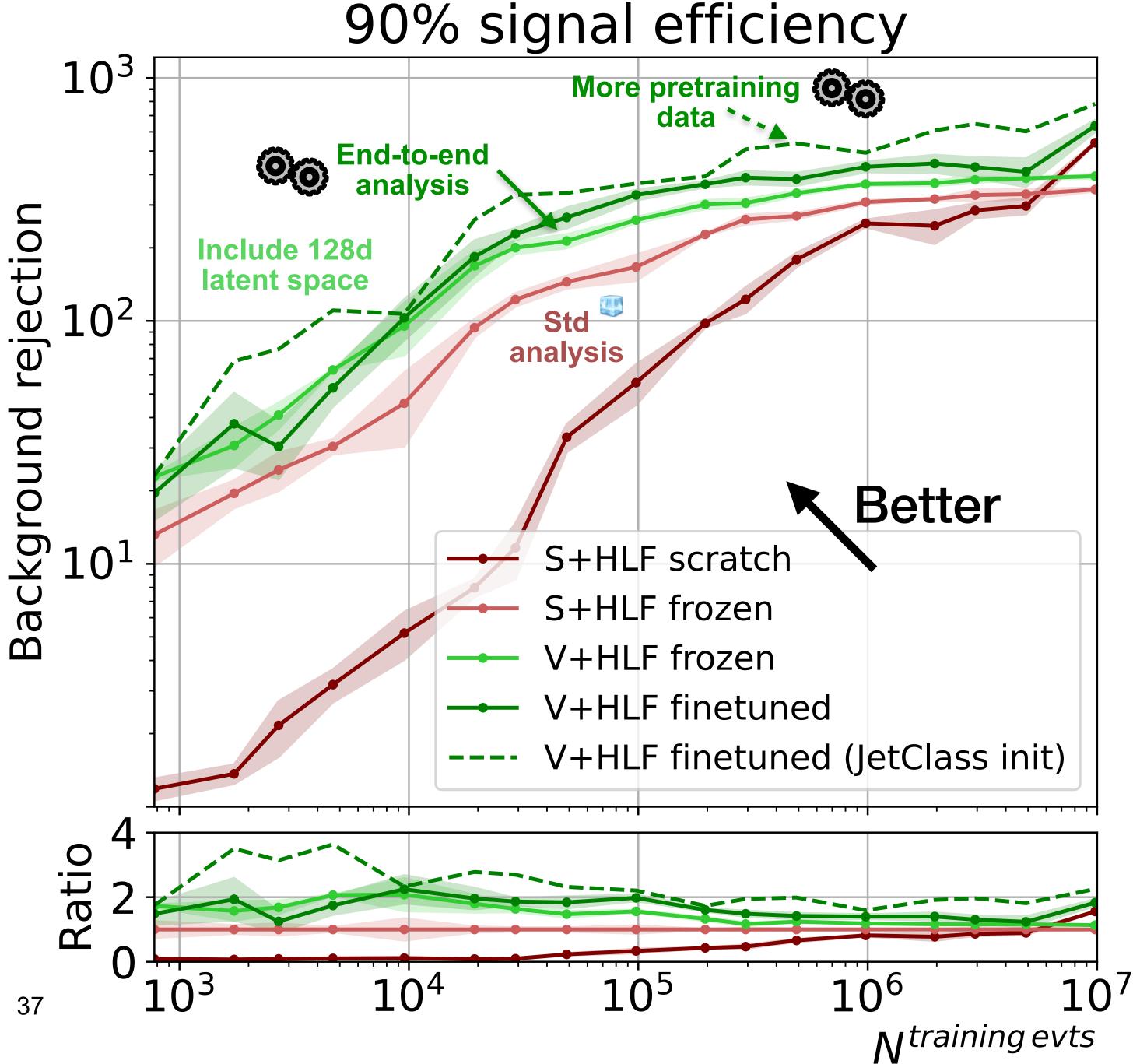
ПШ



Decreases bkg by 2x ...

 S/\sqrt{B} : increases significance by 40%!

ТШП

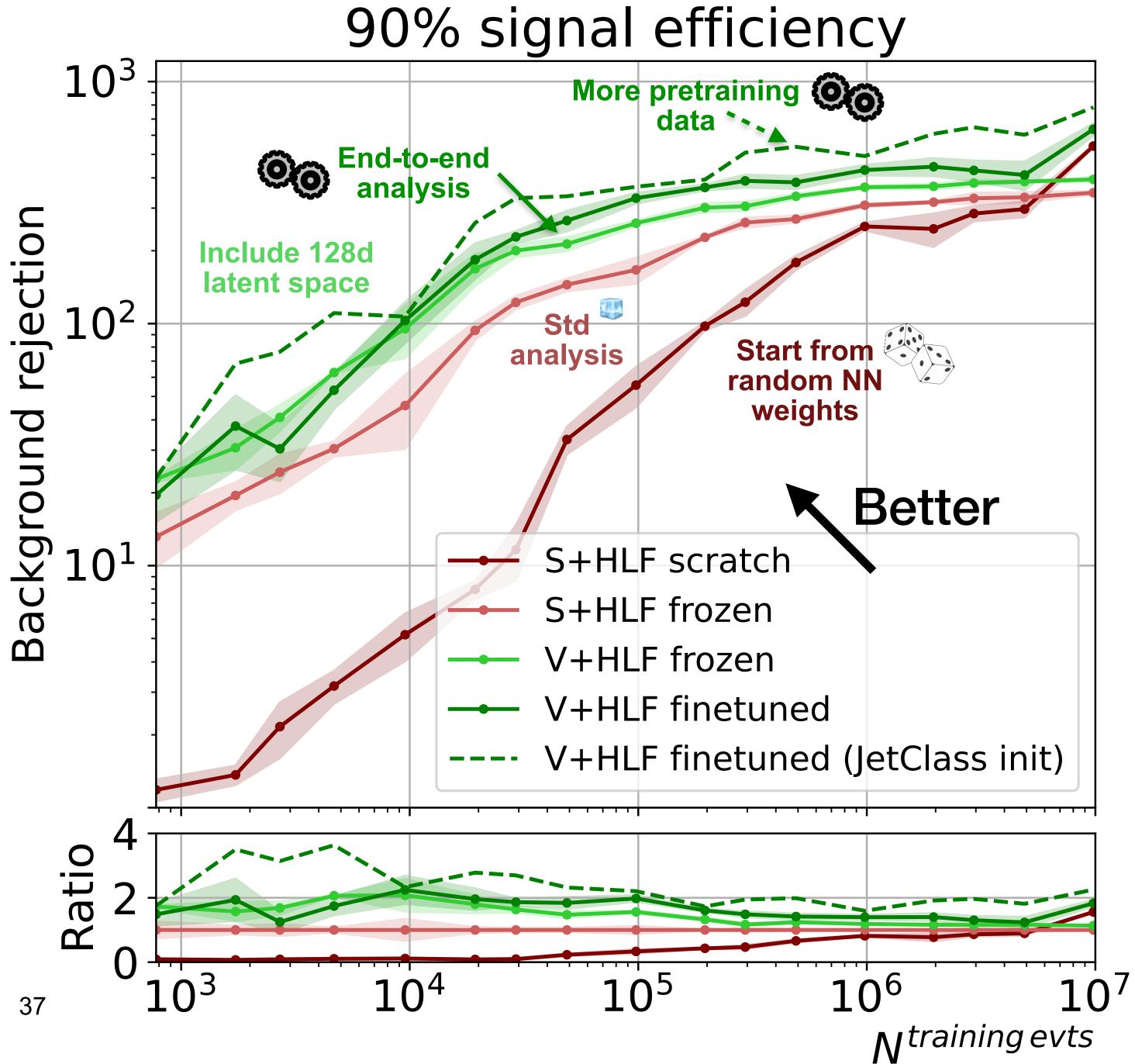


A better Higgs tagger helps analysis performance

Decreases bkg by 2x ...

 S/\sqrt{B} : increases significance by 40%!

ТШП



A better Higgs tagger helps analysis performance

But training from scratch, with enough data, will surpass traditional analyses.

Decreases bkg by 2x ...

 S/\sqrt{B} : increases significance by 40%!

Tasks

How to build this FM: Pretraining

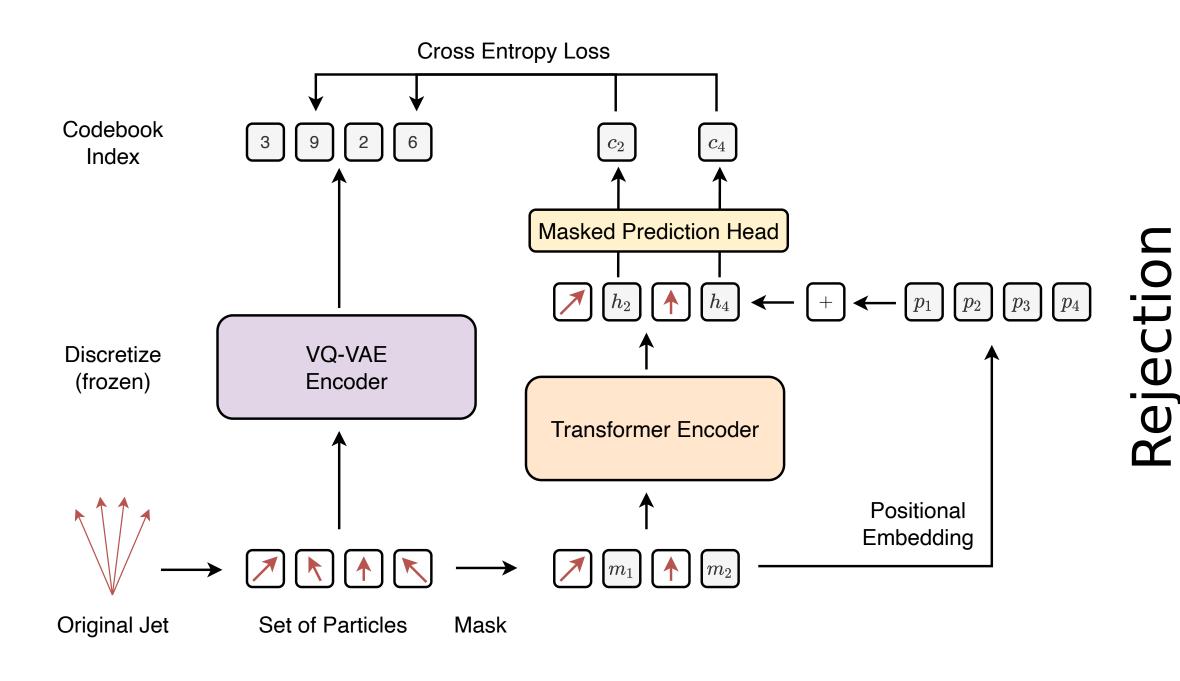
Large Unlabelled Jet Tagging **Dataset** Pileup Charged Mitigation Particle Tracks Calorimeter Track-Calo Fine-tuning Pre-training Clusters Clustering Foundation Calorimeter Model Hits Particle-Flow Reconstruction Muon Tracks **Small Labelled** Event **Dataset** Analysis Multiple modalities Anomaly Detection

More self-supervised training proposals, see M. Kagan & A. Hallin

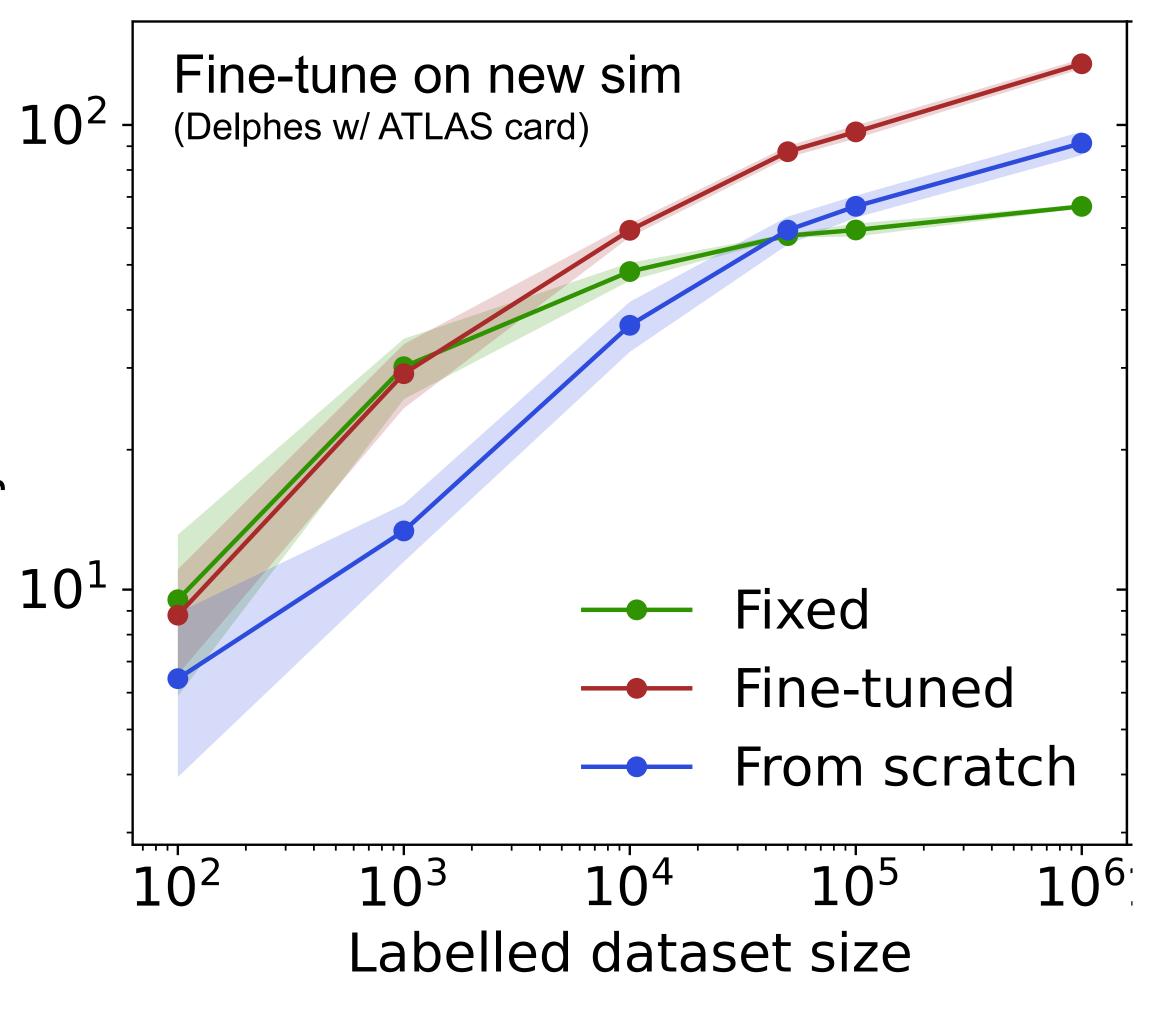
Masked Particle Modelling

Pre-train in an "unsupervised learning" setting, a.k.a, "generative model"

→ Jet completion



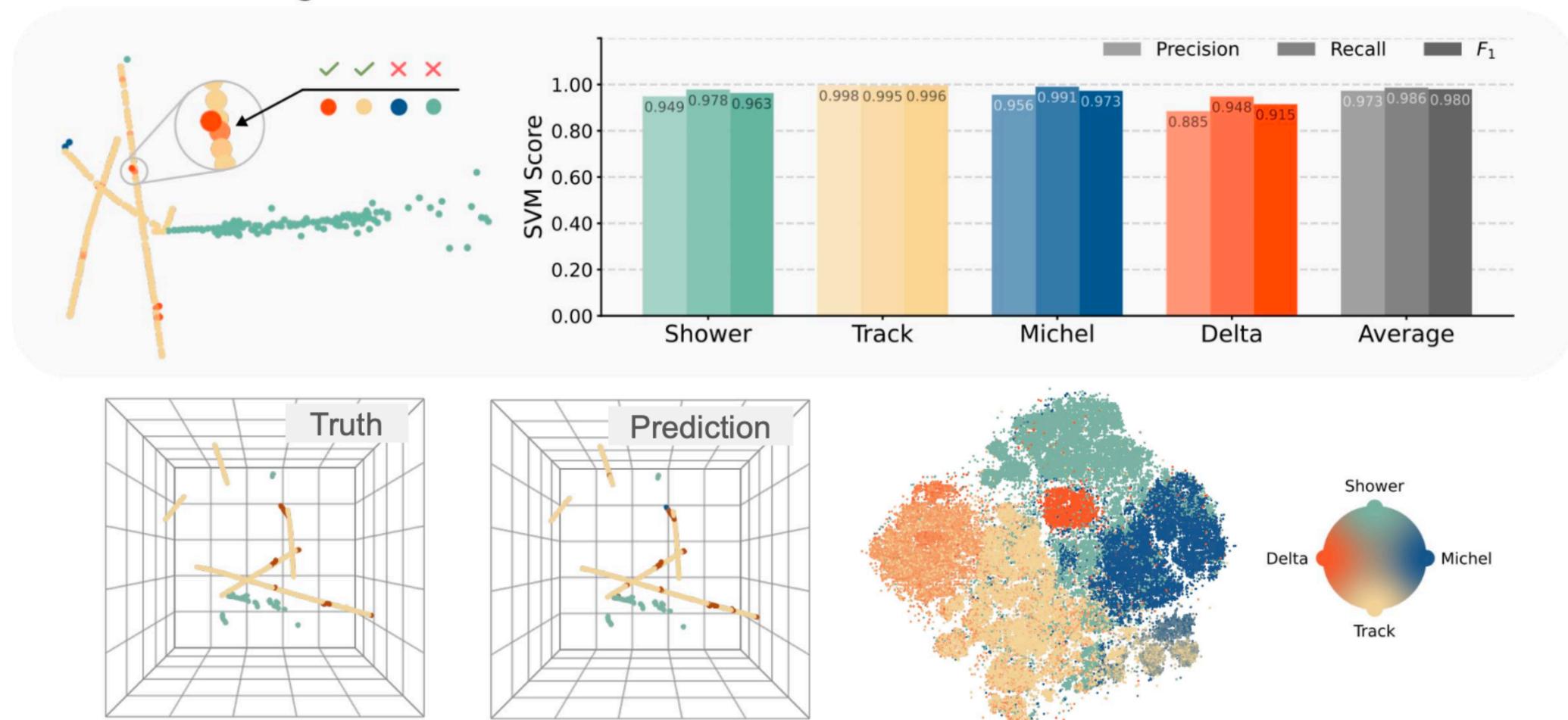
Pretrain Delphes 100m jets.
Ultimately... could be on ATLAS data!



Also in neutrino physics (DUNE)

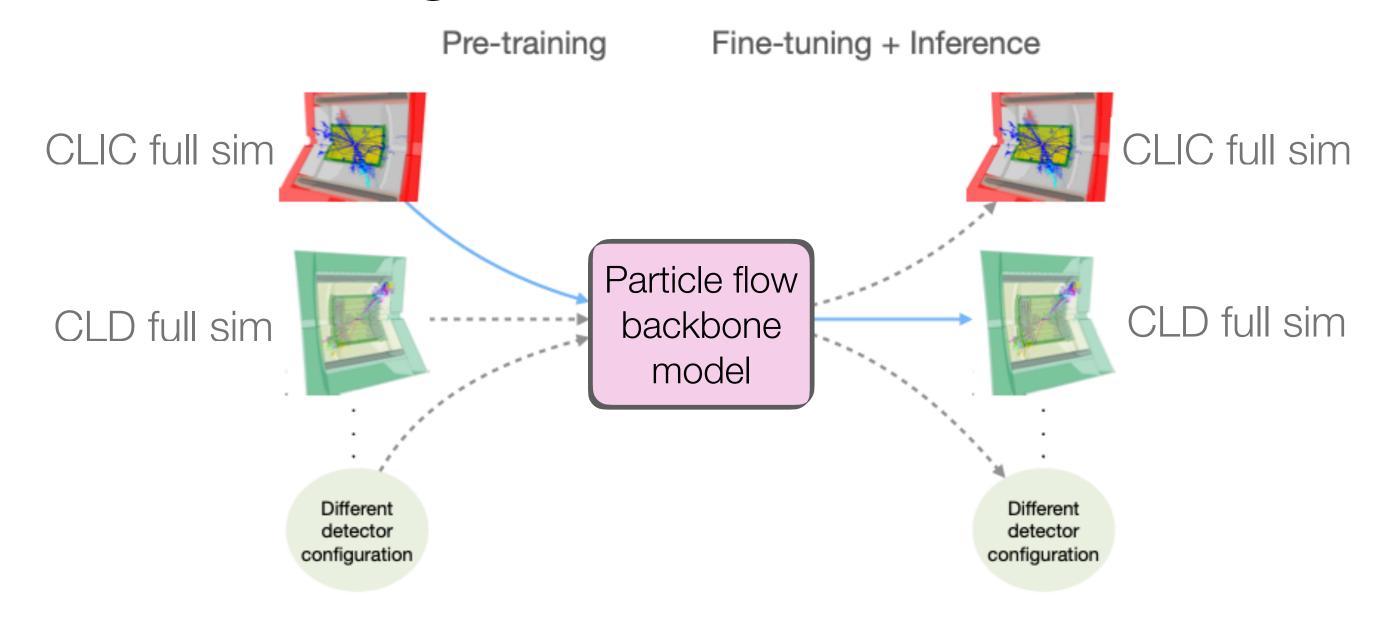
Linear Probing: Semantic Segmentation

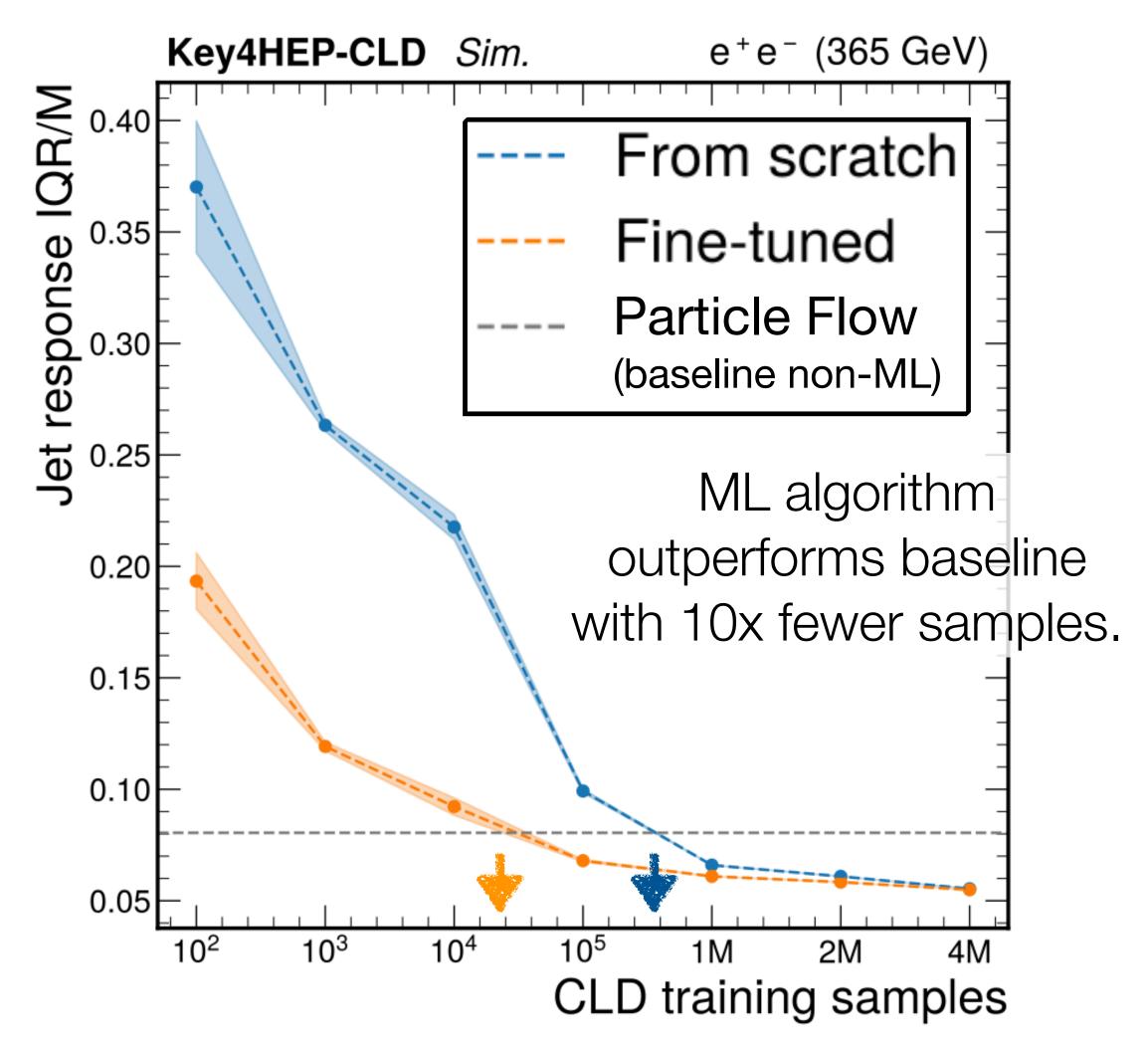
Semantic Segmentation



What does this enable?

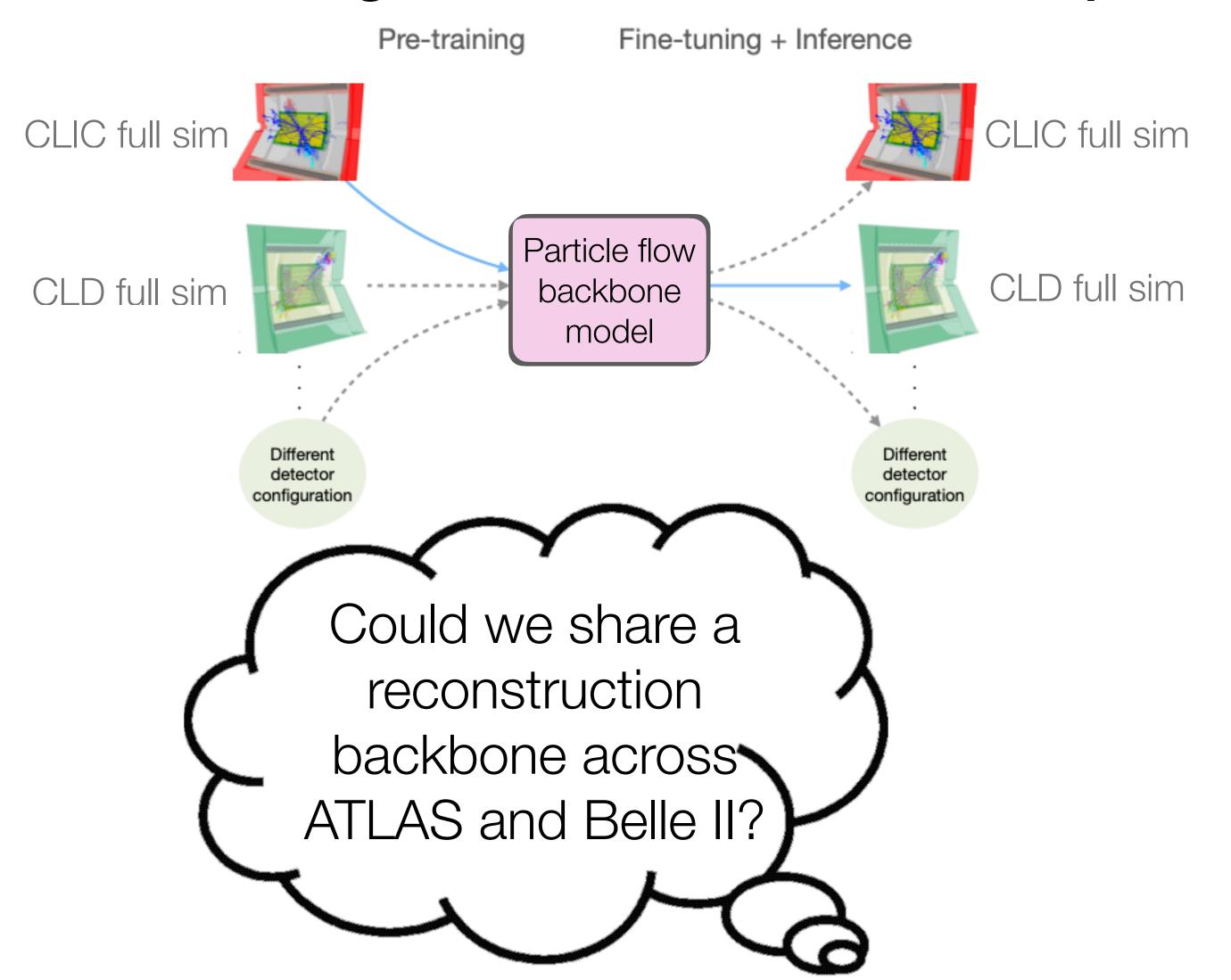
→ Sharing reconstruction across experiments

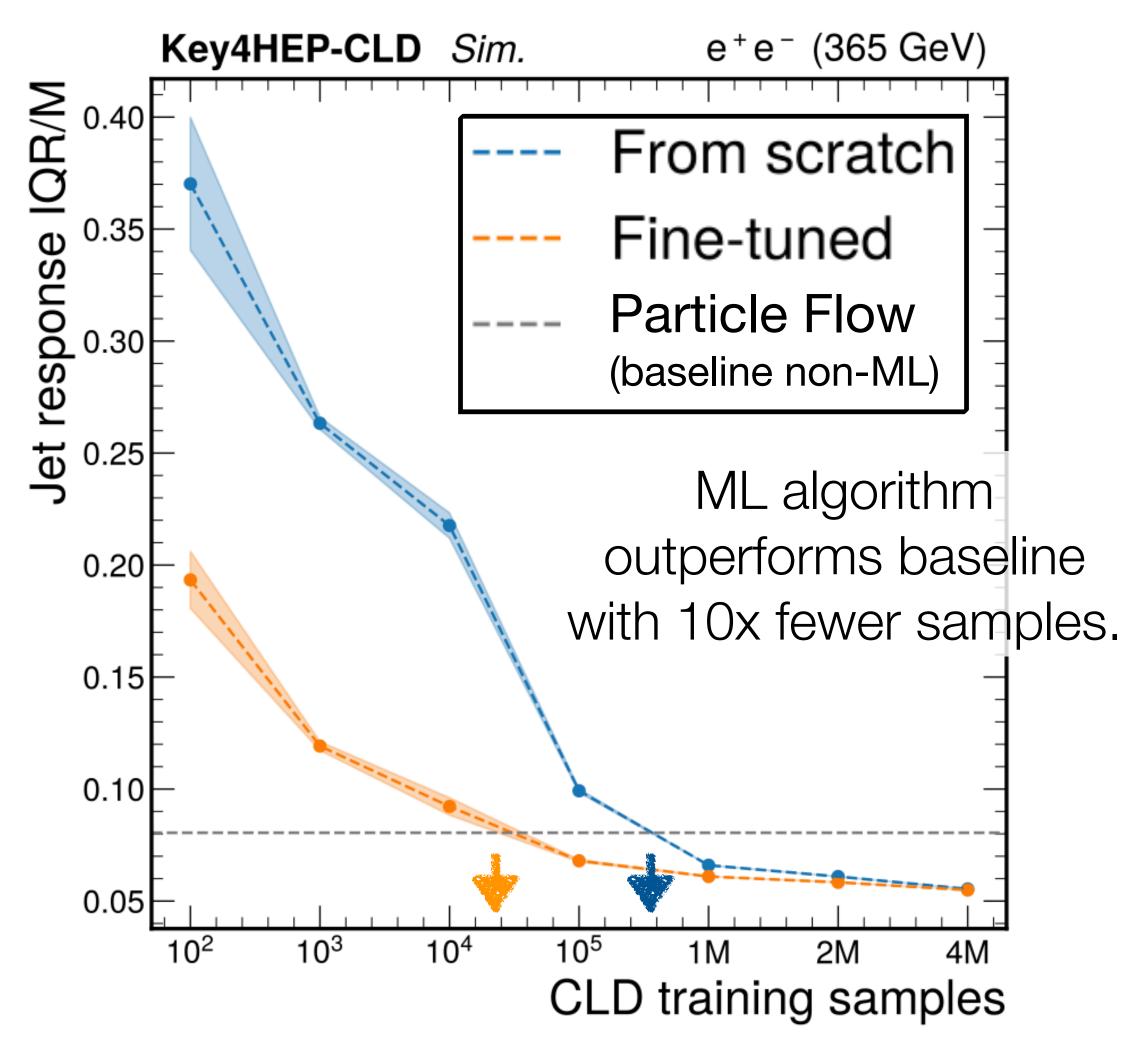




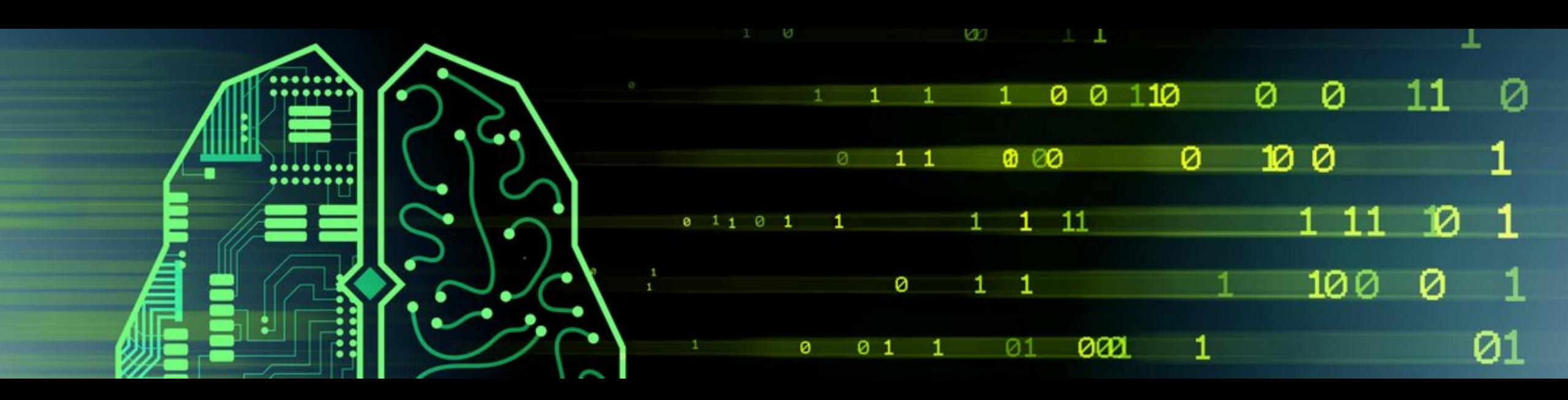
What does this enable?

→ Sharing reconstruction across experiments





The Q: Build big or build smart?



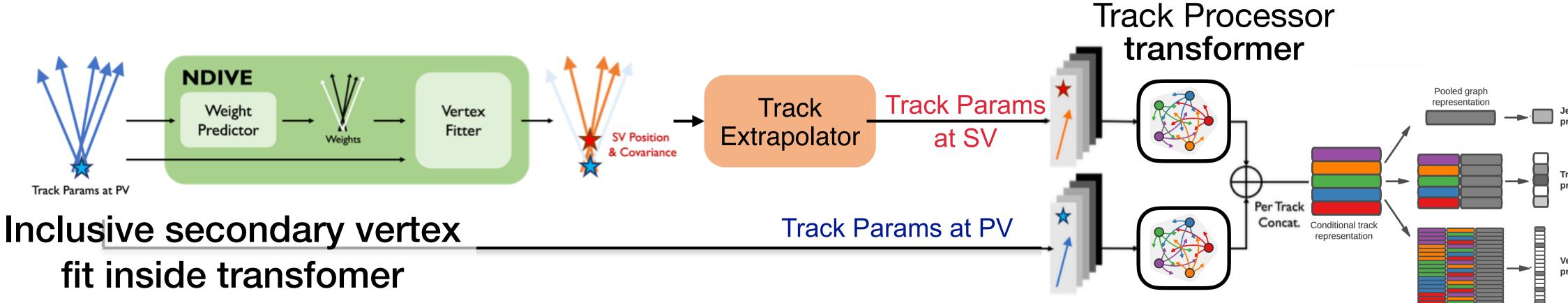
Workshop ongoing in Munich, 25.8 — 19.9

BUILD BIG OR BUILD SMART: EXAMING SCALE AND DOMAIN KNOWLEDGE IN MACHINE LEARNING FOR FUNDAMENTAL PHYSICS

25 August - 19 September 2025
Lukas Heinrich, Michael Kagan, Margarita Osadchy, Tobias Golling, Siddarth Mishra-Sharma

https://www.munich-iapbp.de/activities/activities-2025/machine-learning

Neural Differentiable Vertex Fitter



"Optimization in a loop"

Gradients of fitted vertex to optimize end-to- end trained b-tagger

$$\frac{\mathrm{d}v^*}{\mathrm{d}w} = -\left(\frac{\partial^2\chi^2}{\partial v^2}\right)^{-1}\frac{\partial^2\chi^2}{\partial v\partial w}\bigg|_{v=v^*}$$

$$\frac{\mathsf{Prack}\;\mathsf{I}}{\mathsf{Track}\;\mathsf{I}}$$

$$\frac{\mathsf{q_3,V_3}}{\mathsf{Predicted}\;\mathsf{Secondary}}$$

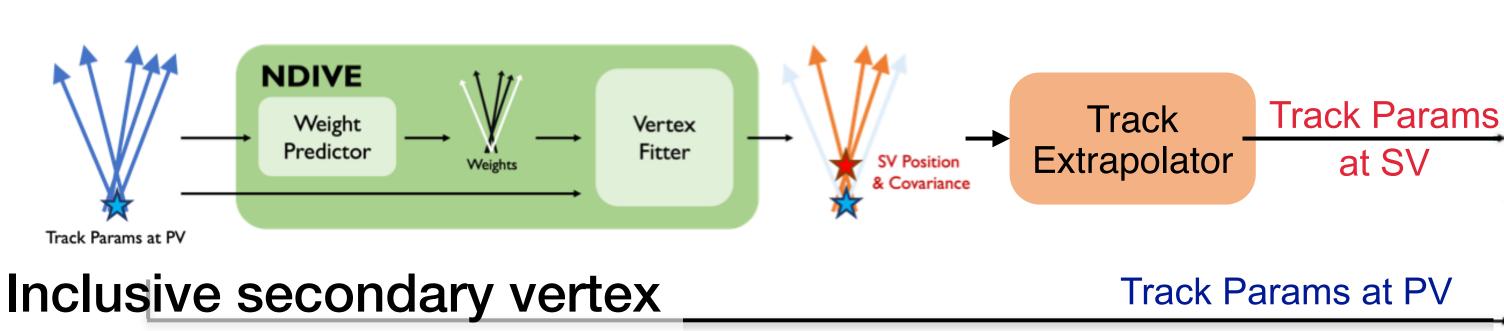
$$\frac{\mathsf{q_2,V_2}}{\mathsf{PV}}$$

$$\mathsf{Preigee\;\mathsf{Parameters\;w.r.t.}}$$

$$\mathsf{Preigee\;\mathsf{Parameters\;w.r.t.}}$$

$$\mathsf{Preigee\;\mathsf{Parameters\;w.r.t.}}$$

Neural Differentiable Vertex Fitter

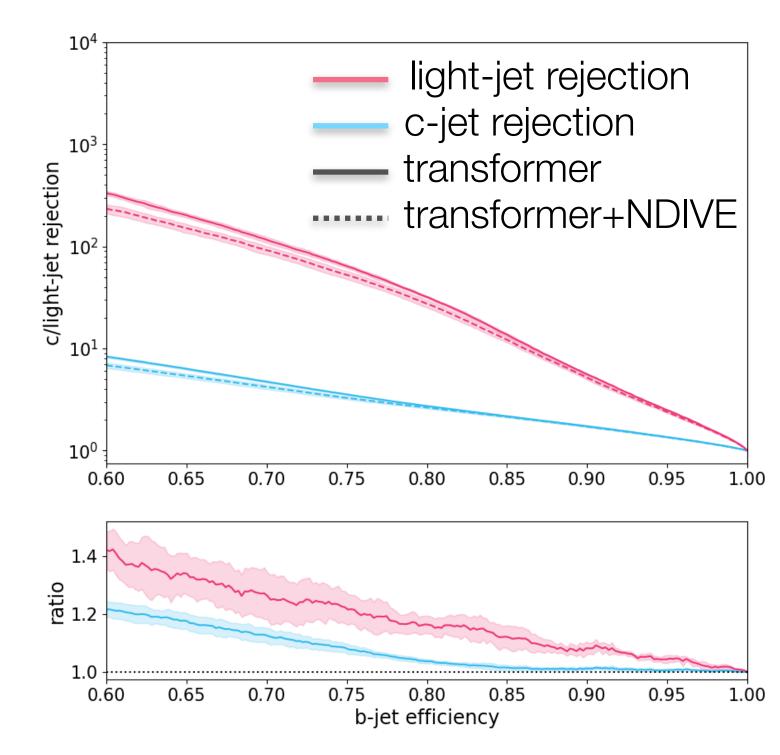


fit inside transfomer

"Optimization in a loop"

Gradients of fitted vertex to optimize end-to- end trained b-tagger

$$\frac{\mathrm{d}v^*}{\mathrm{d}w} = -\left(\frac{\partial^2\chi^2}{\partial v^2}\right)^{-1}\frac{\partial^2\chi^2}{\partial v\partial w}\bigg|_{v=v^*} \\ \text{Track I} \\ \text{Track 2} \\ \text{Predicted Secondary} \\ \text{Vertex} \\ \text{Yertex} \\ \text{Prerigee Parameters w.r.t.} \\ \text{the Primary Vertex} \\ \text{Track 3} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 3} \\ \text{Predicted Secondary} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 5} \\ \text{Track 6} \\ \text{Track 6} \\ \text{Track 7} \\ \text{Predicted Secondary} \\ \text{Track 6} \\ \text{Track 7} \\ \text{Track 7} \\ \text{Track 8} \\ \text{Track 9} \\ \text{Track 9} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 2} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 2} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 6} \\ \text{Track 6} \\ \text{Track 6} \\ \text{Track 9} \\ \text{Trac$$



← 50% improvement in bkg rejection

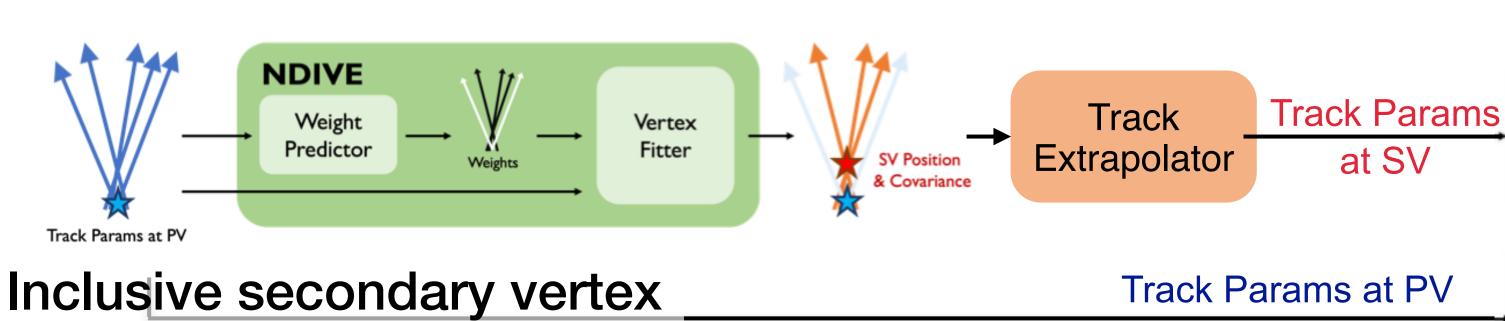
Concat. Conditional track

Track Processor

transformer

(Up to 10x improvement in jet tagging from vtx fit if perfect track selection)

Neural Differentiable Vertex Fitter

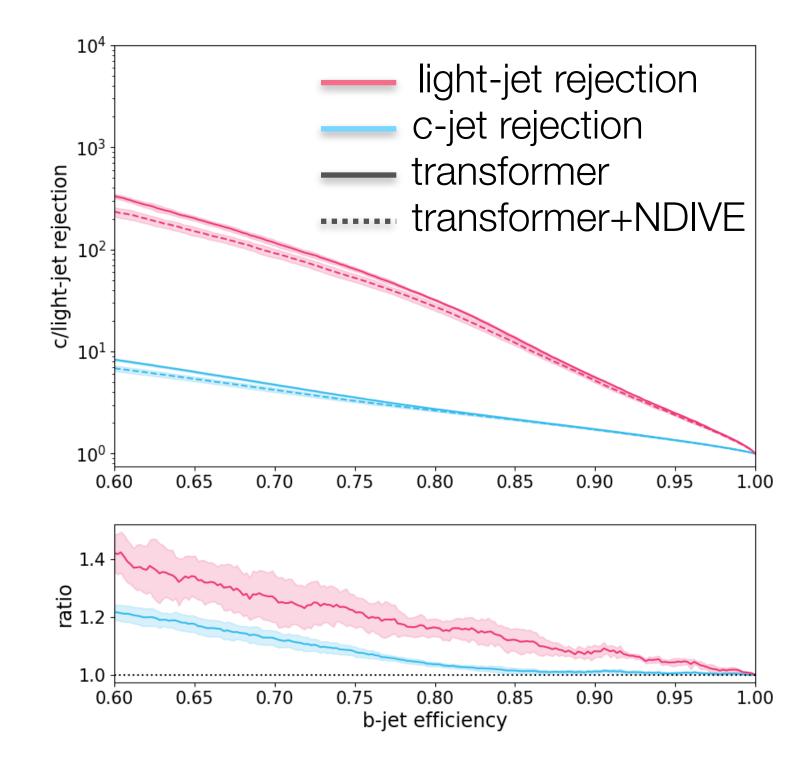


fit inside transfomer

"Optimization in a loop"

Gradients of fitted vertex to optimize end-to- end trained b-tagger

$$\frac{\mathrm{d}v^*}{\mathrm{d}w} = -\left(\frac{\partial^2\chi^2}{\partial v^2}\right)^{-1} \frac{\partial^2\chi^2}{\partial v\partial w} \bigg|_{v=v^*} \\ \text{Track I} \\ \text{Track 2} \\ \text{Track 2} \\ \text{Predicted Secondary} \\ \text{Vertex} \\ \text{Yerigee Parameters w.r.t.} \\ \text{the Primary Vertex} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Predicted Secondary} \\ \text{Yertex} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 5} \\ \text{Predicted Secondary} \\ \text{Yertex} \\ \text{Track 6} \\ \text{Track 7} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 6} \\ \text{Track 7} \\ \text{Track 6} \\ \text{Track 7} \\ \text{Track 8} \\ \text{Track 9} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 1} \\ \text{Track 2} \\ \text{Track 2} \\ \text{Track 3} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 2} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 4} \\ \text{Track 6} \\ \text{Trac$$



← 50% improvement in bkg rejection

Concat. Conditional track

Track Processor

transformer

(Up to 10x improvement in jet tagging from vtx fit if perfect track selection)

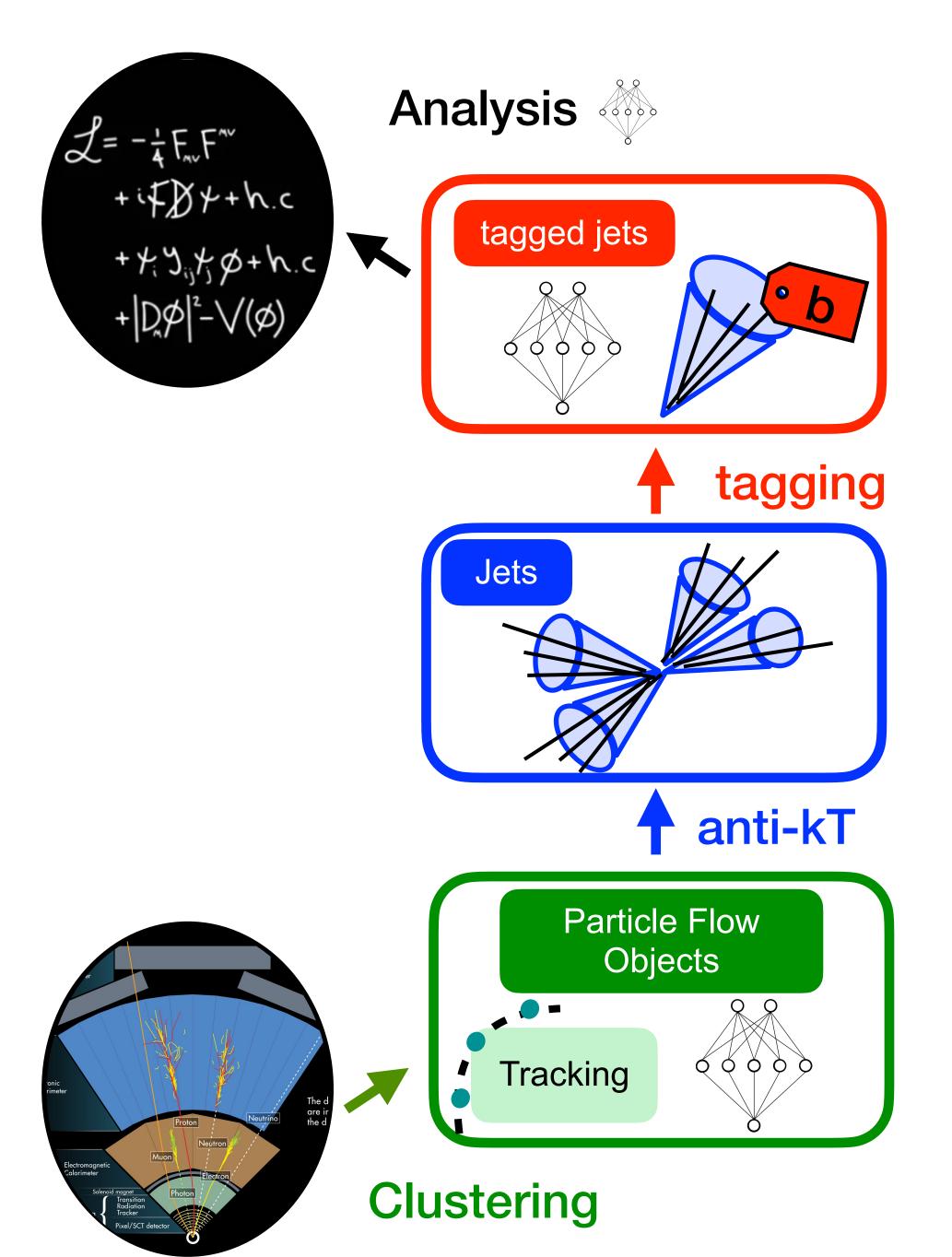
Interesting for Belle 2? 1808.10567

6.2.1. Vertex finding algorithms. The Belle II experiment has deployed three implementations of a vertex fit: KFitter, developed for the Belle experiment, RAVE [59], a standalone package originating from the CMS vertex fitting libraries, and TreeFitter [71], initially conceived by the BaBar collaboration. We use both KFitter and RAVE for kinematic fits and RAVE for geometric fits. TreeFitter is used for the fitting entire decay chains.

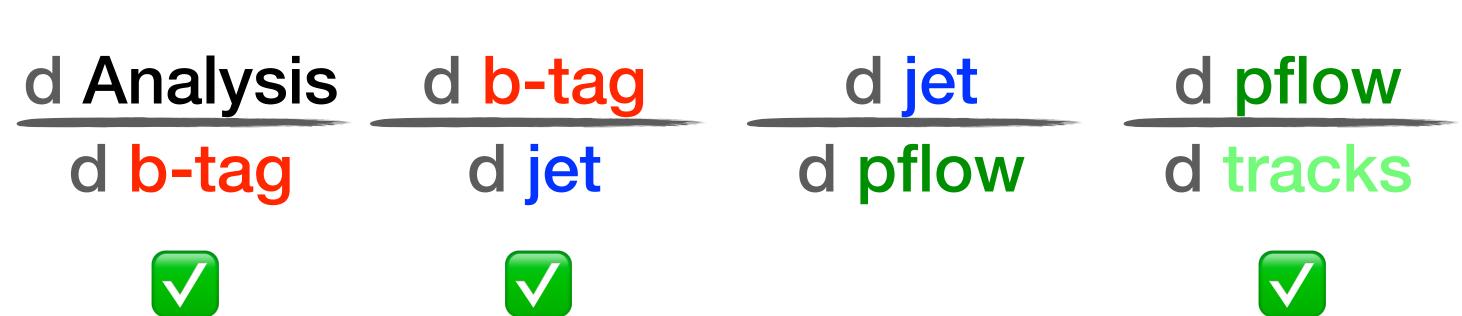
Kinematic fitts. Kinematic fitting uses the known properties of a specific decay chain to improve the measurements of the process. Lagrangian multipliers are used in order to impose the kinematic constraints to the fit. Given the measurements, $\mathbf{q} = (q_1, ..., q_n)$, with a covariance matrix, V, and kinematic constraints, $\mathbf{h}(\mathbf{q})$, the function to be minimised in terms of the most suitable vertex is:

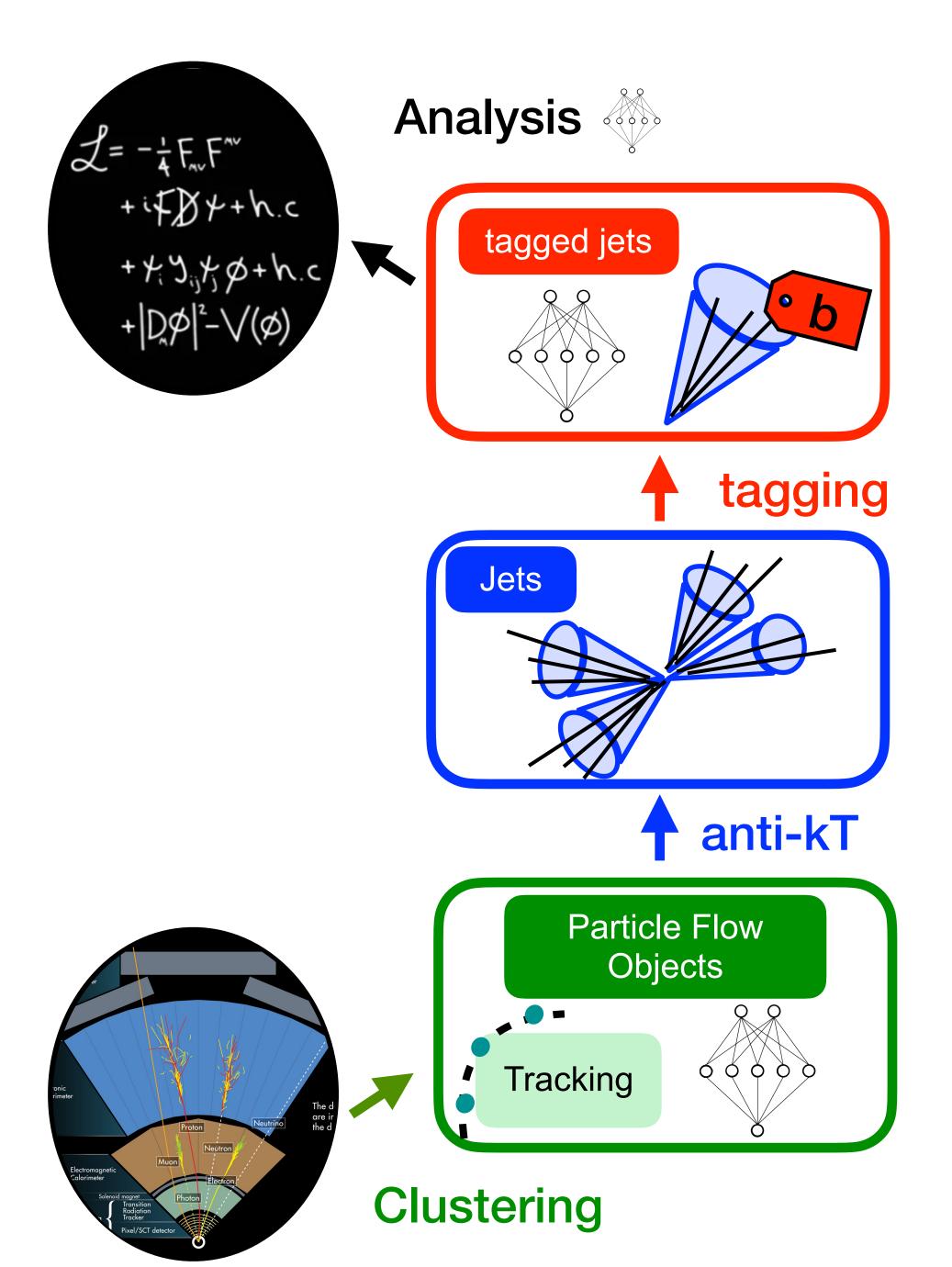
$$\chi^{2} = (\mathbf{q} - \bar{\mathbf{q}})^{T} V^{-1} (\mathbf{q} - \bar{\mathbf{q}}) + 2\lambda^{T} (\mathbf{D}\delta \mathbf{y} + \mathbf{h}(boldsymbolq))$$
(15)

where λ is the Lagrange multiplier, $h(\bar{\mathbf{q}}) = 0$ and $\mathbf{D} = \partial \mathbf{h}/\partial \mathbf{y}$. Here $\bar{\mathbf{q}}$ represents the improved measurements, \mathbf{d} is the kinematic constraint at the starting value and $\delta \mathbf{y}$ is the difference between the improved measurement and the starting value.

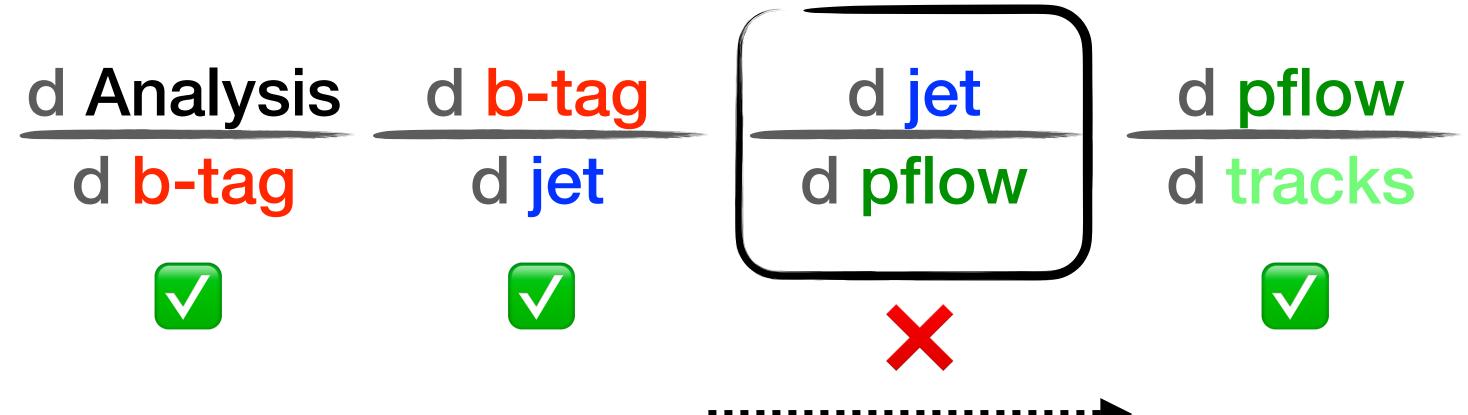


How to use gradients MORE?

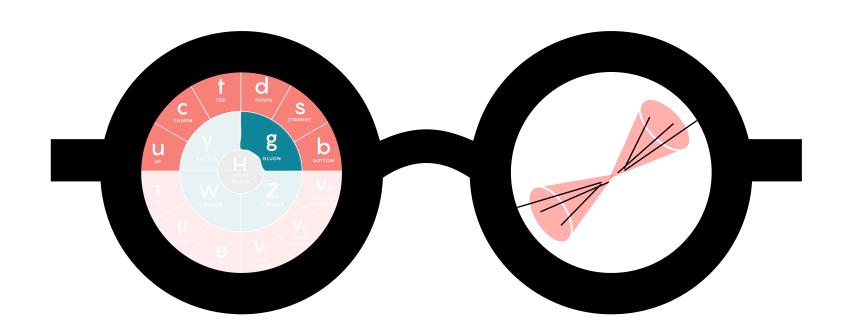




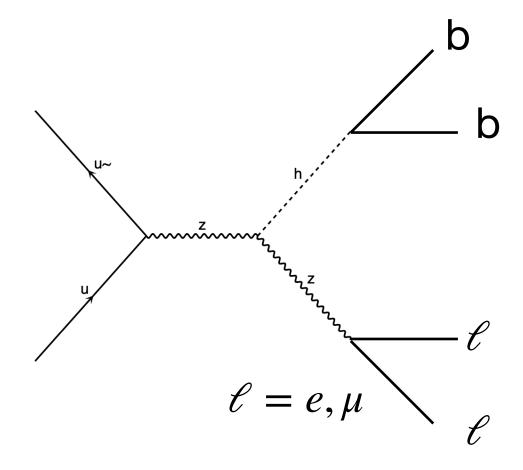
How to use gradients MORE?



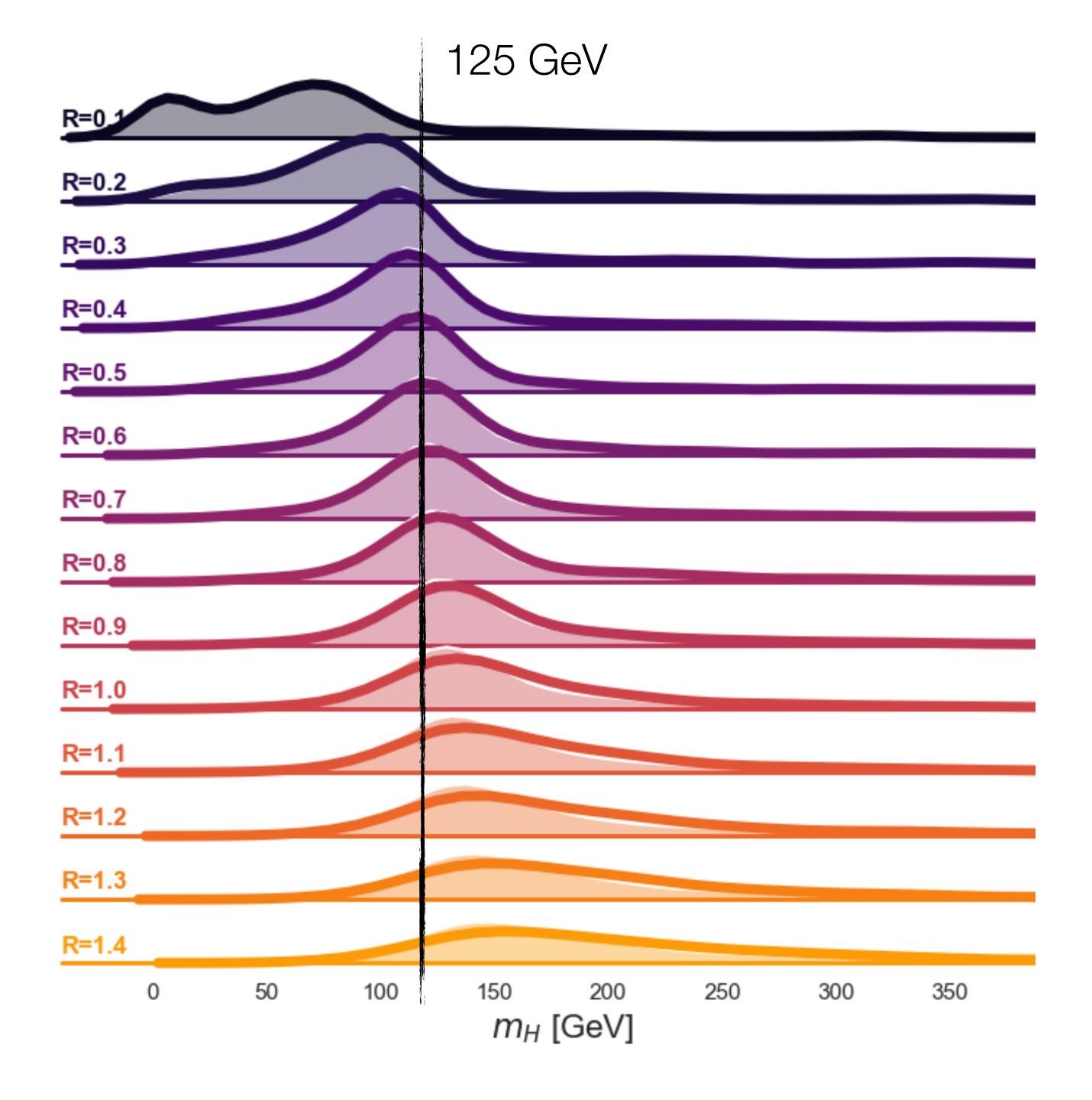
Clustering *not* differentiable... how to fully exploint the richness of our high dim data (tracks / hits).



How to see a Higgs: wearing the right glasses

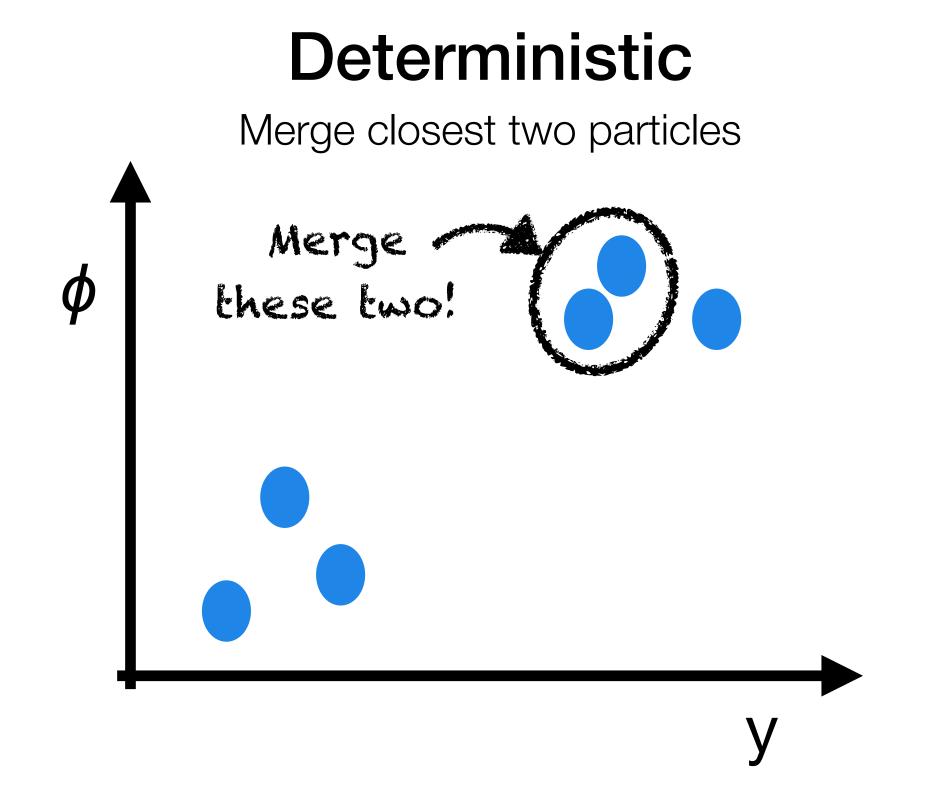


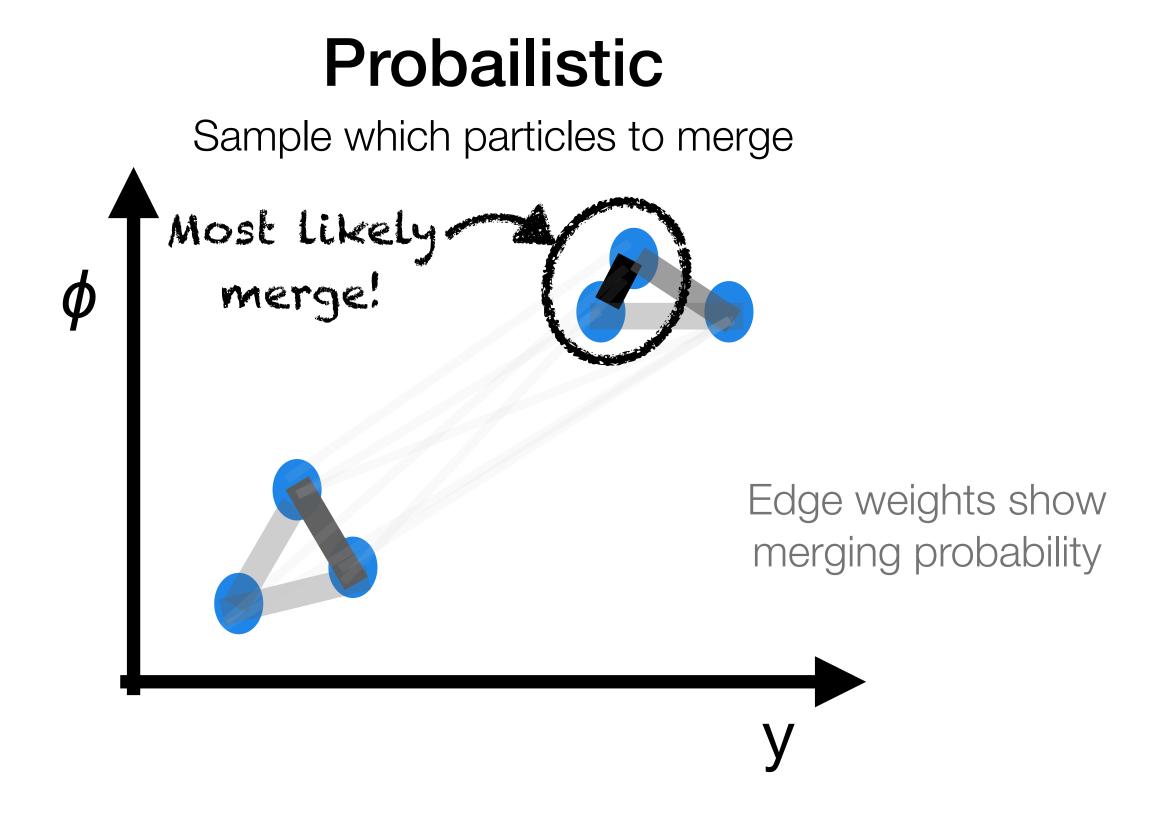
Madgraph + pythia, pp collisions @ 14 TeV



Step 1: Interpret clustering decision probabilistically

Qanti-kT, M. Schwartz + collab: <u>1201.1914</u>, <u>1304.2394</u>





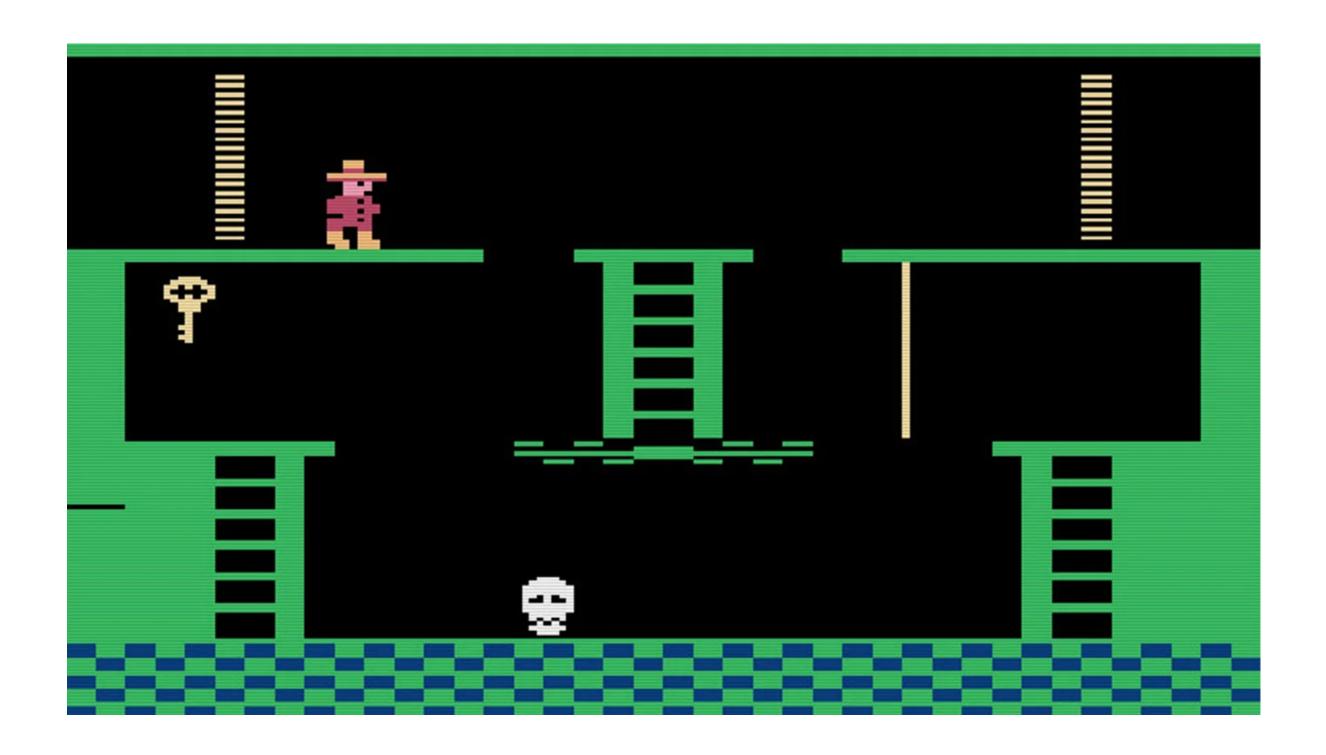
Step 1: Interpret clustering decision probabilistically

Step 2: Gradient with score based estimate

As in <u>2308.16680</u>

$$\nabla_{\theta} \mathbb{E}_{x \sim p(\theta)}[f(x)]$$

$$= \mathbb{E}_{x \sim p(\theta)}[f(x) \nabla_{\theta} \log p(\theta)]$$

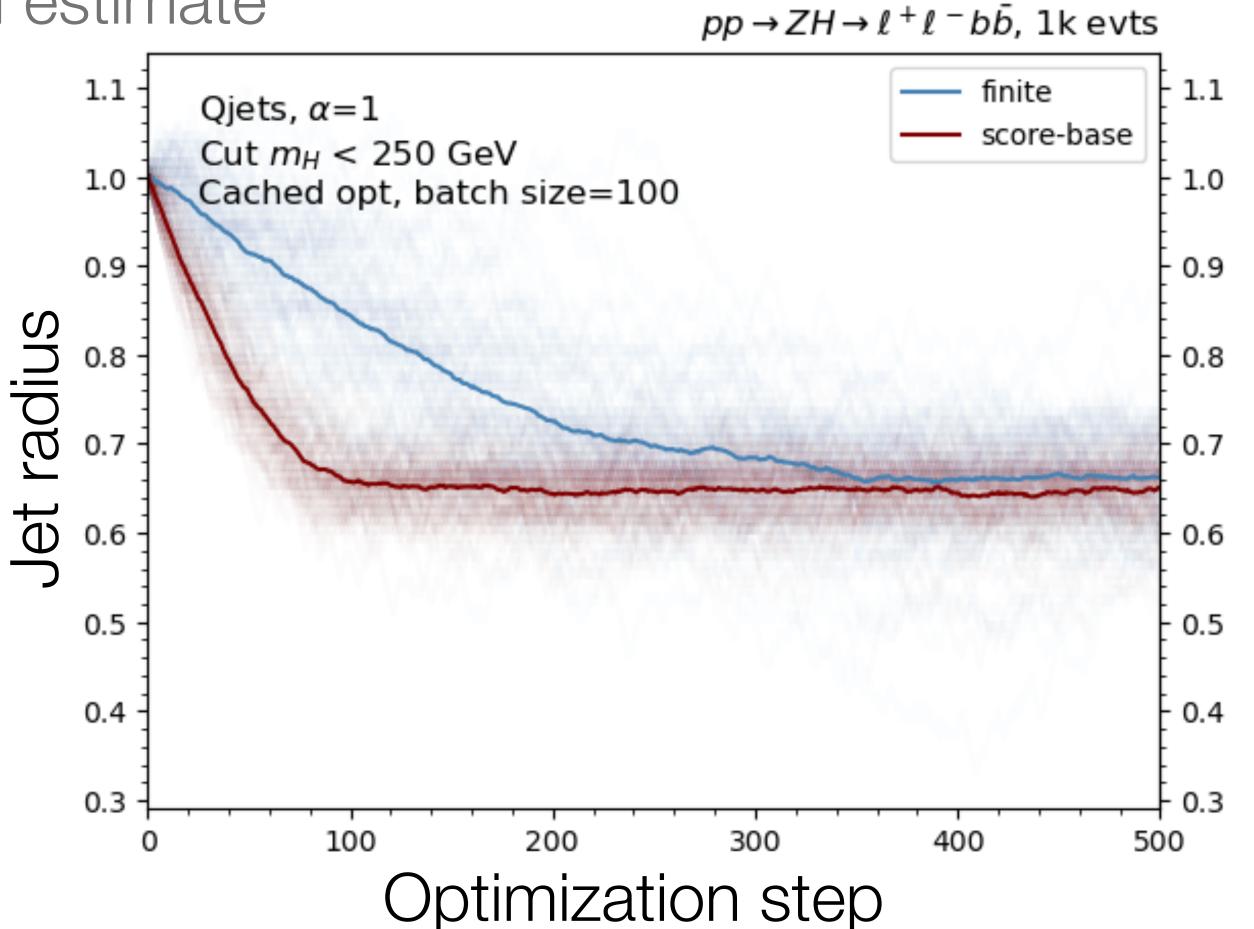


Step 1: Interpret clustering decision probabilistically

Step 2: Gradient with score based estimate

DIE IN

Step 3: Gradient based opt for clustering radius parameter

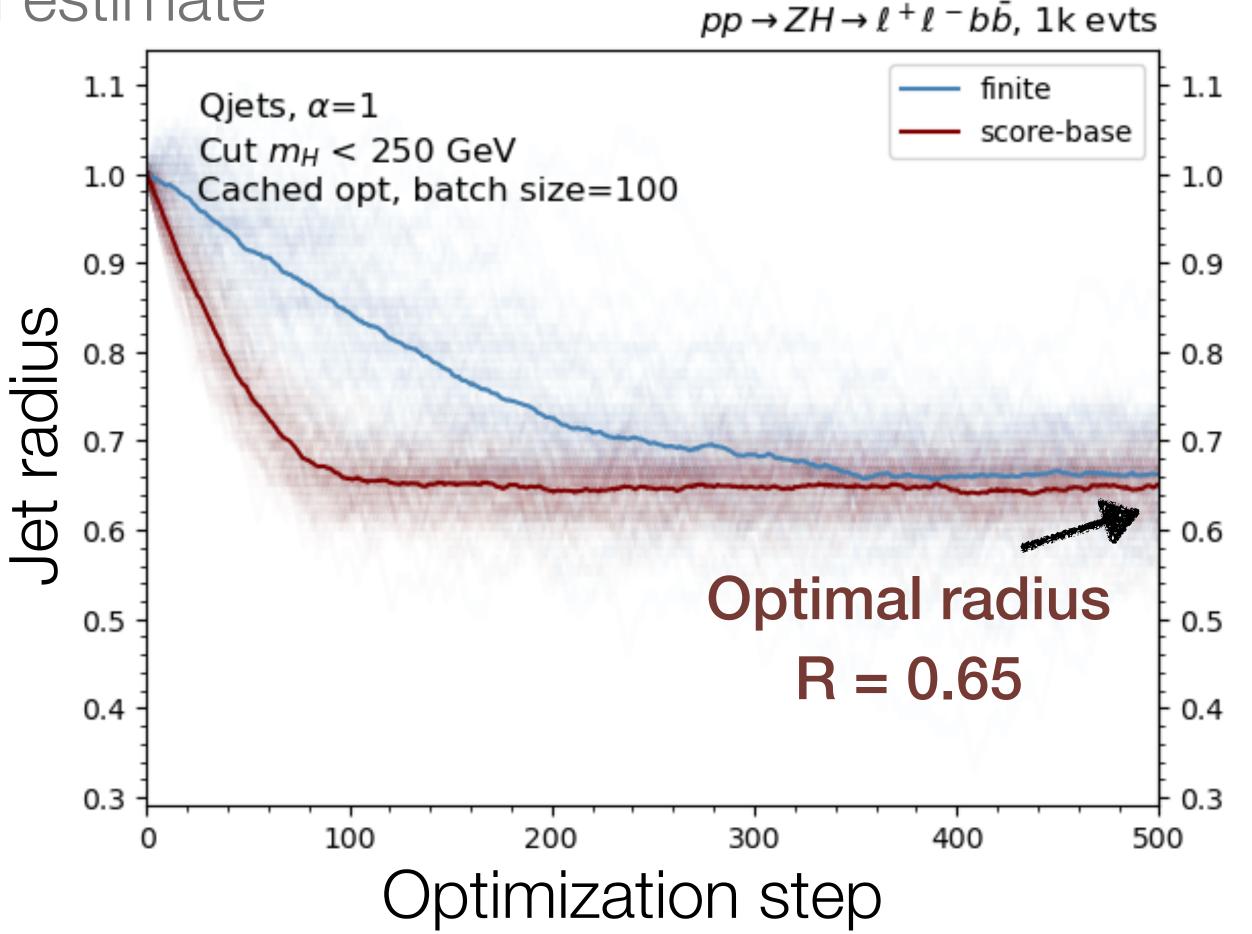


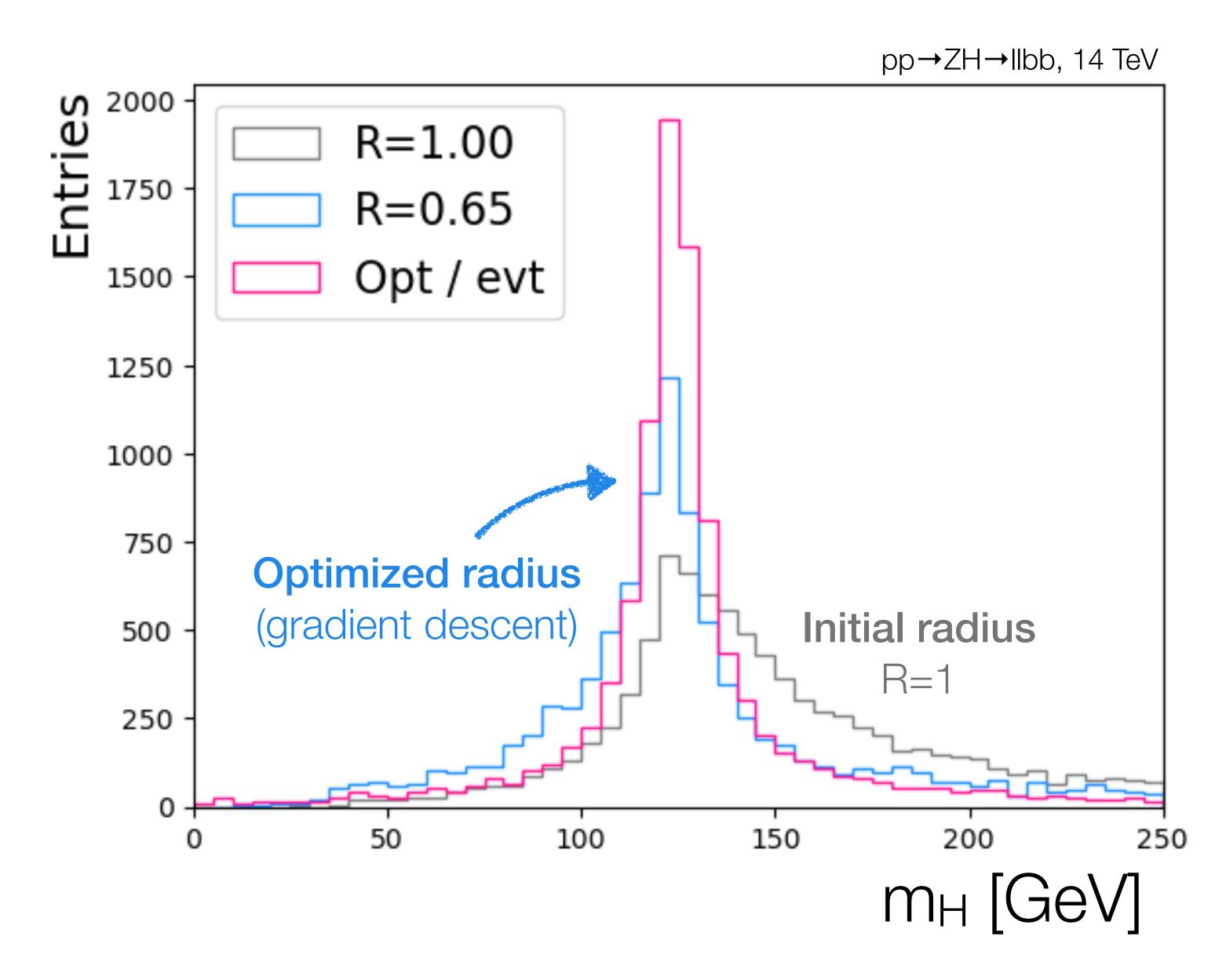
Step 1: Interpret clustering decision probabilistically

Step 2: Gradient with score based estimate

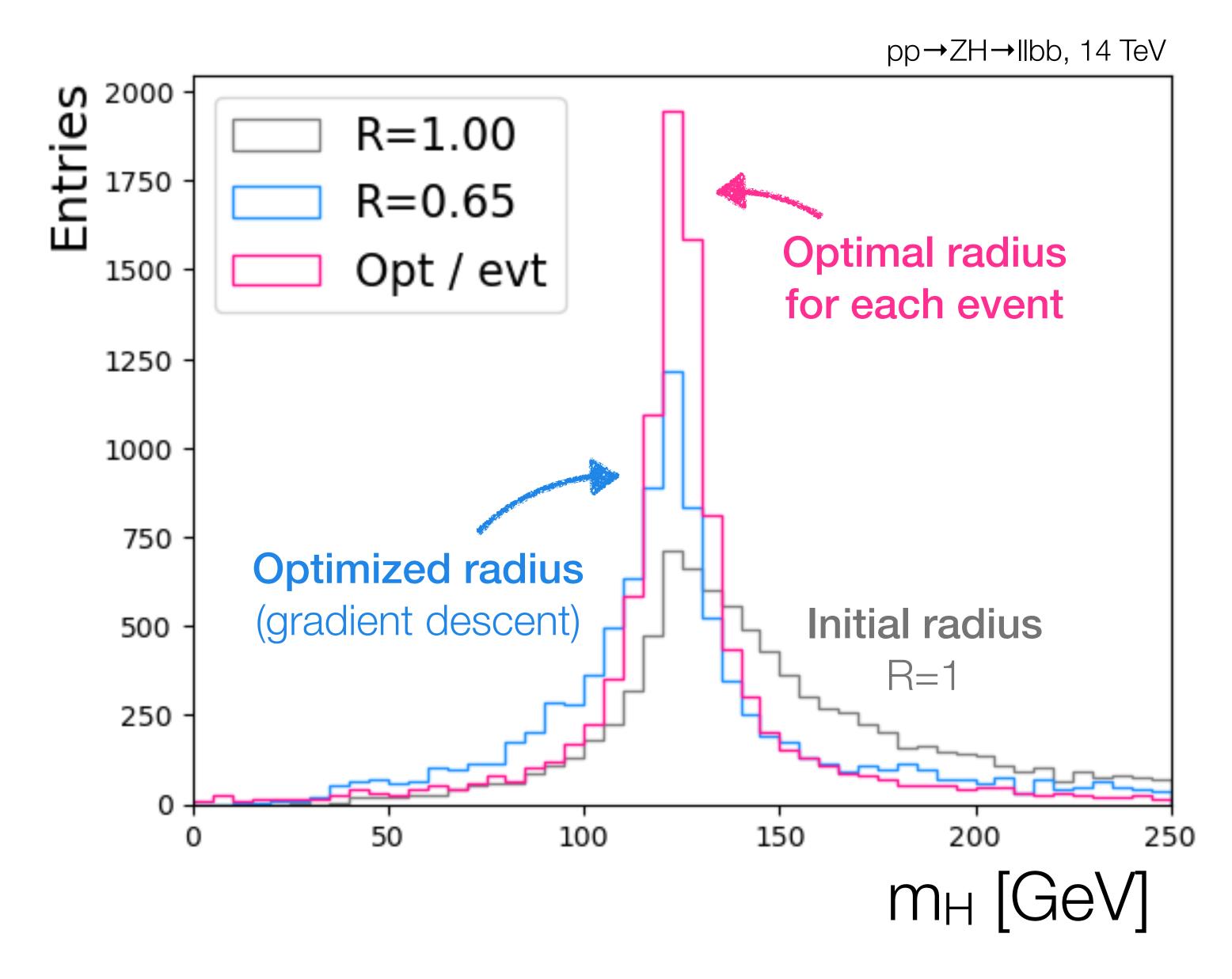
BIBILIO

Step 3: Gradient based opt for clustering radius parameter

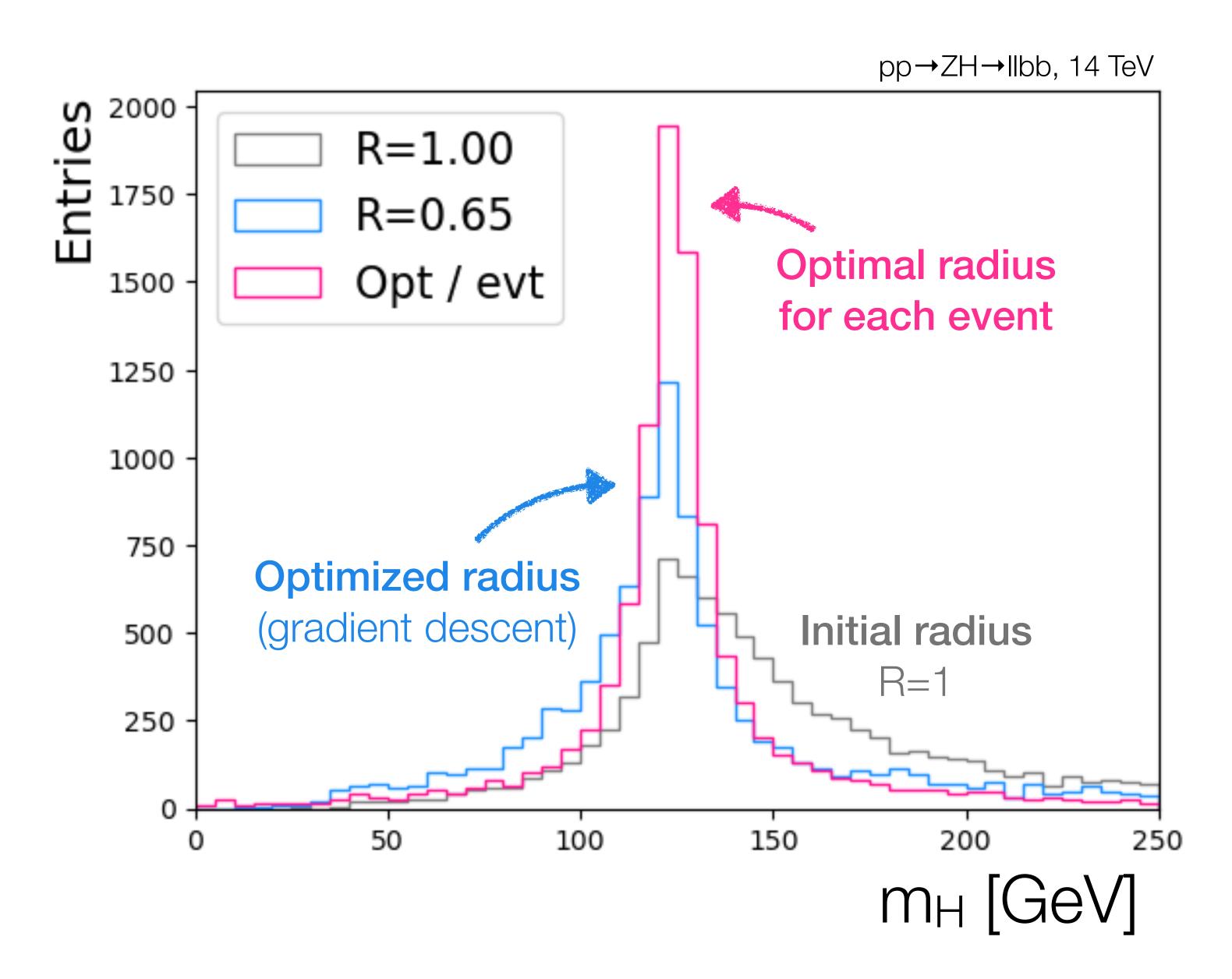


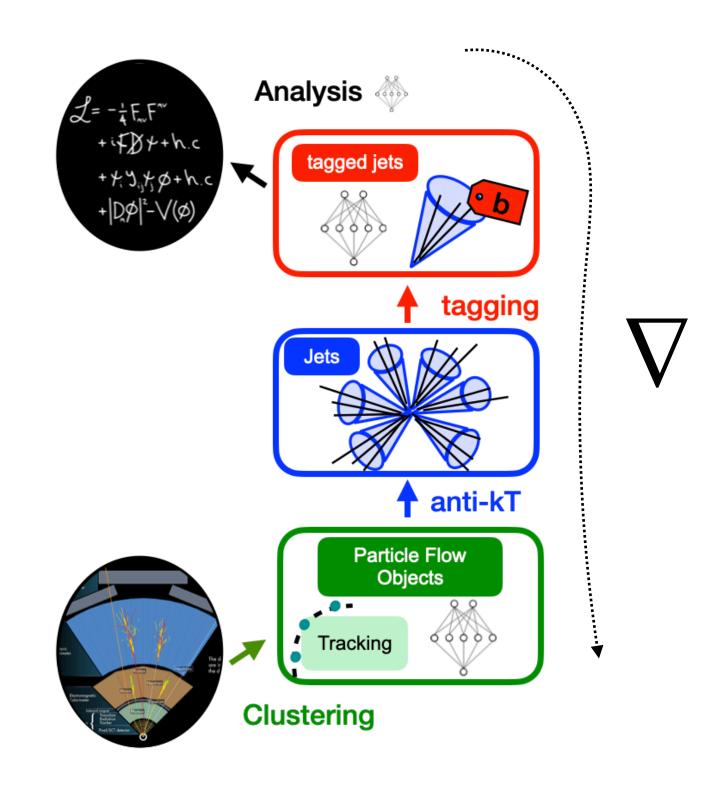


With A. Kofler, M. Kagan and L. Heinrich, in prep



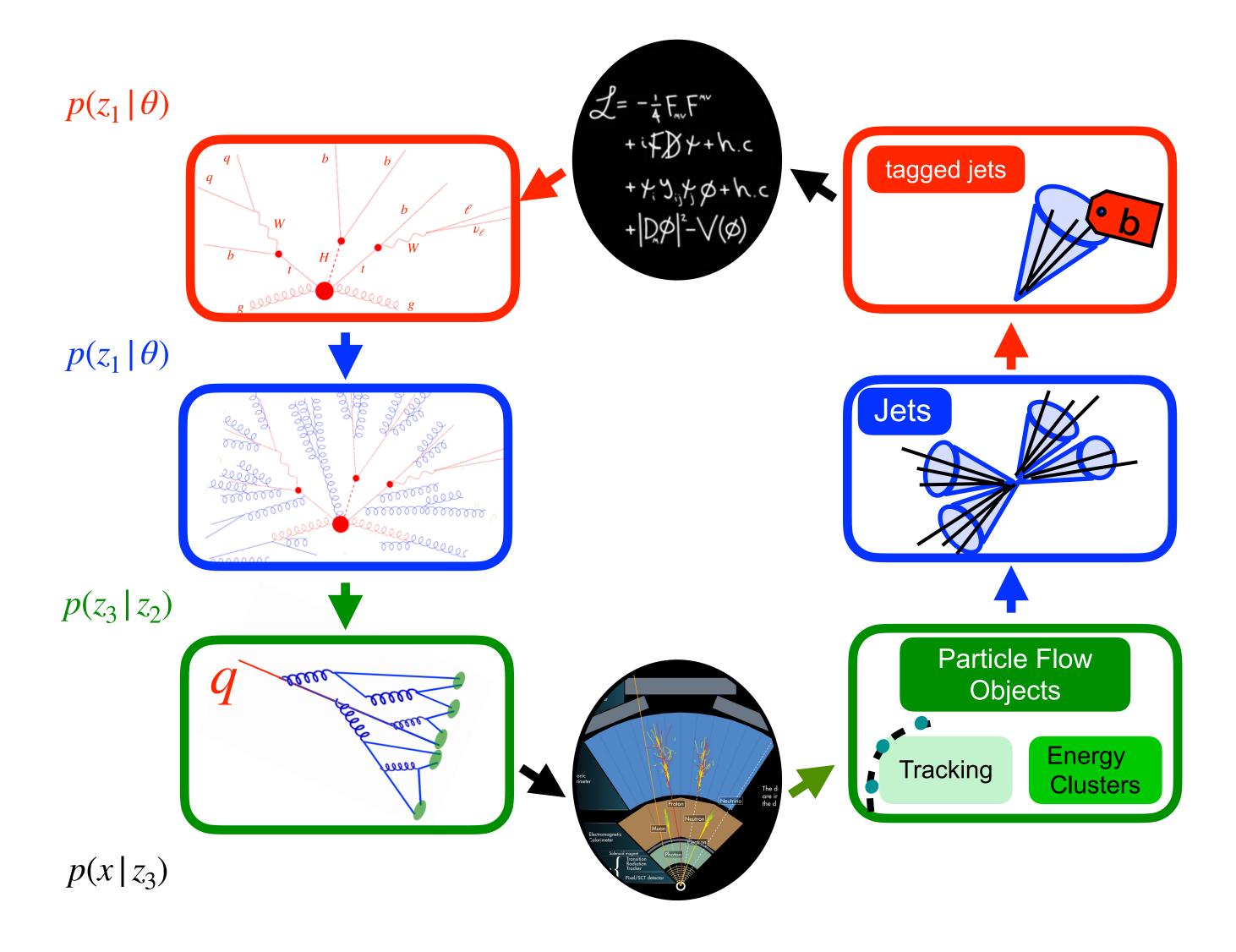
With A. Kofler, M. Kagan and L. Heinrich, in prep

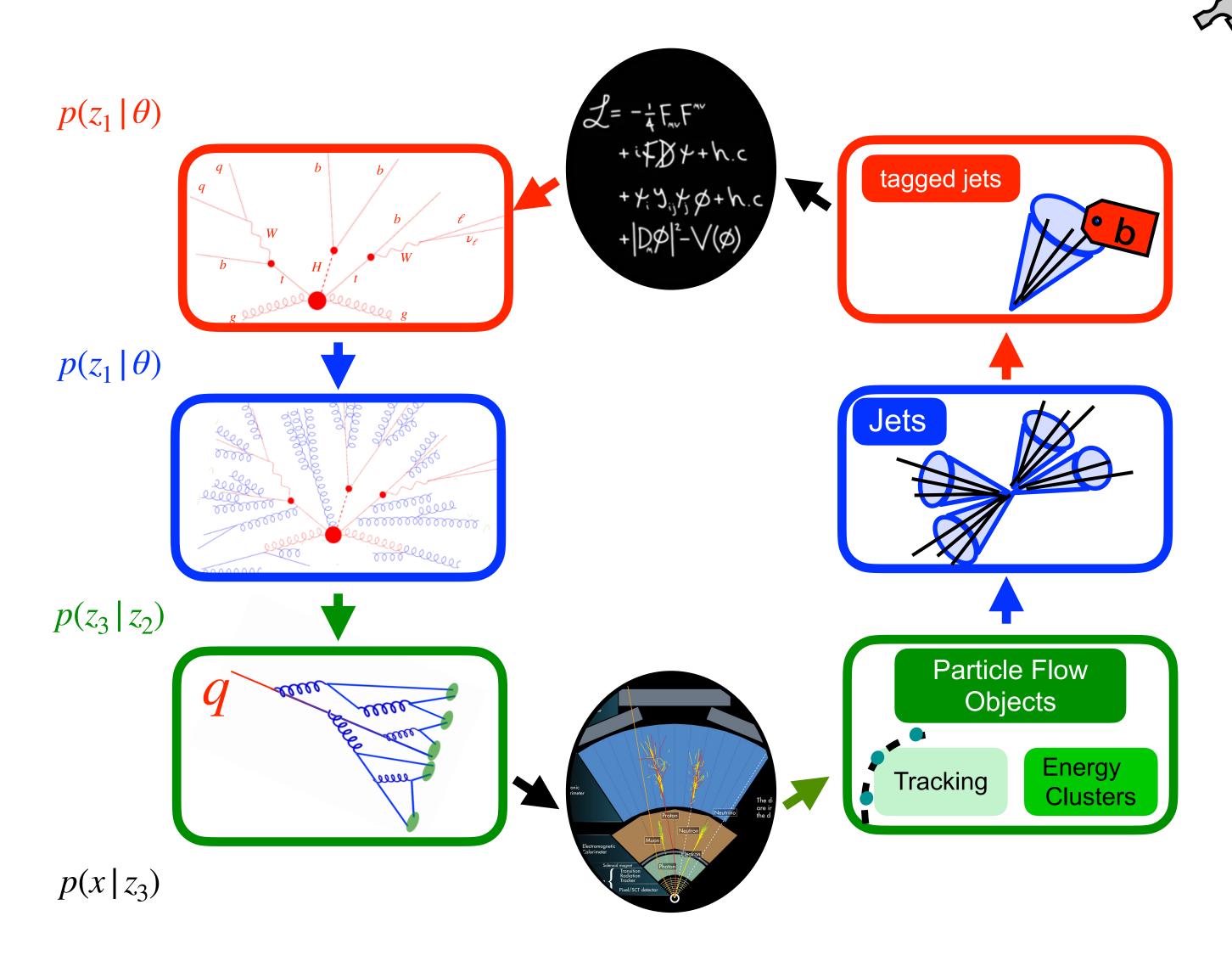


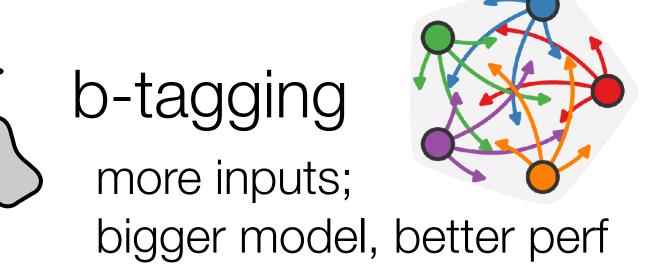


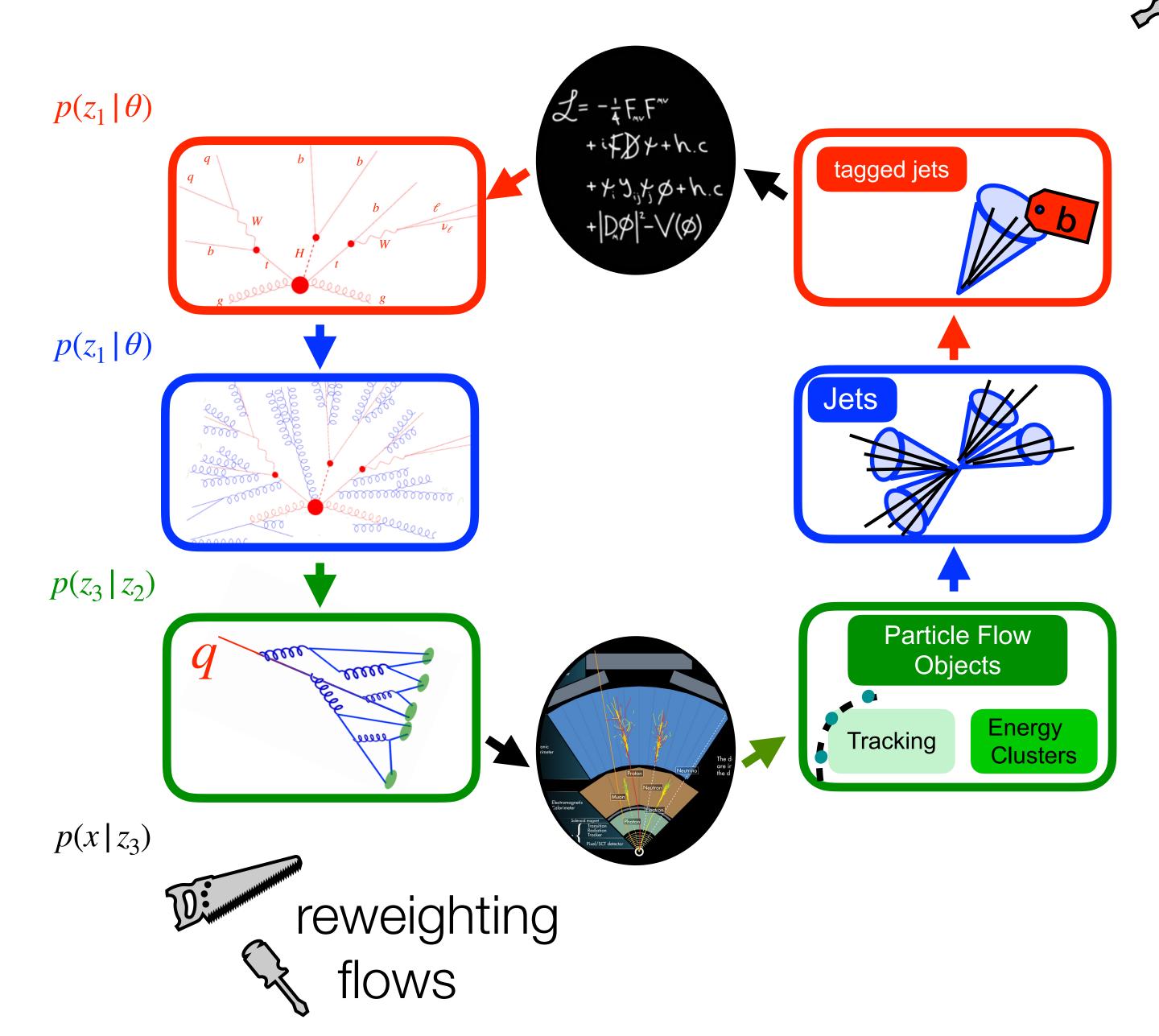
Promising for end-to-end models... holy grail for physics data anlaysis!

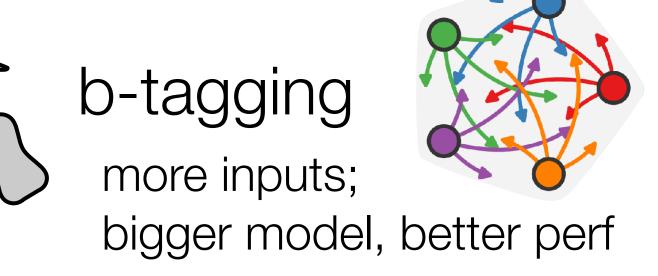
With A. Kofler, M. Kagan and L. Heinrich, in prep

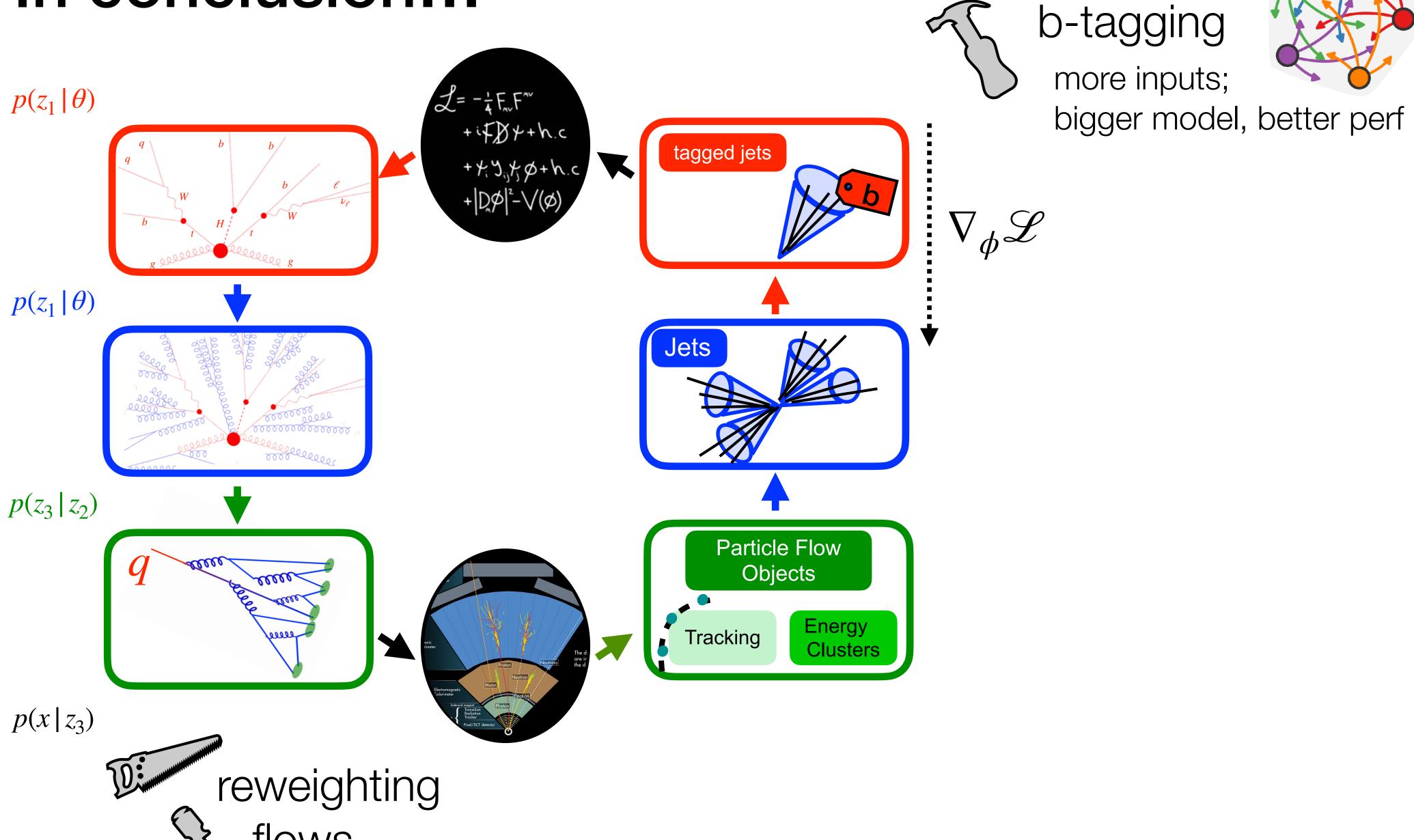












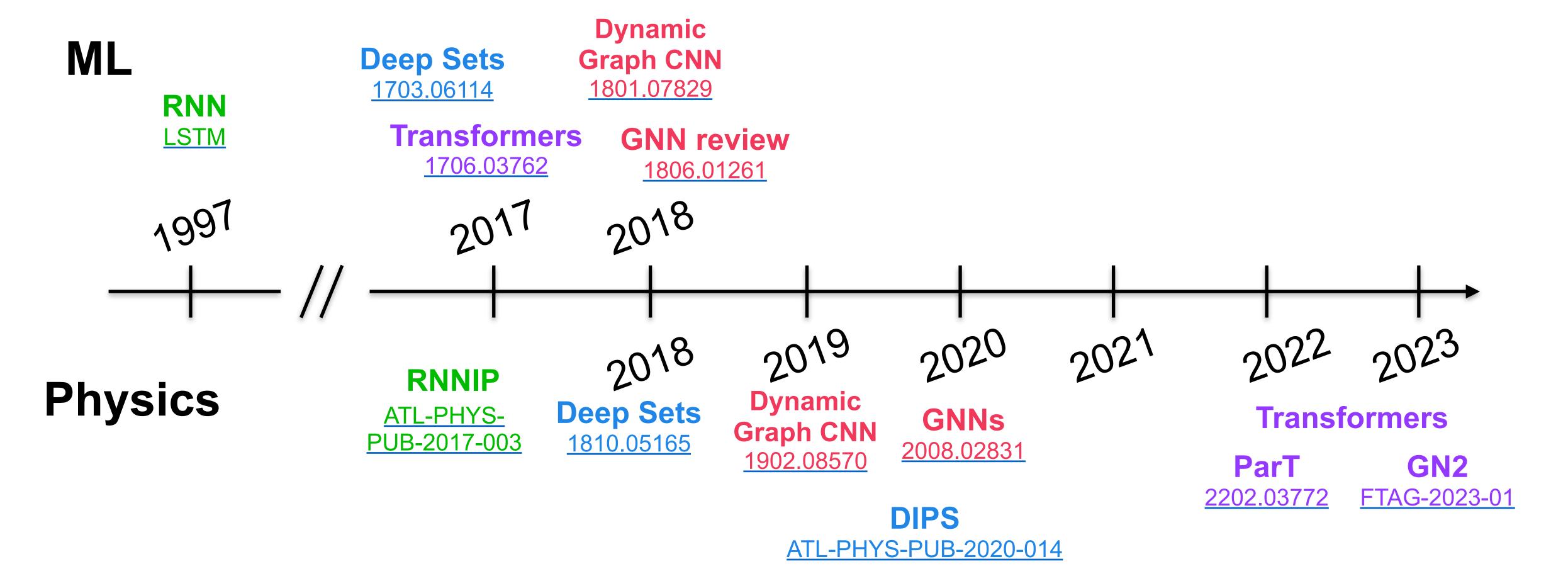
In conclusion... b-tagging more inputs; $p(z_1 | \theta)$ Z= -+ Fn-Fnbigger model, better perf + iFD++h.c tagged jets + + + y , + p + h.c $+\left|\sum_{n}\phi\right|^{2}-\bigvee(\phi)$ $p(z_1 | \theta)$ Jets $p(z_3 \mid z_2)$ Particle Flow Objects Energy Tracking Clusters $p(x | z_3)$ reweighting flows

In conclusion... b-tagging more inputs; $p(z_1 | \theta)$ Z= - + F_~ F^~ bigger model, better perf + i F) + + h.c tagged jets + + + y , + p + h.c $+\left|\sum_{n}\phi\right|^{2}-\bigvee(\phi)$ $p(z_1 | \theta)$ Jets $\nabla_{\theta} \log p(x; \theta)$ $p(z_3 \mid z_2)$ Particle Flow Objects Energy Tracking Clusters $p(x \mid z_3)$ reweighting flows 50

In conclusion... b-tagging more inputs; $p(z_1 | \theta)$ Z= - + F_~ F^~ bigger model, better perf + i F) + + h.c tagged jets + + + y , + p + h.c $+\left|\sum_{n}\phi\right|^{2}-\bigvee(\phi)$ $p(z_1 | \theta)$ Jets $\nabla_{\theta} \log p(x; \theta)$ $p(z_3 \mid z_2)$ Particle Flow Objects Energy Tracking Clusters Thanks. $p(x \mid z_3)$ reweighting flows 50

Backup

Timeline



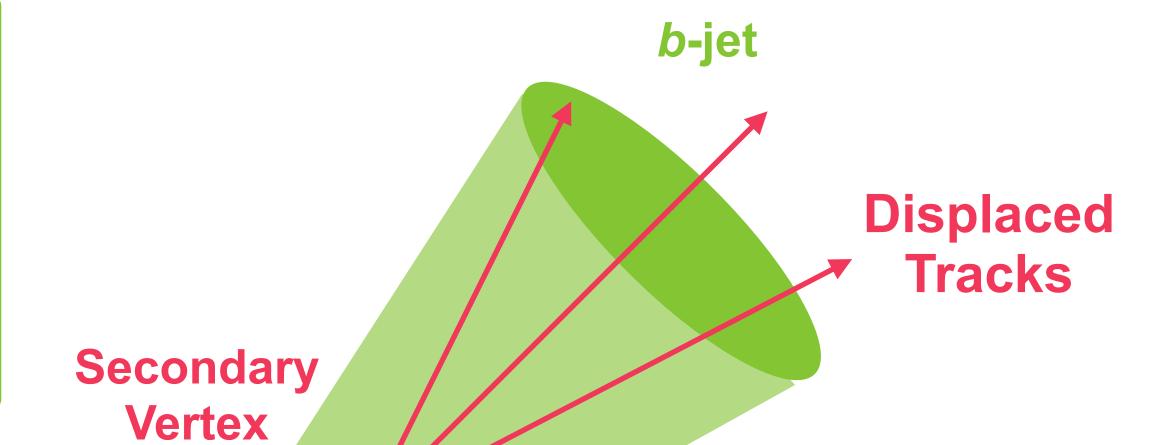




✓ "Long" lifetime: τ = 1.2 ps

✓ Many (≈ 5) displaced tracks

Jet

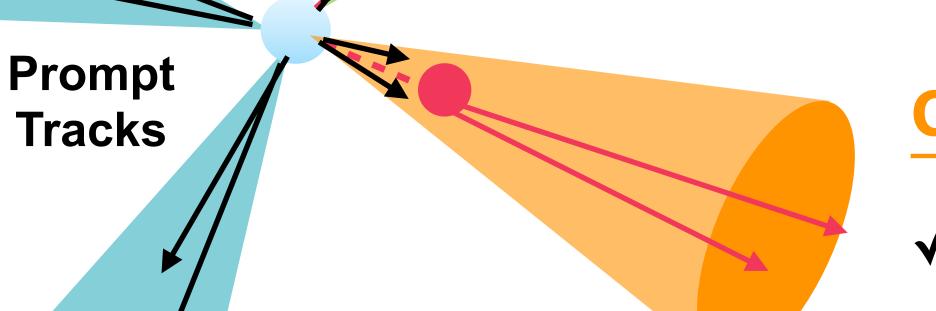


light jet

u,d,s g

✓ Most tracks originating from the PV

√ Few displaced tracks



Primary

Vertex

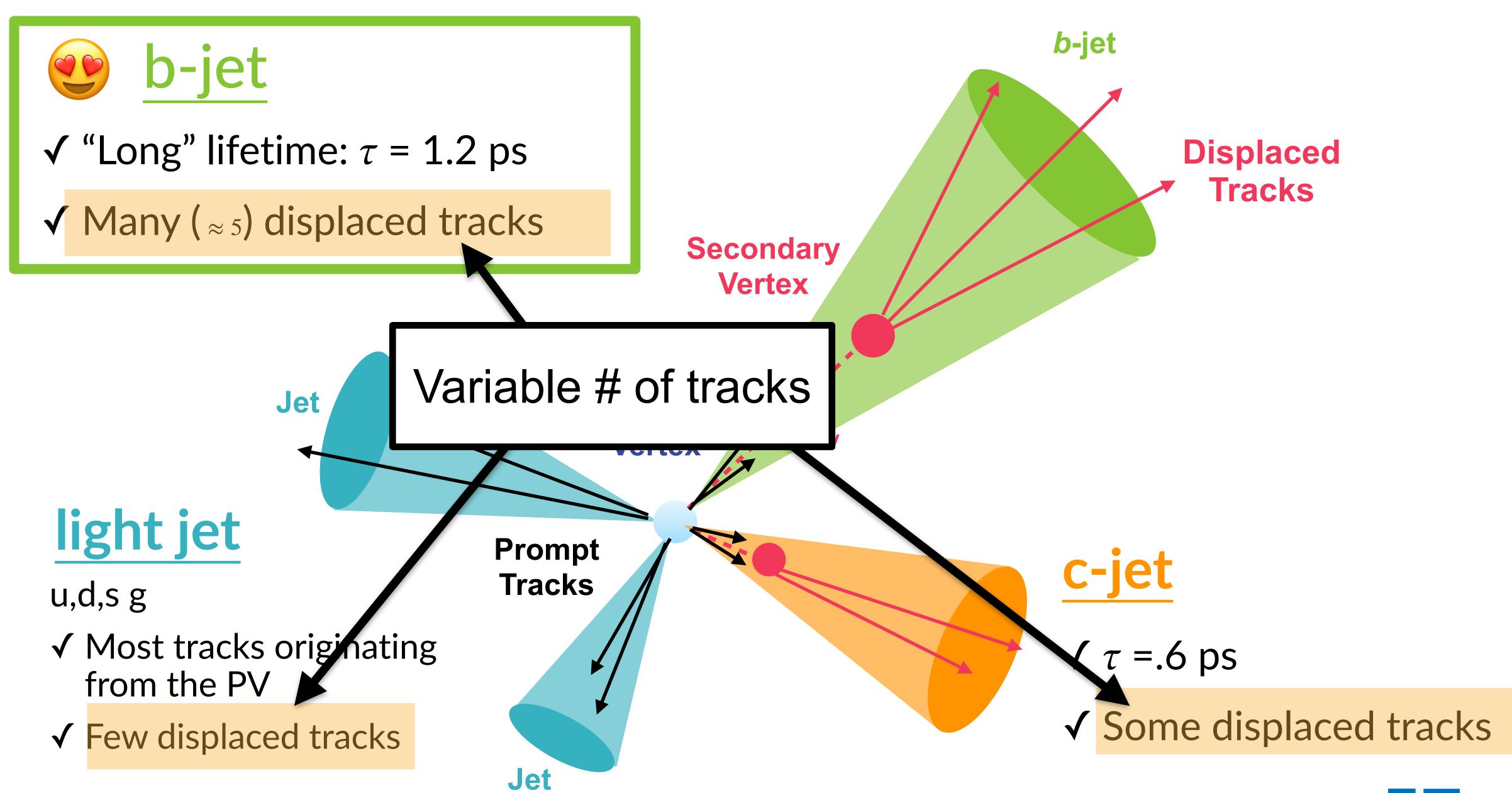
Jet

c-jet

$$\sqrt{\tau}$$
 = .6 ps

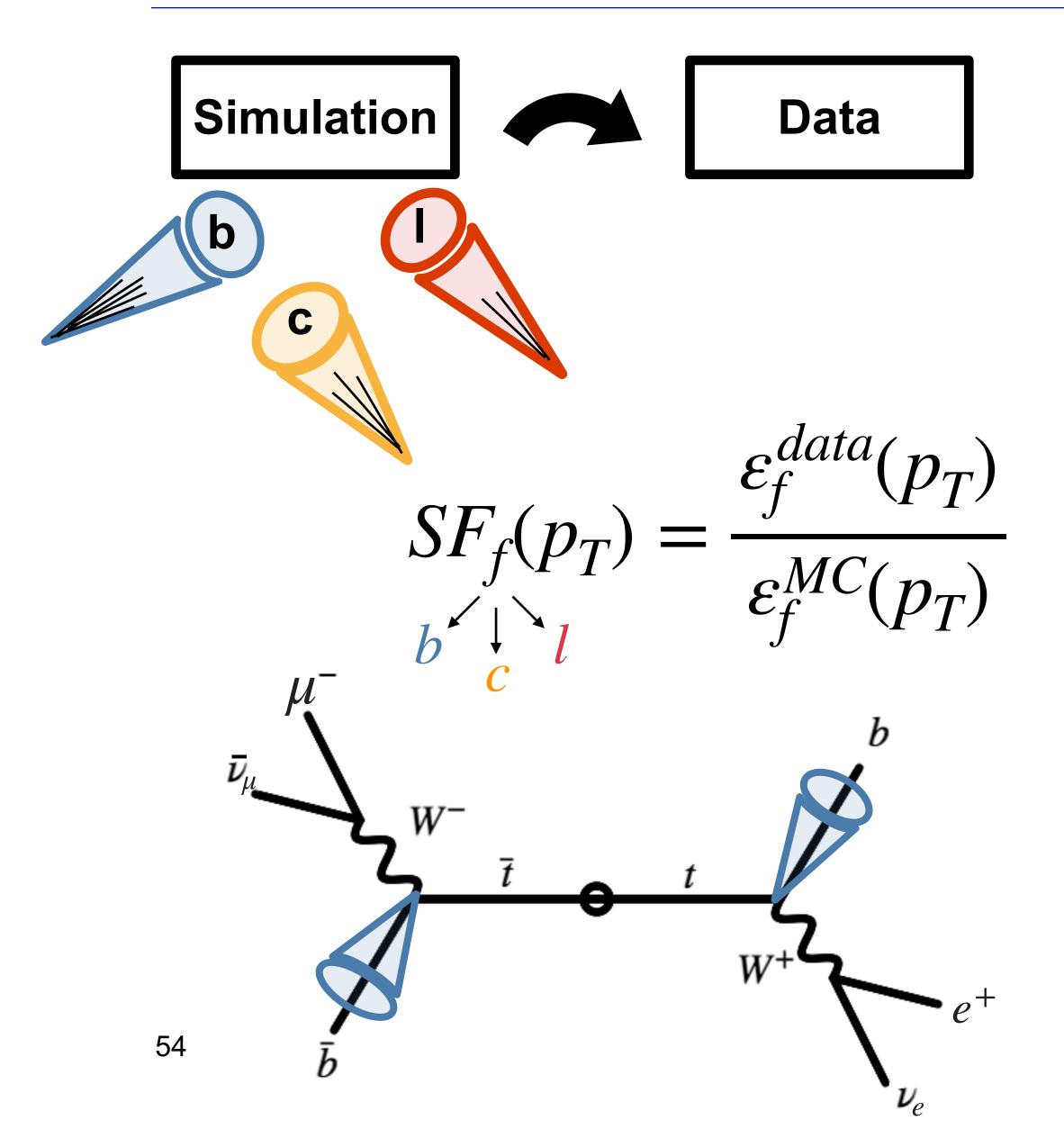
√ Some displaced tracks

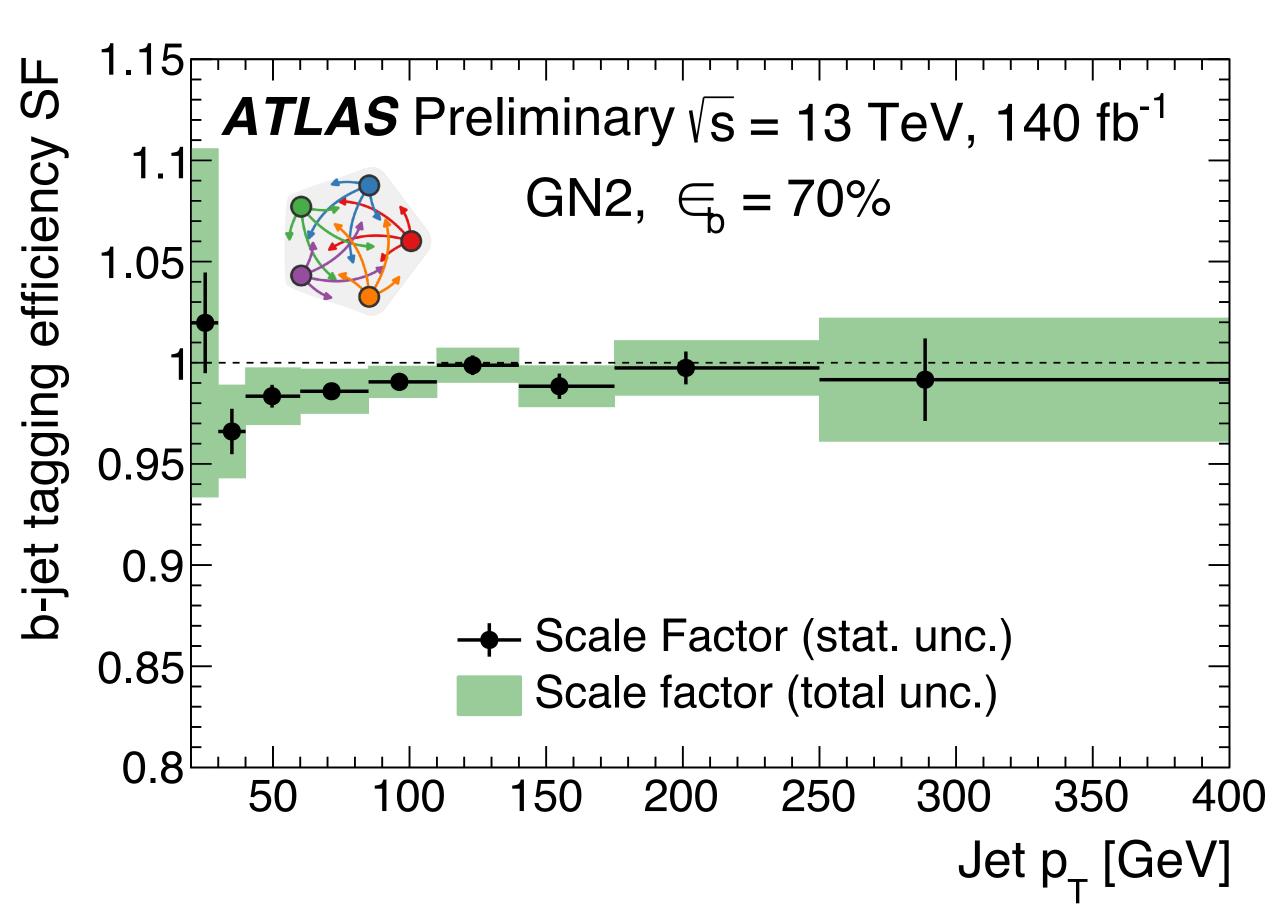






Calibration



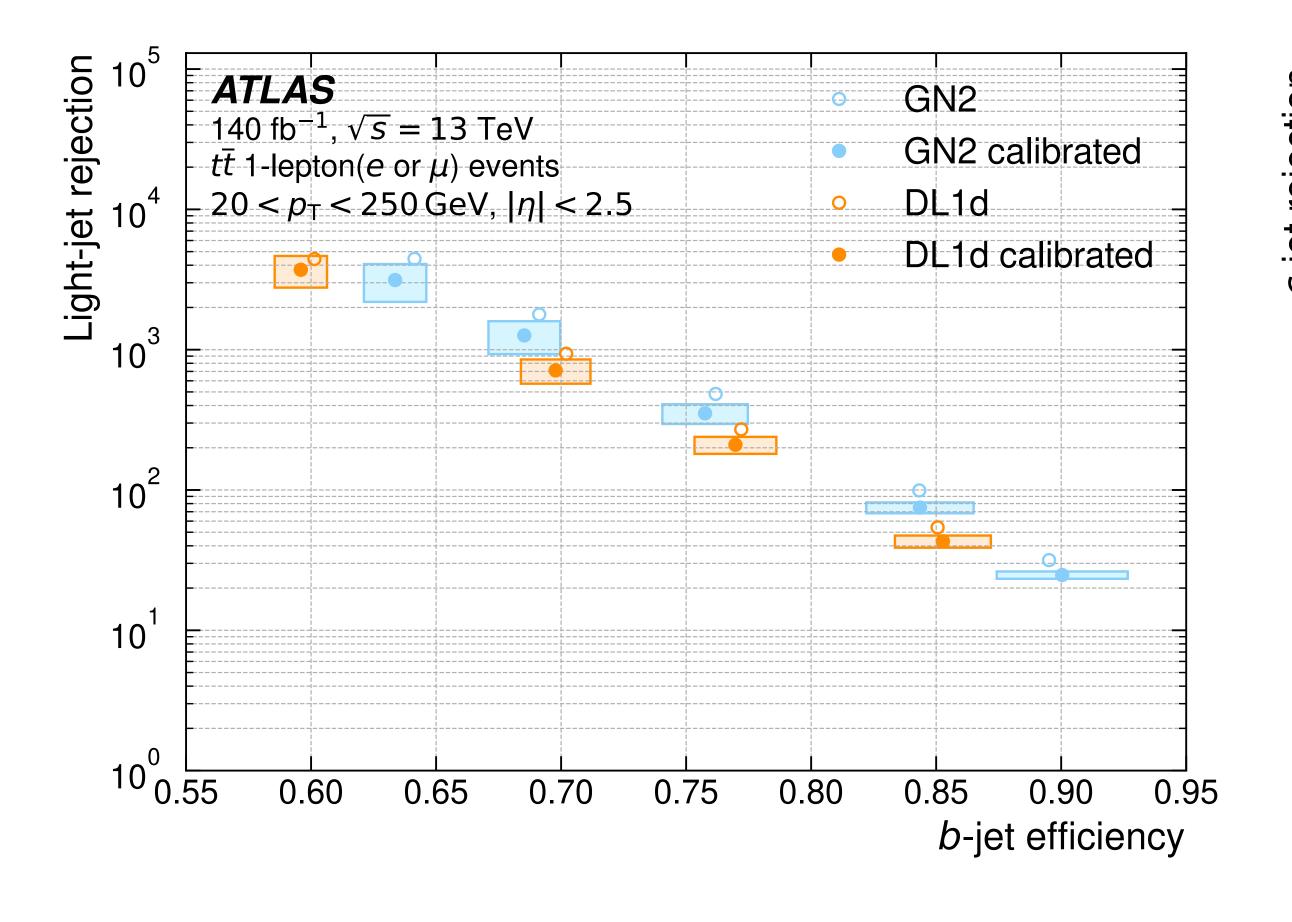


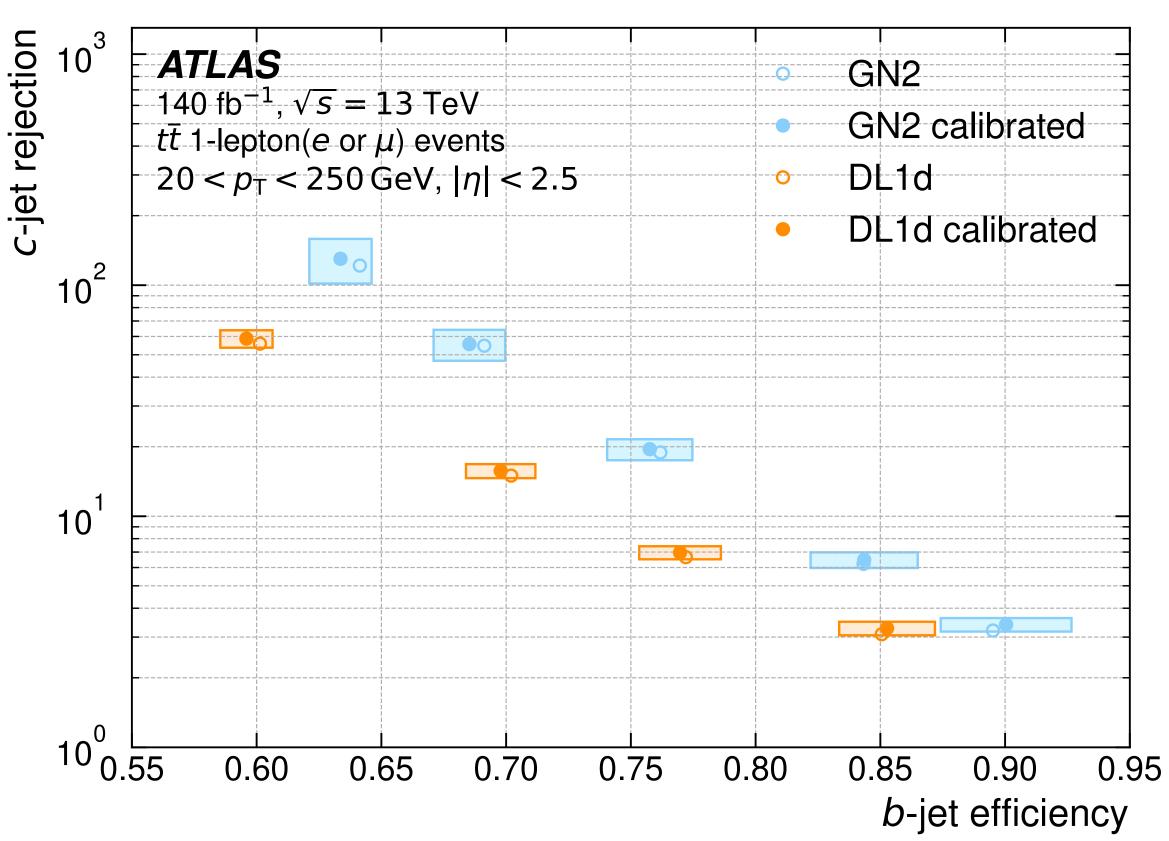


% level non-closure



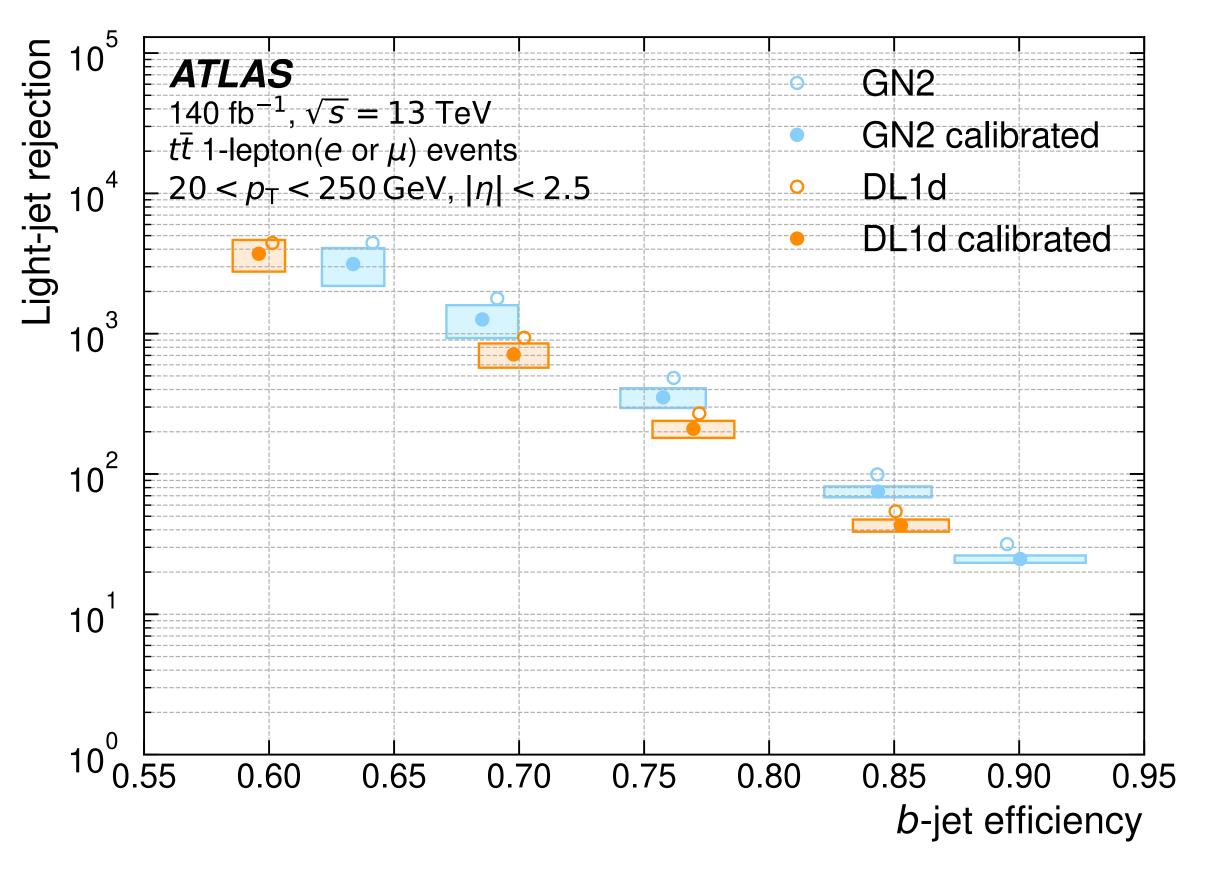
And... it translates to the physics (!)

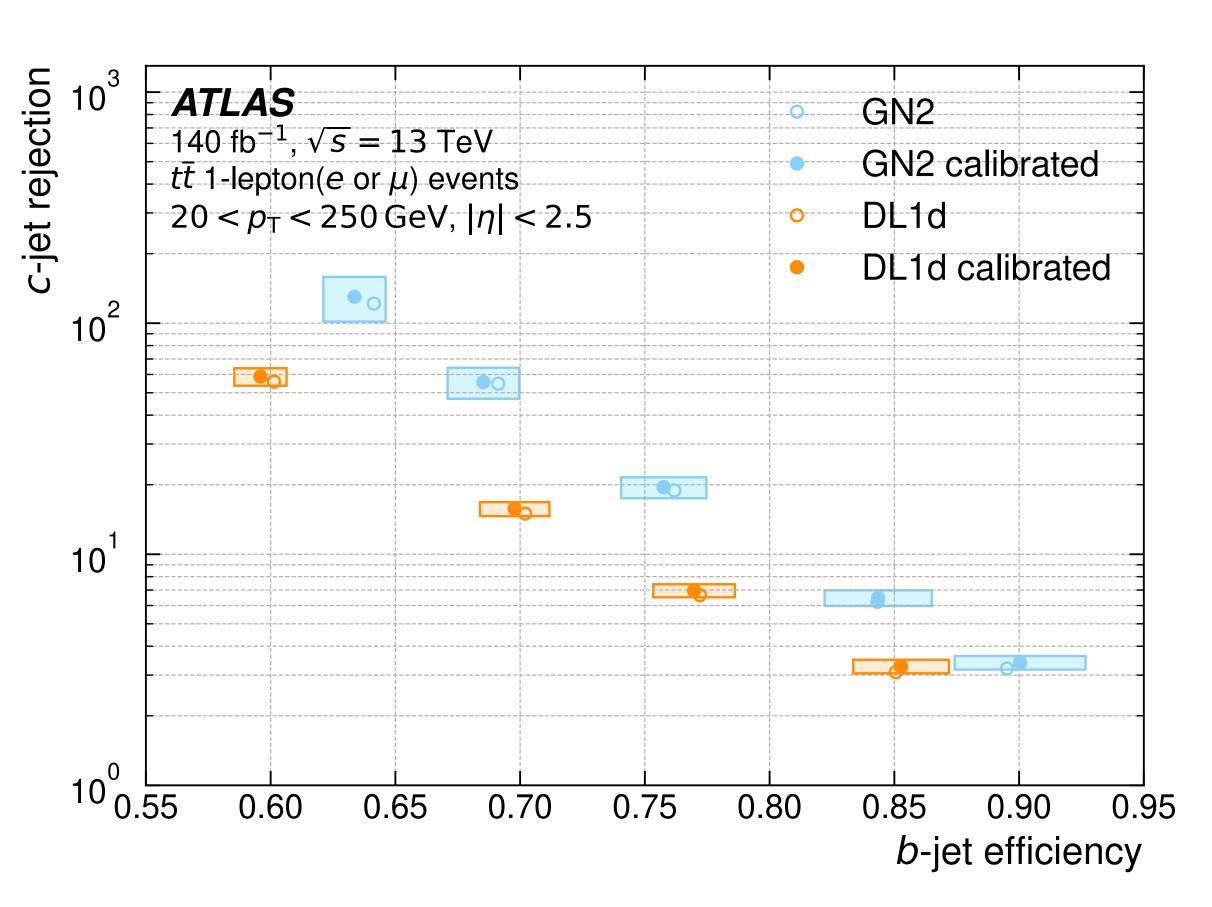






And... it translates to the physics (!)



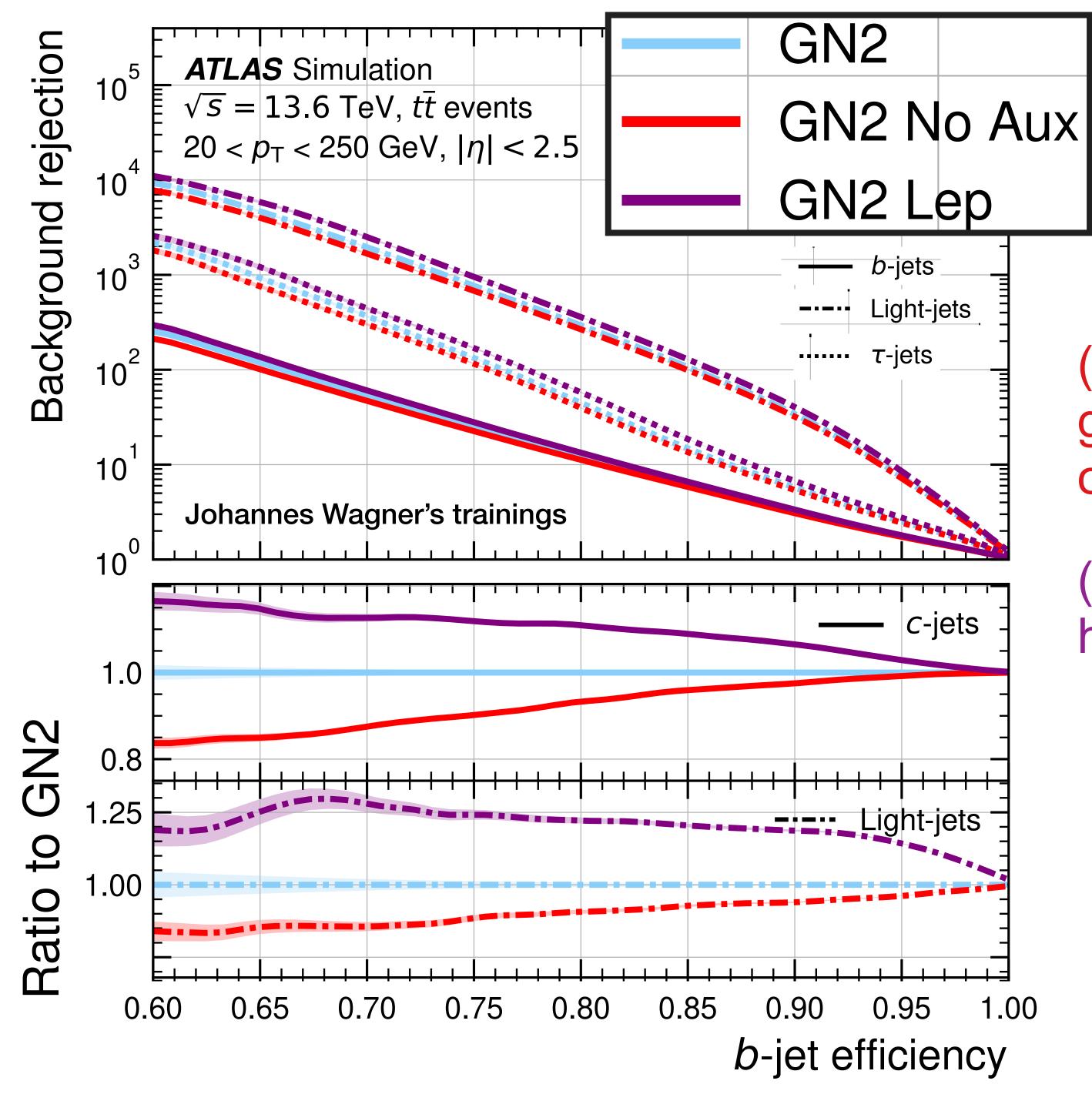




GN2 used in *new* $HH \rightarrow bb\gamma\gamma$ analysis (with 2022 — 2024 dataset)

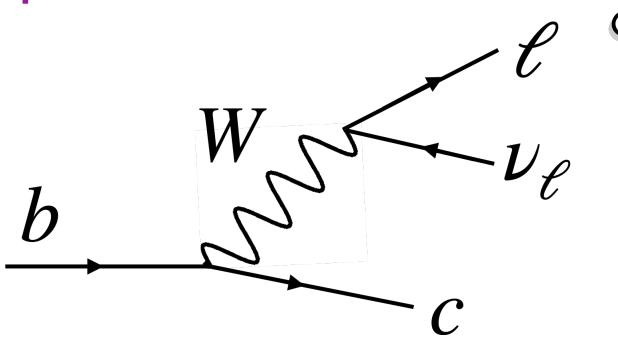


GN3



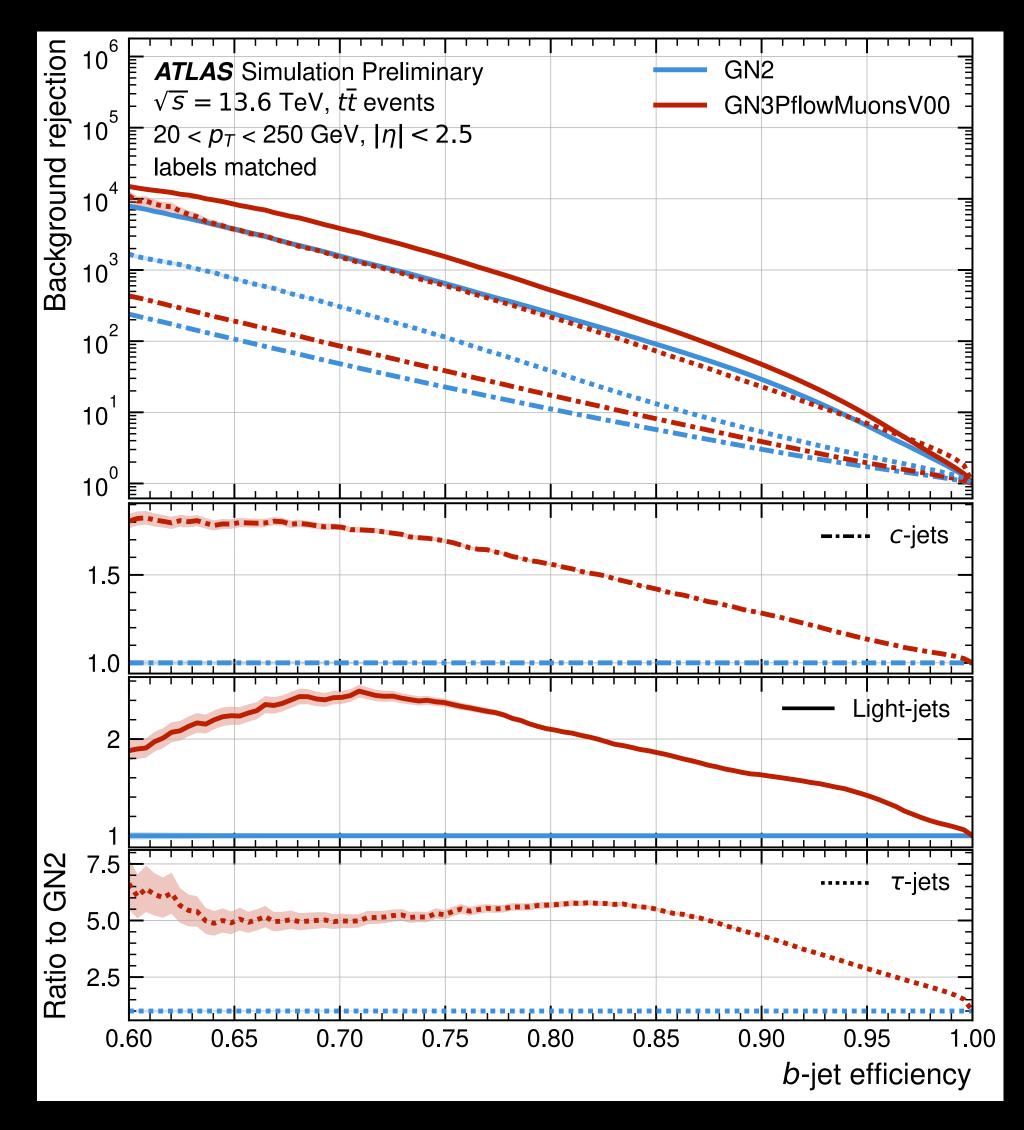
(1) The auxiliary tasks guide the model to more optimal solutions

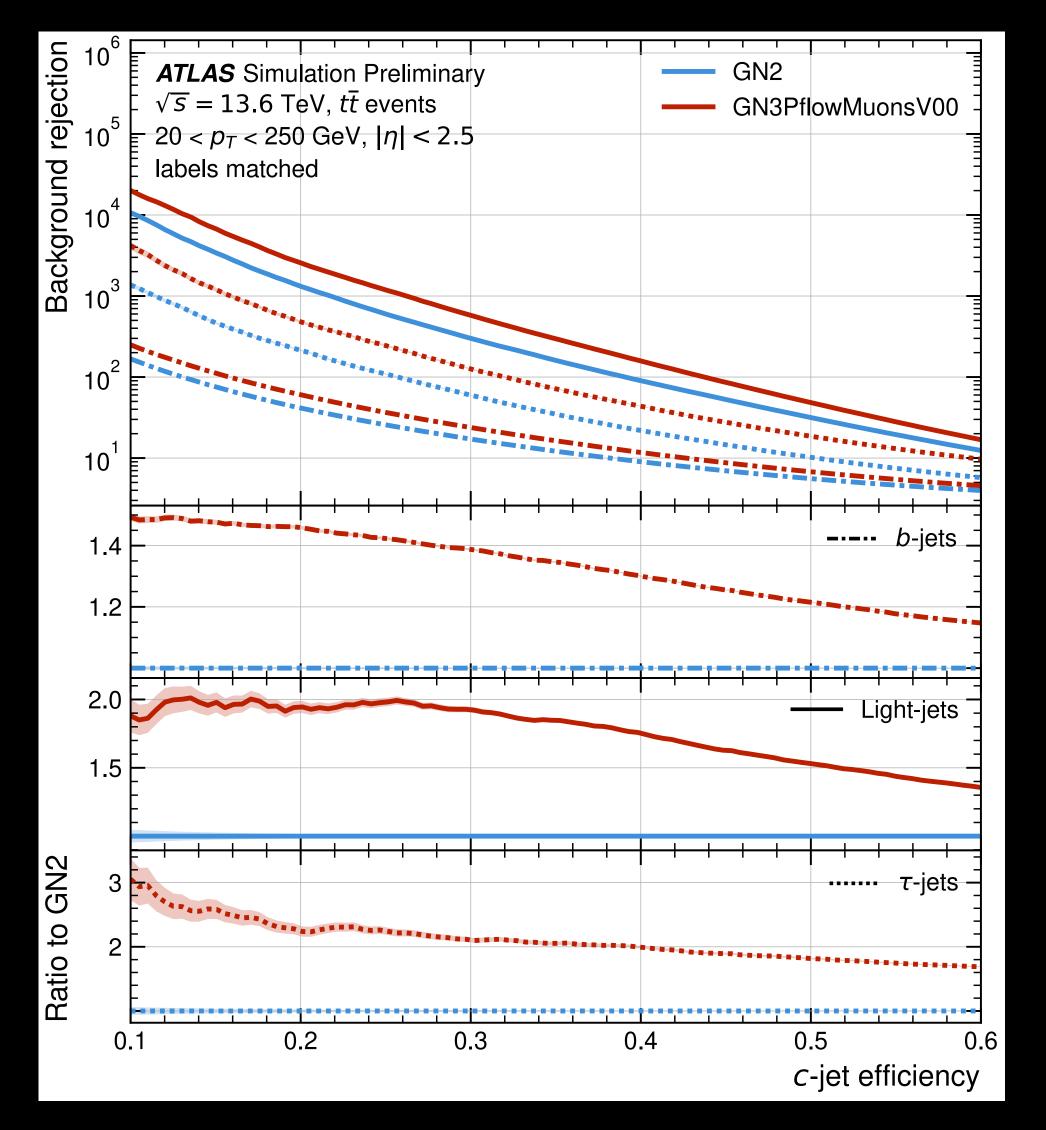
(2) Extra physics info helps more





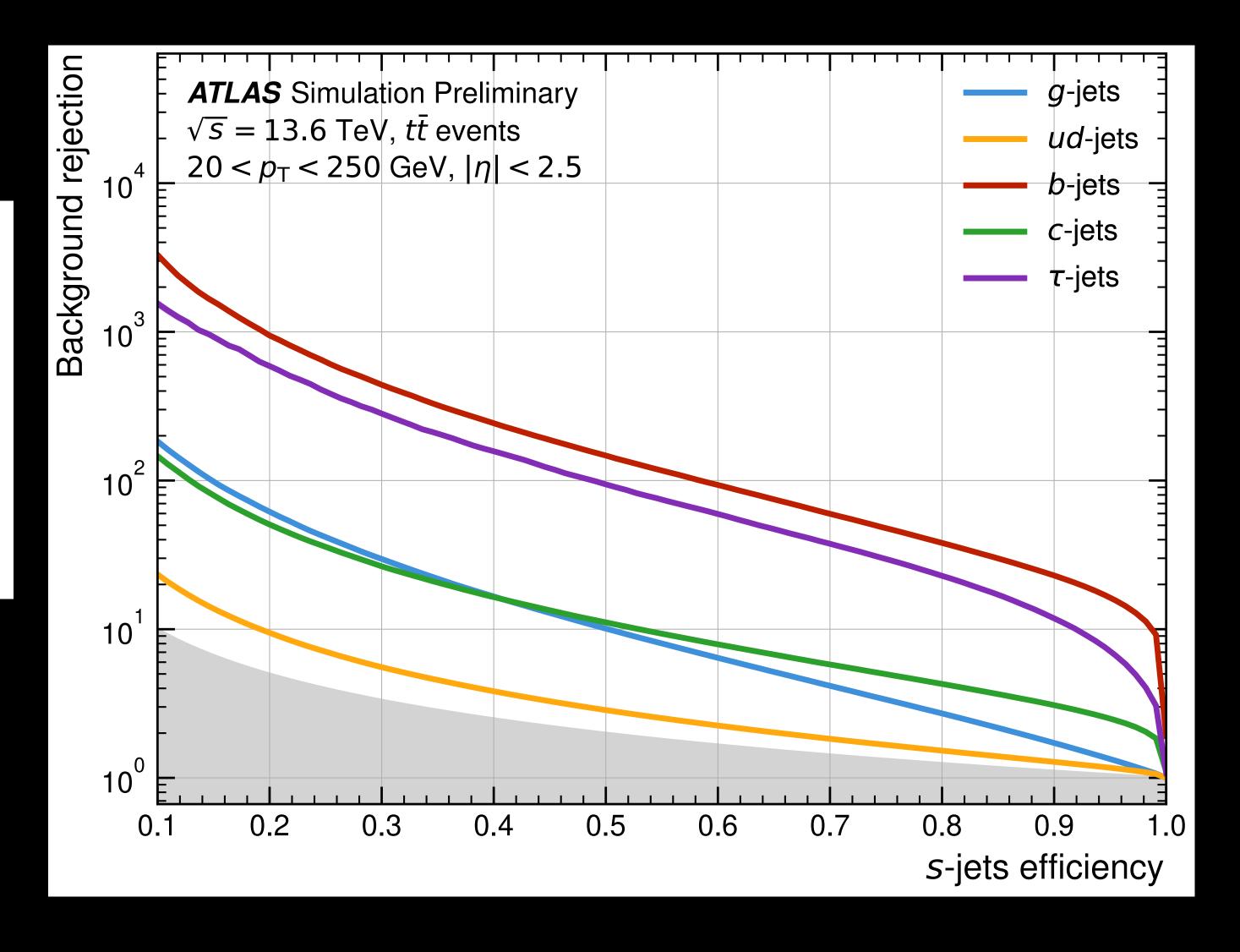
GN3: Multi-task, multi-modal b-tagger



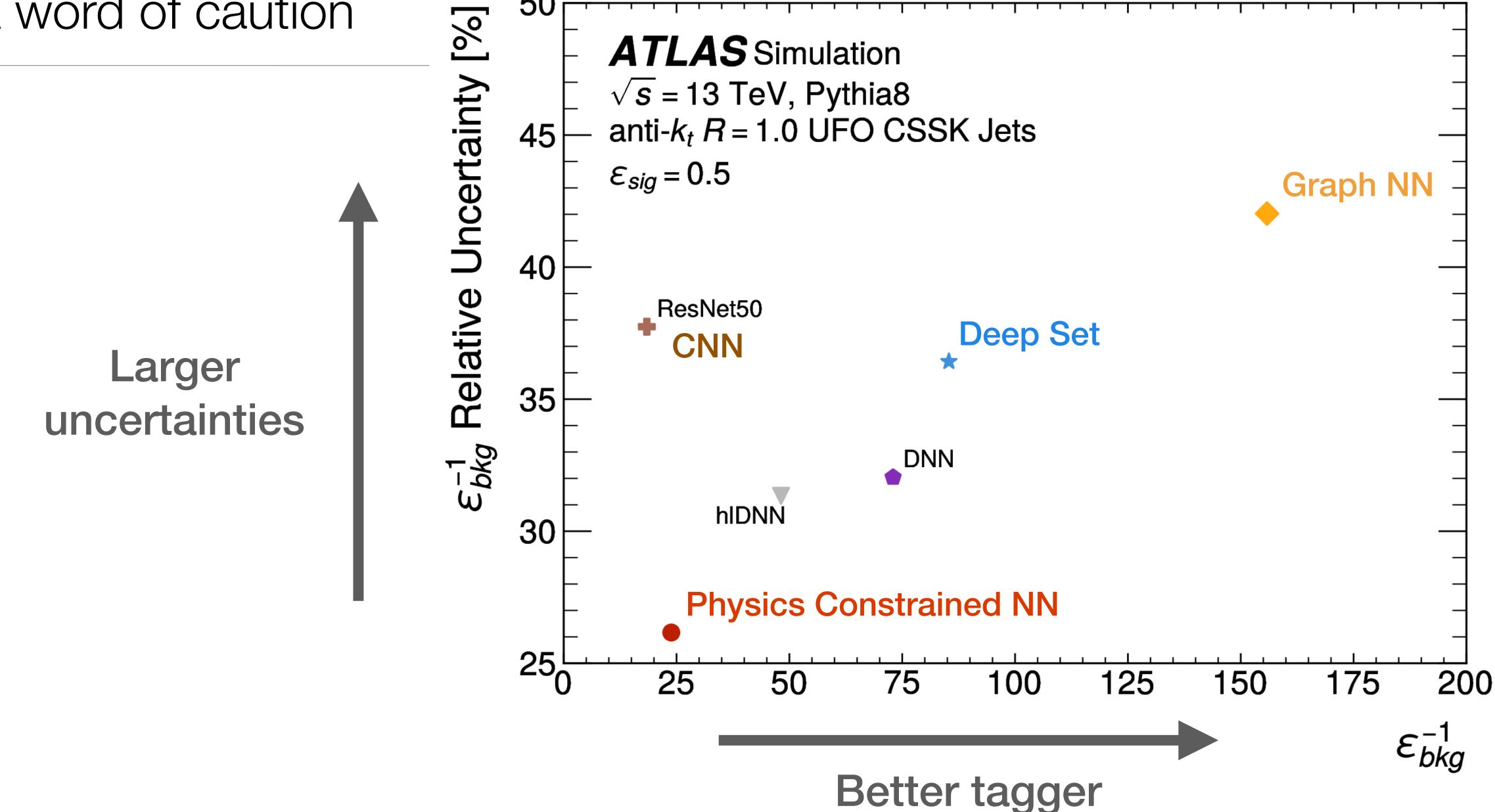


GN3: Multi-task, multi-modal b-tagger

Component	GN2	GN3
Track-jet association	ΔR -based association	Ghost-association
Track selection	$d_0 < 3.5 \text{mm}, 40 \text{tracks}$	d_0 < 5 mm, 50 tracks
Inputs	Jets and tracks	Jets, ghost-associated tracks, soft muons
		PFlow objects
Activation function	ReLU	SiLU
Initialisation layers	256	512
Transformer encoder	4 scaled dot-product atten-	4 Flash Attention, Gated Linear Units
	tion	8 register tokens
Embedding dimension	256	512
Jet classification	4 classes: b, c, ℓ, τ	6 classes: b, c, ud, s, g, τ
Loss balancing	Fixed weights per task	Geometric mean of losses (GLS)
Optimiser	AdamW	Lion







50

2019: The bitter lesson (Richard Sutton)

"The biggest lesson that can be read from 70 years of Al research is that general methods that leverage computation are ultimately the most

effective, and by a large margin."

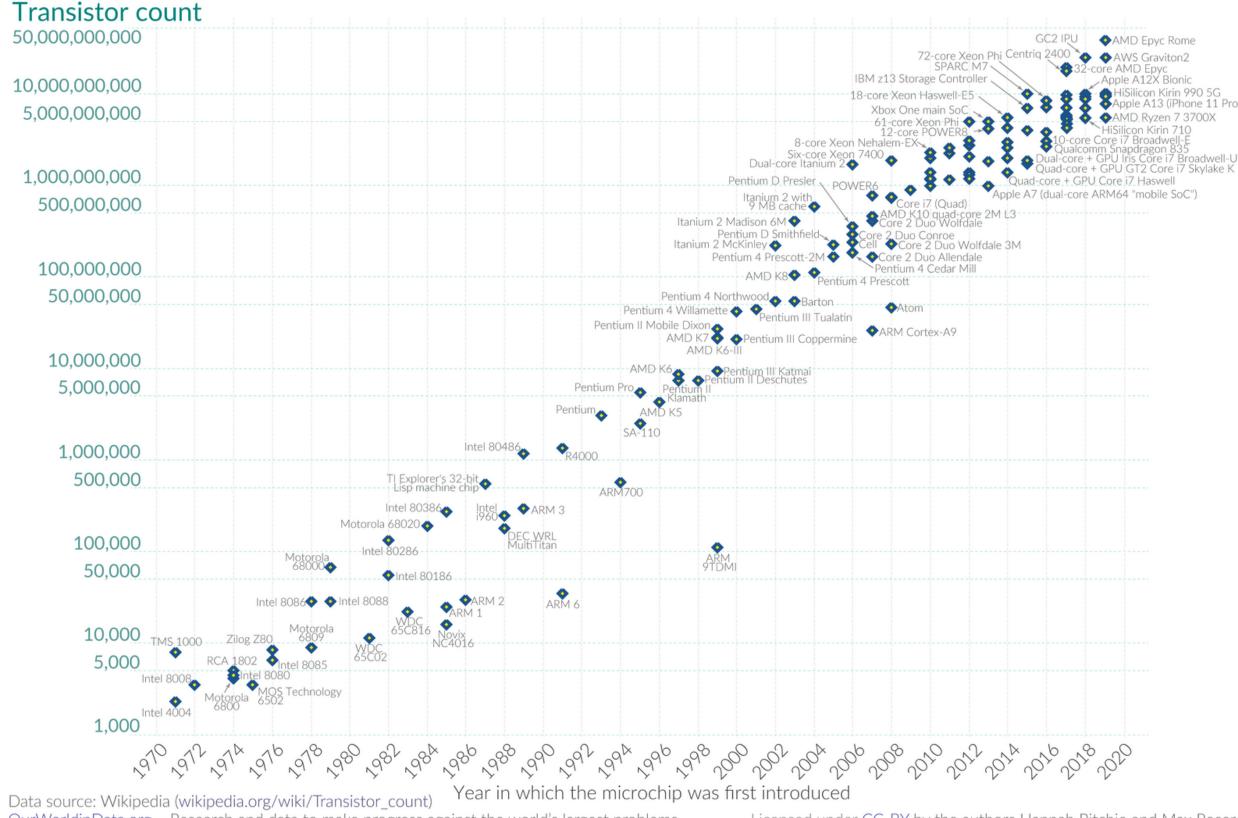
Exs: Image Maps Input **Fully Connected** Convolutions Subsampling

Illustration of LeCun et al. 1998 from CS231n 2017 Lecture

Moore's Law: The number of transistors on microchips doubles every two years Our World

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers

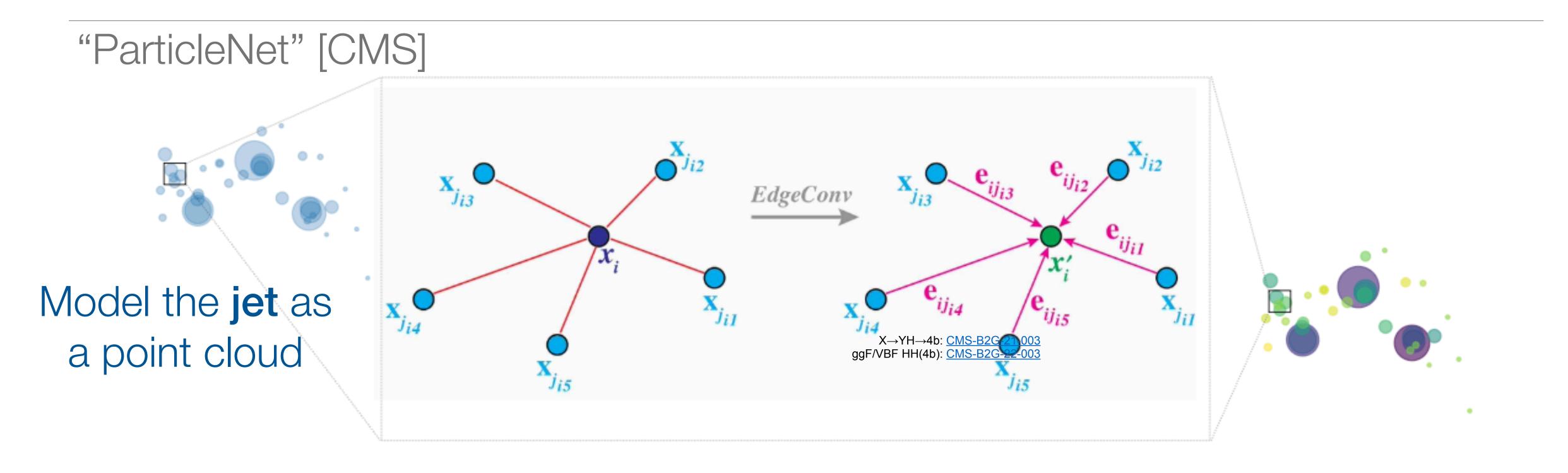




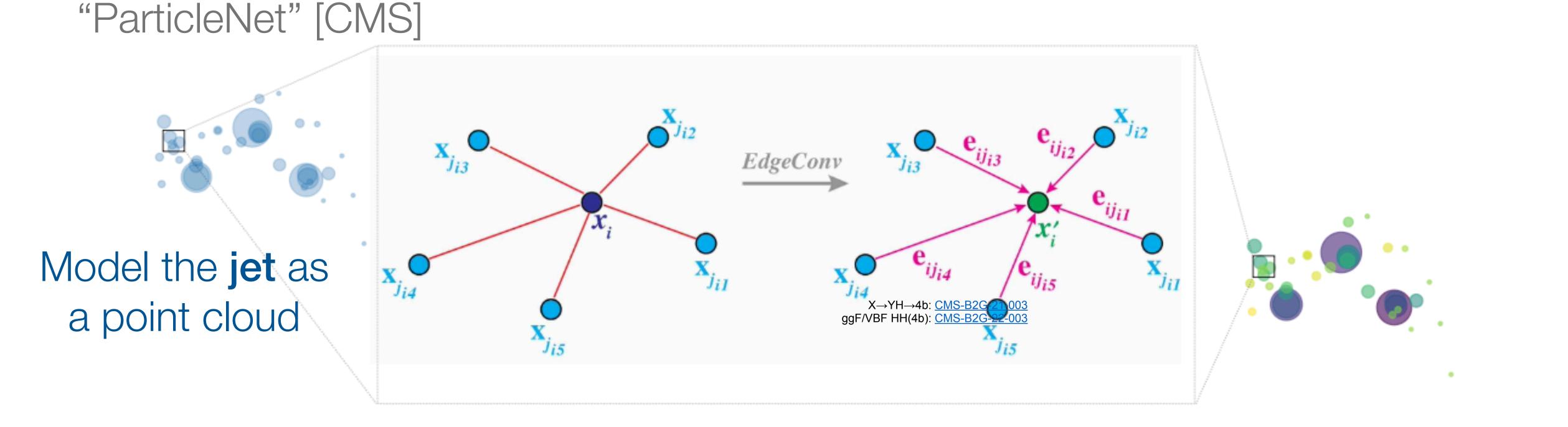
OurWorldinData.org - Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

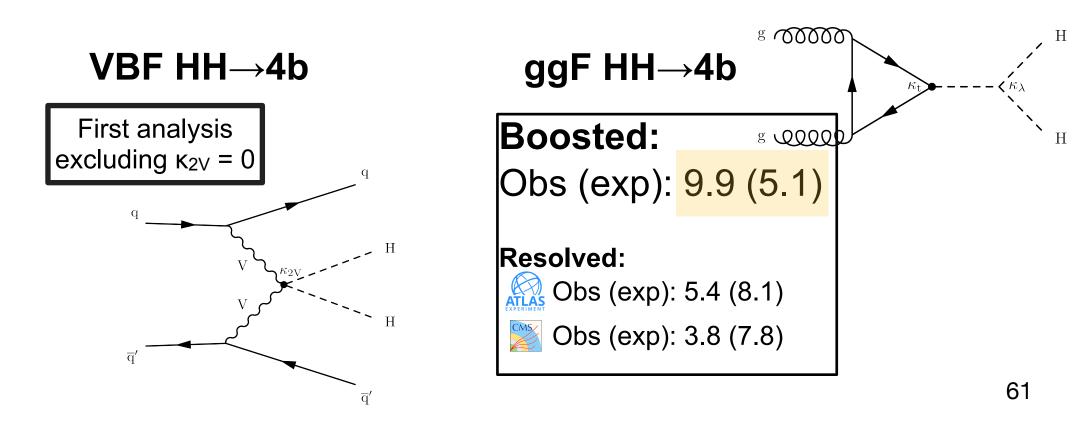
Dynamic Graph CNN



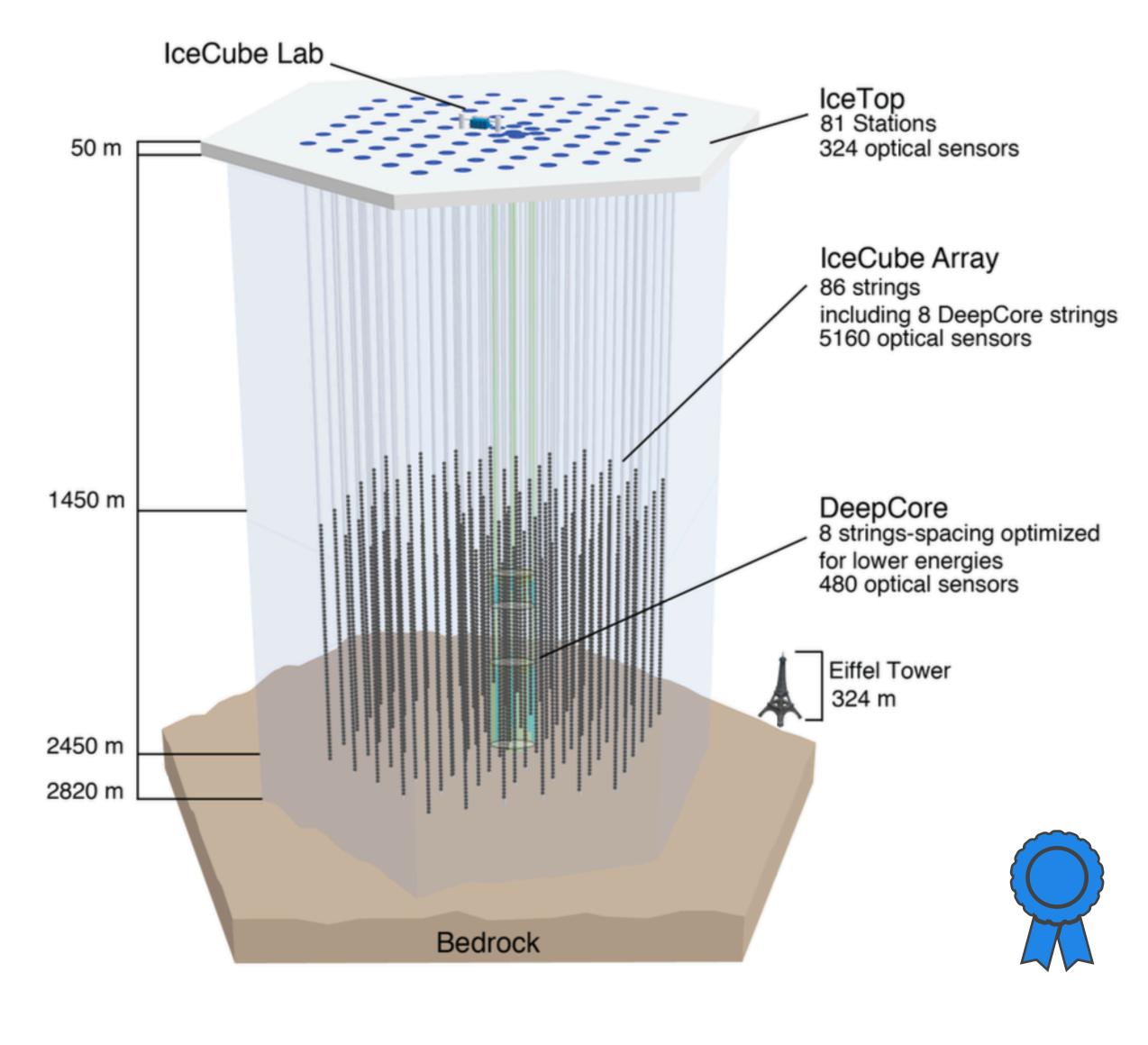
Dynamic Graph CNN



Very impressive physics results
But was the graph representation needed?









3.1.4 Edge Selection

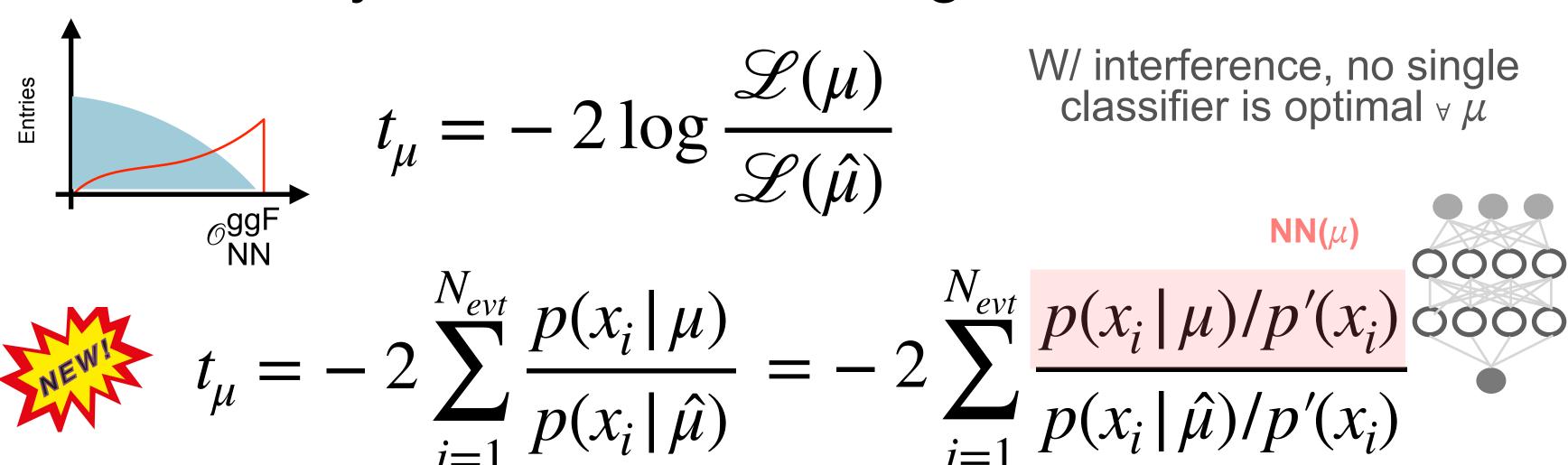
In the original EdgeConv implementation, edges are calculated in each layer dynamically by k-Nearest Neighbors (kNN). However, this edge selection scheme is not differentiable in itself and therefore does not have gradients. This would still work well for the segmentation task in the original paper, as points in the same segment are trained to be close in the latent space. However, the situation is different in this task, where it would not make sense to dynamically select edges. Therefore, the edges used in EdgeConv are calculated from the input features only once in our model.

All three of the winning solutions used a transformer architecture.

Also in jet tagging (ATLAS and CMS), transformers outperform GNN architectues.

Simulation based inference

Traditionally, train MVA & histogram-ize

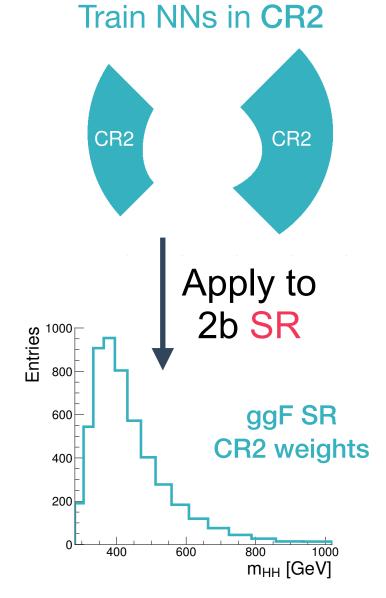


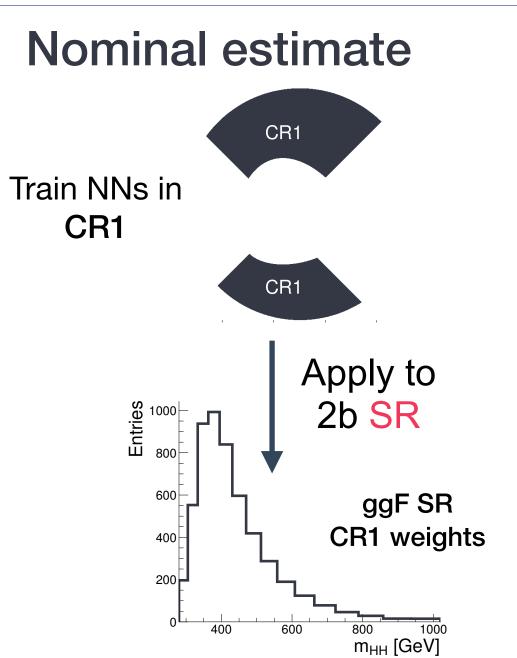
Trick: don't need $NN(\mu)$

$$\frac{p(x \mid \mu)}{p'(x)} = \frac{1}{\sigma(\mu)} \left(\mu \sigma_s \frac{p_S(x)}{p'(x)} + \sqrt{\mu} \sigma_I \frac{p_I(x)}{p'(x)} + \sigma_B \frac{p_B(x)}{p'(x)} \right)$$
Neural Networks

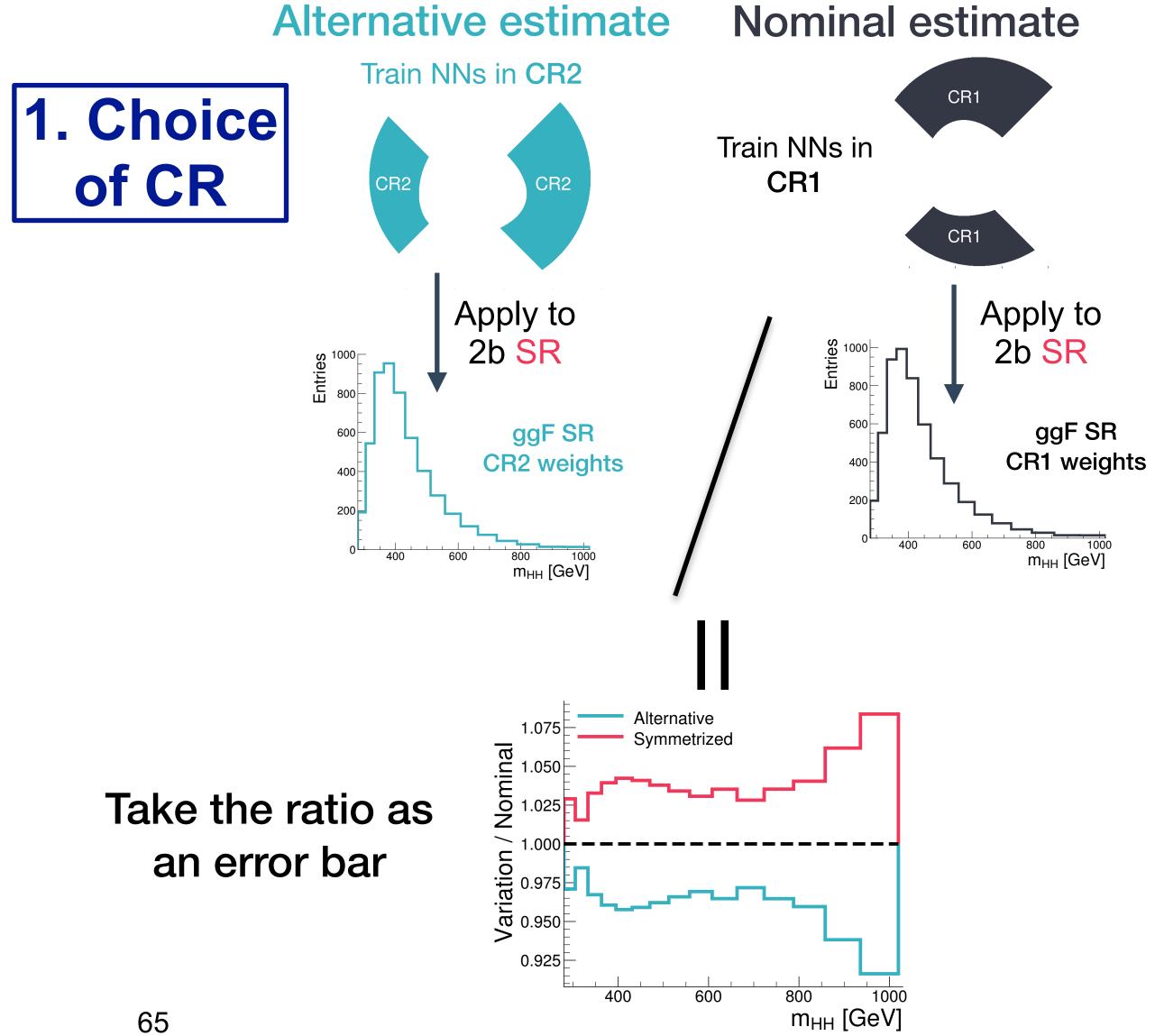
Alternative estimate

1. Choice of CR





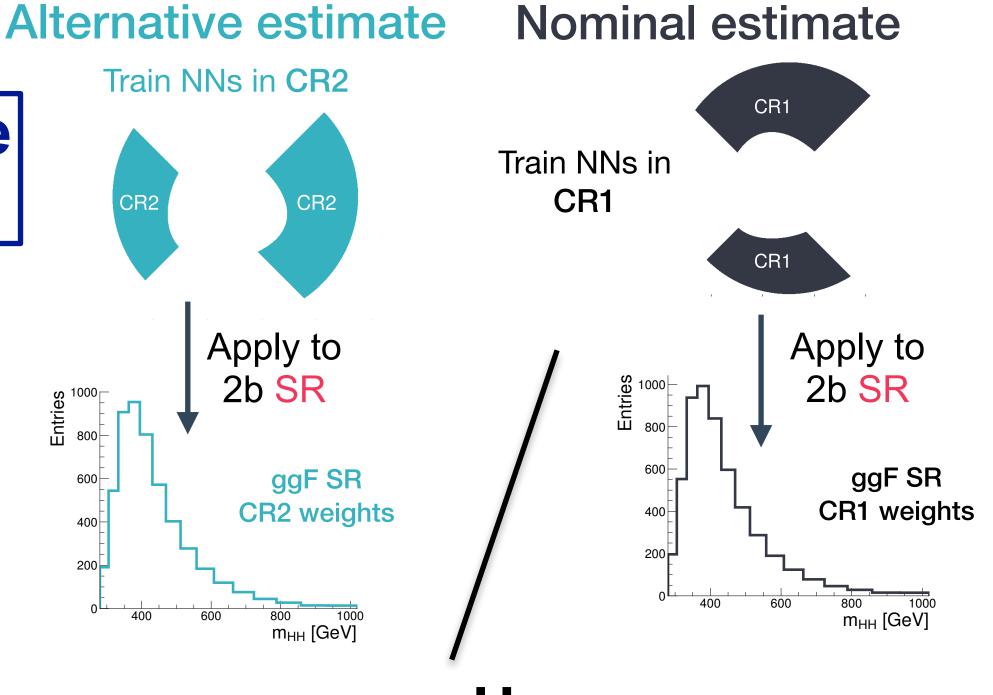




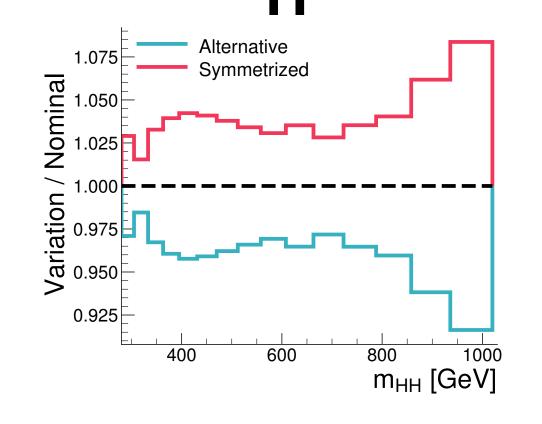




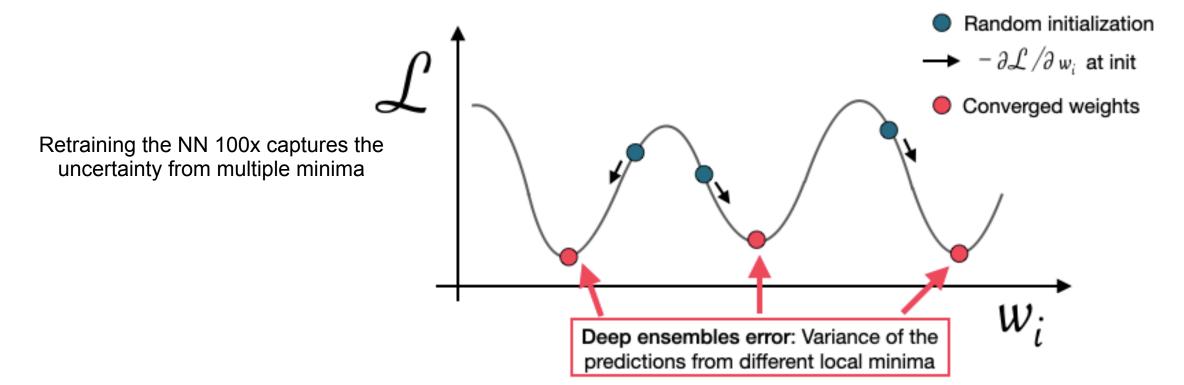




Take the ratio as an error bar

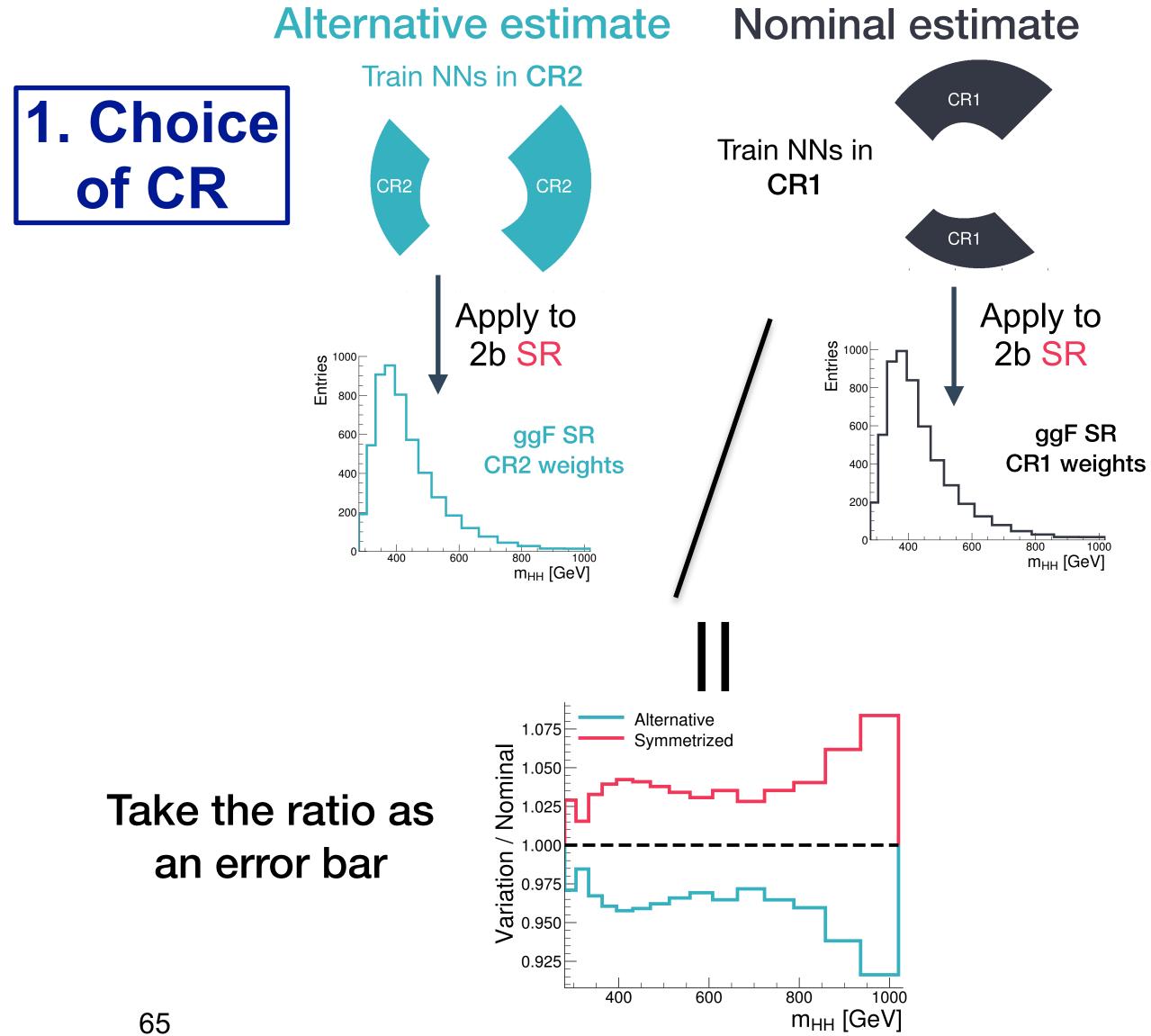


2. NN initialization

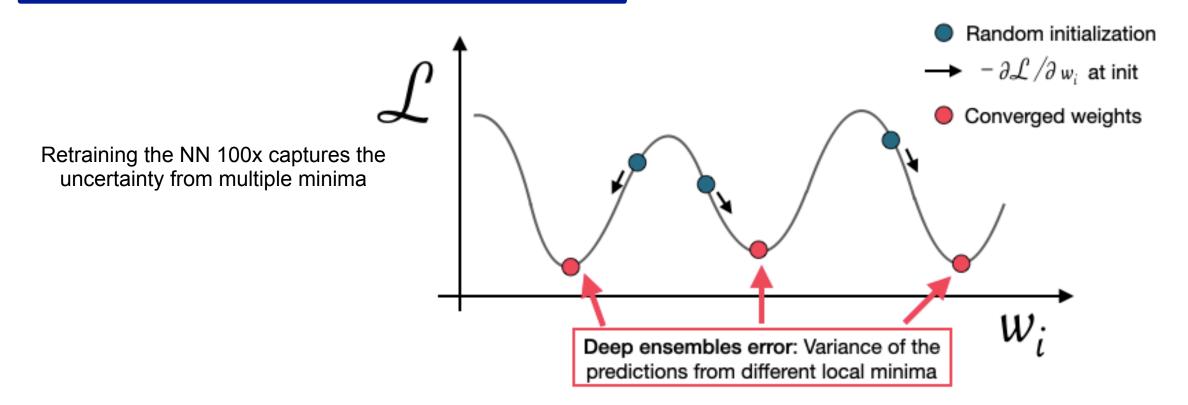


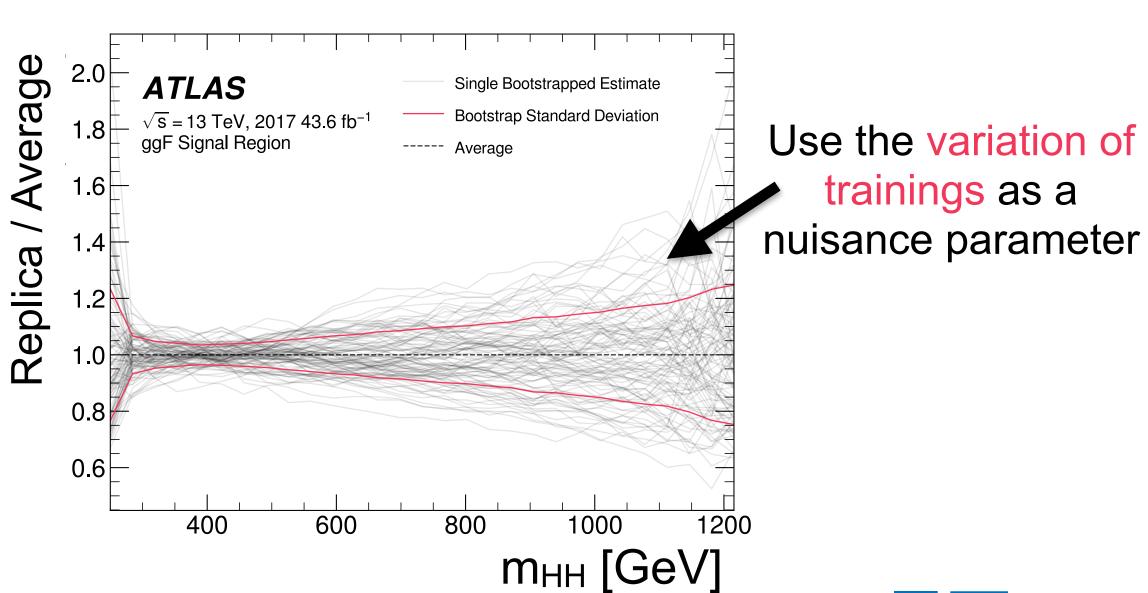


2301.03212



2. NN initialization







Flows for high dimensional interpolation

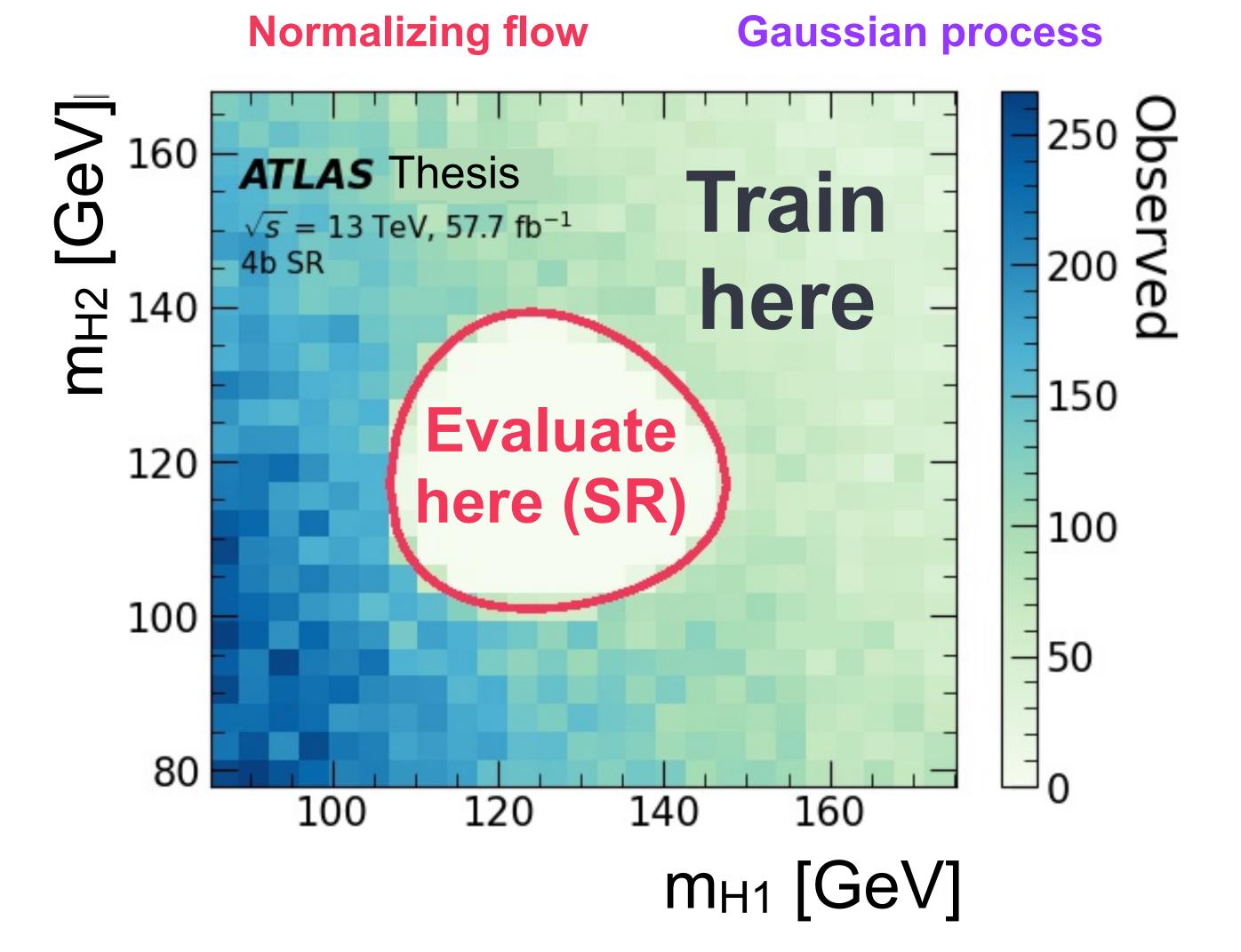
Hierarchical $p(x, m_{H1}, m_{H2}) = p(x|m_{H1}, m_{H2}) p(m_{H1}, m_{H2})$ podel:

x: Event kinematics

 $p_{T,H1}$, $p_{T,H2}$, η_{H1} , η_{H2} , $\Delta \phi_{HH}$, X_{Wt} [top veto]

Use the smoothly varying (m_{H1}, m_{H2}) to predict SR kinematics.

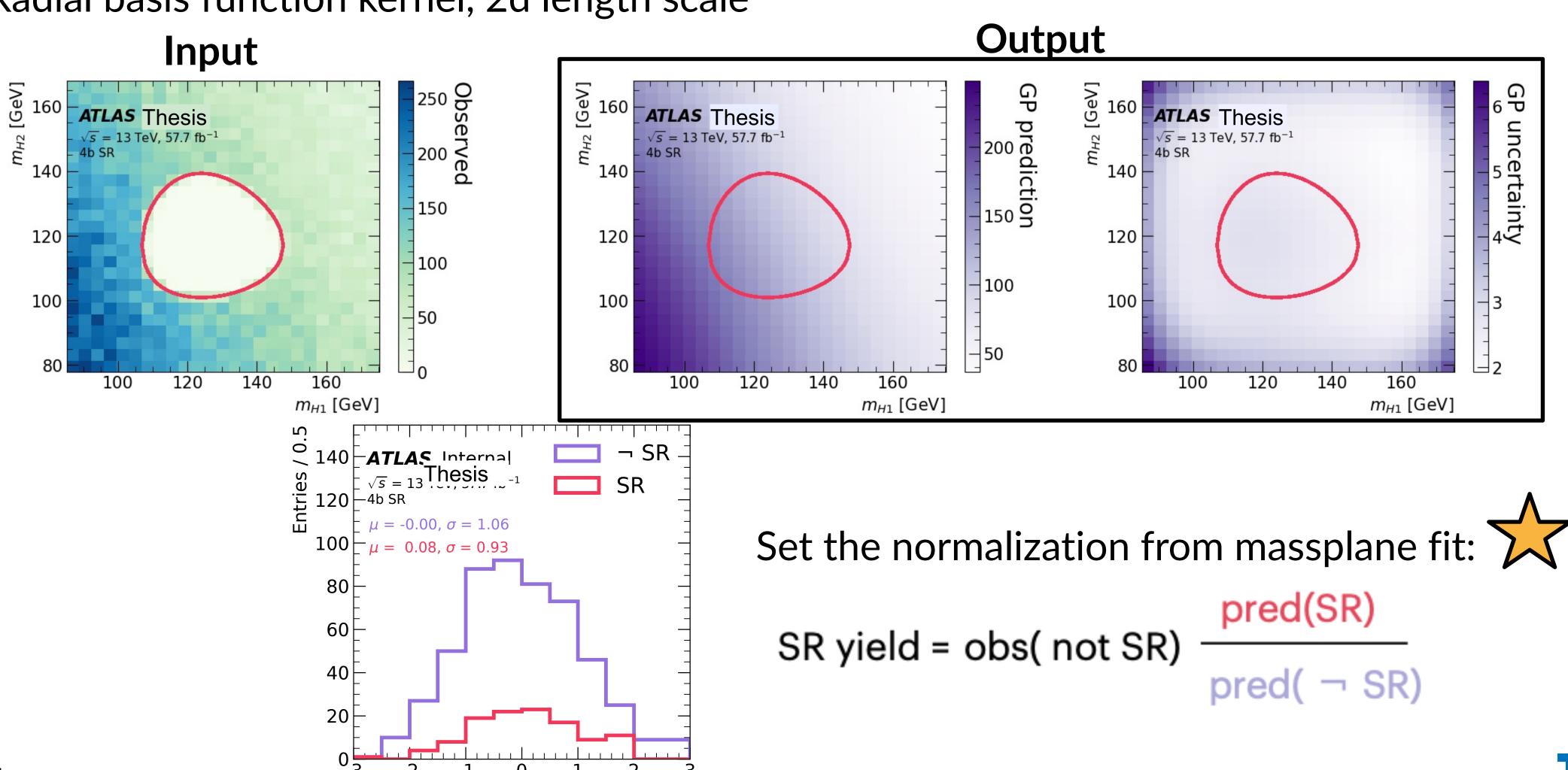
Conceptually identical to CATHODE



1) GP fits

 $p(x, m_{H1}, m_{H2}) = p(x | m_{H1}, m_{H2}) p(m_{H1}, m_{H2})$

Fit a GP to the 2d (m_{H1} , m_{H2}) histogram Radial basis function kernel, 2d length scale



 $\frac{\mathsf{GP}-\mathsf{obs}}{\sqrt{\mathsf{obs}}}$

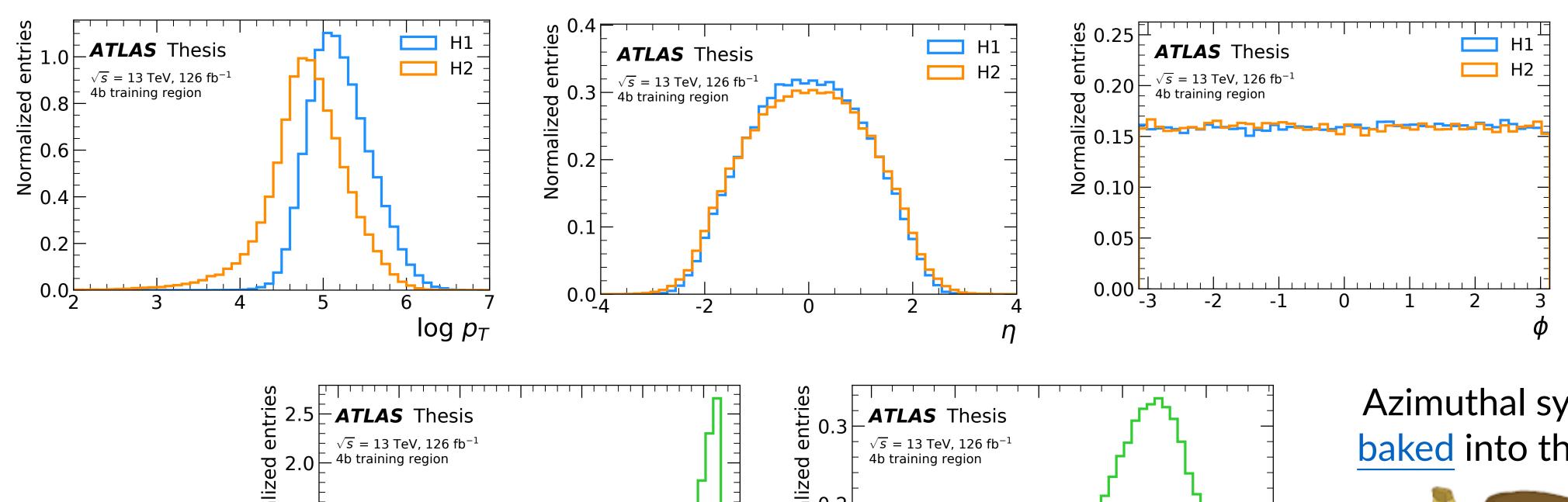


Input processing: HH -> 4b background modeling

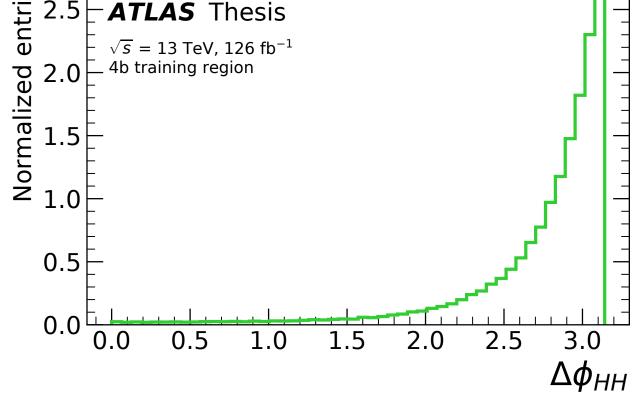


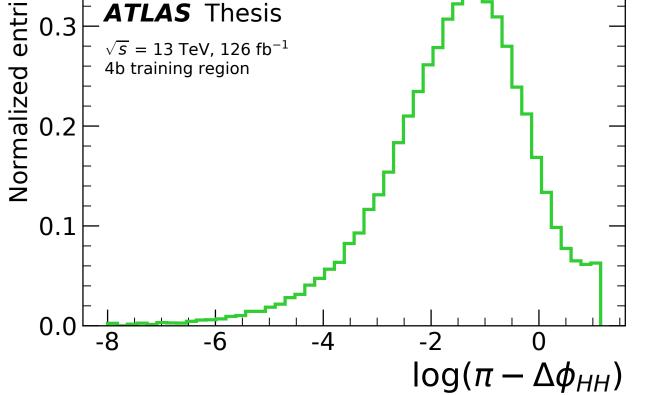


Concerned about modeling b/c no way to know that $-\pi \rightarrow \pi$.

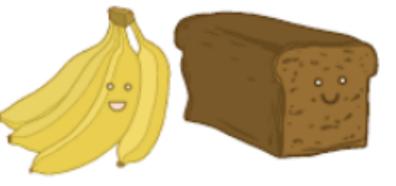


Constant $\Delta\Phi_{HH}$ - will give the same m_{HH}





Azimuthal symmetry baked into the model

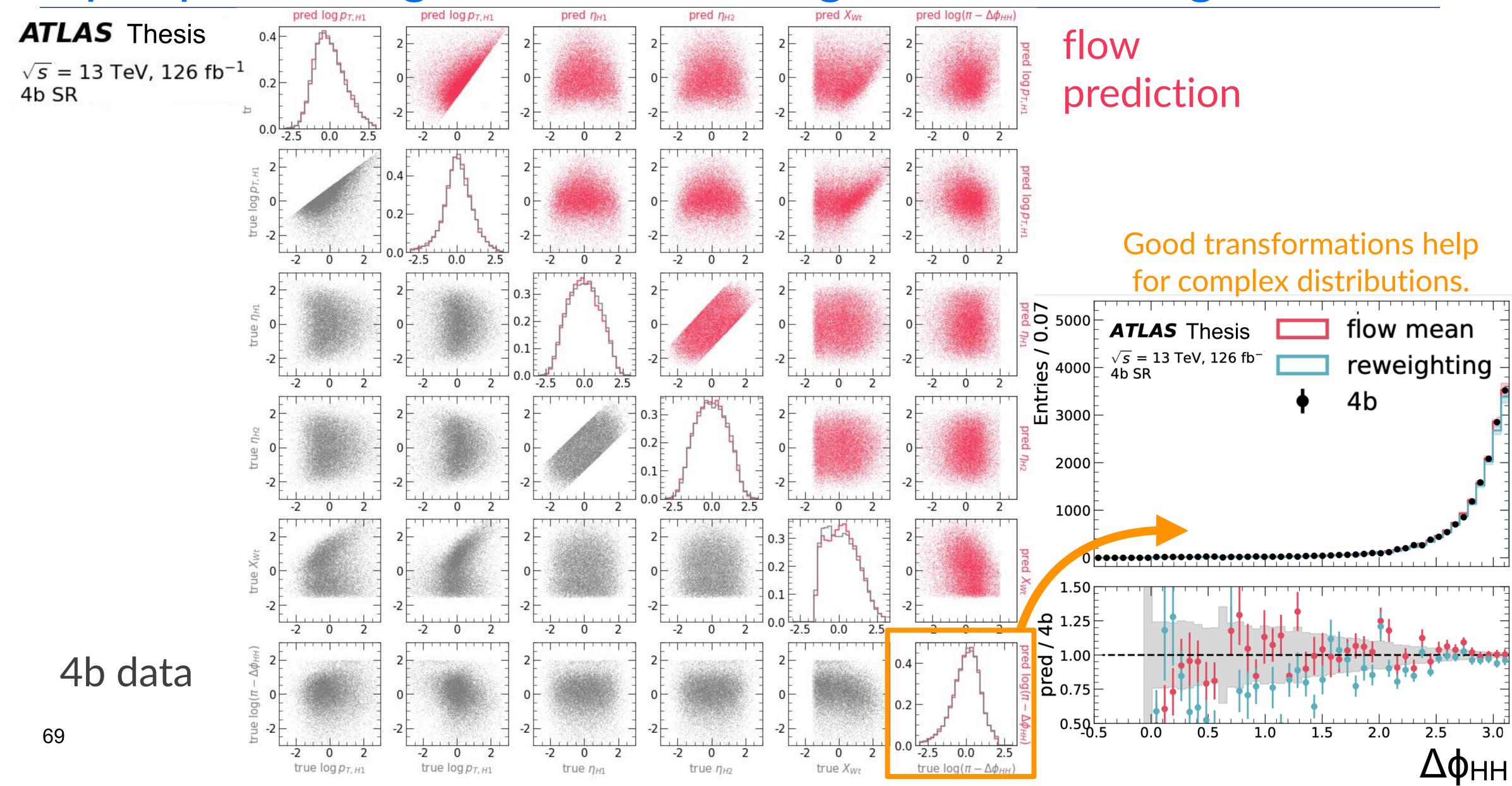


Symmetry

model with symmetry

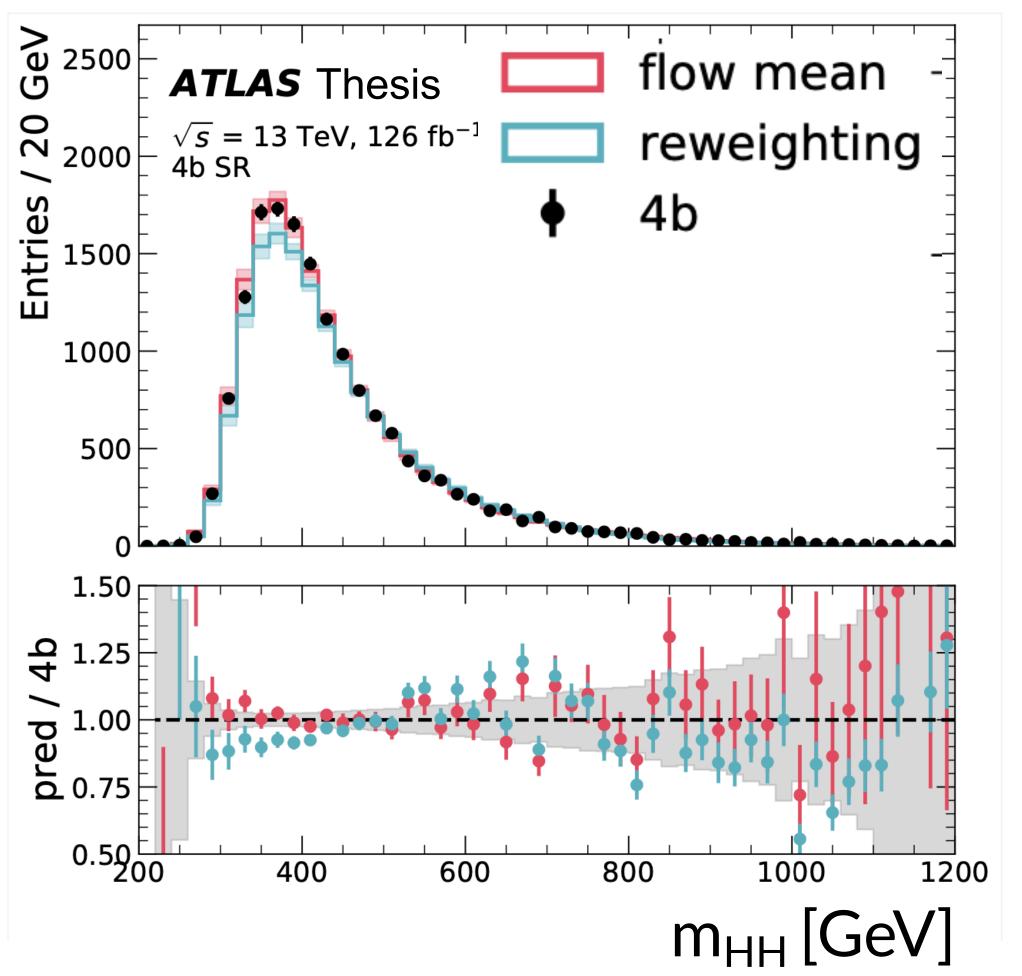


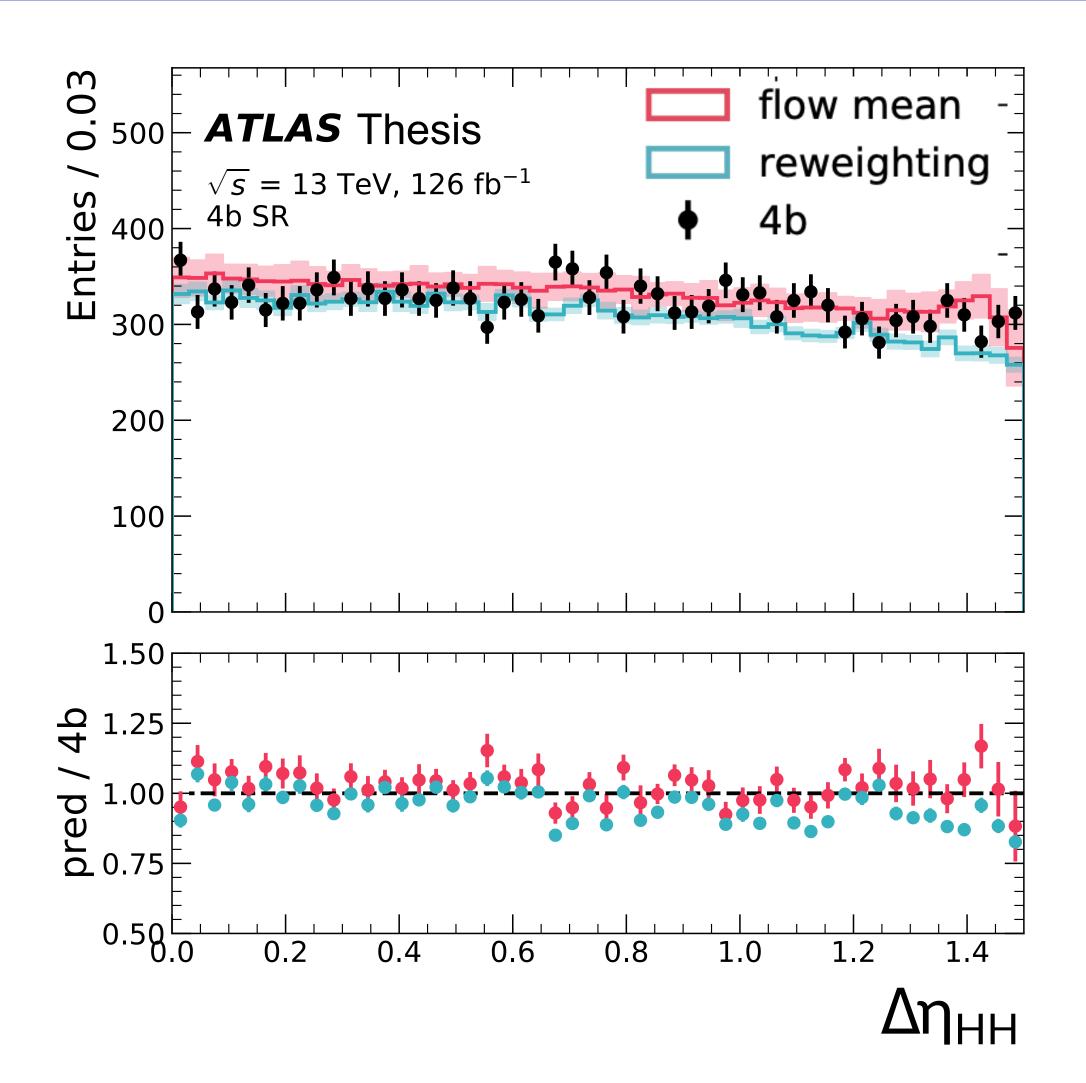
Input processing: HH -> 4b background modeling



Input processing : HH → 4b background modeling

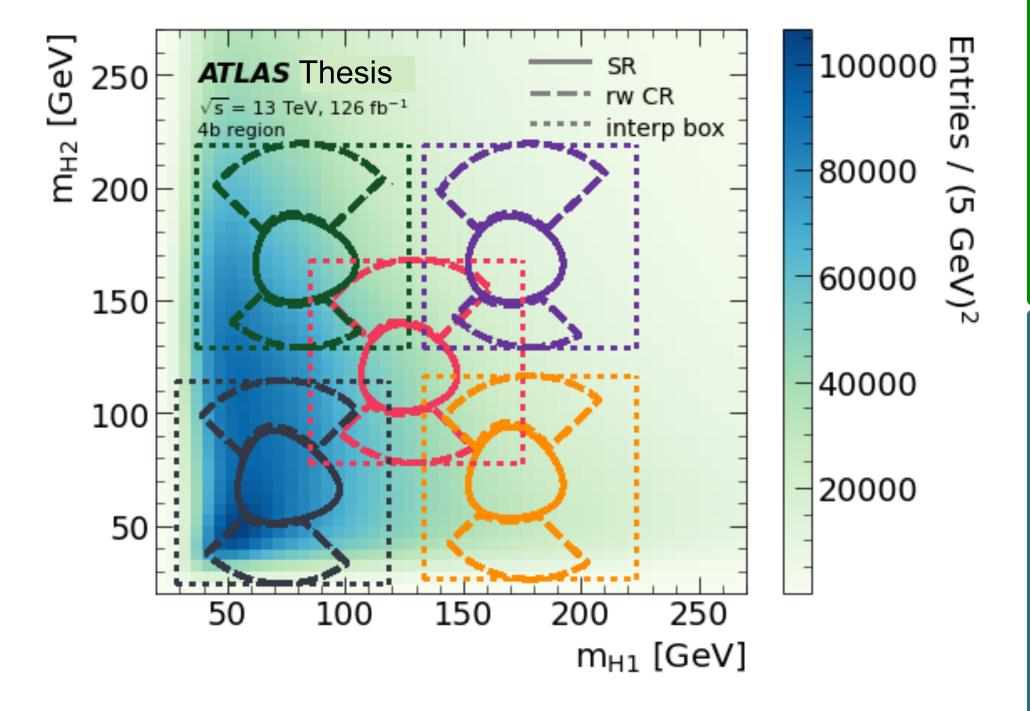
Variable transformations

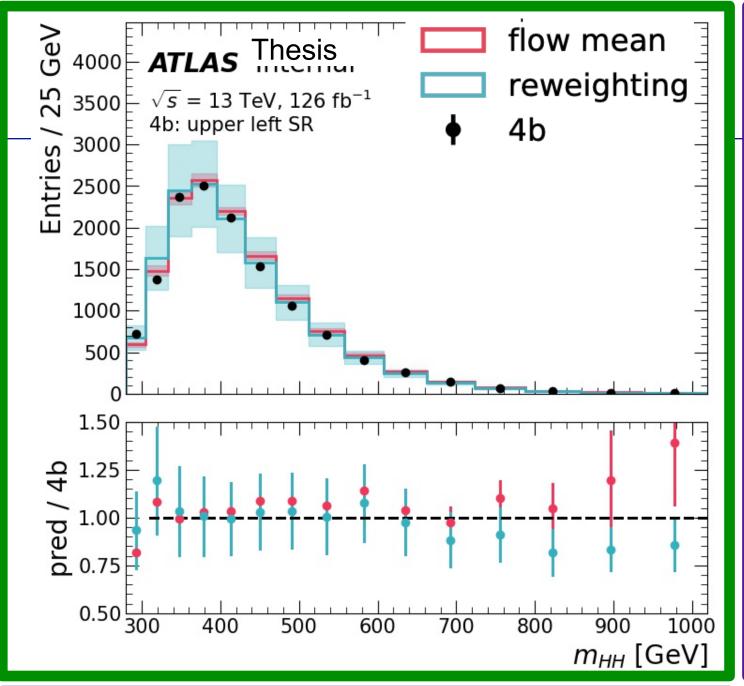


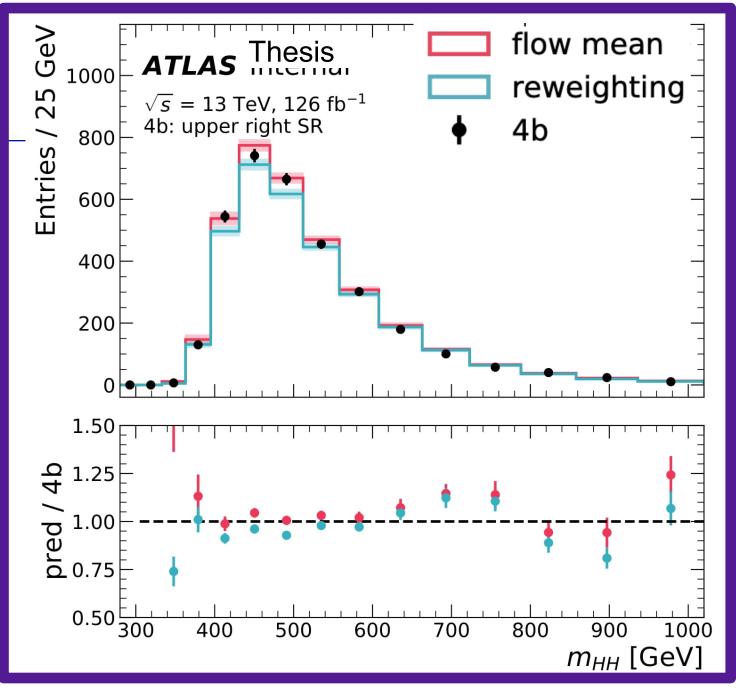


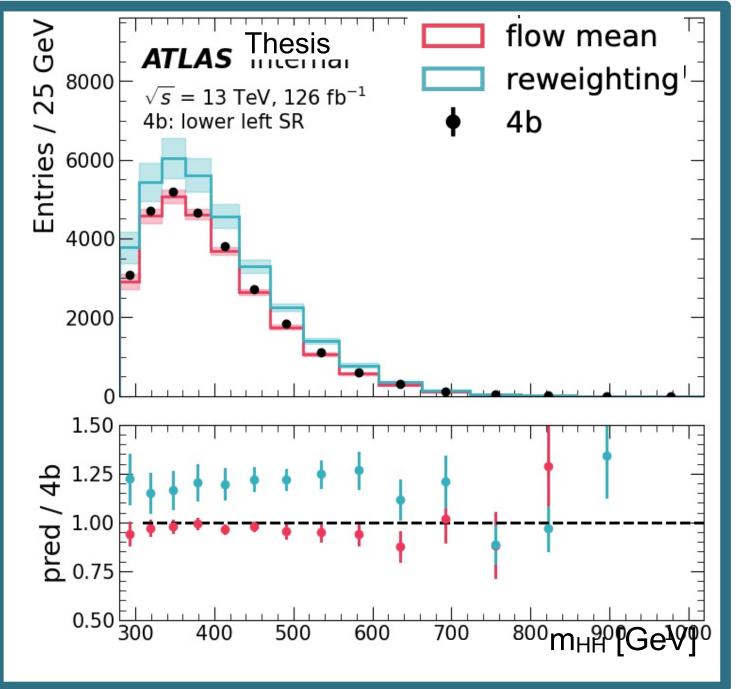


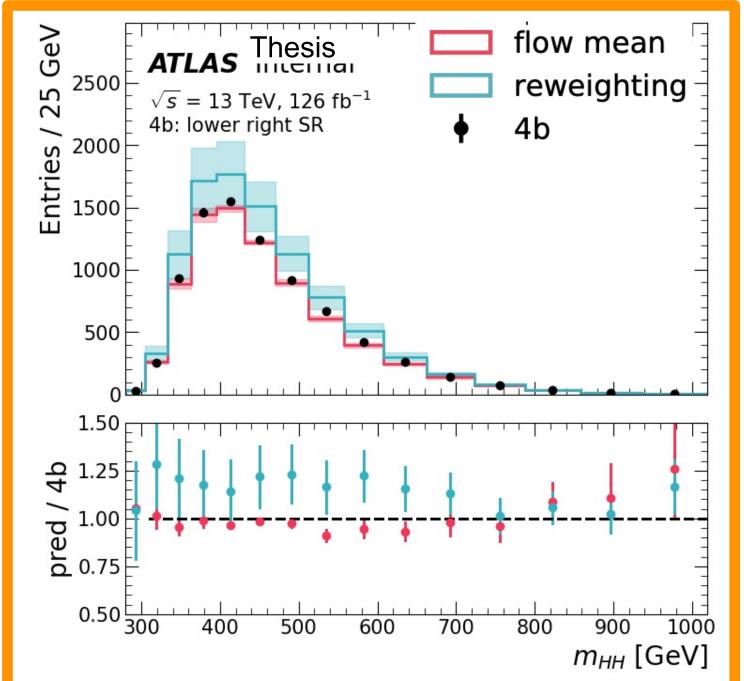
Shifted regions







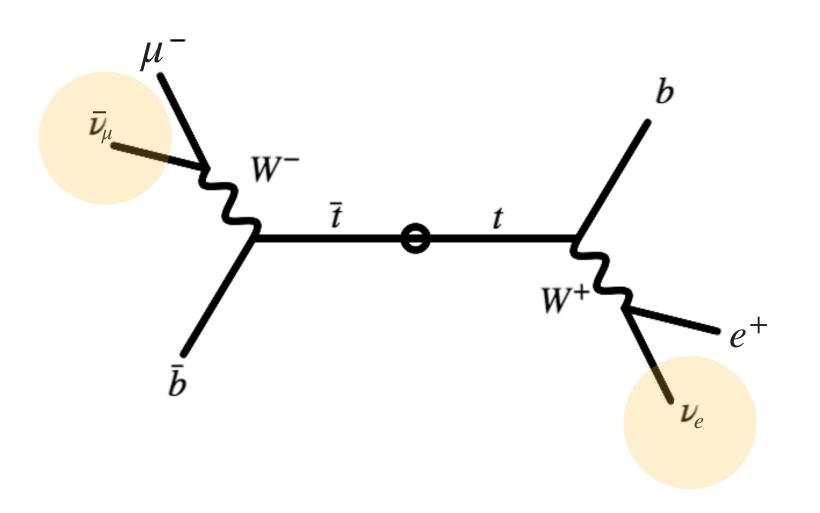




2207.00664 2307.02405

TOWN THE STATE OF THE STATE OF

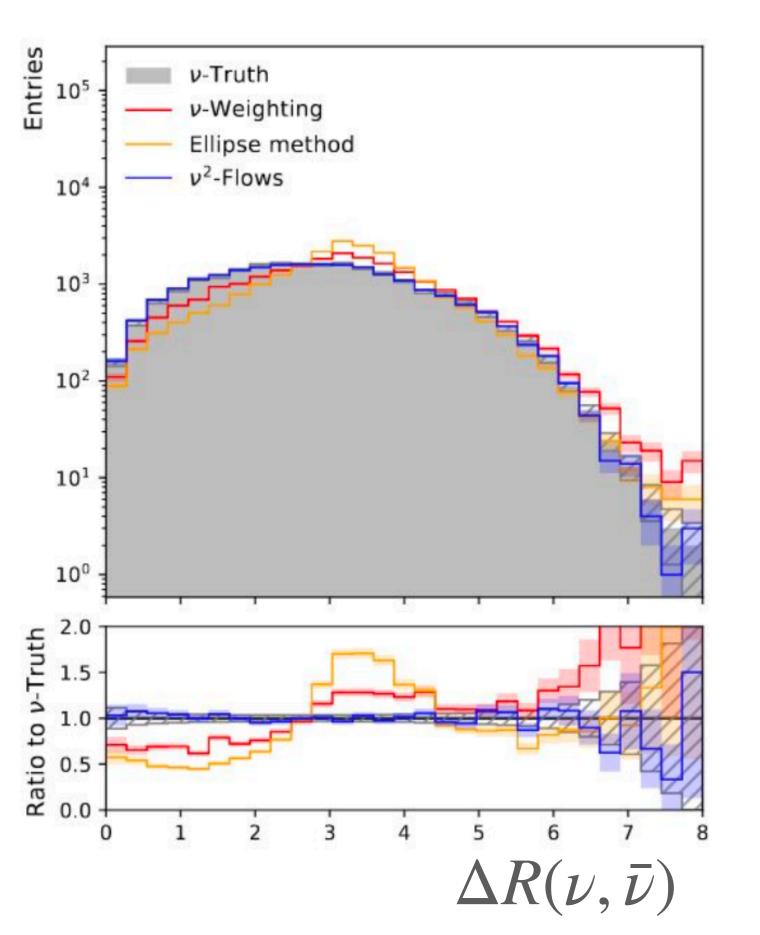
Probailistic reconstruction of two neutrinos

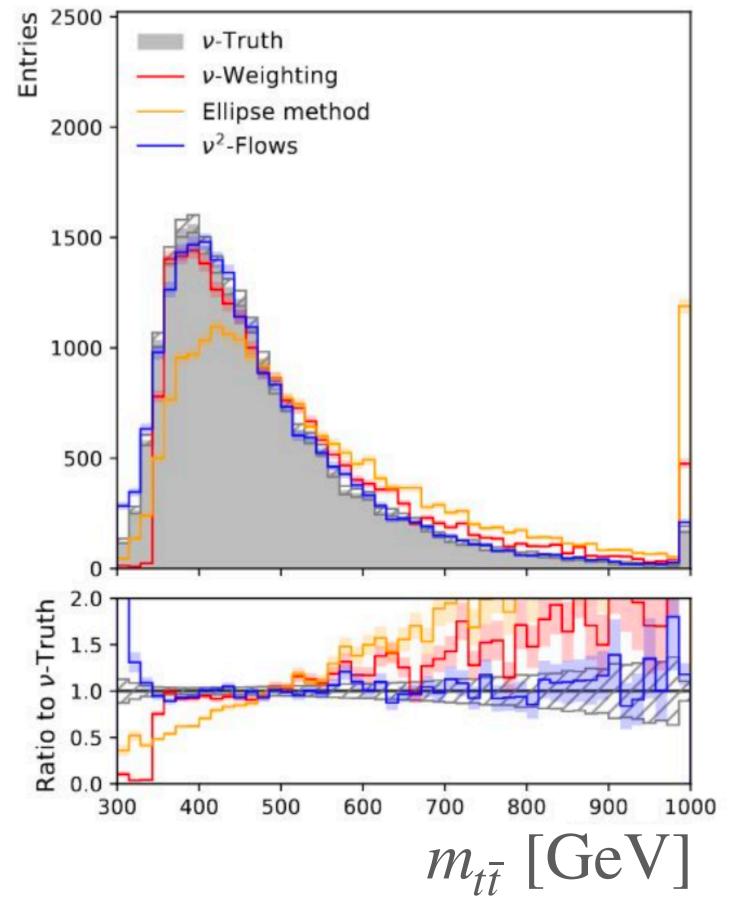


Useful for Belle 2? 1808.10567

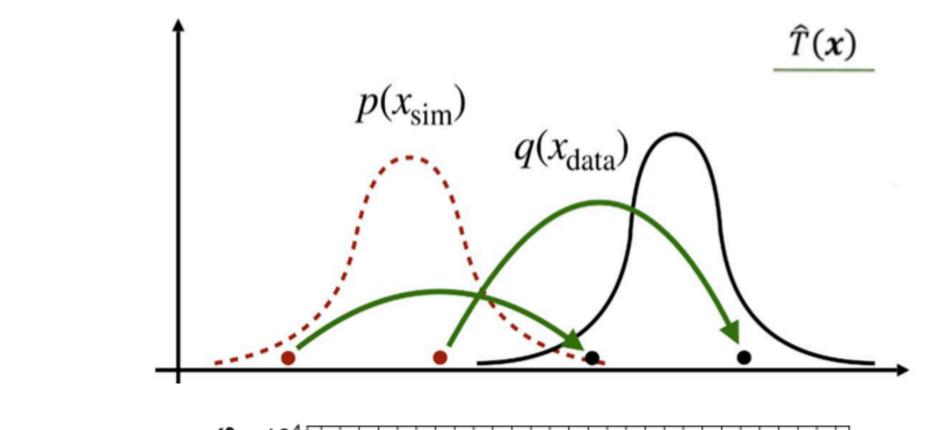
Table 5: Expected errors on several selected observables in radiative and electroweak penguin B decays. Note that 50 ab⁻¹ projections for B_s decays are not provided as we do not expect to collect such a large $\Upsilon(5S)$ data set.

Observables	Belle	Belle II	
	(2017)	$5~{ m ab}^{-1}$	$50~{ m ab^{-1}}$
$\mathcal{B}(B \to K^{*+} \nu \overline{\nu})$	$< 40 \times 10^{-6}$	25%	9%
$\mathcal{B}(B \to K^+ \nu \overline{\nu})$	$< 19 \times 10^{-6}$	30%	11%
· · · · · · · · · · · · · · · · · · ·	Λ-		





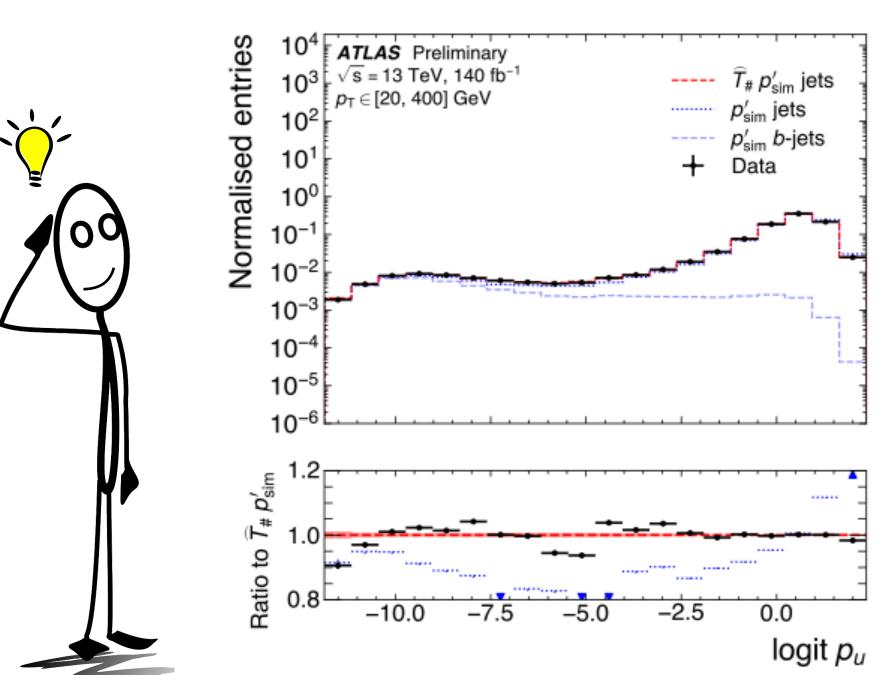
Calibration... in a differentiable way

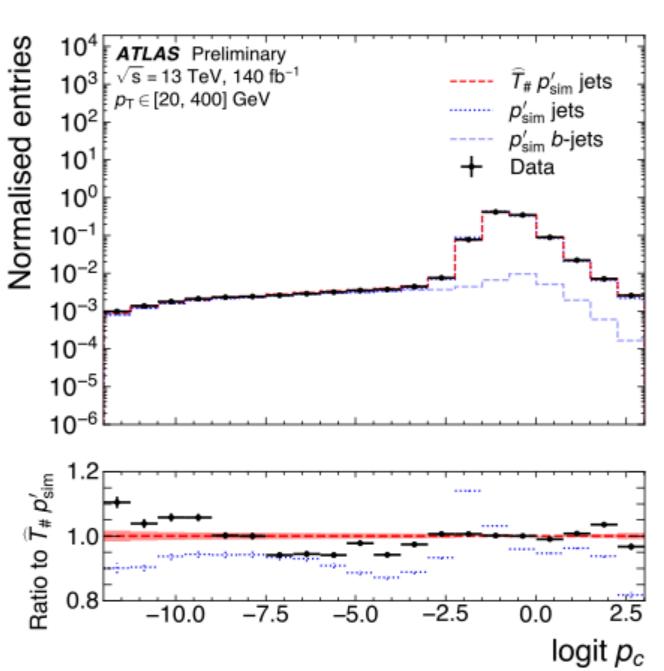


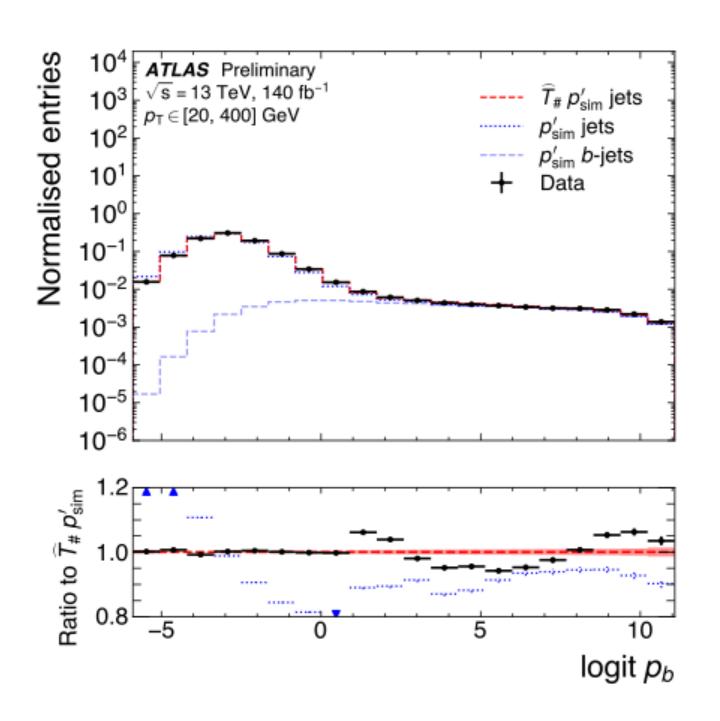
Idea: How do you have a mapping from

 $p_b^{MC} \rightarrow p_b^{data}$

Architecture: Normalizing flow with a constraint to ensure the transport map is minimal.









HEP Reconstruction

Backbone

Reconstruction

Foundation Model
Pretext tasks
(next word prediction)

<u>HEP</u>

Reconstruction
Reconstruction closure

Backbone

Reconstruction

Foundation Model

Pretext tasks

(next word prediction)

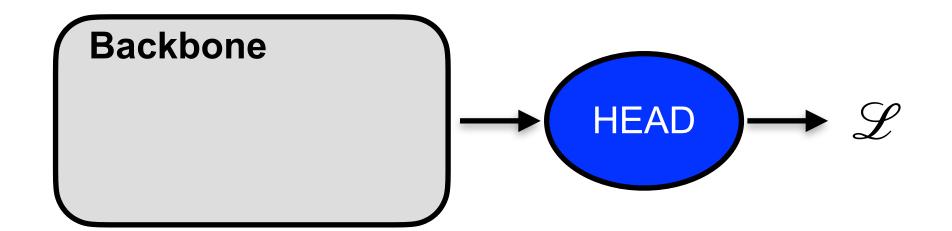
Downstream head

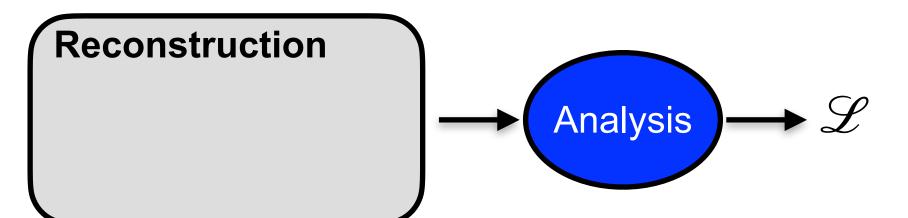
<u>HEP</u>

Reconstruction

Reconstruction closure

Analysis





Foundation Model

Pretext tasks

(next word prediction)

Downstream head

Finetuning

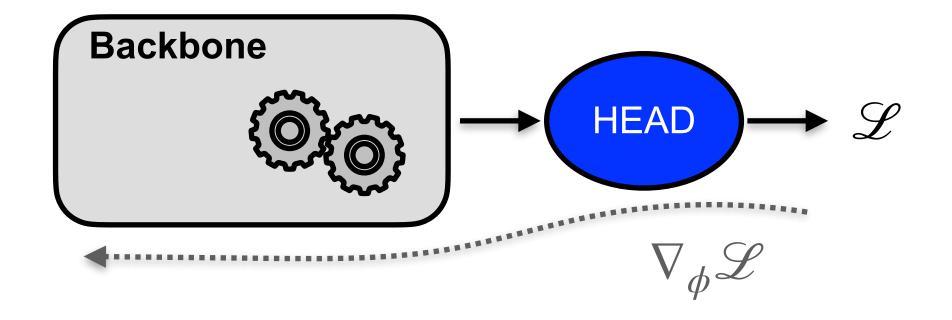
<u>HEP</u>

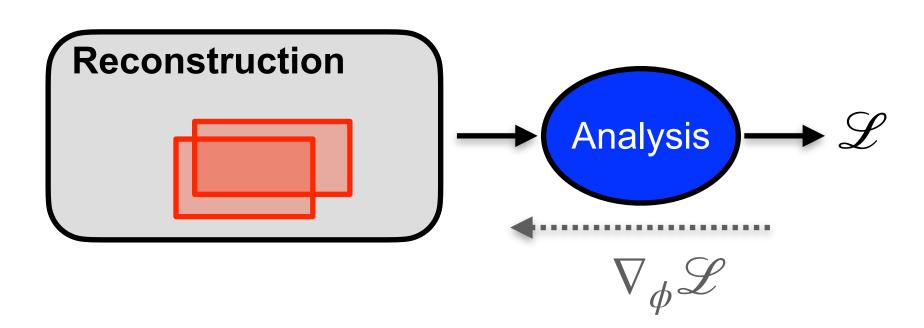
Reconstruction

Reconstruction closure

Analysis

Analysis specific reconstruction, e.g. operating points





Foundation Model

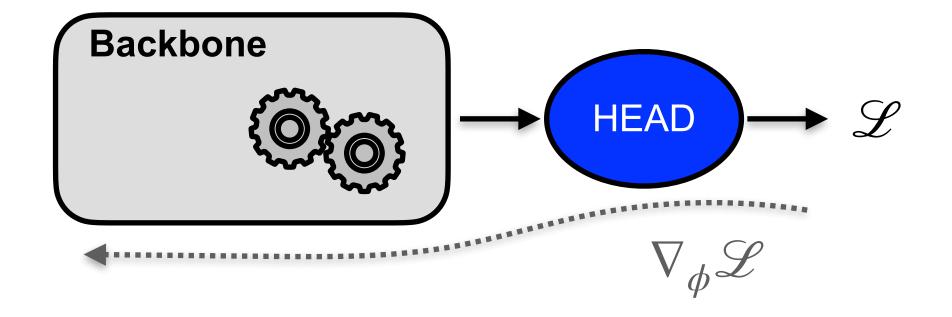
Pretext tasks

(next word prediction)

Downstream head

Finetuning

Embedding



<u>HEP</u>

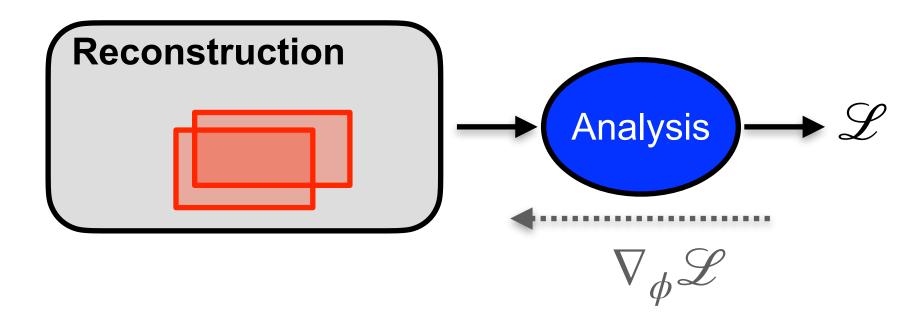
Reconstruction

Reconstruction closure

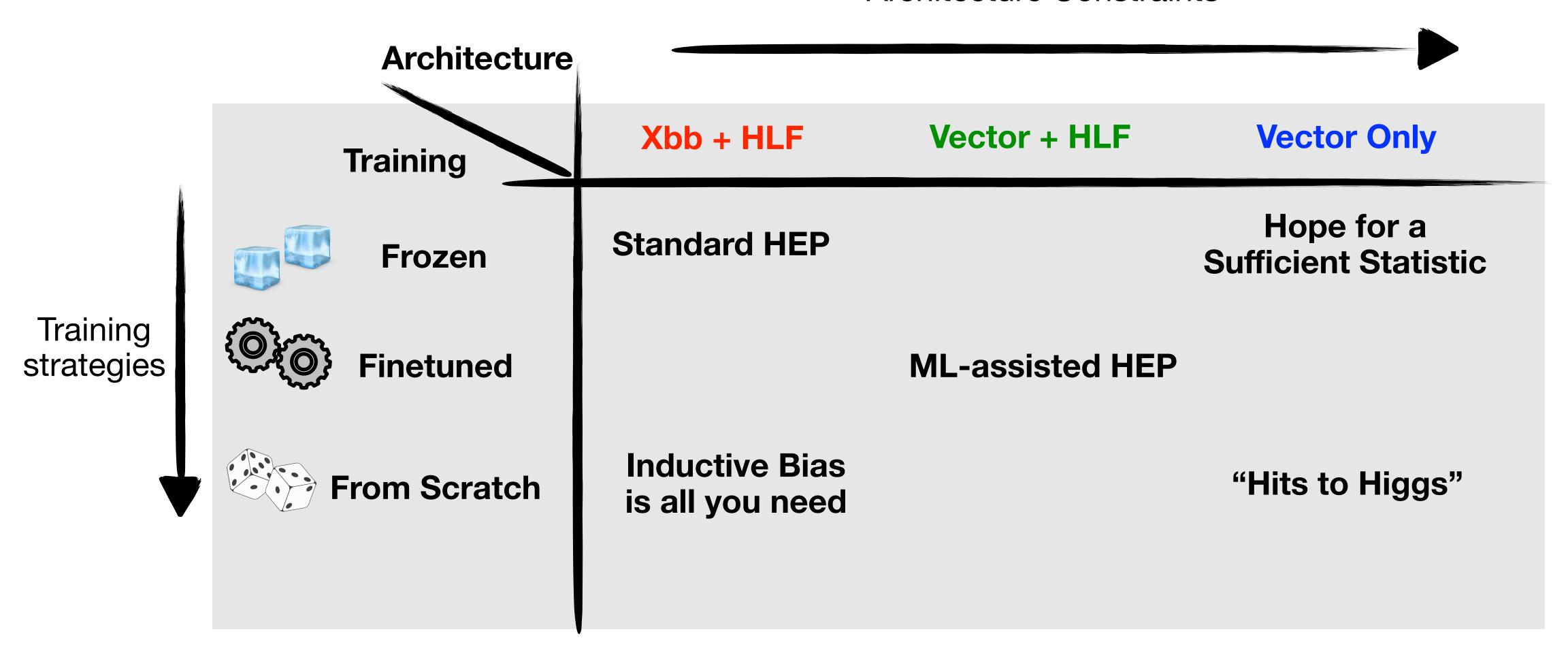
Analysis

Analysis specific reconstruction, e.g. operating points

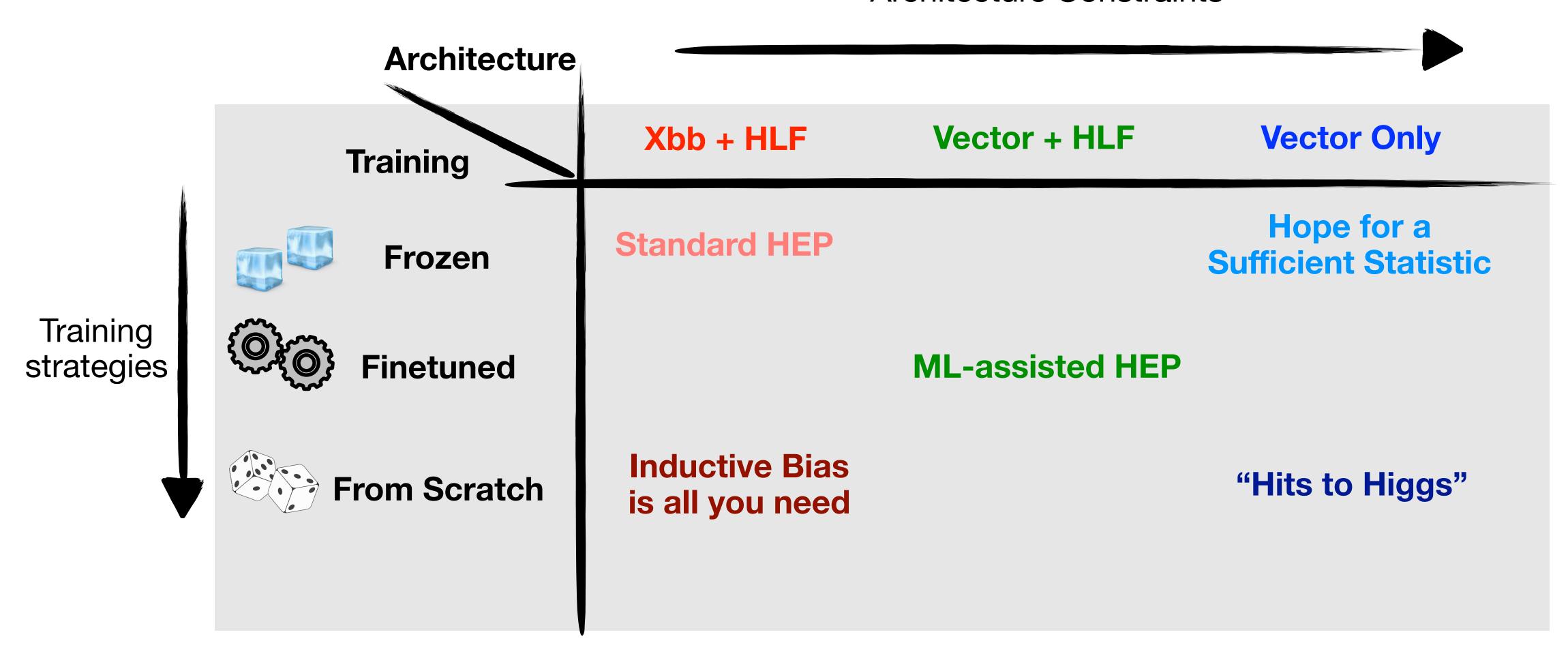
Object observables (jets)



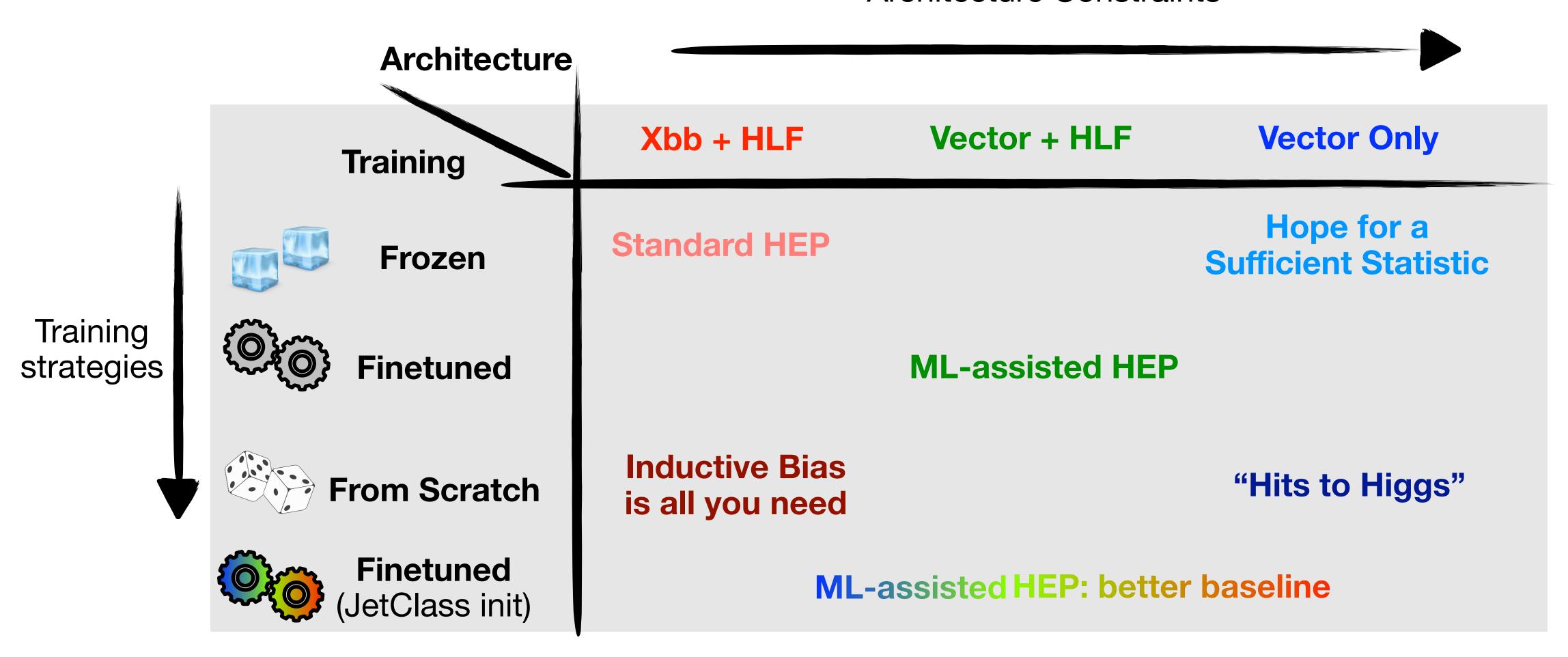
Architecture Constraints



Architecture Constraints



Architecture Constraints

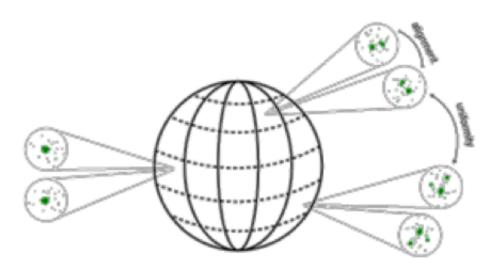


M. Kagan's slide

& A. Hallin for <u>up-to-date list</u>

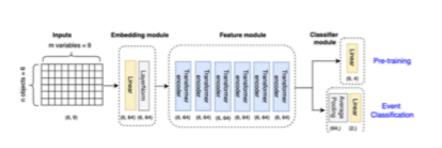
How to build this HEP foundation model

Contrastive Learning: Symmetry Augmentation



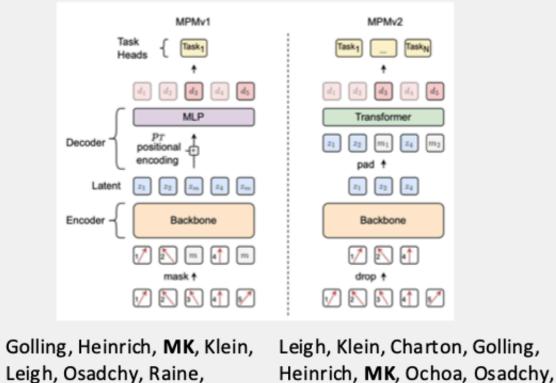
Dillon, Kasieczka, Olischlager Plehn, Sorrenson, Vogel, 2108.04253

Masked Particle Type Prediction

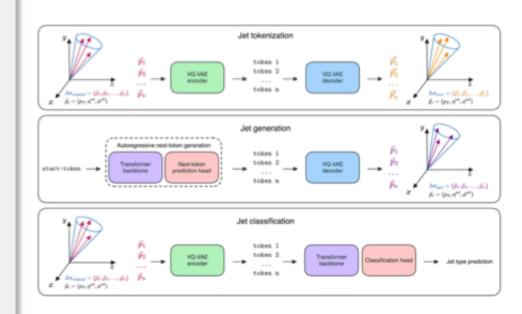


Kishimoto, Morinaga, Saito Tanaka, <u>2312.06909</u>

Masked Particle Modeling

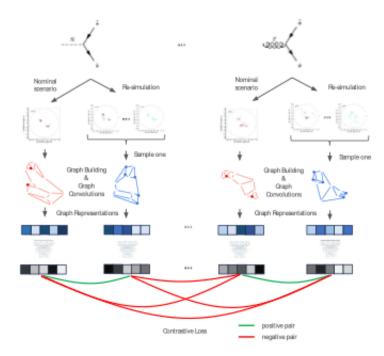


Next Token Predictoin



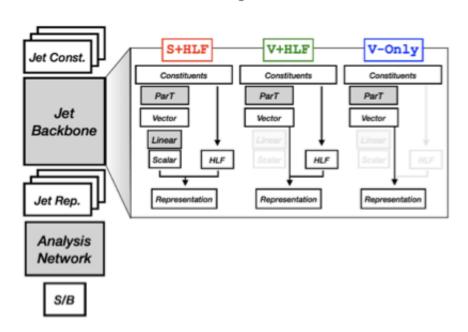
Birk, Hallin, Kasieczka, 2403.05618

Contrastive Learning: Re-Simulation



Harris, MK, Krupa, Maier, Woodward, 2403.07066

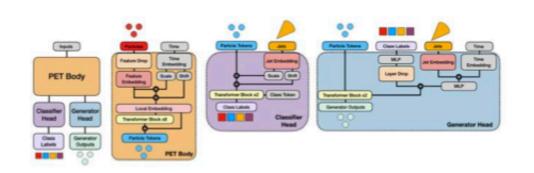
Supervised Pre-training and Joint Optimization



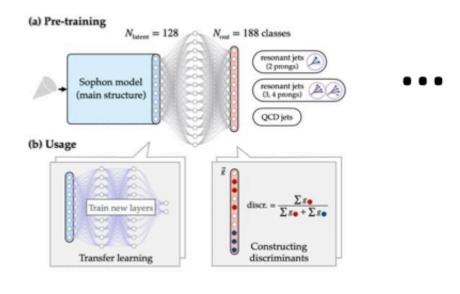
Supervised Classification and Generation

2409.12589

2401.13537



Large-Scale Fine-Grained Classification

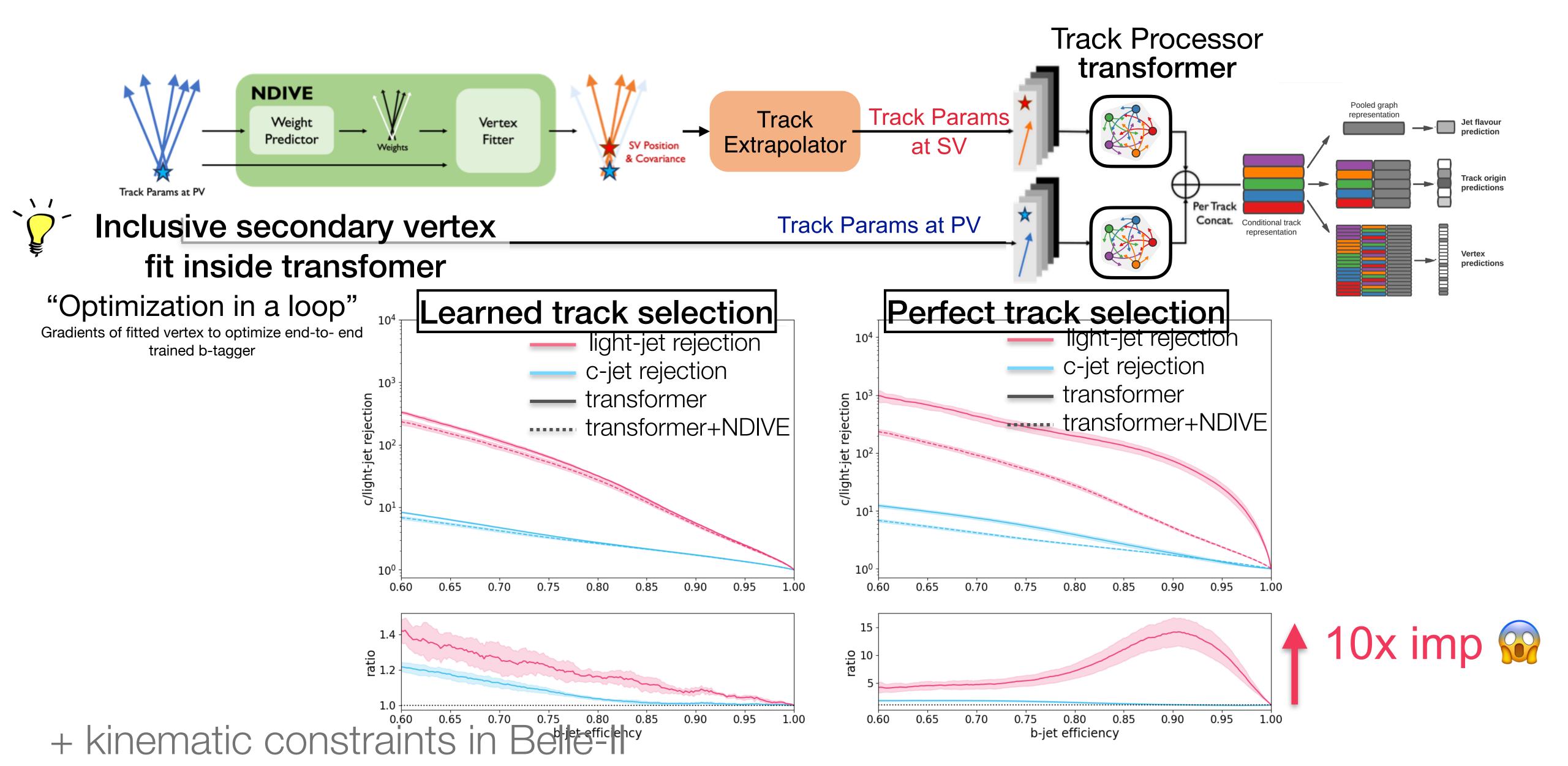


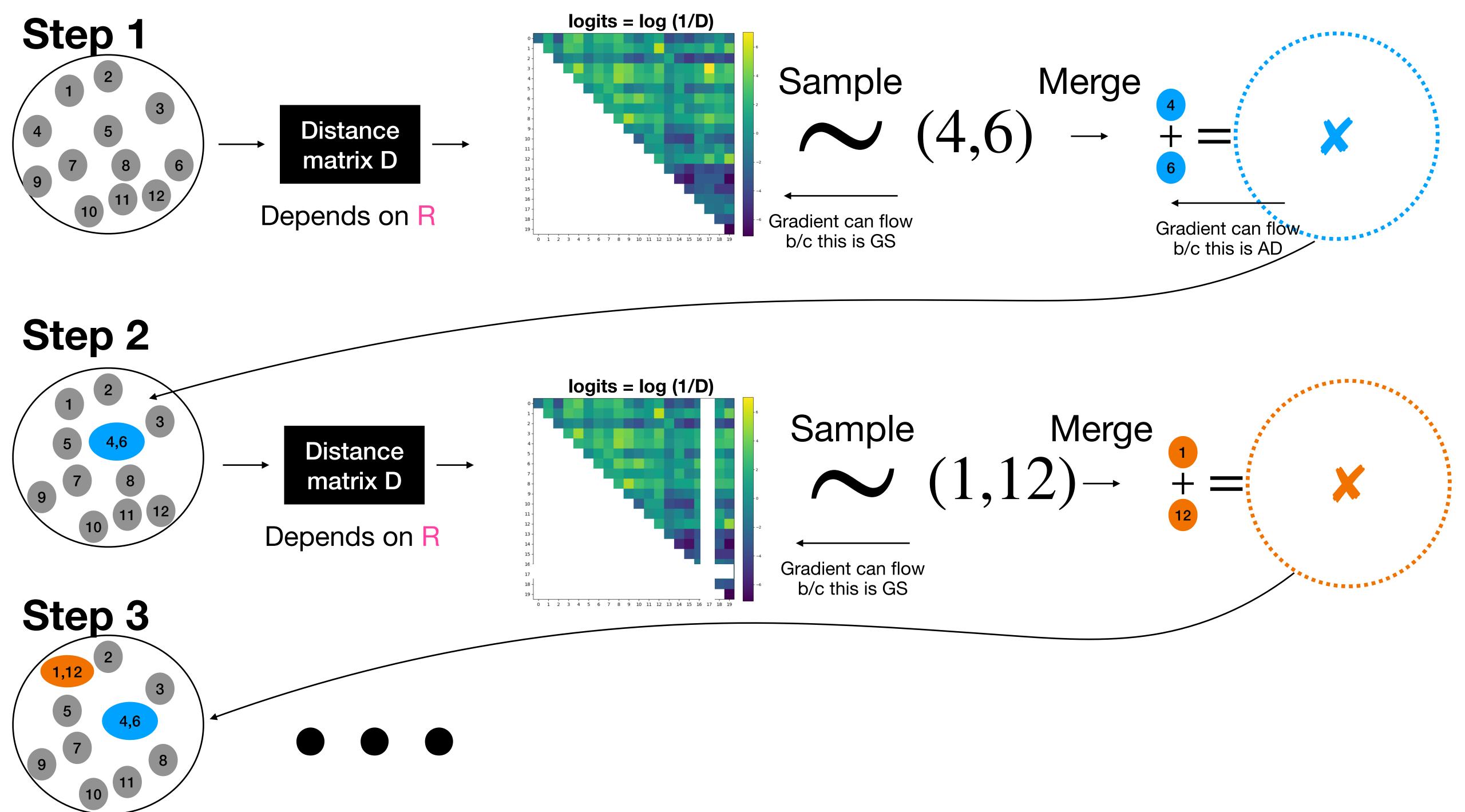
Vigl, Hartman, Heinrich, 2401.13536

Mikuni, Nachman <u>2404.16091</u>

Li, Li, et al. <u>2405.12972</u>

Neural Differentiable Vertex Fitter





Differentiating clustering

Step 1: Interpret clustering decision probabilistically

Step 2: Gradient with score based estimate

As demonstrated by M. Kagan and L. Heinrich for particle interactions in material: 2308.16680

In general:

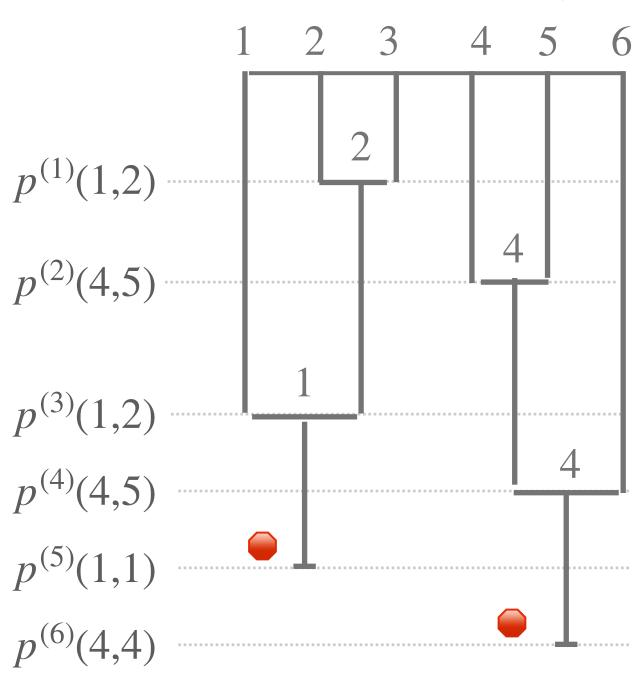
$$\nabla_{\theta} \mathbb{E}_{x \sim p(\theta)}[f(x)] = \mathbb{E}_{x \sim p(\theta)}[f(x) \nabla_{\theta} \log p(\theta)]$$

For jet clustering:

$$\nabla_R \mathbb{E}_{p(R)}[m_H] = \mathbb{E}_{p(R)}[m_H \nabla_R \log p(R)]$$

Ex: sample i through event w/ 6 p'cles

Particles "pseudo-jets"



$$m_{H}^{(i)} \left(\begin{array}{c} \nabla p^{(1)}(1,2) + \nabla p^{(2)}(4,5) + \nabla p^{(3)}(1,2) \\ + \nabla p^{(4)}(4,5) + \nabla p^{(5)}(1,1) + \nabla p^{(6)}(4,4) \end{array} \right)$$