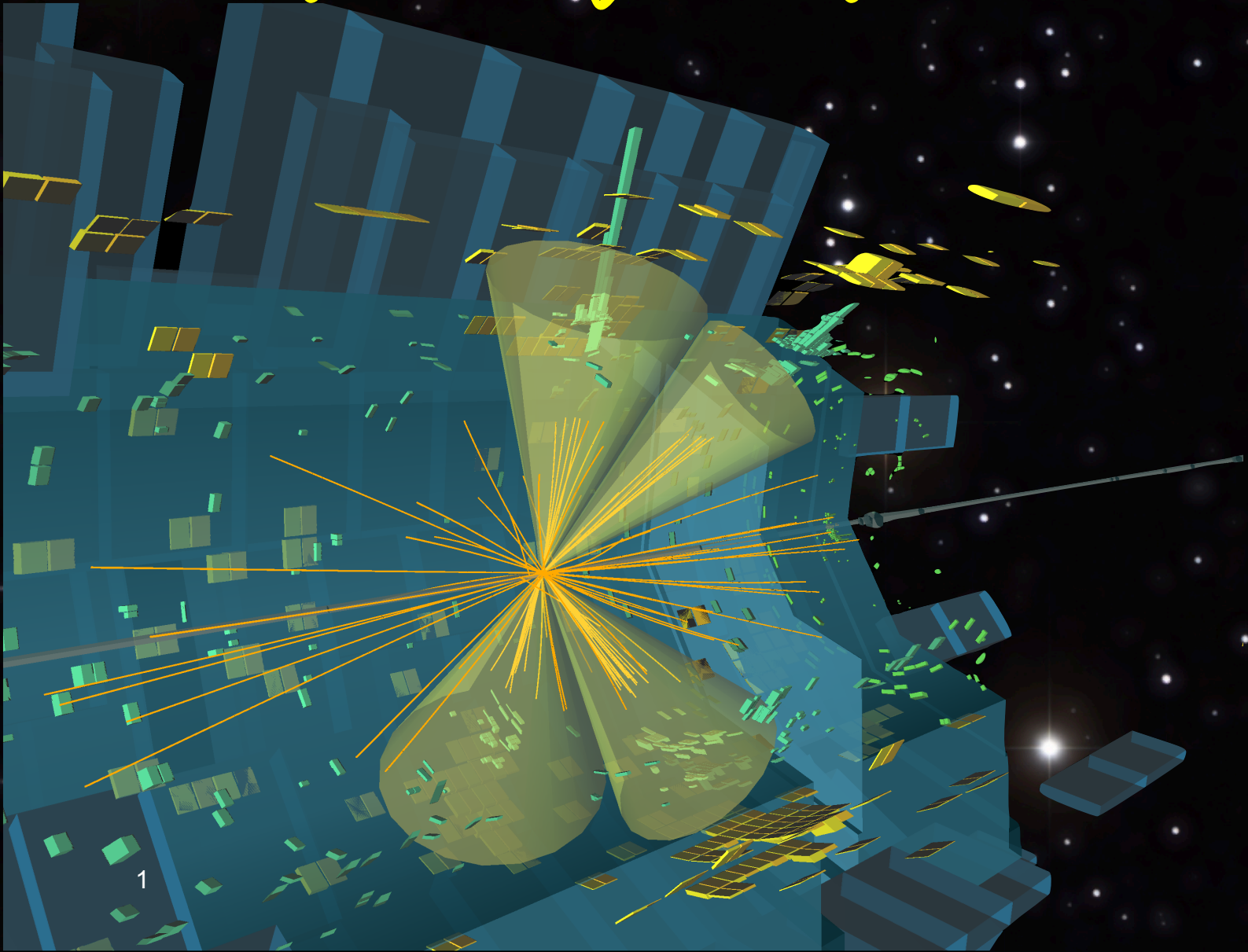


Transforming Jet flavour tagging on ATLAS



Nicole Hartman
nicole.hartman@tum.de

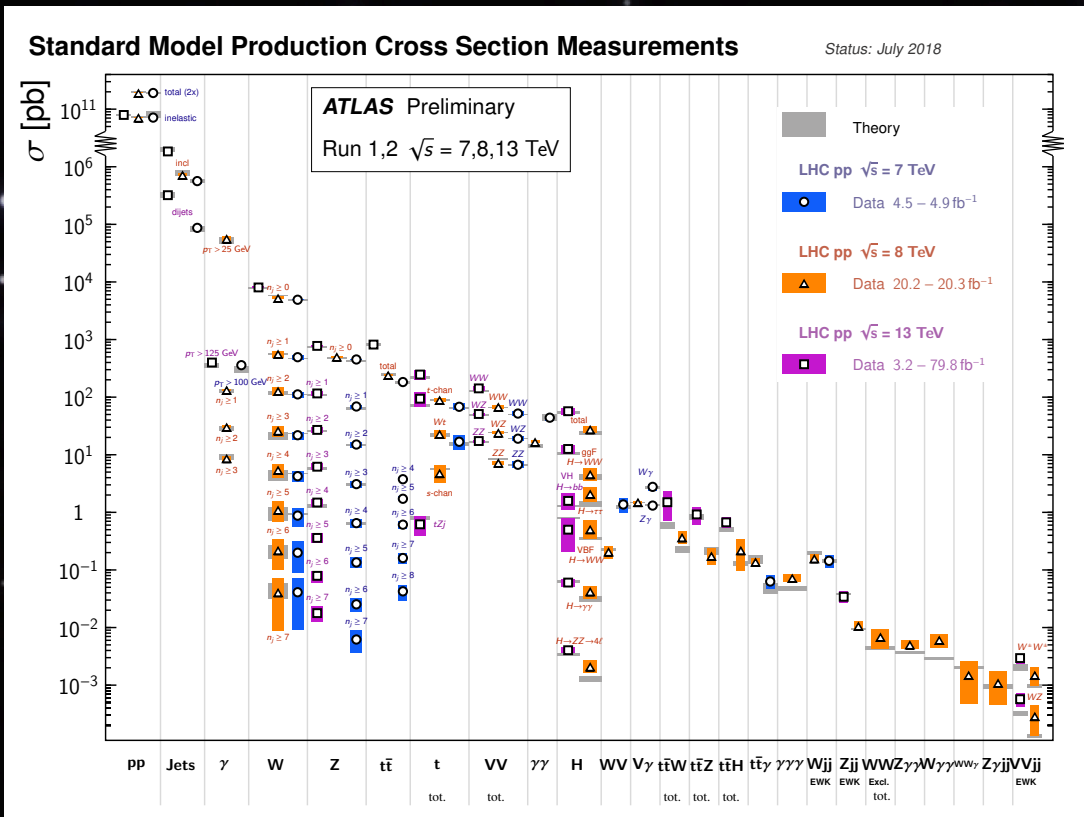
Siegen CPPS seminar
26th Nov 2024



Run: 362619
Event: 524614423
2018-10-03 08:06:34 CEST

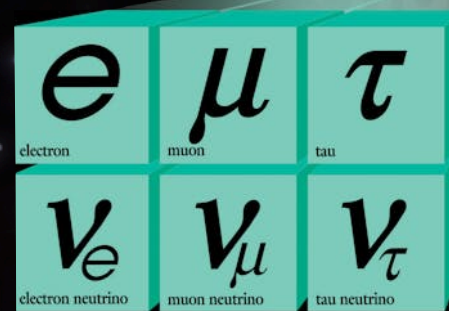
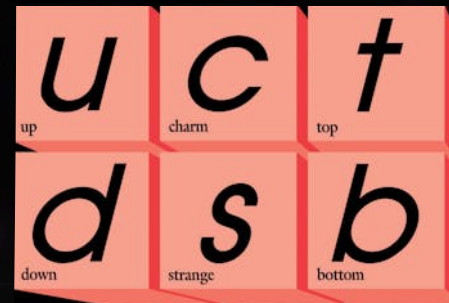


$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$



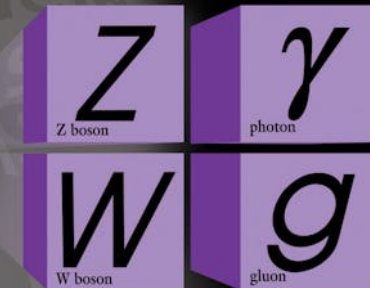
ATLAS SM summary plots

Quarks



Leptons

Forces

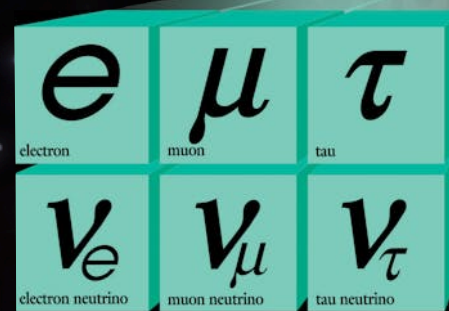
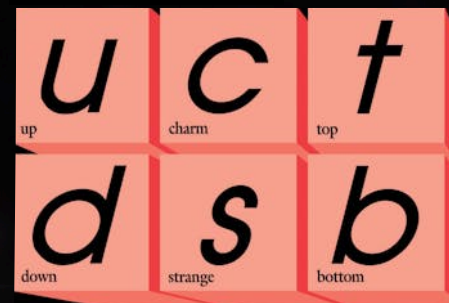


$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$

$$\begin{pmatrix} |V_{ud}| & |V_{us}| & |V_{ub}| \\ |V_{cd}| & |V_{cs}| & |V_{cb}| \\ |V_{td}| & |V_{ts}| & |V_{tb}| \end{pmatrix}$$

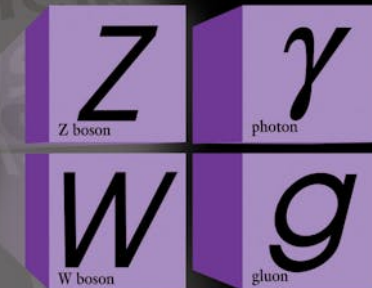
$$= \begin{pmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{pmatrix}$$

Quarks

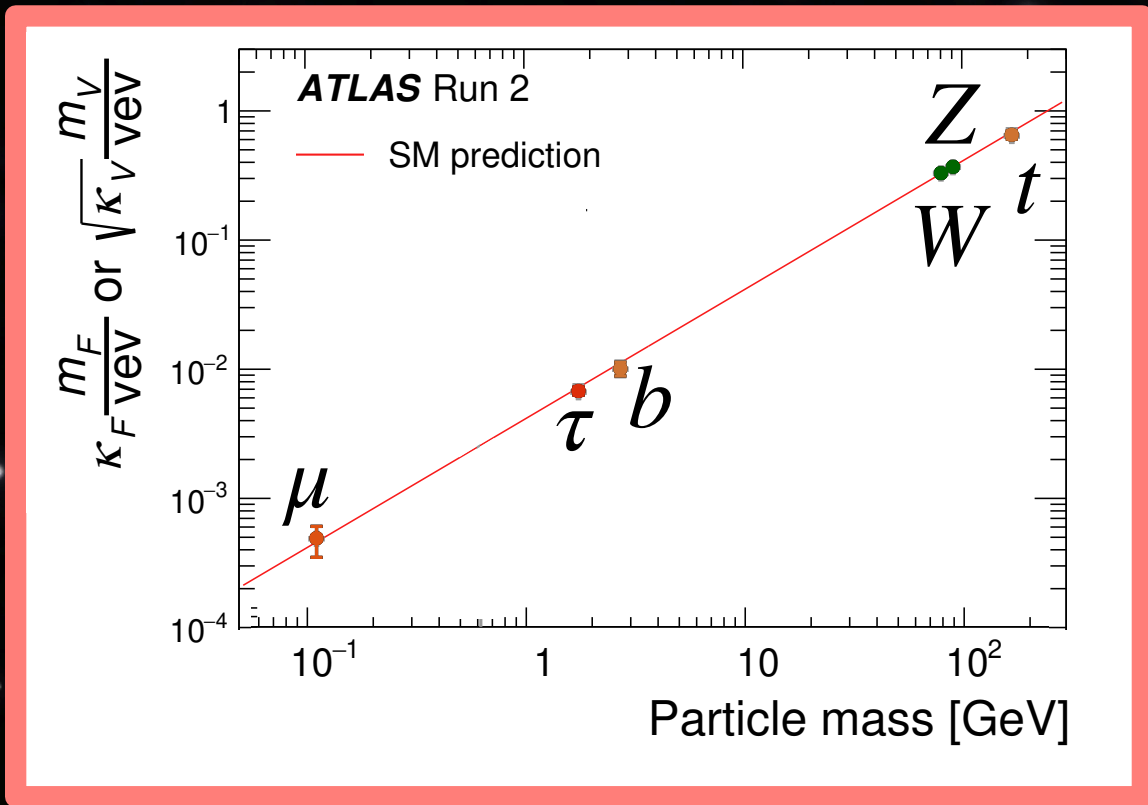


Leptons

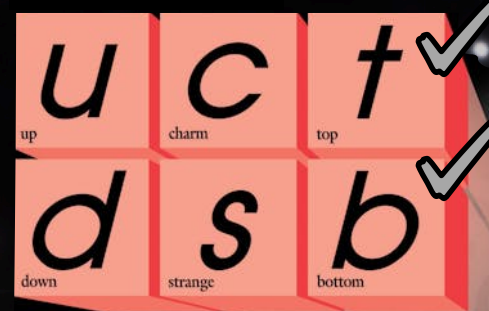
Forces



$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$



Quarks

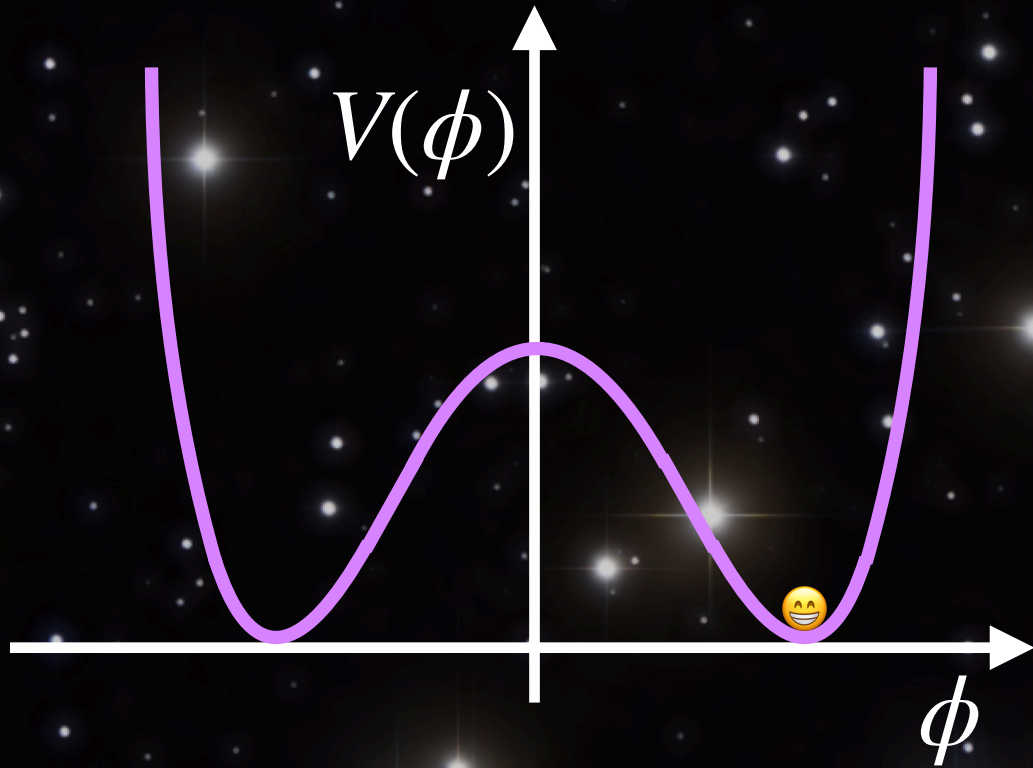


Leptons

Forces



$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$



Quarks

u up	c charm	t top
d down	s strange	b bottom

e electron	μ muon	τ tau
ν_e electron neutrino	ν_{μ} muon neutrino	ν_{τ} tau neutrino

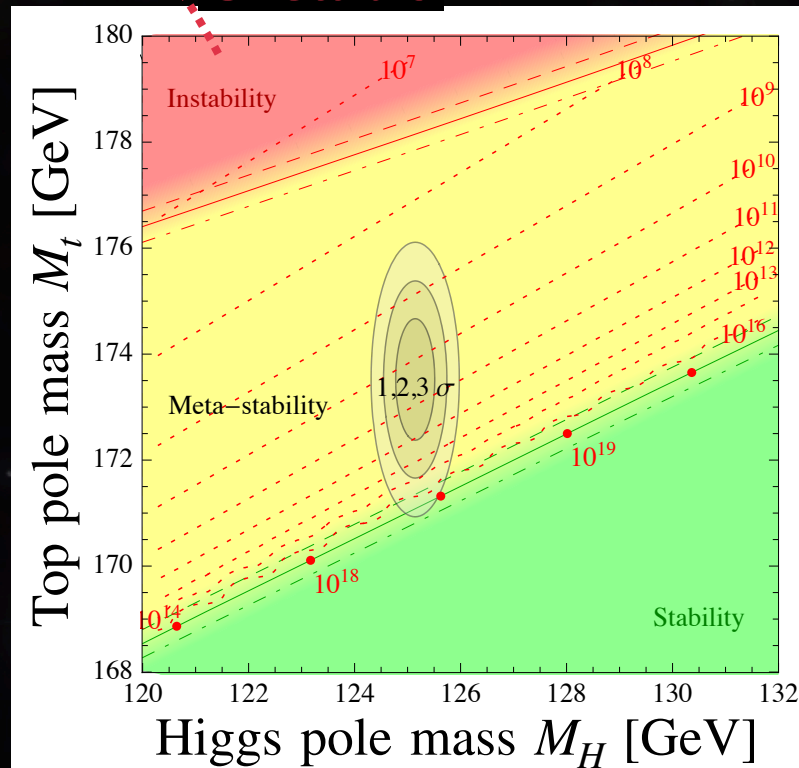
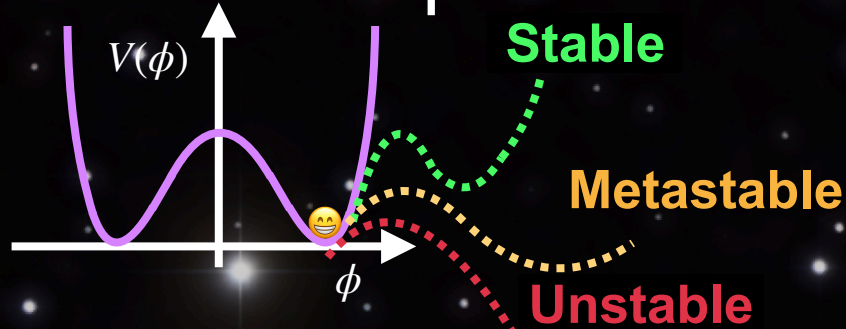
Leptons

Forces

Z Z boson	γ photon
W W boson	g gluon

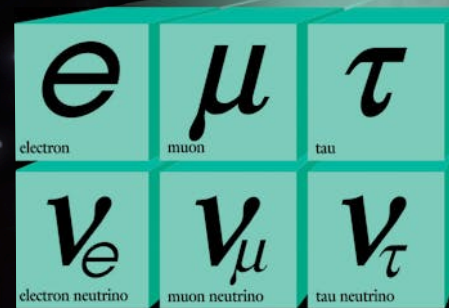
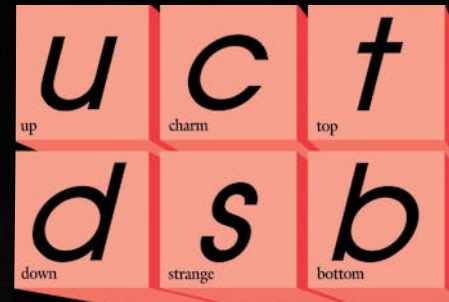


$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$



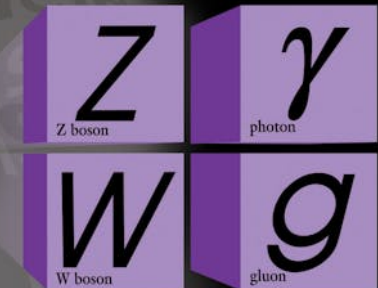
Physics 8, 108
1307.3536

Quarks



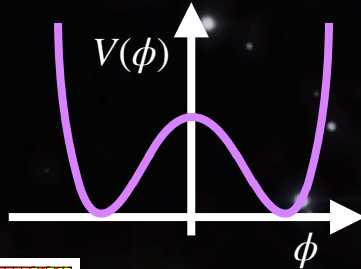
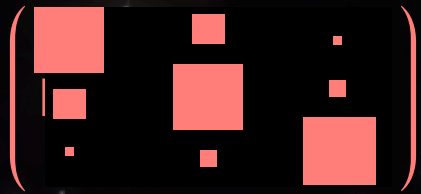
Leptons

Forces



$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$

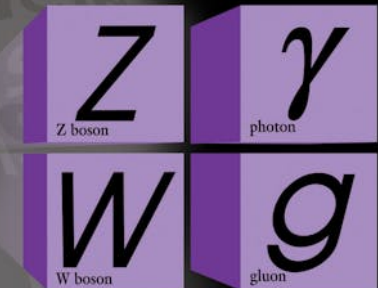
Key questions:



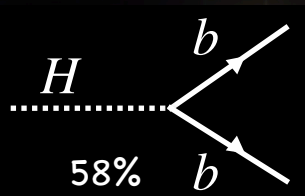
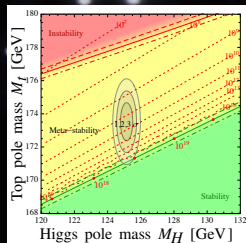
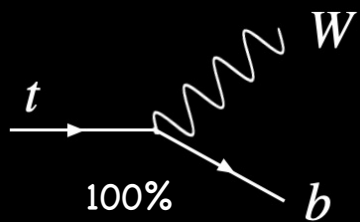
Quarks



Forces

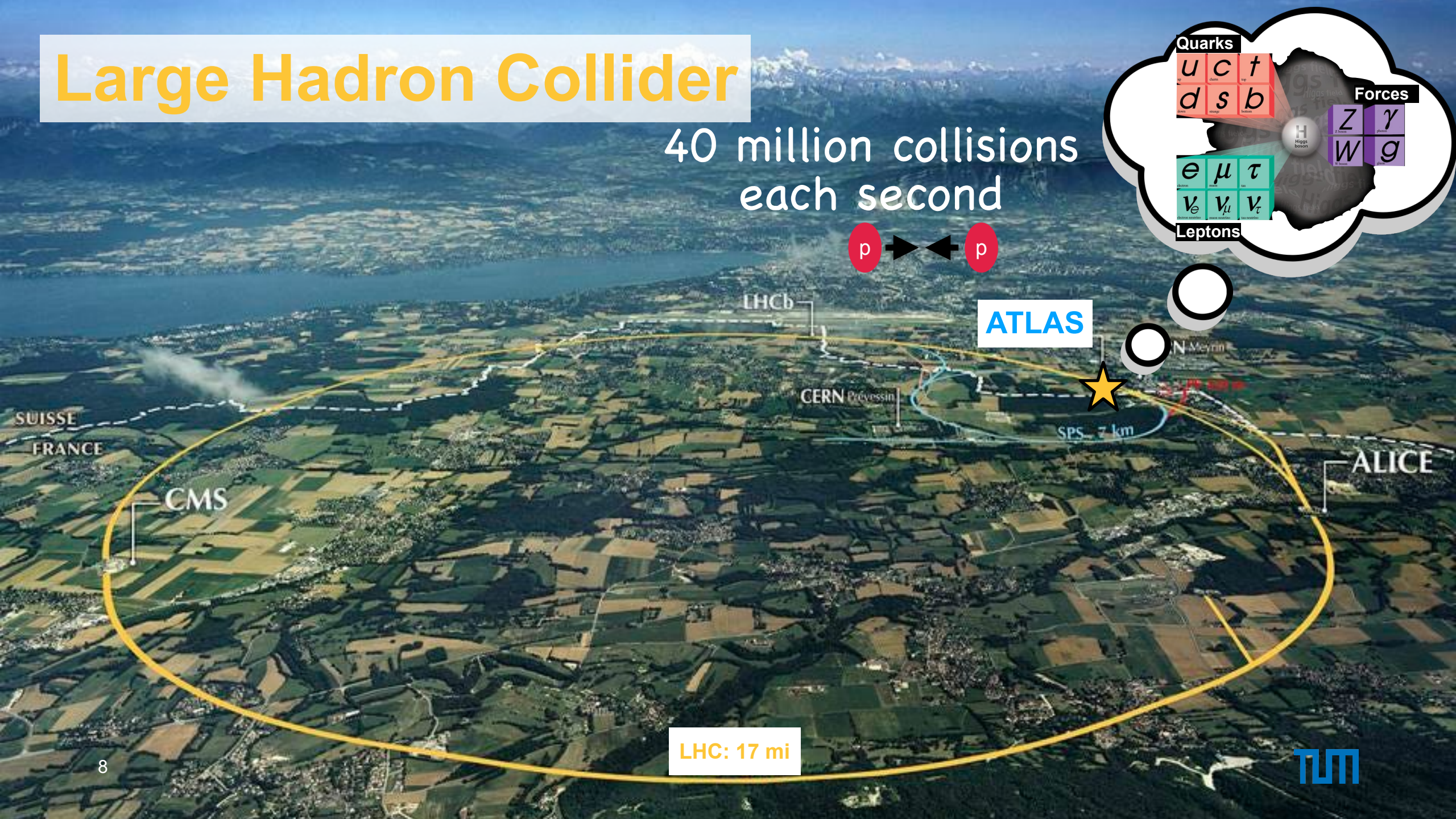
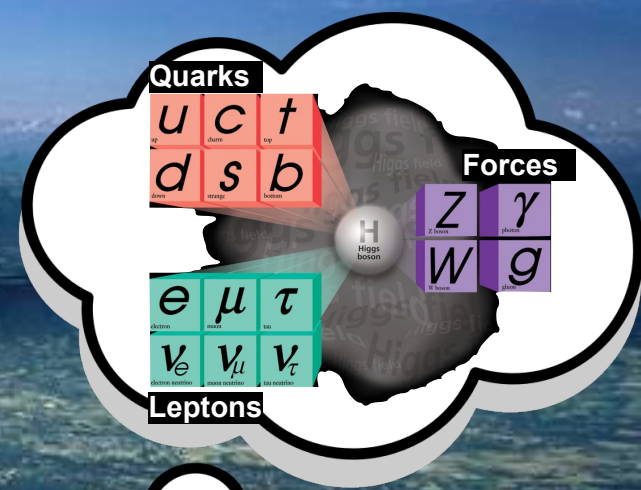


Leptons



Large Hadron Collider

40 million collisions
each second



ATLAS

LHCb

CERN Provenance

SPS, 7 km

ALICE

CMS

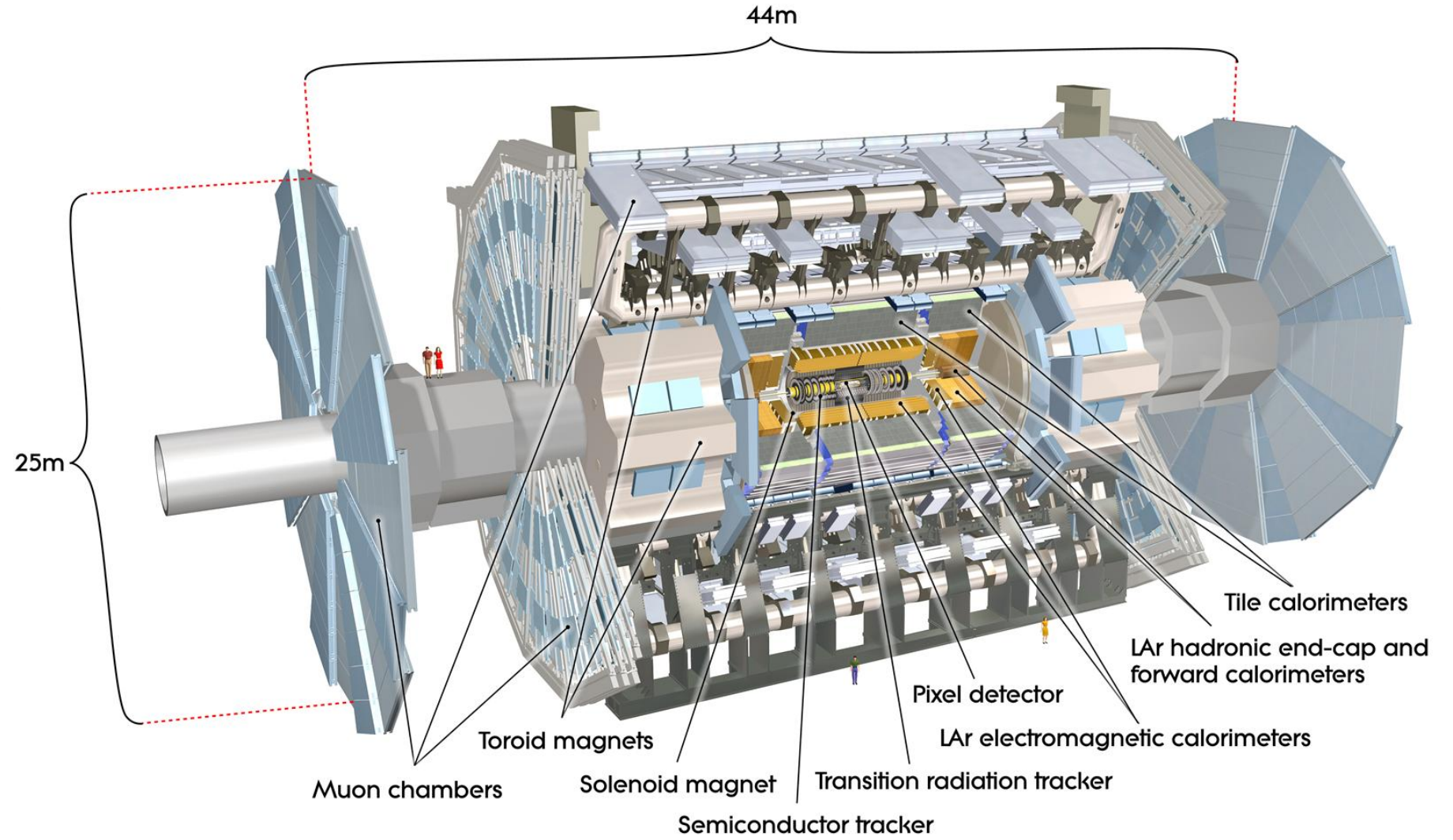
SUISSE
FRANCE

LHC: 17 mi

ATLAS



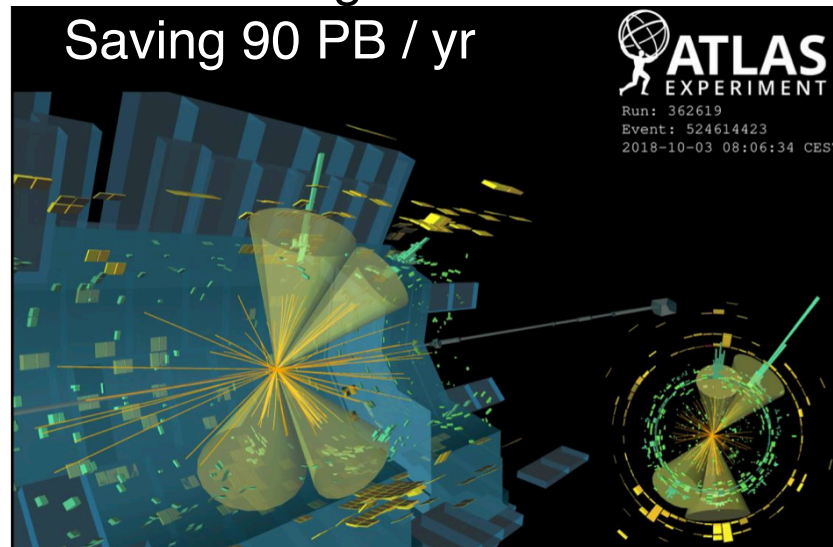
100 million
read out
channels!



What's exciting about data analysis now?

Increasingly large datasets!!

LHC : Collecting PB/ sec

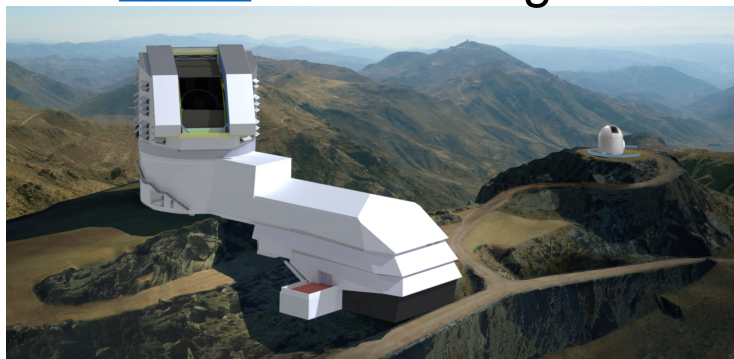


Saving 90 PB / yr

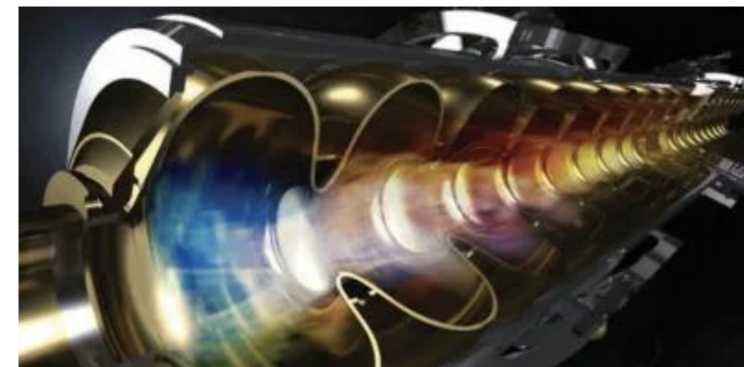
SKA : TB / s



LSST : 20 TB / night



LCLS : Tb / s



What's exciting about methods now?

The New York Times

Physicists Find Elusive Particle Seen as Key to Universe

Share full article | 122



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](#)

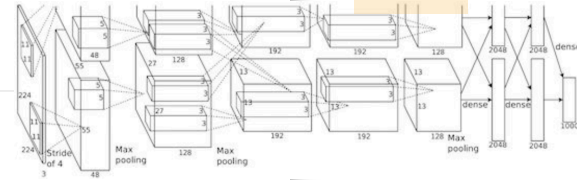


Physics, 2014


The New York Times

Scientists See Promise in Deep-Learning Programs

Share full article



“Alex Net”
Krizhevsky, Sutskever, Hinton, 2012



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

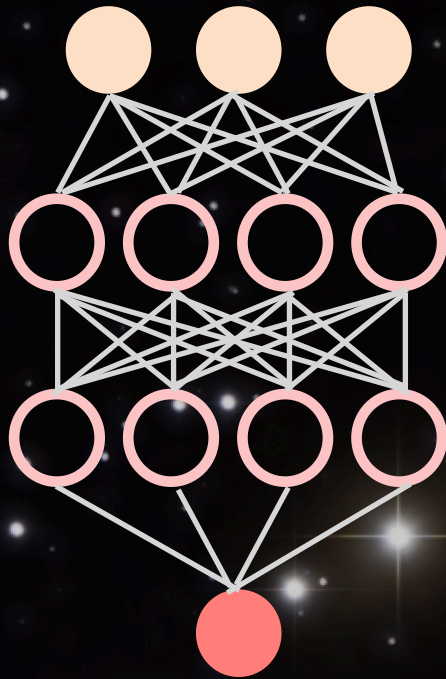
[article](#)



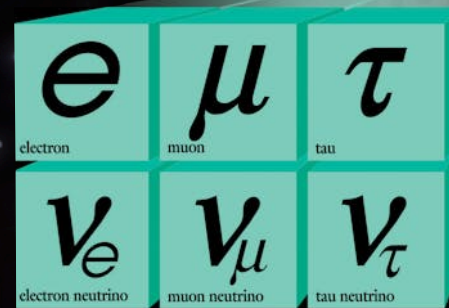
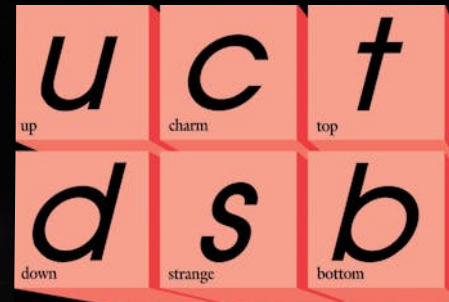
Physics, 2024



$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\gamma_{\mu}D^{\mu}\psi + |D_{\mu}\phi|^2 - V(\phi) + (y_{ij}\bar{\psi}_i\psi_j + \text{h.c.})$$

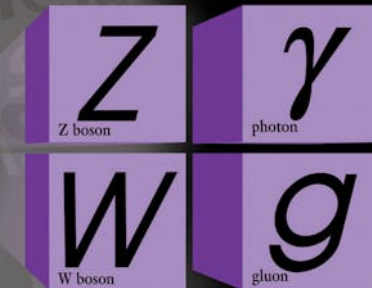


Quarks



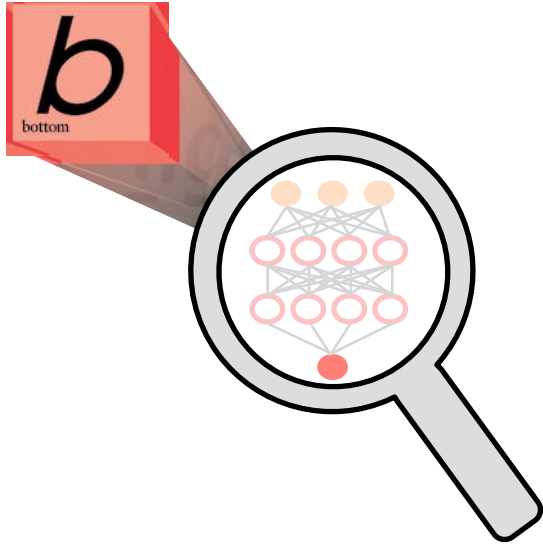
Leptons

Forces



AI for science: How to gain the *most* out of our physics datasets

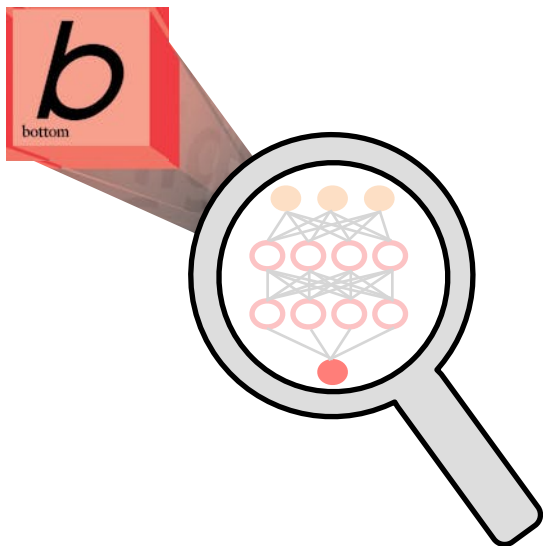
Future Outlook



Transformer-era

Deep Learning in FTAG

Future Outlook

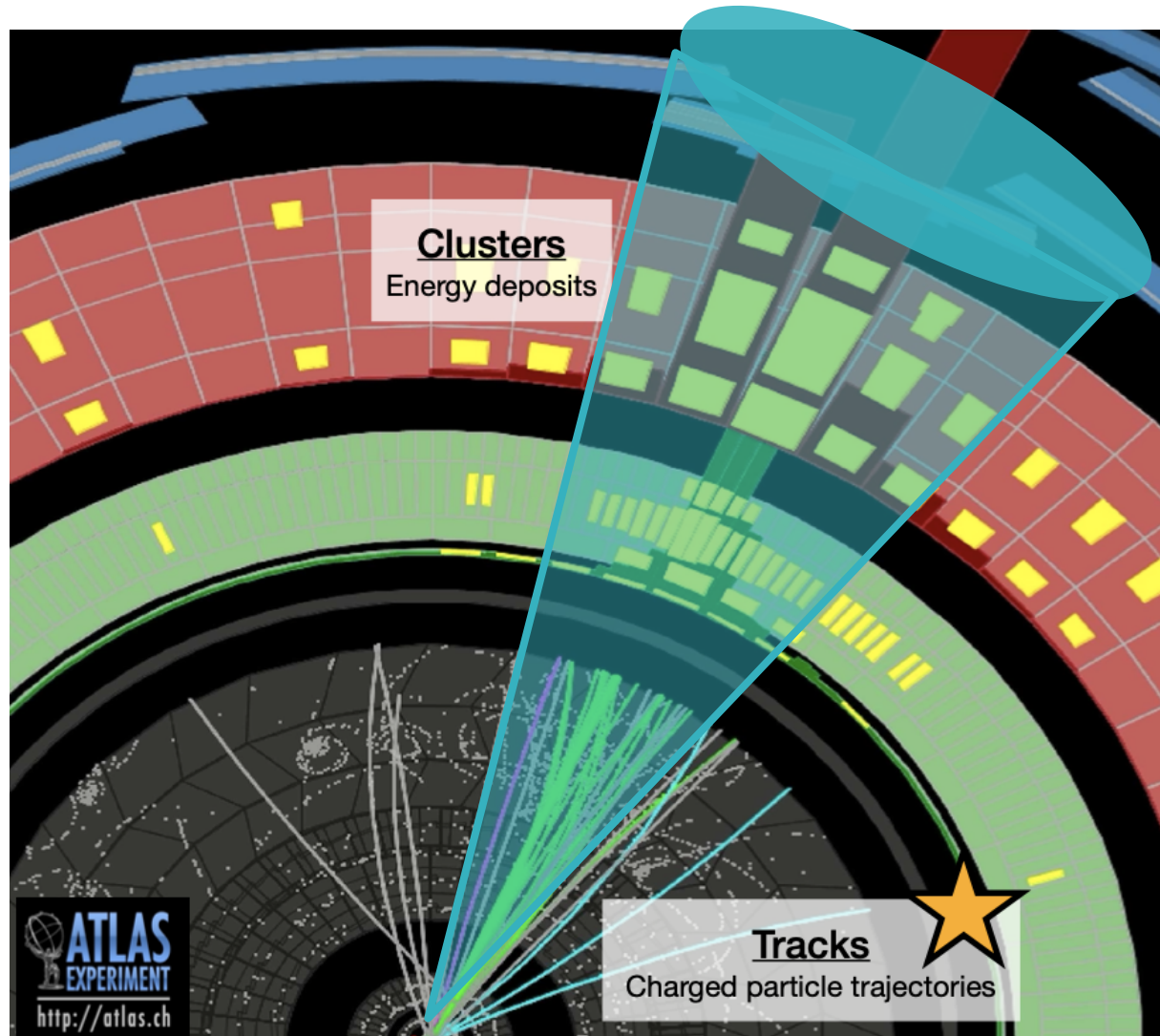


Transformer-era

Deep Learning in FTAG

- b-jets on ATLAS
- Why Deep Learning
 - RNNIP
 - Physics impact
- DIPS

Quark signature



Quark → reconstructed as collimated spray of particles

Jet : Unsupervised clustering algorithm

- ▶ Proxy for the **quark**
- ▶ Cluster with anti-kT
 - ▶ $R=0.4$ for b-jets ✨
 - ▶ $R=1.0$ for Higgs-jets



b-jet

✓ “Long” lifetime: $\tau = 1.2$ ps

✓ Many (≈ 5) displaced tracks

Variable # of tracks

b-jet

Displaced Tracks

Secondary Vertex

Jet

light jet

u,d,s g

✓ Most tracks originating from the PV

✓ Few displaced tracks

Prompt Tracks

Jet

c-jet

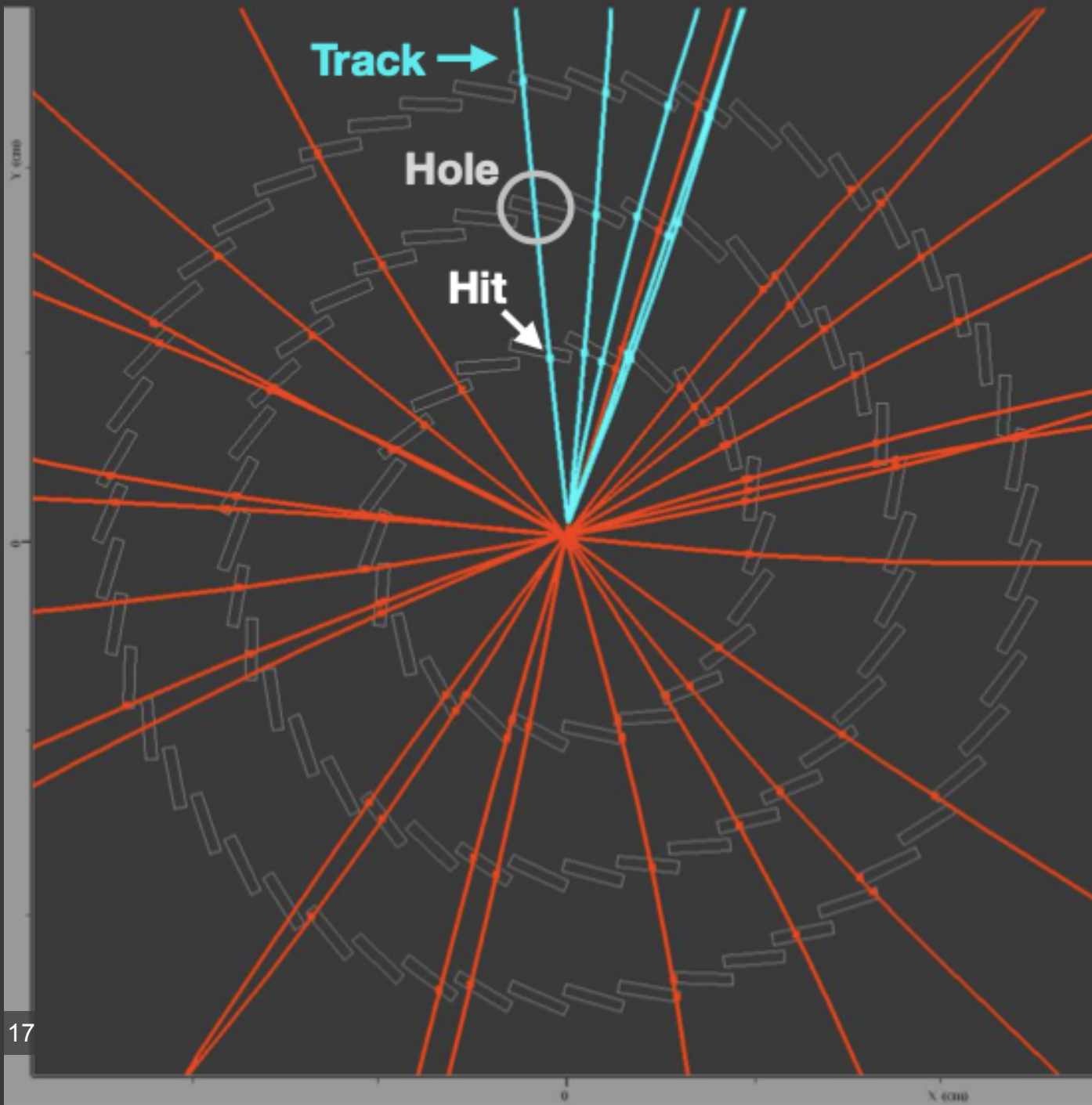
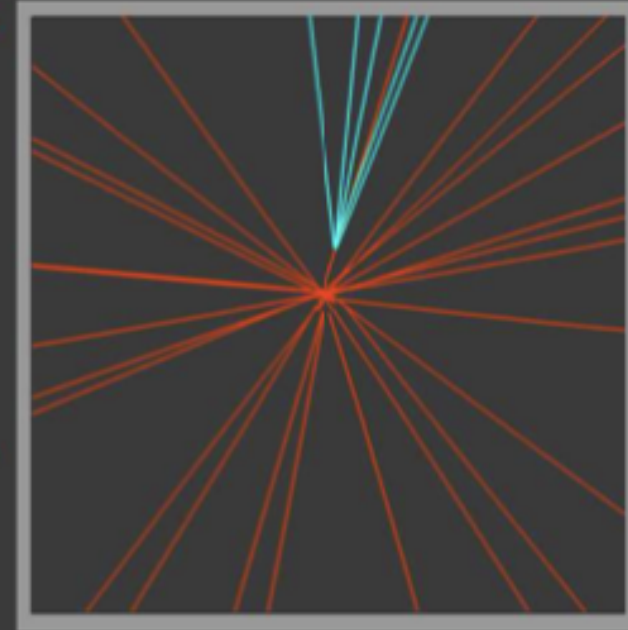
$\tau = .6$ ps

✓ Some displaced tracks



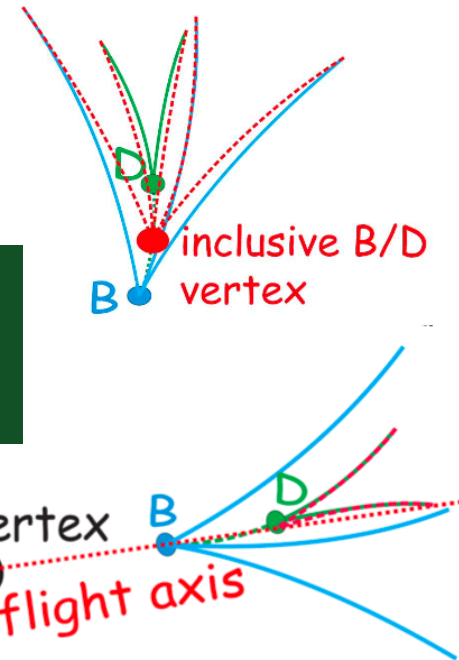
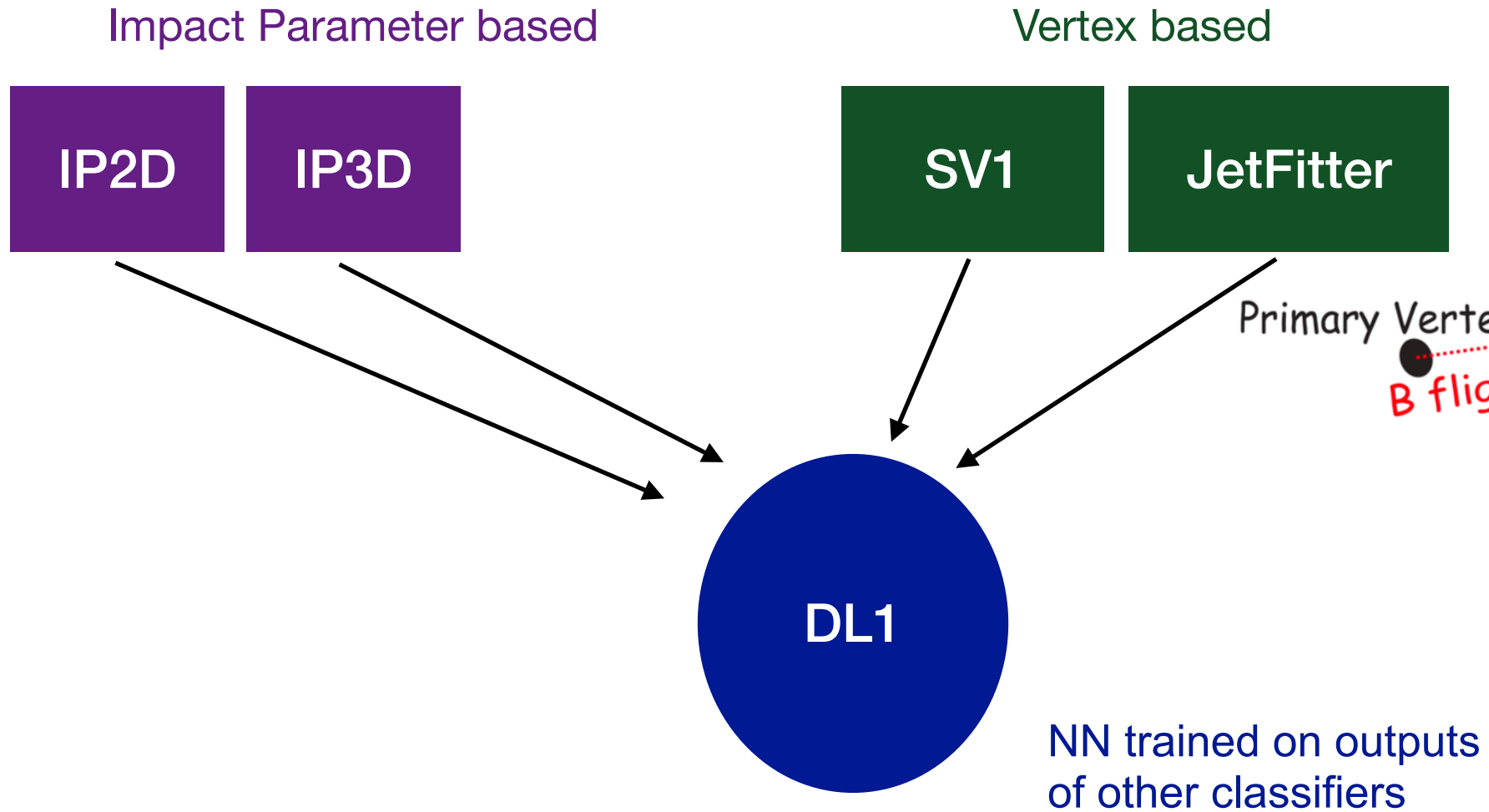
Run 142195, Event 284154

Decay length = 3.7 mm
Decay length significance = 22
Lifetime = 3.1 ps
Vertex mass = 2.5 GeV
Number of tracks = 5



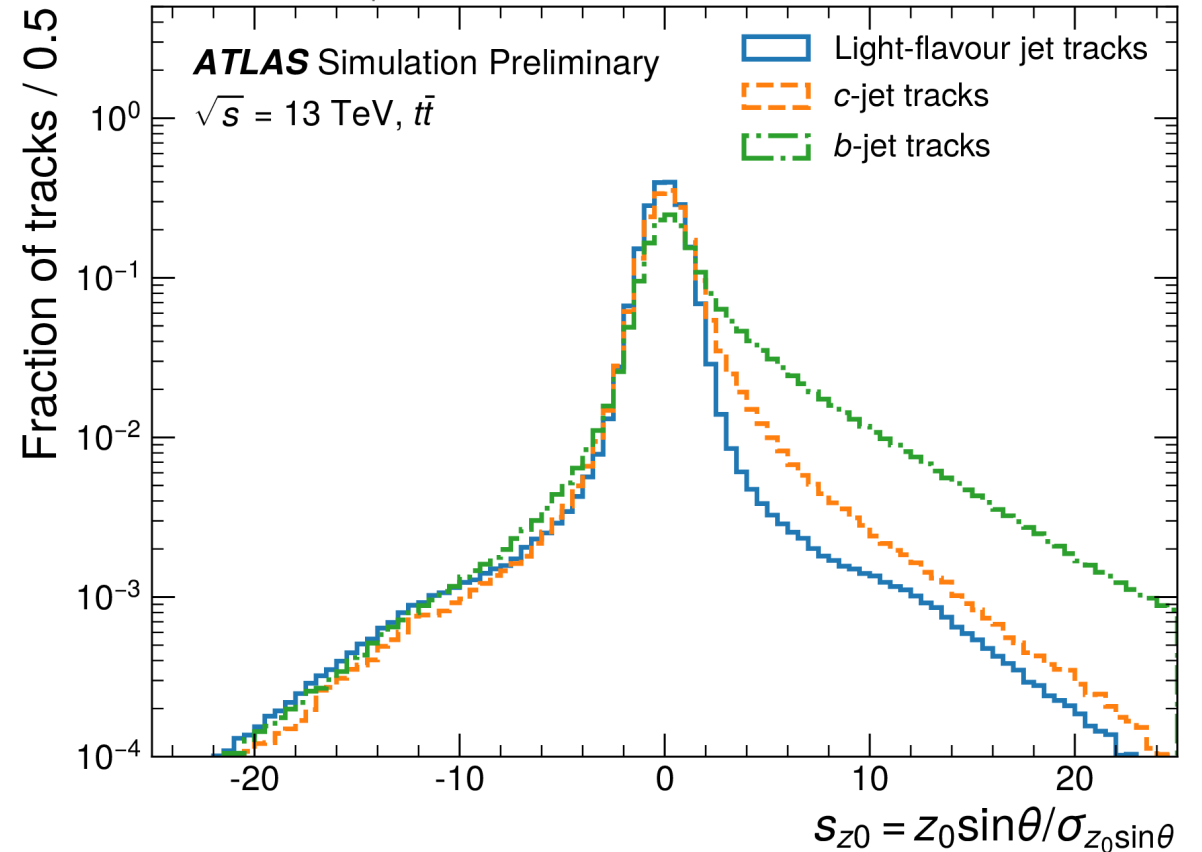
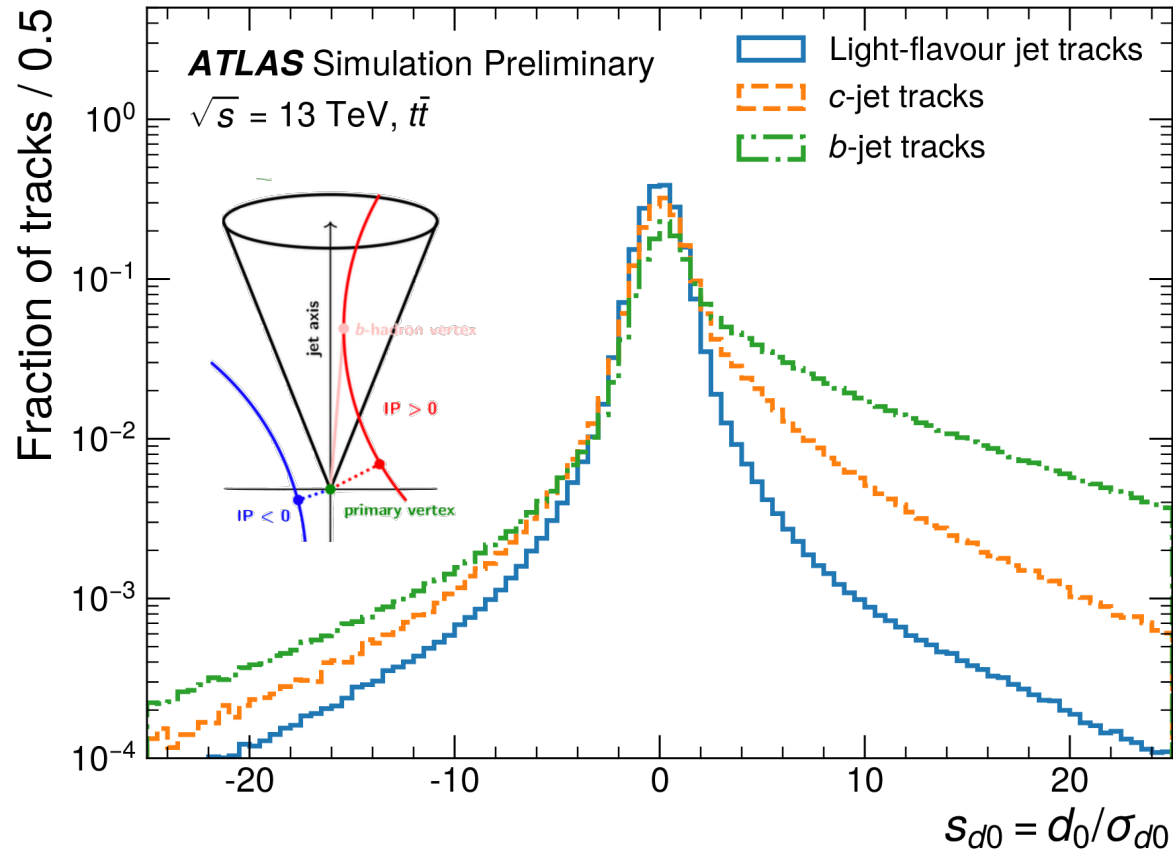
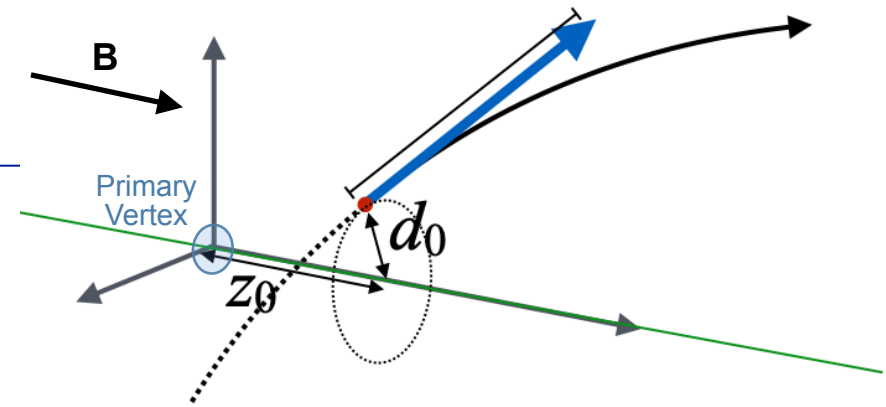
ATLAS b-jet classifiers

Low level
High level



b-tagging

Key variable: impact parameter

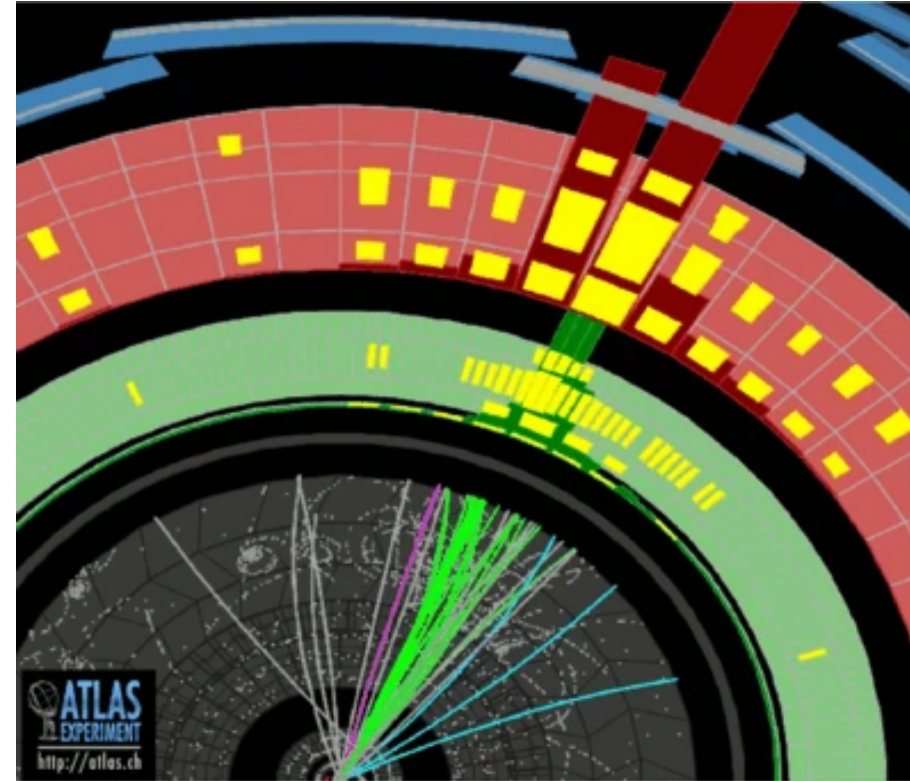


How do we aggregate this information?

What we have:

- Collection of tracks
 - $X_i : i = \{1, \dots, \underline{n}\}$
🔑 variable # of tracks
- Each track has features
 - $X_i \in \mathbb{R}^m$ ← E.g, impact parameters
momenta, quality
- Jet has labels
 - $Y: \{\mathbf{b}, \mathbf{c}, \mathbf{light}\}$ – or –
 - $Y: \{\mathbf{bb}, \mathbf{cc}, \mathbf{top}, \mathbf{QCD}\}$

What we want: $p(Y | X_1, \dots, X_n)$



High dimensional problem

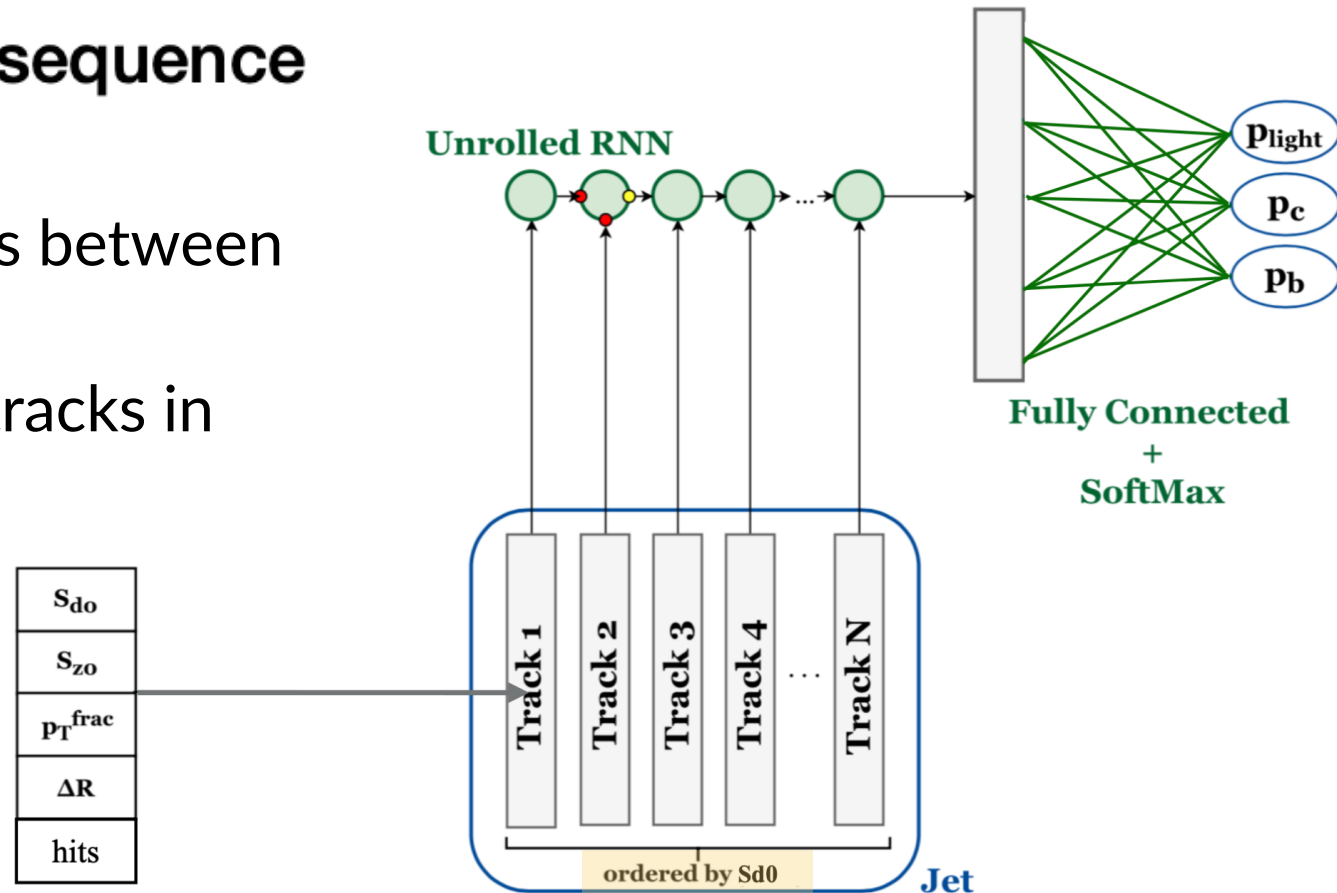
$$n \cdot m \sim \mathcal{O}(10^3)$$

Recurrent Neural Network

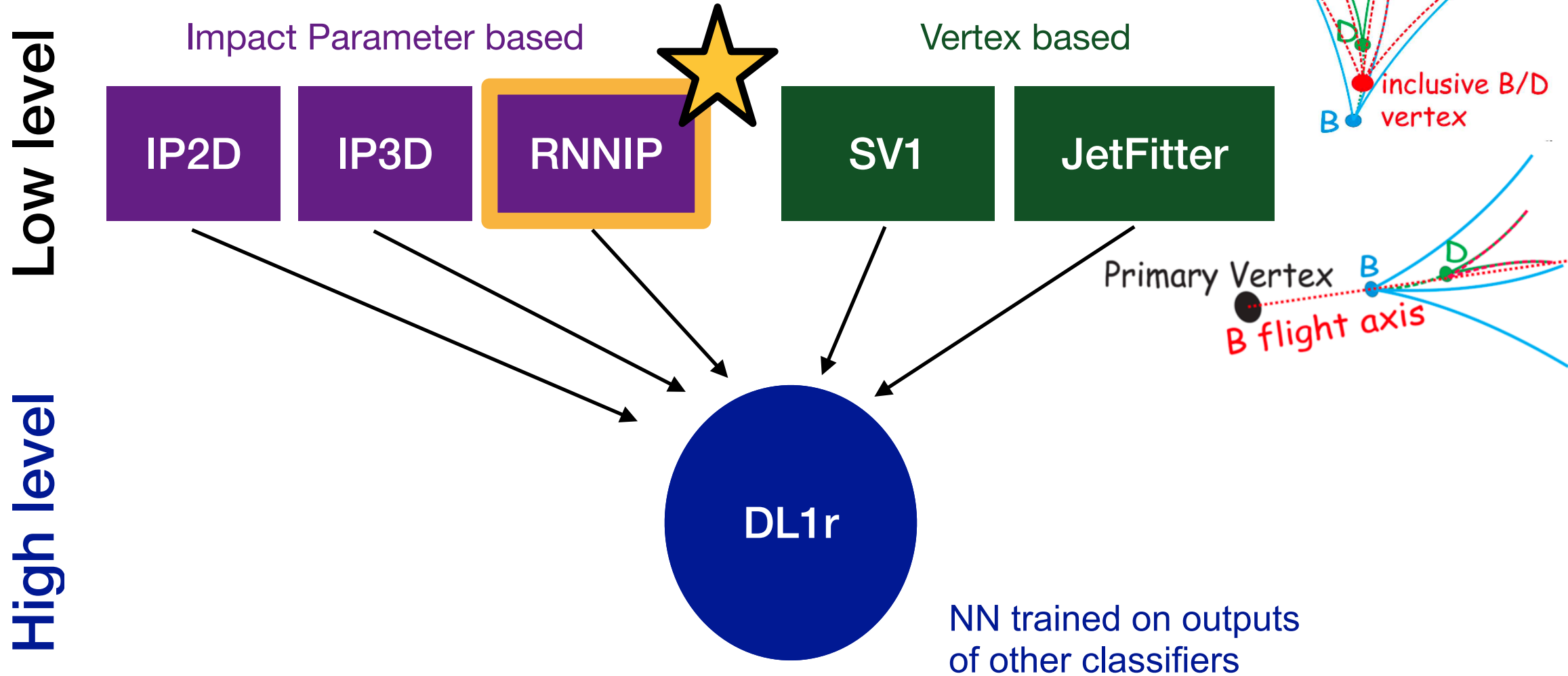


Model the jet as a sequence

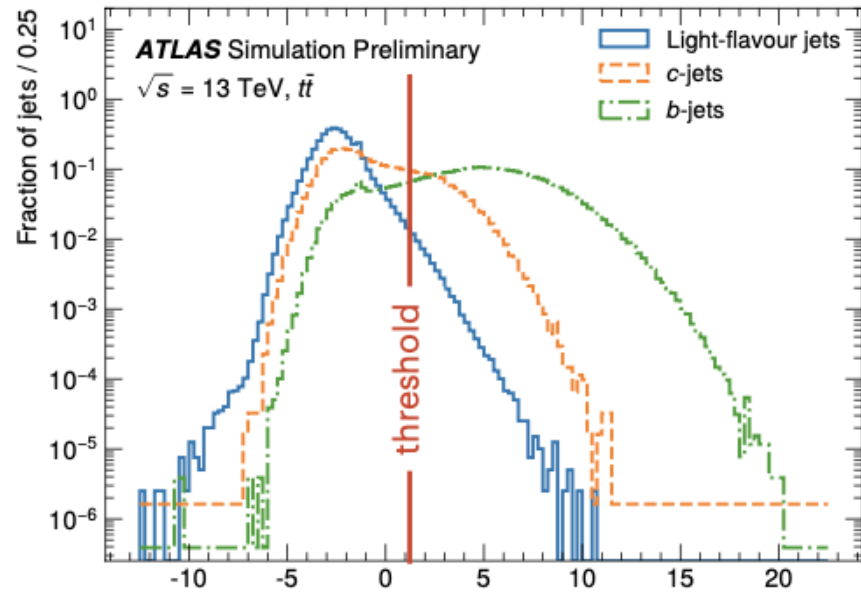
- ☑ Account for correlations between tracks
- ☑ Allow for variable # of tracks in the jet
- ☑ Avoids curse of dimensionality - add more features



ATLAS b-jet classifiers



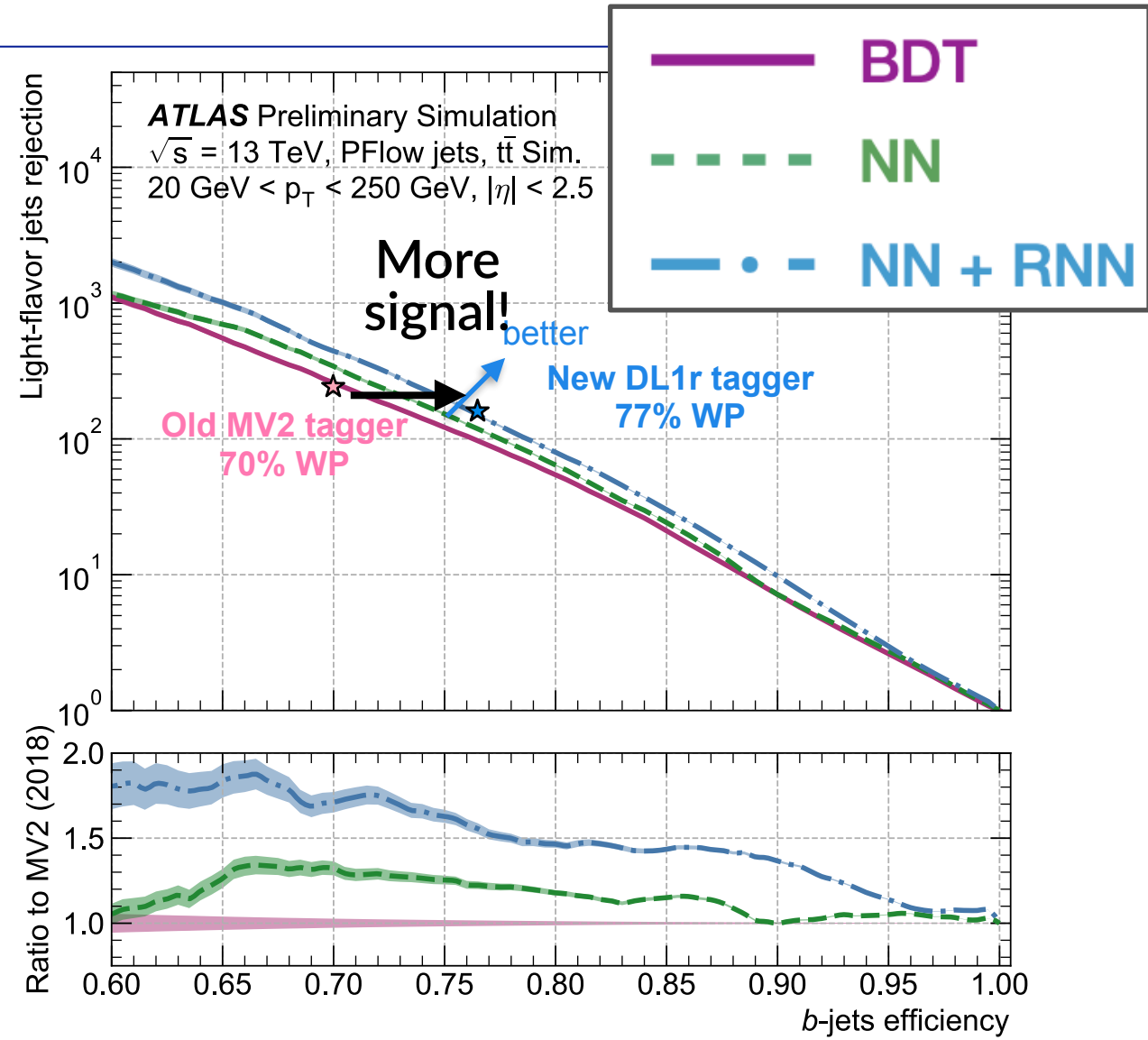
DL1r classifier



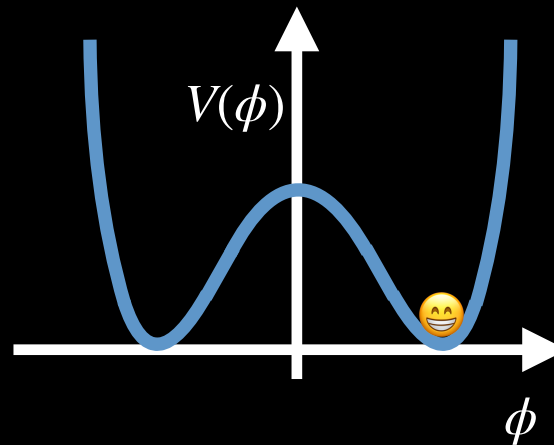
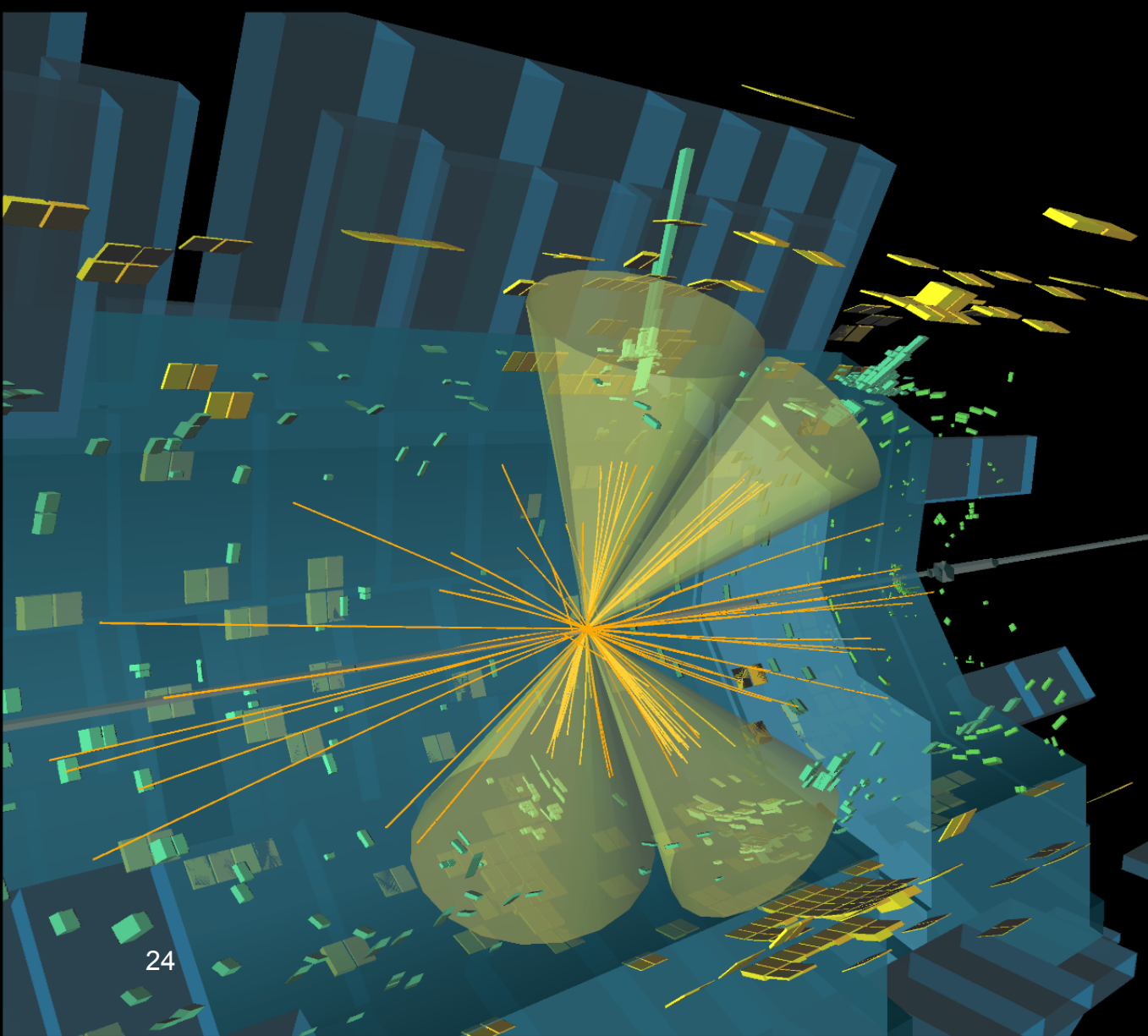
$$D_b = \log \frac{p_b}{f_c p_c + (1 - f_c) p_l}$$

$$\text{b-jet efficiency} = \frac{\# \text{ b-jets} > \text{threshold}}{\# \text{ b-jets}}$$

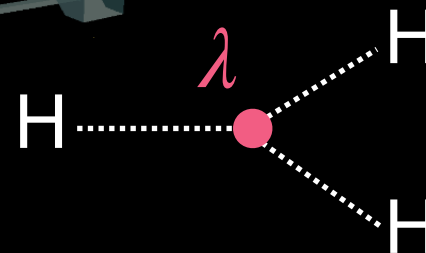
$$\text{Background rejection} = \frac{1}{\text{Background efficiency}}$$



How did **b-tagging** improvements help our HH analyses?



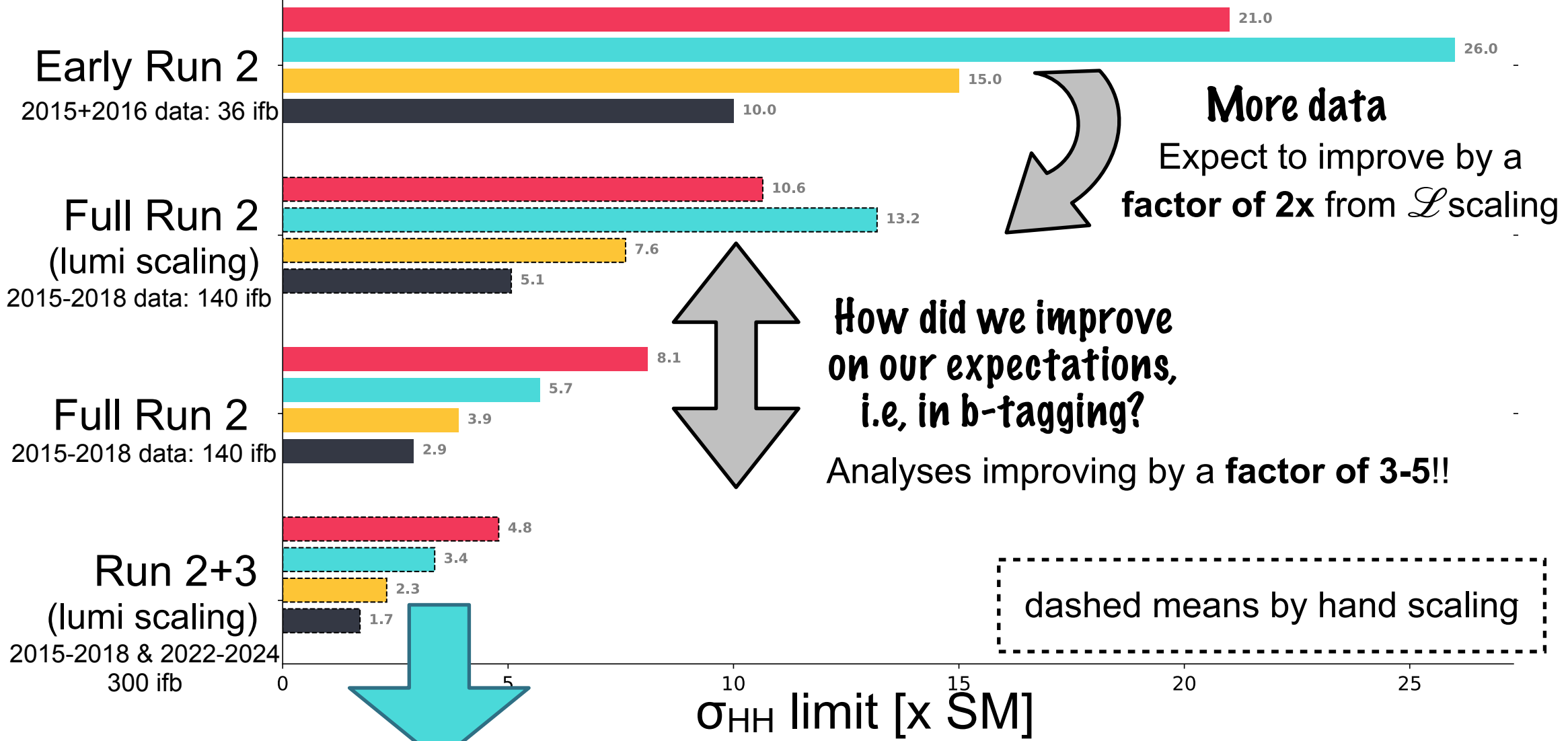
$$V(\phi) = \mu^2 h(x)^2 + \lambda v h(x)^3 + \frac{1}{4} \lambda h(x)^4$$



signal
background $\sim 10^{-13}$

Main HH channels *4b*, *bbγγ*, *bbττ*,
all need *b*-tagging.

4b $b\bar{b}\gamma\gamma$ $b\bar{b}\tau\tau$ Combination



More data
Expect to improve by a **factor of 2x** from \mathcal{L} scaling

How did we improve on our expectations, i.e, in b-tagging?
Analyses improving by a **factor of 3-5!!**

dashed means by hand scaling

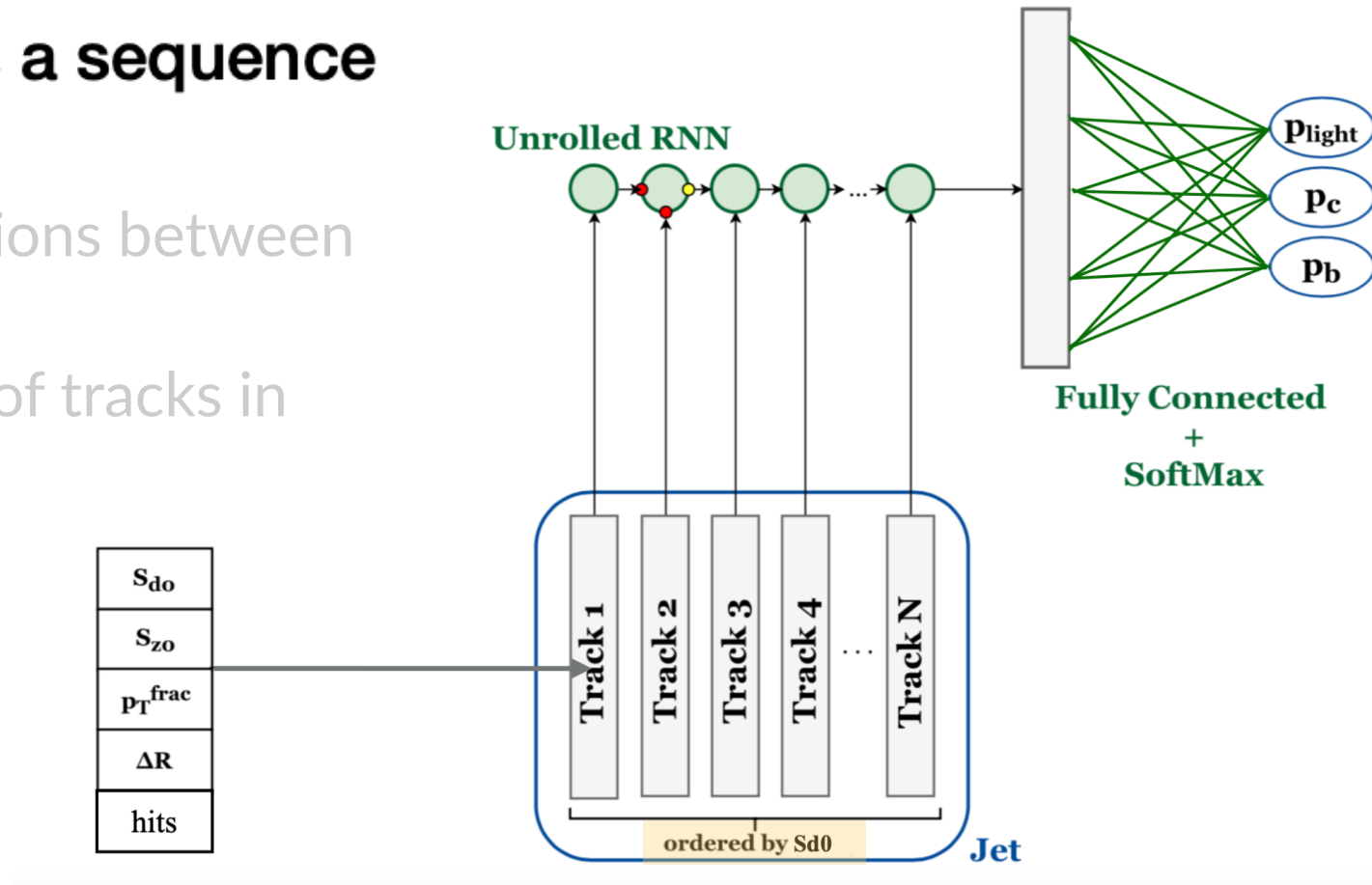
How will b-tagging improvements help $b\bar{b}\gamma\gamma$ again exceed expectations?

Issue: ordering



Model the jet as a sequence

- ✓ Account for correlations between tracks
- ✓ Allow for variable # of tracks in the jet
- ✓ Avoids curse of dimensionality - add more features



Challenge: How to order??

When does order matter (?)

Natural language

(1) Mary likes John



(2) John likes Mary

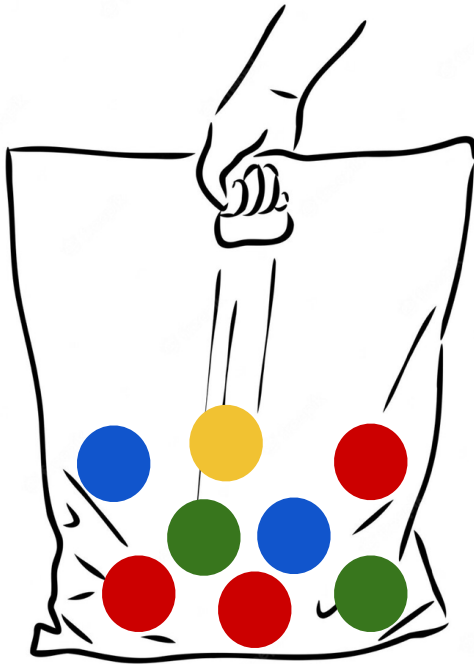


Same words... the order changes the meaning.

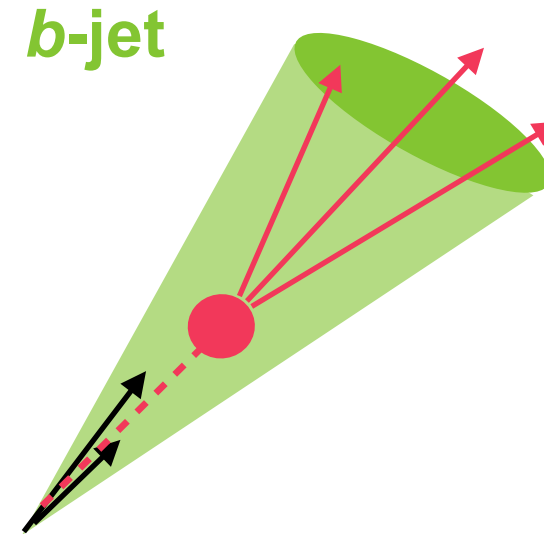
Permutation invariance

Set: Collection of objects without any specified order

Ex 1: # of colored balls in a bag



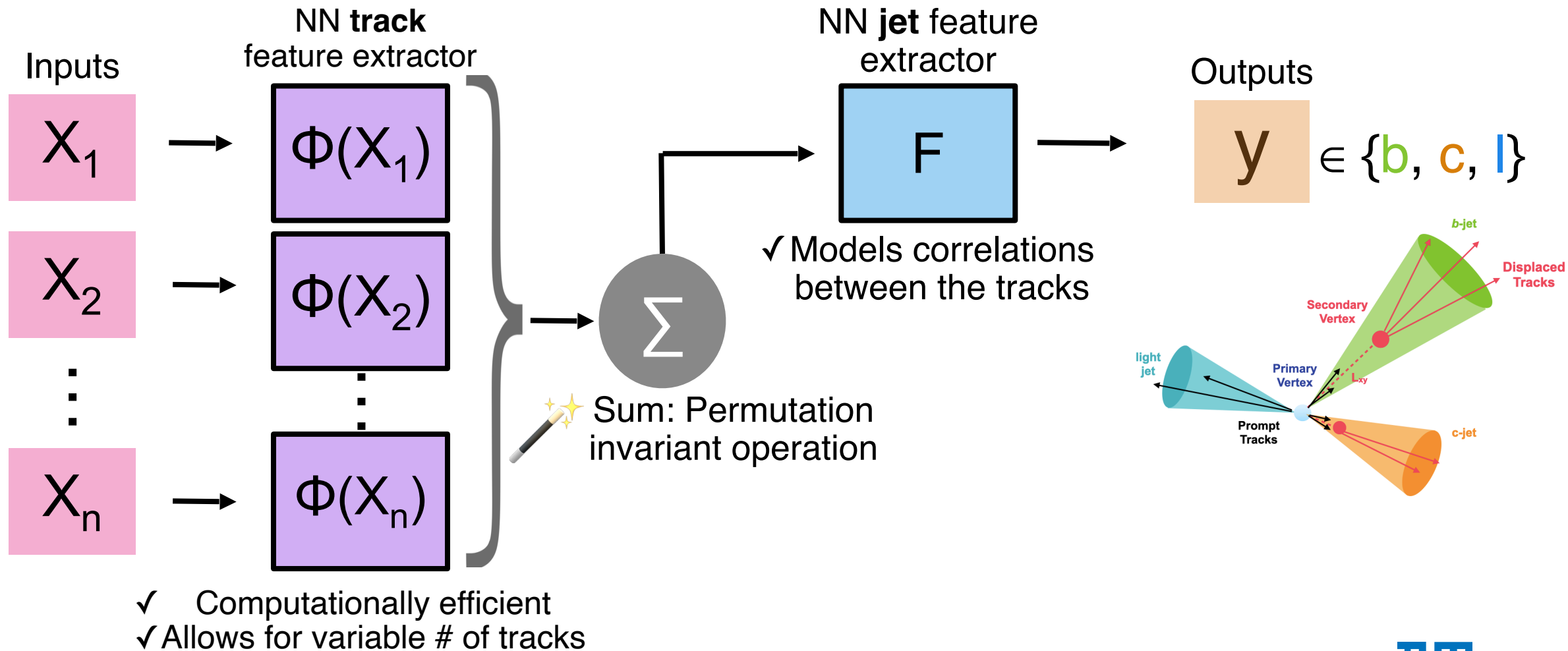
Ex 2: tracks in heavy flavour decay



Deep Set: Neural Network designed to operate on sets

Deep Sets

Network of networks!!



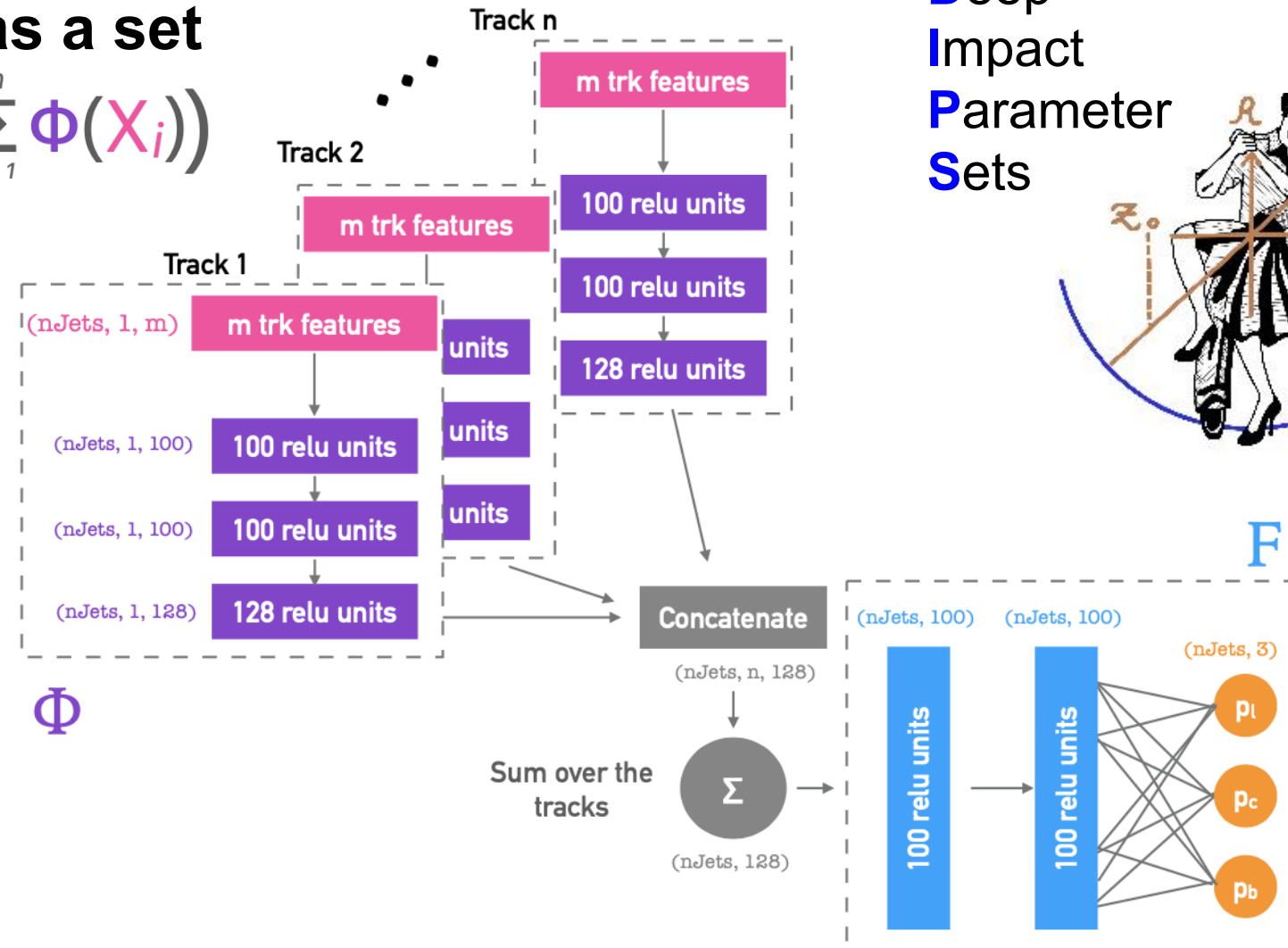
DIPS



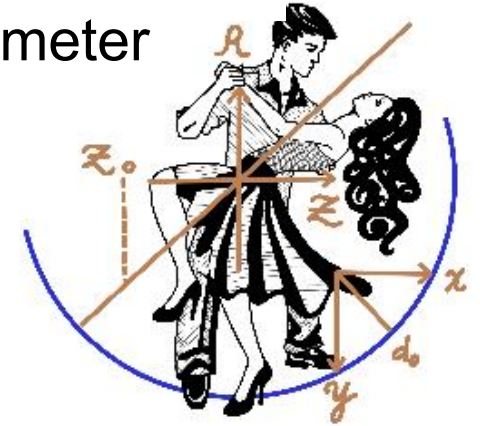
Model the jet as a set

$$O(\{X_1, \dots, X_n\}) = F\left(\sum_{i=1}^n \Phi(X_i)\right)$$

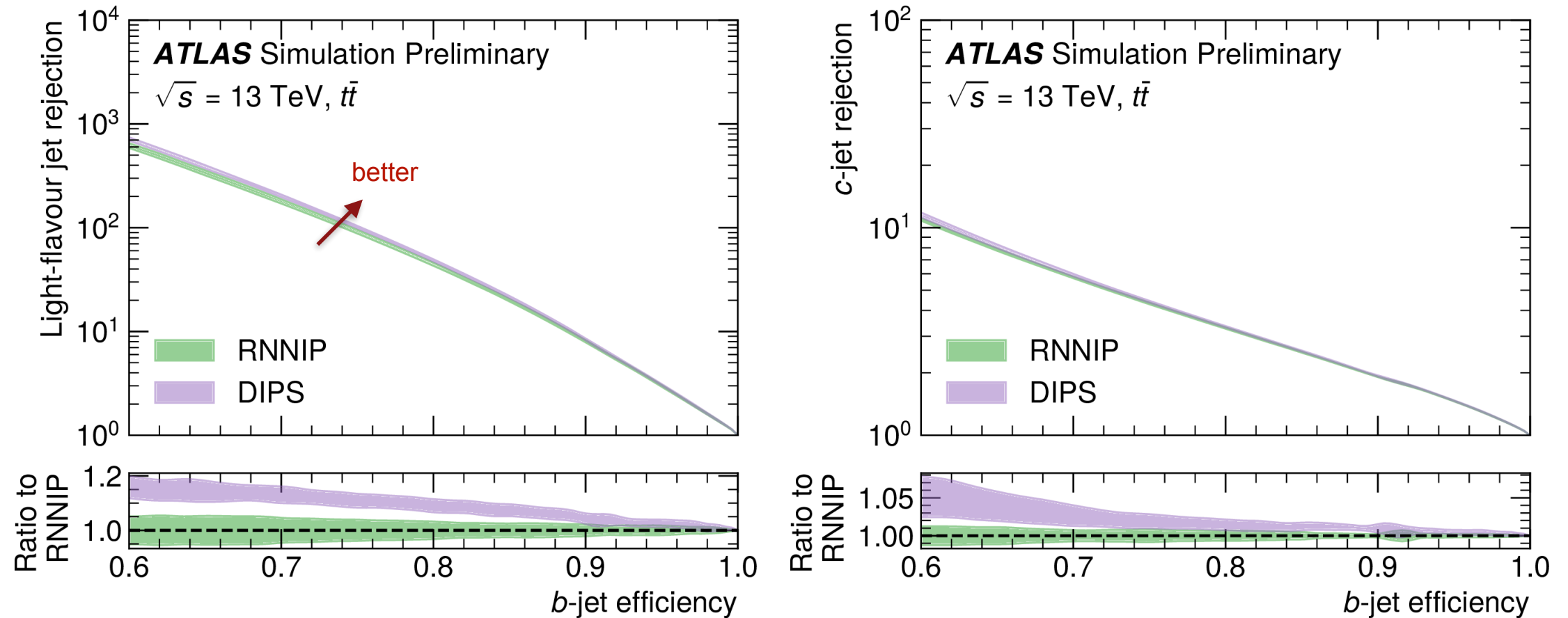
Same inputs as b-tagging RNN



Deep
Impact
Parameter
Sets



DIPS



Similar performance with the same inputs



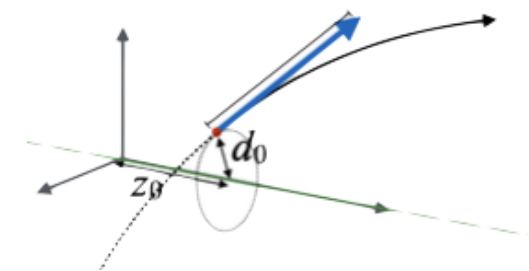
4x speed-up in the training time!!!

Faster turn around time for physics optimizations

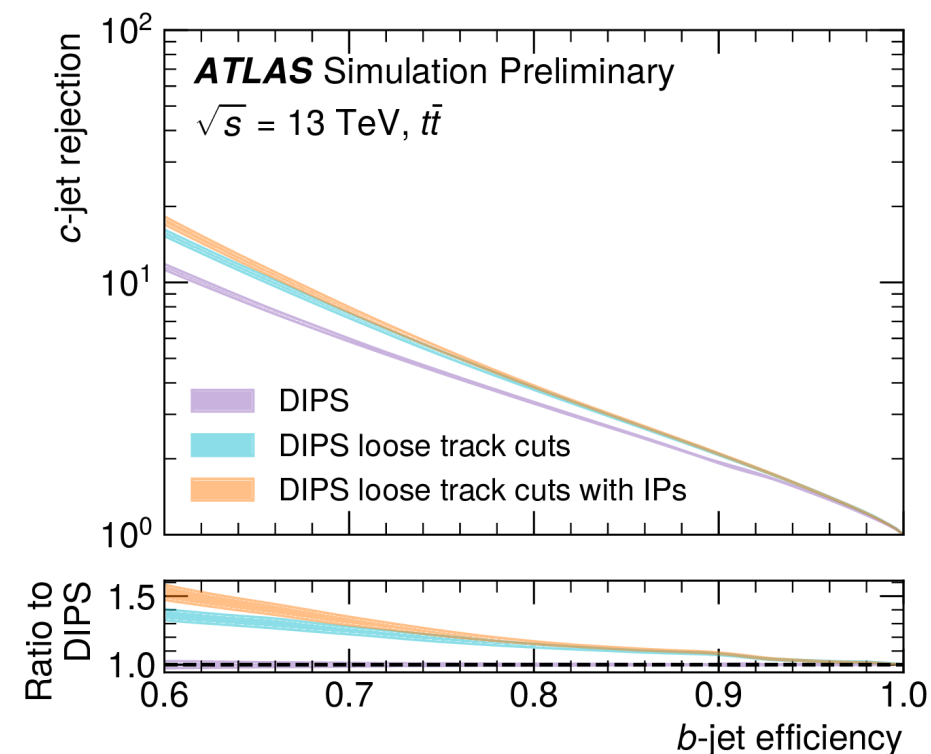
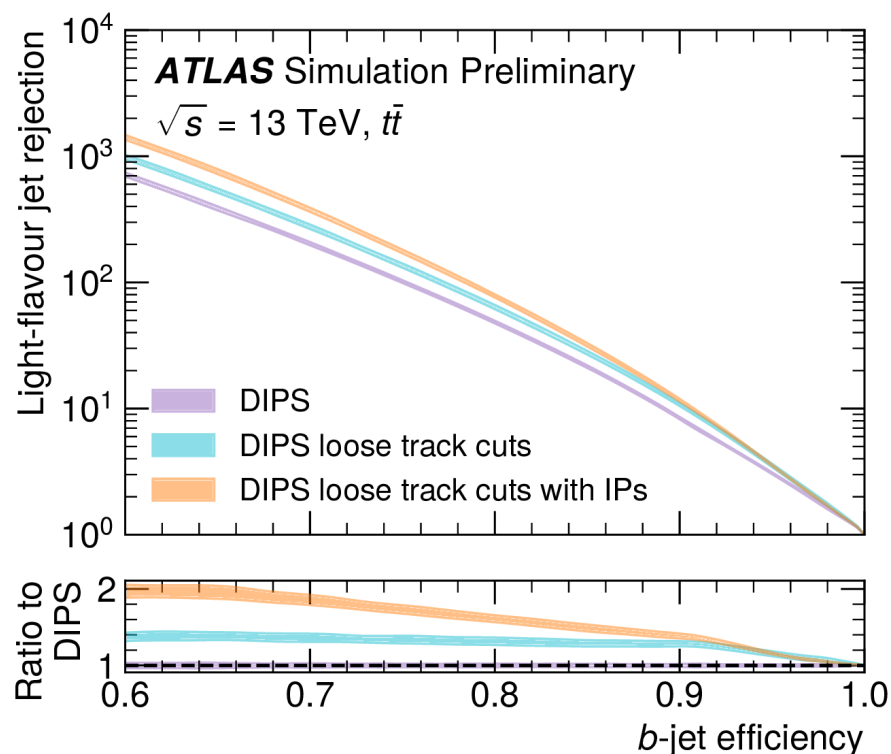
nominal: from before ($p_T > 1$ GeV, $|d_0| < 1$ mm, $|z_0 \sin \theta| < 1.5$ mm)

loose: poorer quality ($p_T > 1$ GeV, $|d_0| < 1$ mm, $|z_0 \sin \theta| < 1.5$ mm)

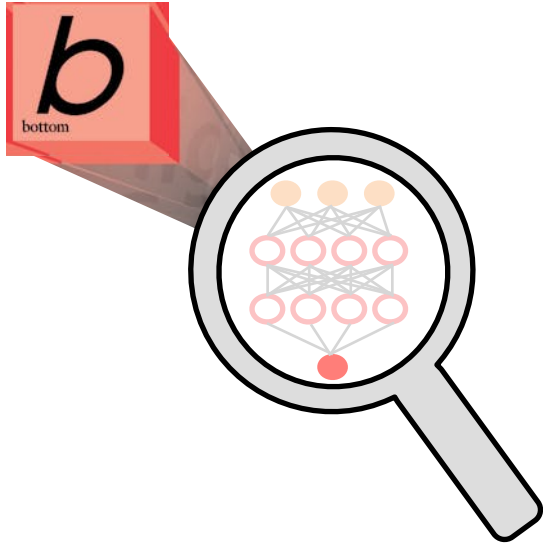
Optimized DIPS: loose selection + new inputs ($d_0, z_0 \sin \theta$)



DIPS more performant for a high efficiency track selection



Future Outlook



Transformer-era

- GN1
- Vertexing plots
- Performance in data
- Other similar applications
 - pT regression
 - Higgs tagging

Deep Learning in FTAG



Transformers are transforming our day-to-day lives!

What can I help with?

Message ChatGPT

- Create image
- Brainstorm
- Summarize text
- Code
- More

• Math

how do you integrate a multidimensional gaussian with a source term

To integrate a multidimensional Gaussian with a source term, you are typically dealing with an integral of the form:

$$I = \int d^n \mathbf{x} \exp\left(-\frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{J}^T \mathbf{x}\right)$$

where:

- \mathbf{x} is an n -dimensional vector,
- A is an $n \times n$ symmetric positive-definite matrix (the covariance or precision matrix),
- \mathbf{J} is an n -dimensional source term (linear term),
- $d^n \mathbf{x}$ is the measure for the n -dimensional integral.

This integral represents a Gaussian function in n -dimensions with a linear source term \mathbf{J} . The steps to solve this are outlined below:

Step 1: Completing the square

First, you need to complete the square in the exponent. The exponent is:

$$-\frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{J}^T \mathbf{x}$$

You can rewrite it by factoring terms:

$$-\frac{1}{2} (\mathbf{x}^T A \mathbf{x} - 2 \mathbf{J}^T \mathbf{x})$$

Next, complete the square by introducing a new shifted variable $\mathbf{y} = \mathbf{x} - A^{-1} \mathbf{J}$, which removes the linear term from the exponent.

$$-\frac{1}{2} (\mathbf{x} - A^{-1} \mathbf{J})^T A (\mathbf{x} - A^{-1} \mathbf{J}) + \frac{1}{2} \mathbf{J}^T A^{-1} \mathbf{J}$$

Now, the exponent is:

$$-\frac{1}{2} \mathbf{y}^T \downarrow + \frac{1}{2} \mathbf{J}^T A^{-1} \mathbf{J}$$



Transformers are transforming our day-to-day lives!

What can I help with?

Message ChatGPT

🔗 🌐 🔊

- Create image
- Brainstorm
- Summarize text
- Code
- More

- **Math**
- **Poetry**

Write me a song for the ATLAS flavour tagging group.

(Verse)
Impact parameter, a tiny shift,
Lifetime of beauty, a precious gift.
Secondary vertices, a hidden clue,
Decays revealing what particles flew.
Machine learning, a powerful friend,
Boosting performance, till the very end.
(Chorus)
Flavour tagging, that's our game,
Sorting through the chaos, calling out the name.
b-jets, c-jets, light quark too,
Reconstructing truth in all we do.
Algorithms humming, a digital art,
Unraveling secrets, playing our part.

Curtesy of Diptaparna
* Gemini 1.5 Pro 002



Transformers are transforming our day-to-day lives!

What can I help with?

Message ChatGPT



Create image

Brainstorm

Summarize text

Code

More

Could you draw me a 16x9 image of Siegen in impressionist style

- Math
- Poetry
- **Art**



Here is an impressionist-style illustration of Siegen, Germany, showcasing its landmarks and natural beauty with soft, textured brushstrokes. Let me know if you need further adjustments!





Transformers for collecting your Nobel prize

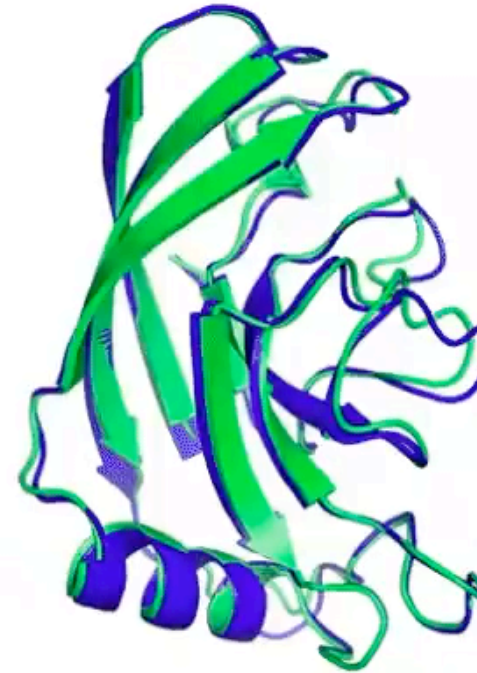


Chemistry, 2024

AlphaFold2



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)



T1049 / 6y4f
93.3 GDT
(adhesin tip)

- Experimental result
- Computational prediction



Transformers 101

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

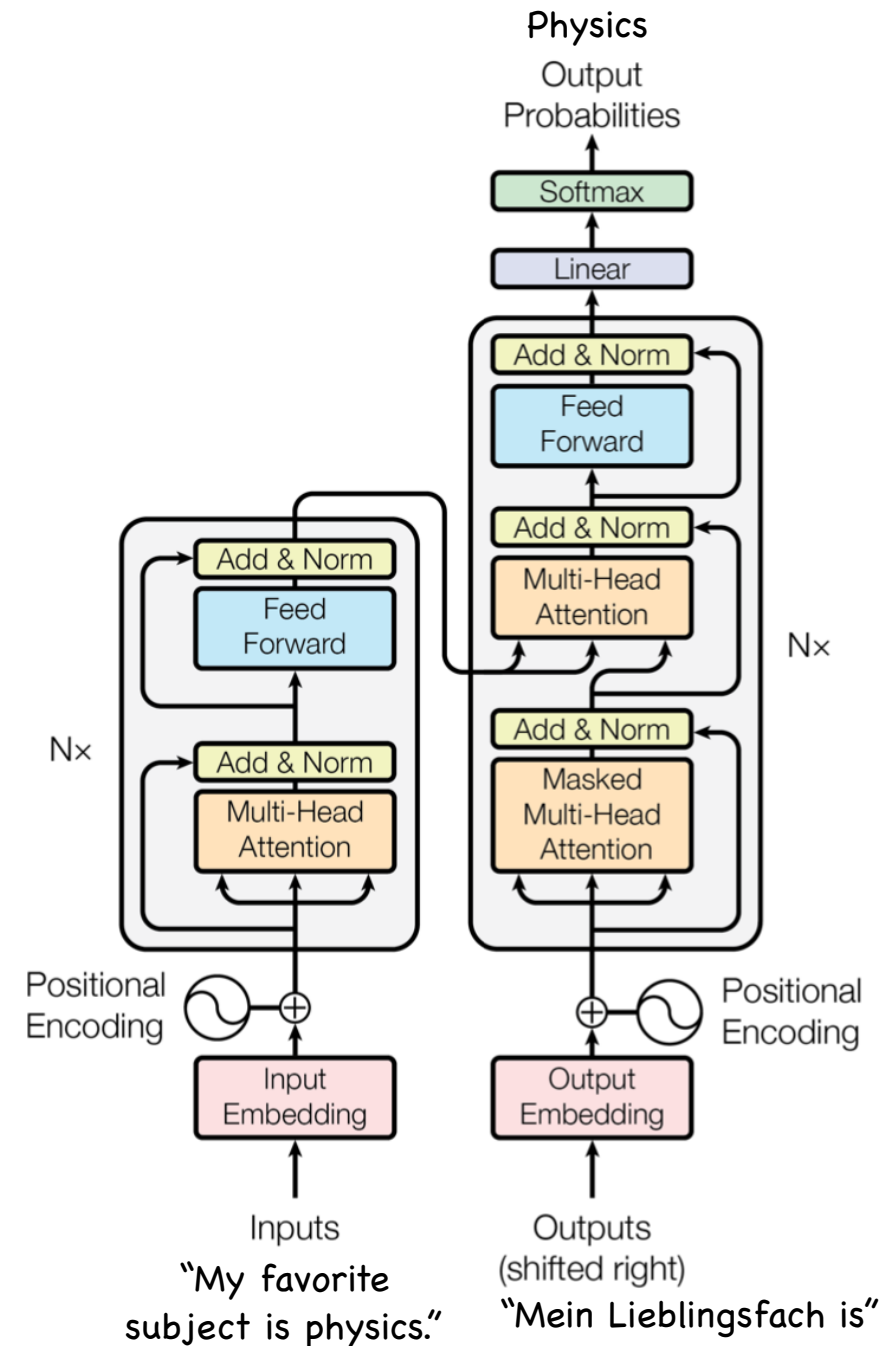
Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

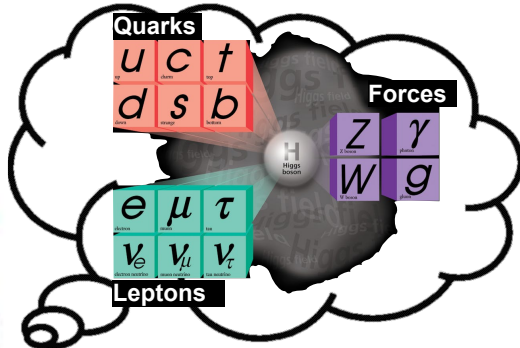
142k
citations





Transformers 101

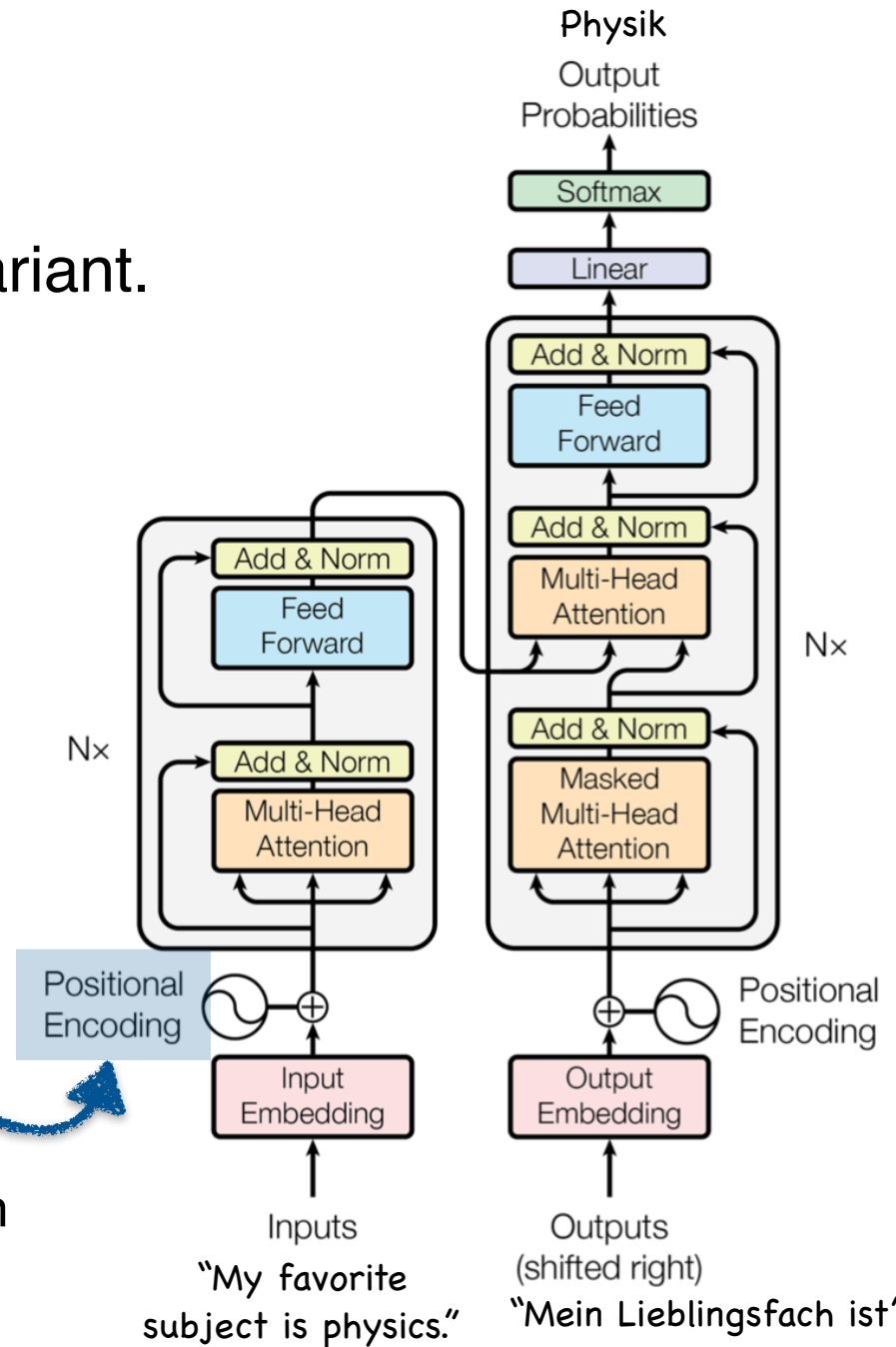
Architecture natively permutation invariant.



"My favorite subject is physics."

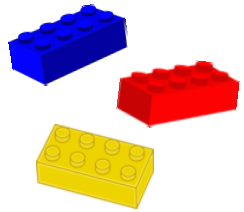
<START> 0 1 2 3 4 <END>

Sequence information encoded via an additional input





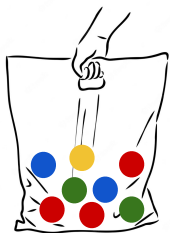
Transformers for FTAG



And do this lots of times!

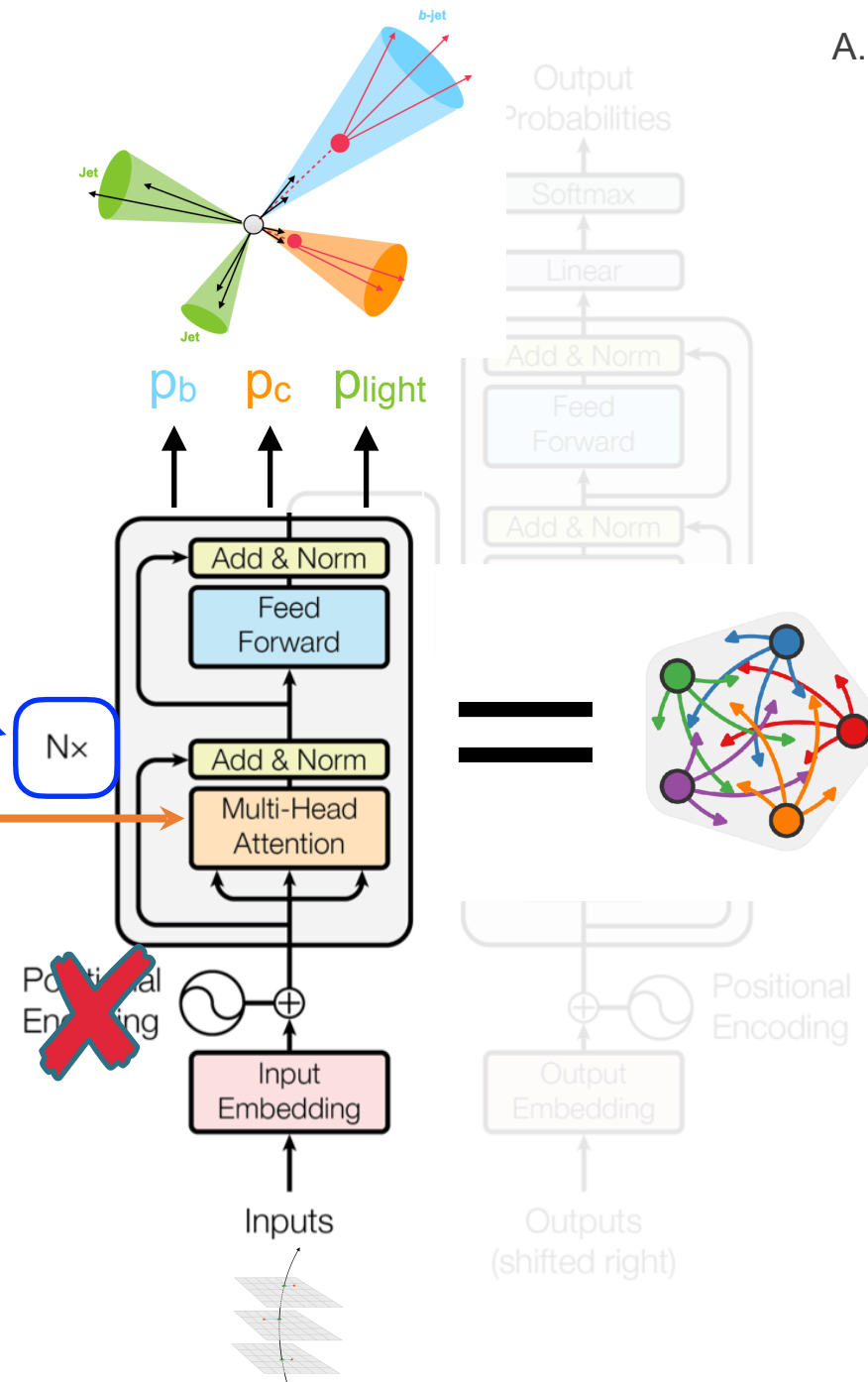


Replace the Deep Sets Sum with a **weighted sum** (still permutation invariant)



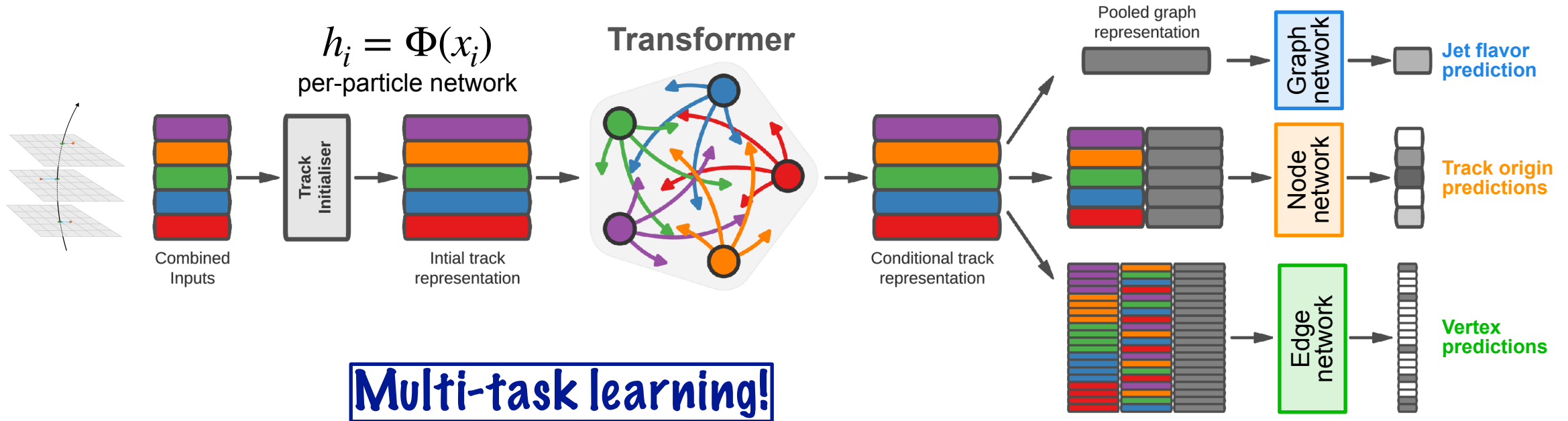
No need for positional encoding
(permutation invariance over the tracks)

40



GN2: architecture

Transformer-based flavour tagger

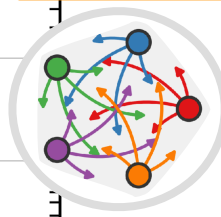
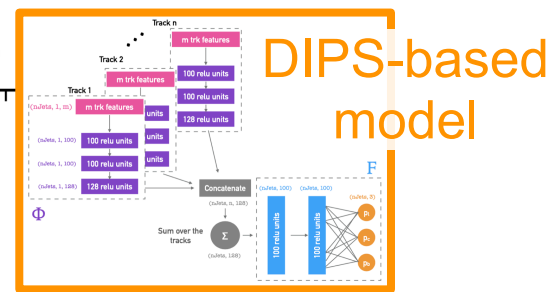
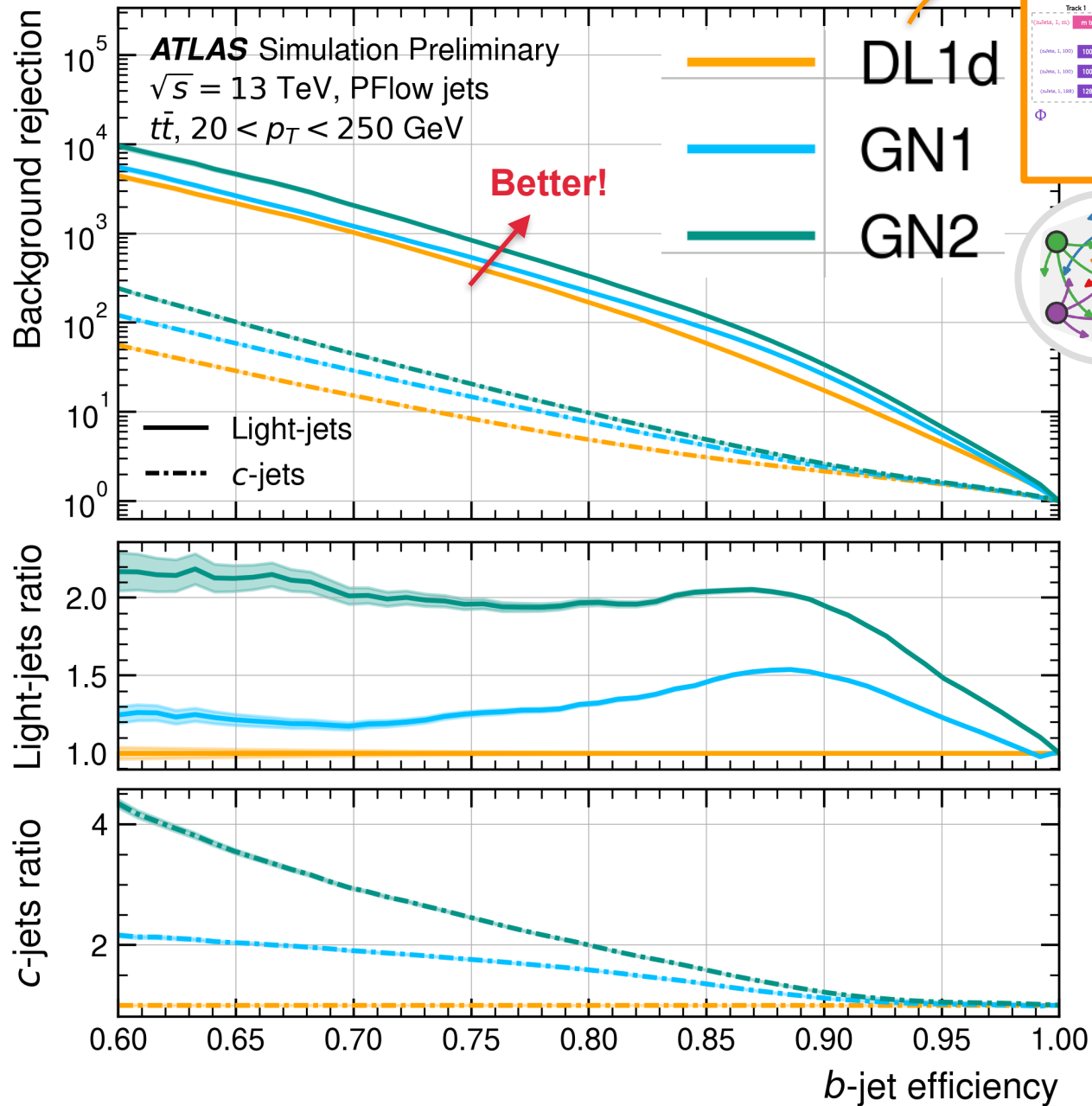


$$\mathcal{L}_{tot} = \mathcal{L}_{jet} + \alpha \mathcal{L}_{trk} + \beta \mathcal{L}_{vtx}$$

$\alpha = 0.5$ $\beta = 1.5$



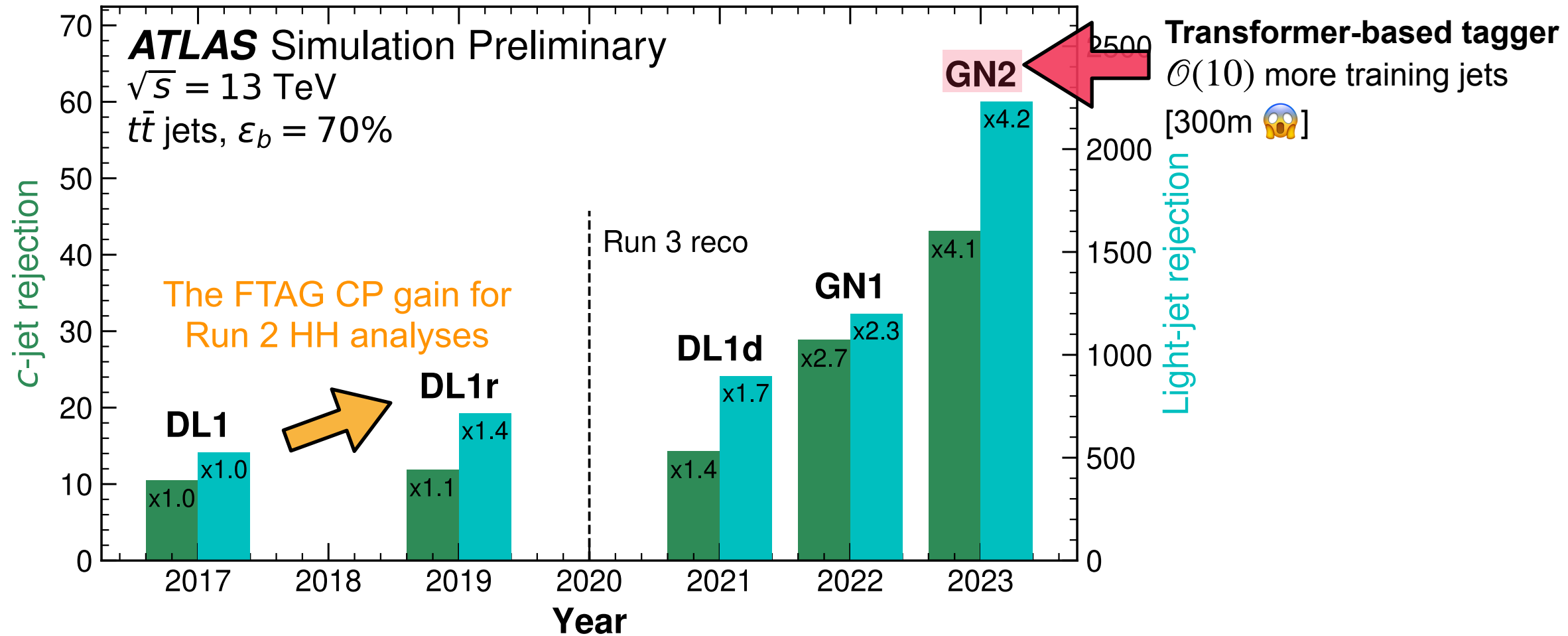
GN1: graph neural network
 GN2: transformer



2x increase in light rejection

4x increase in charm rejection

FTAG over time



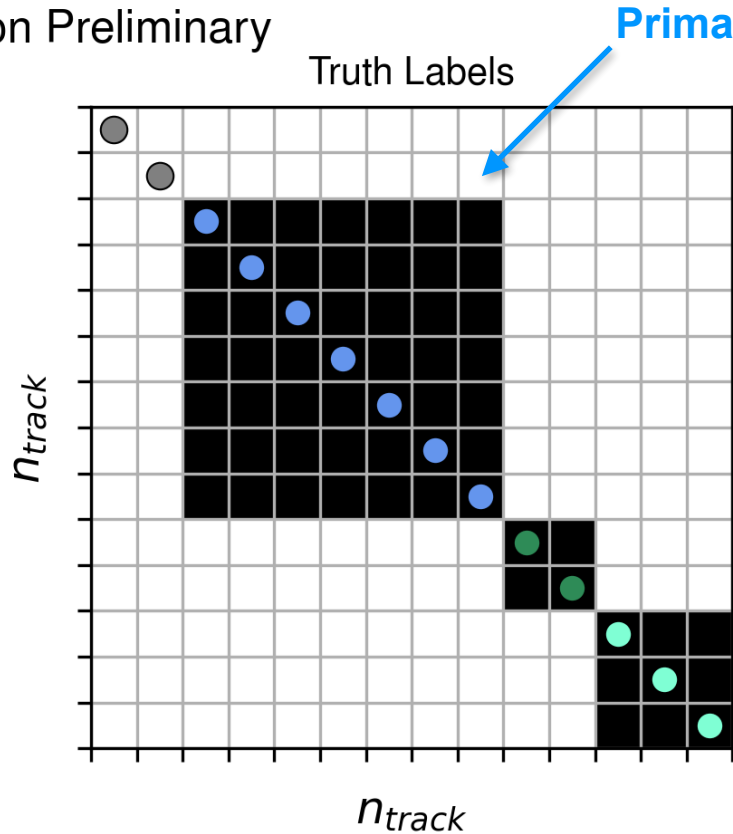
Vertex finding

Group tracks with **pair-wise compatibility** > 0.5

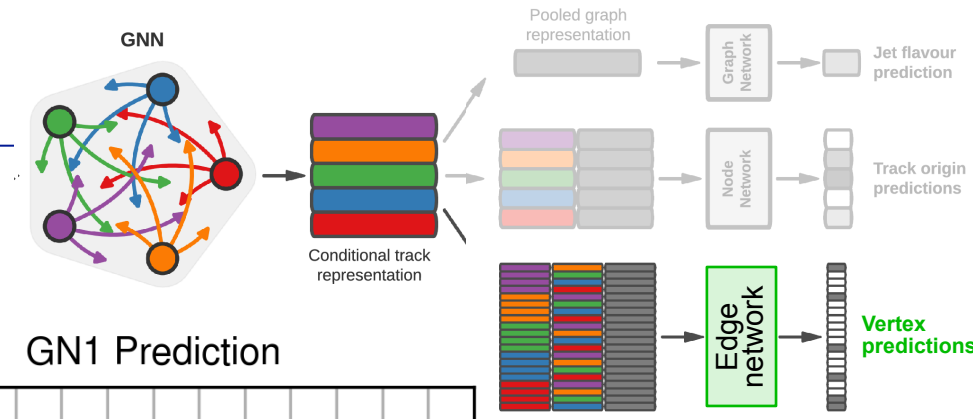
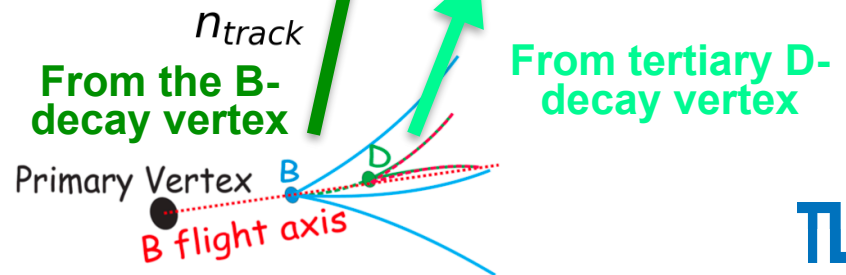
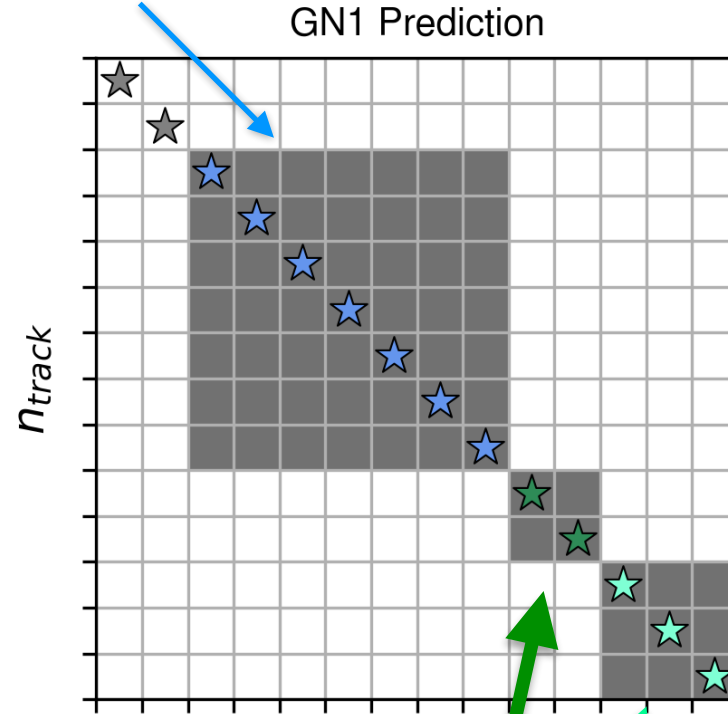
ATLAS Simulation Preliminary
 $\sqrt{s} = 13$ TeV
 $t\bar{t}$ jets

Truth b -jet
 $p_T = 134.1$ GeV

$p_b = 0.995$
 $p_c = 0.005$
 $p_u = 0.000$

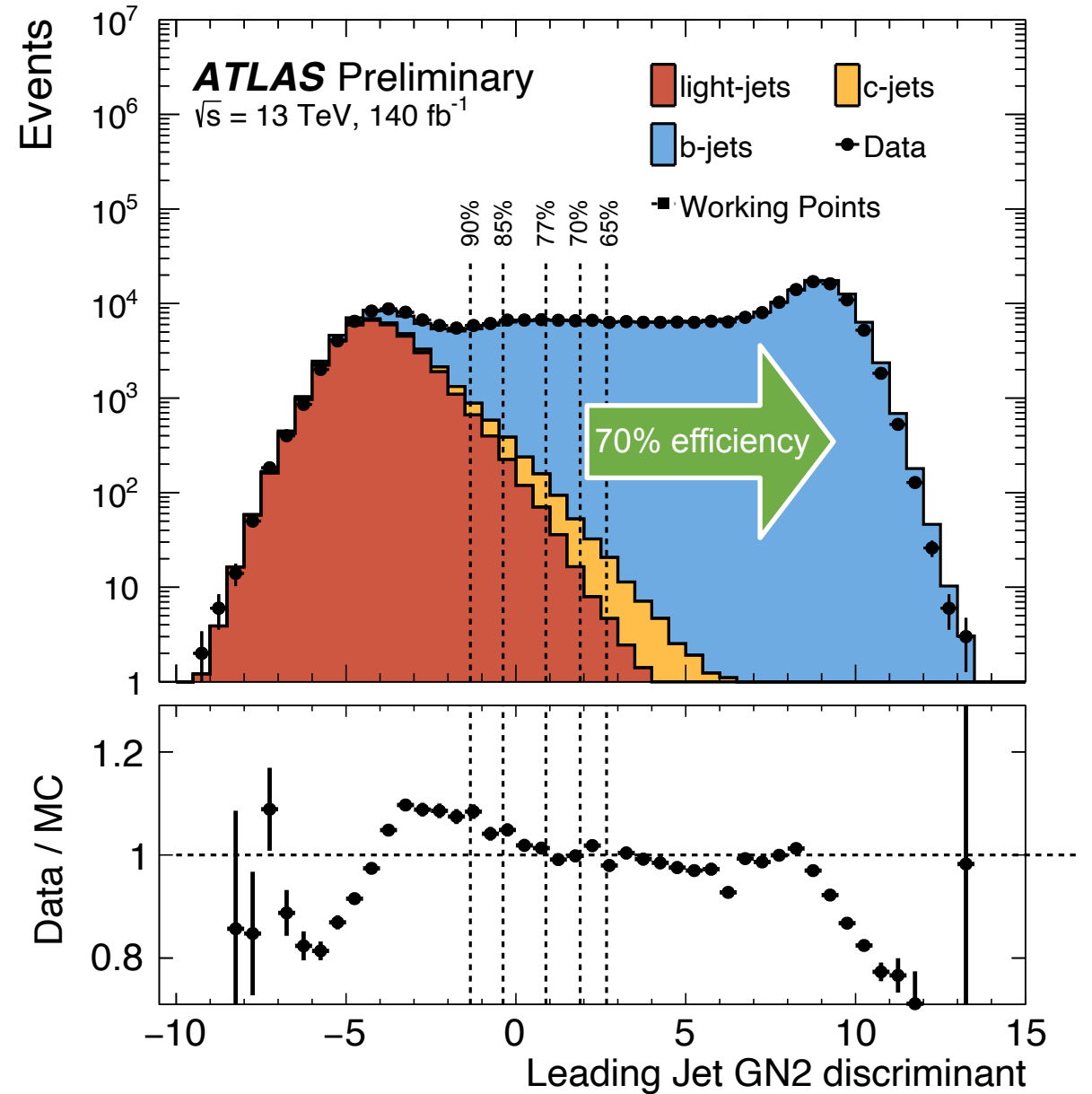
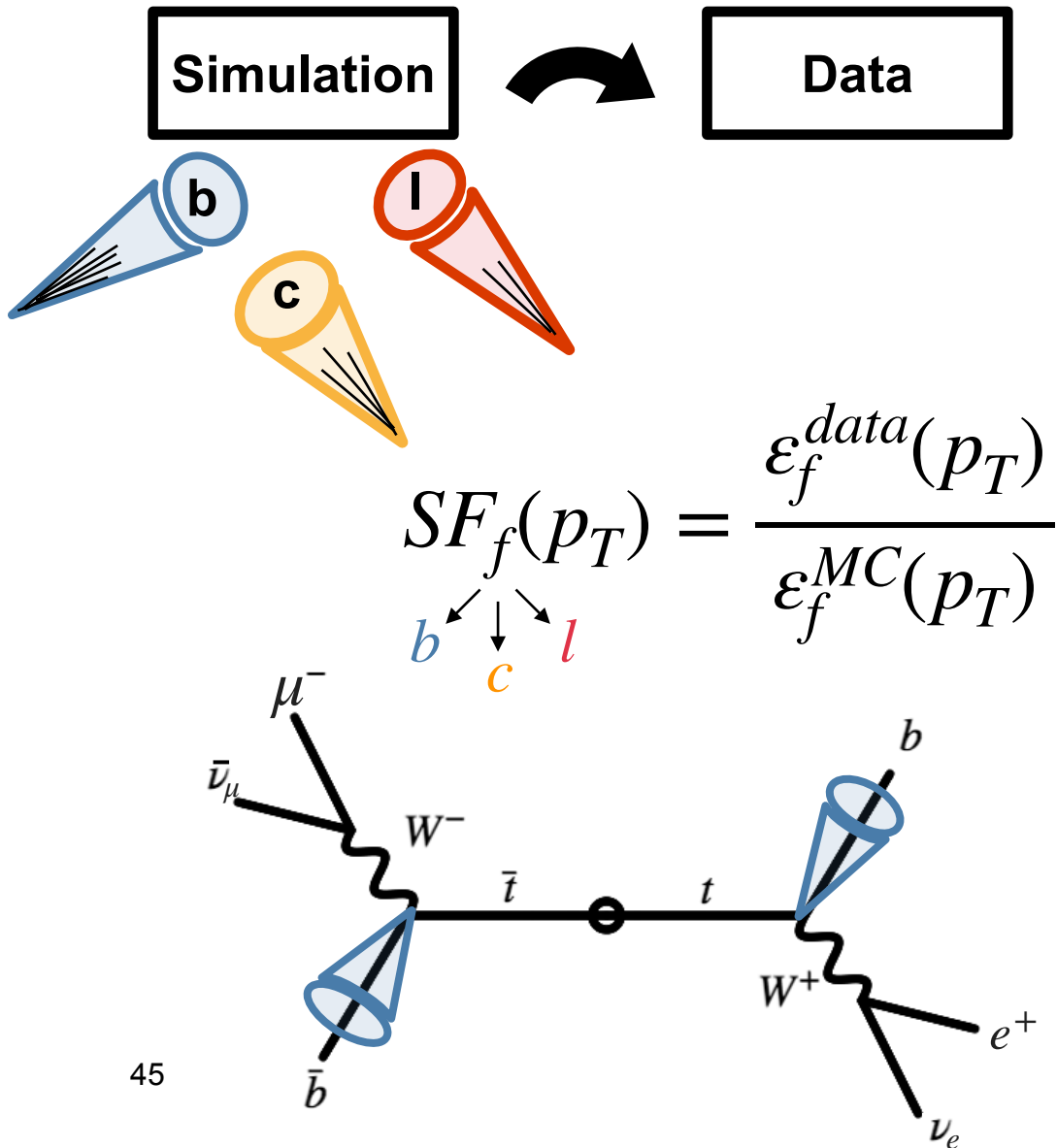


Primary vertex

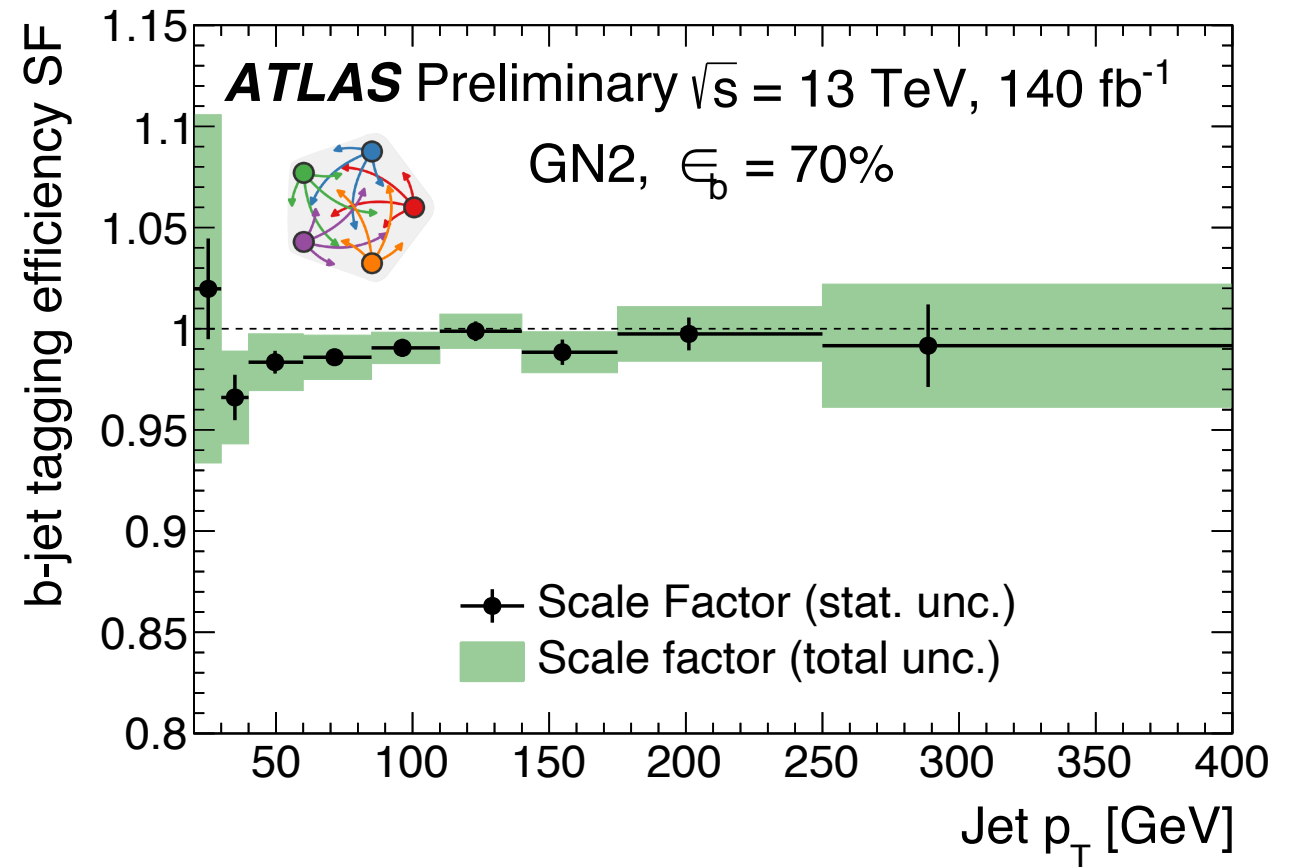
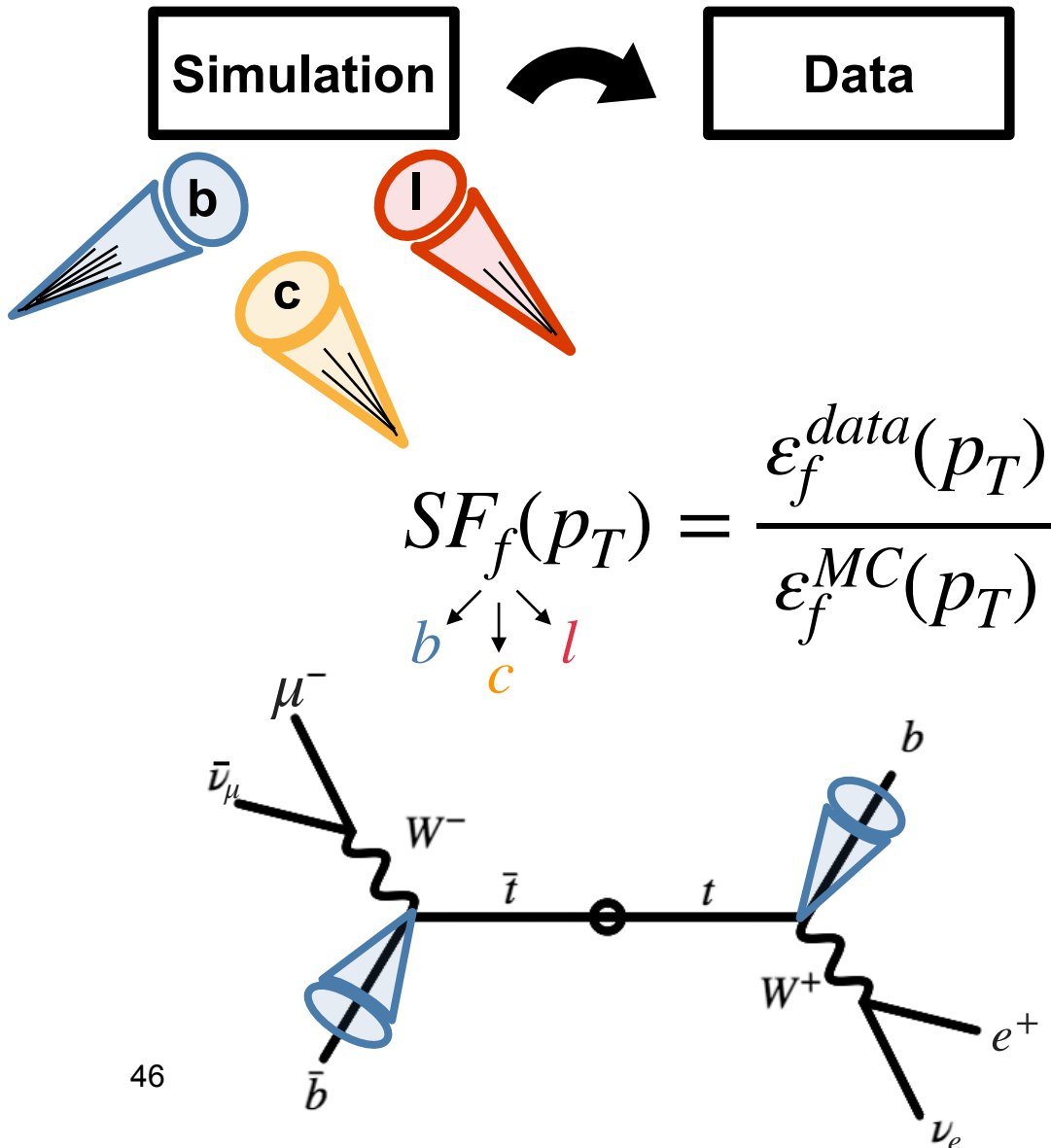


- Pileup
- Primary
- From B
- From BC

Calibration

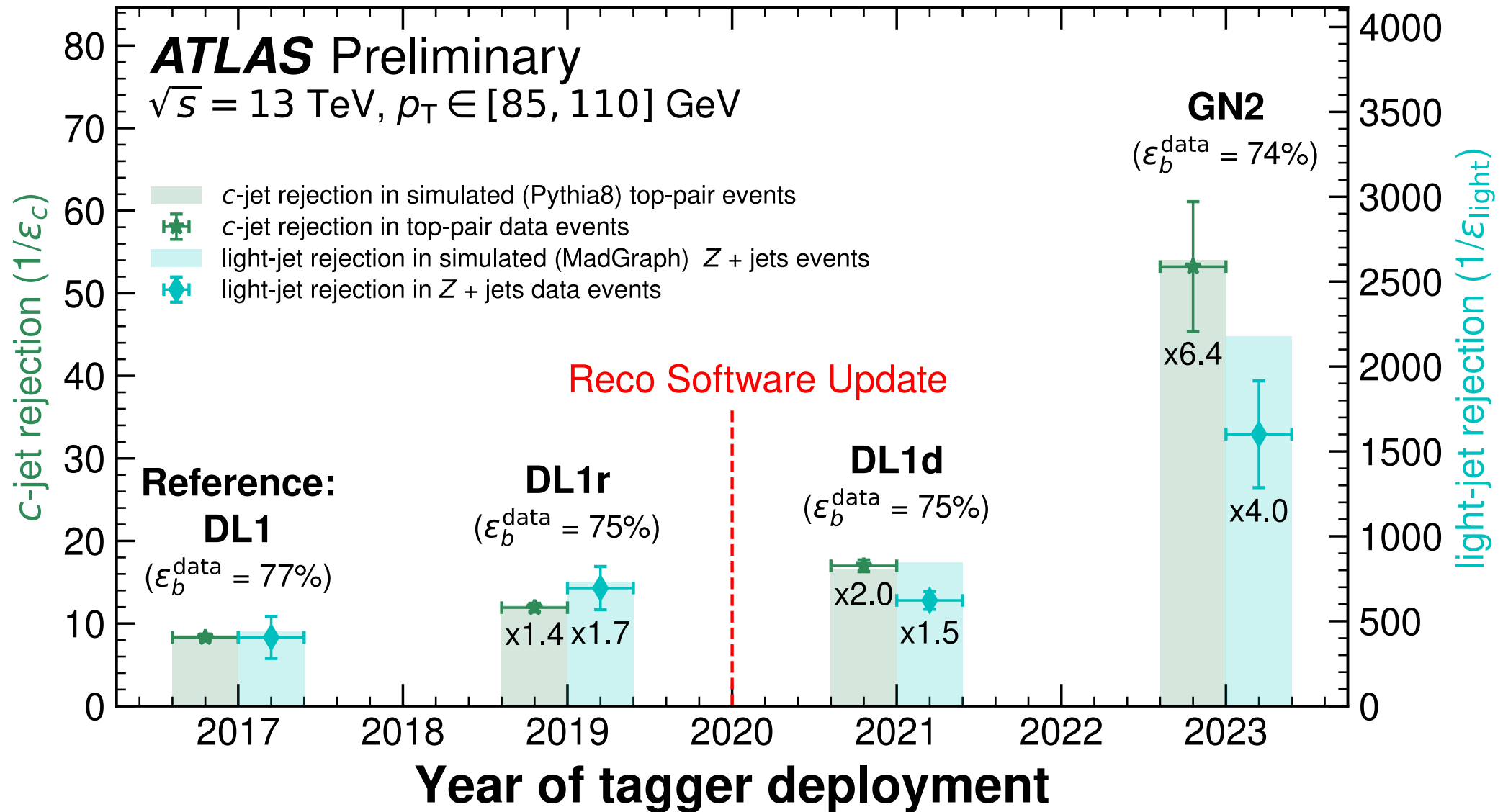


Calibration



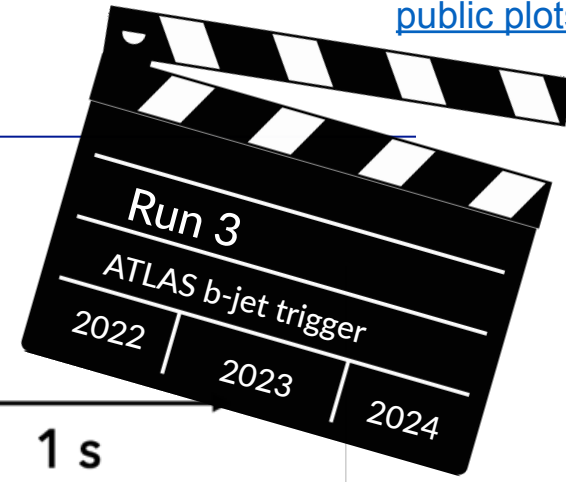
% level non-closure

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

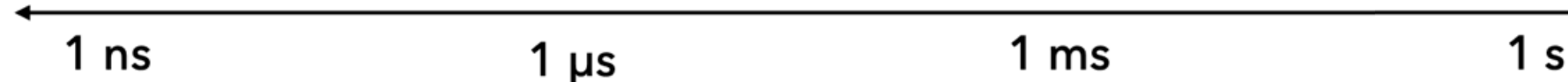


Impacting the physics...

Which events we save?



Compute Latency



40 MHz



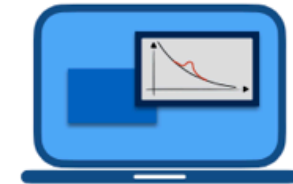
L1 Trigger

100 kHz



High-Level Trigger

1 kHz



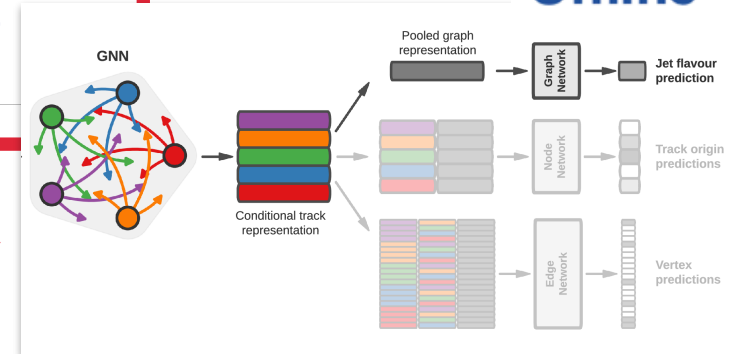
Offline

2022: DL1d

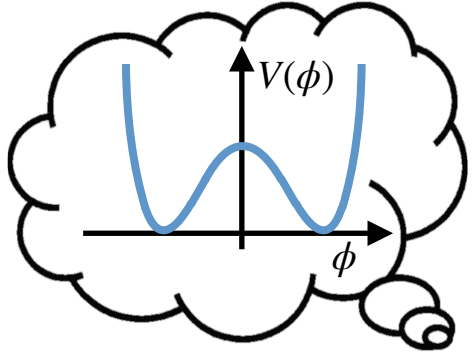
2023: GN1

2024: GN2

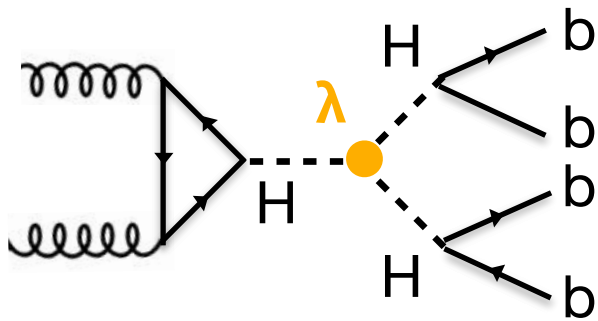
New b-taggers in trigger each year of Run 3!



FTAG in trigger



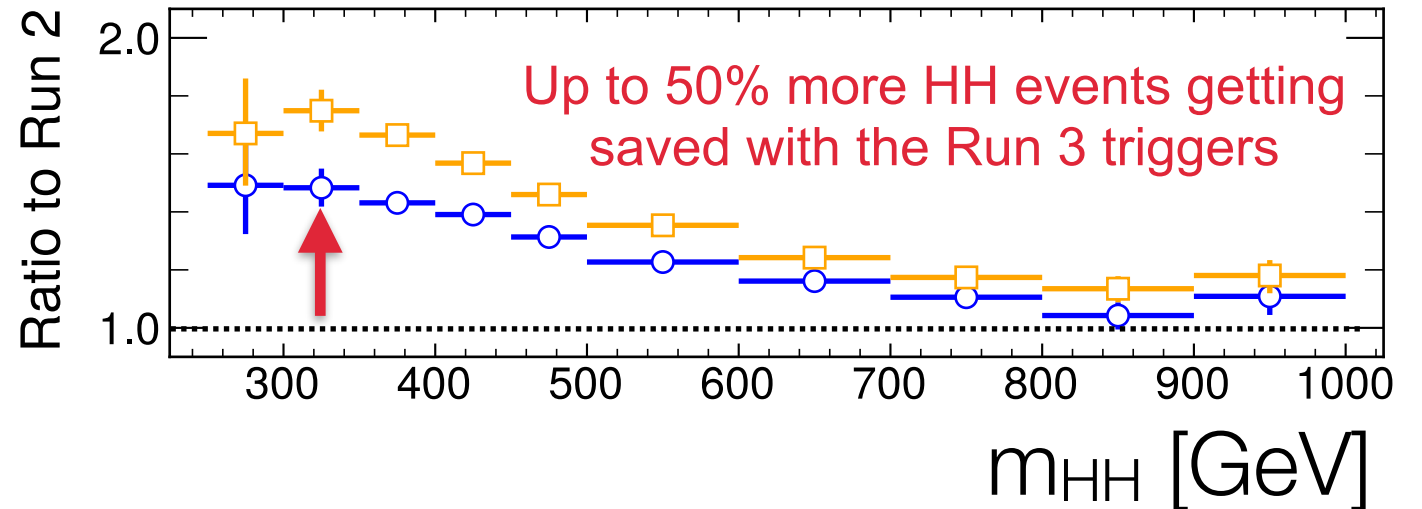
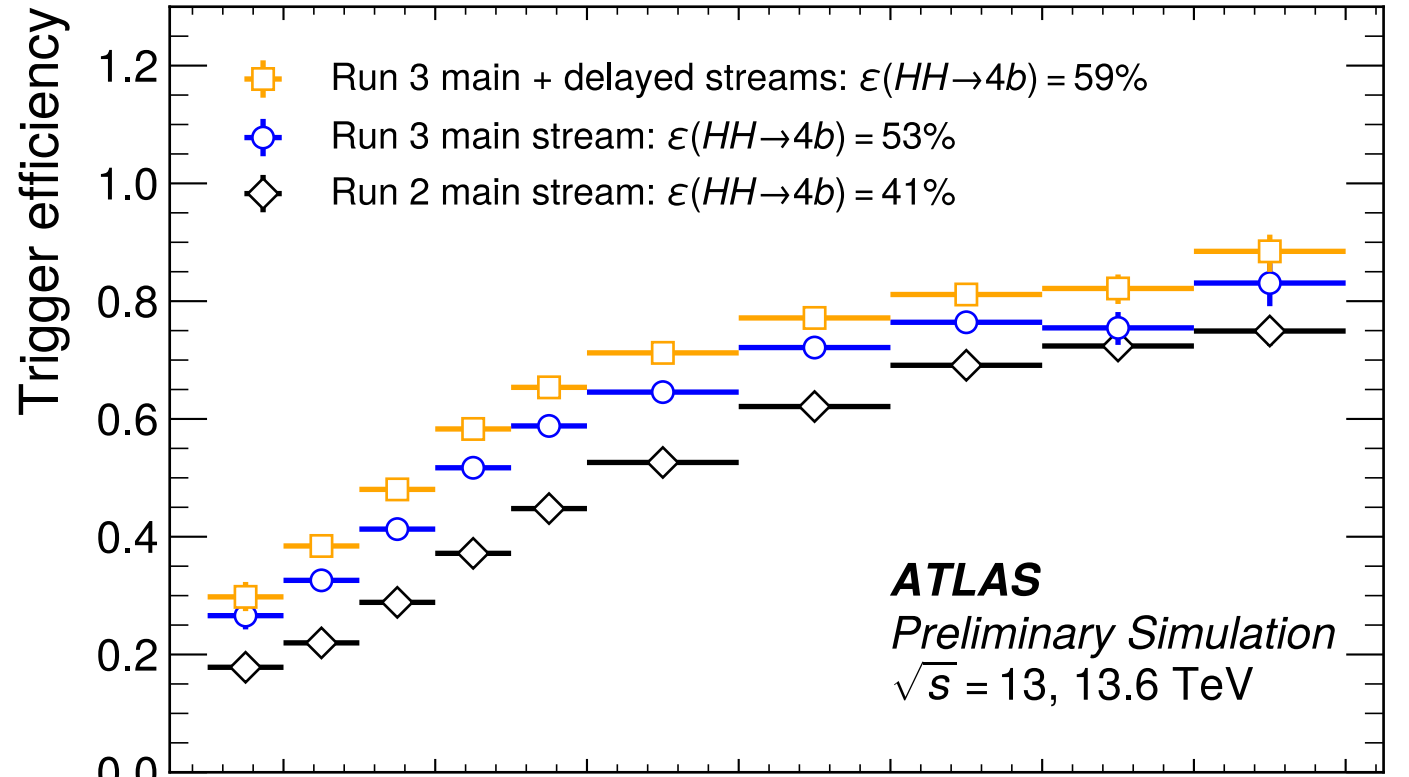
Ex: HH4b efficiency



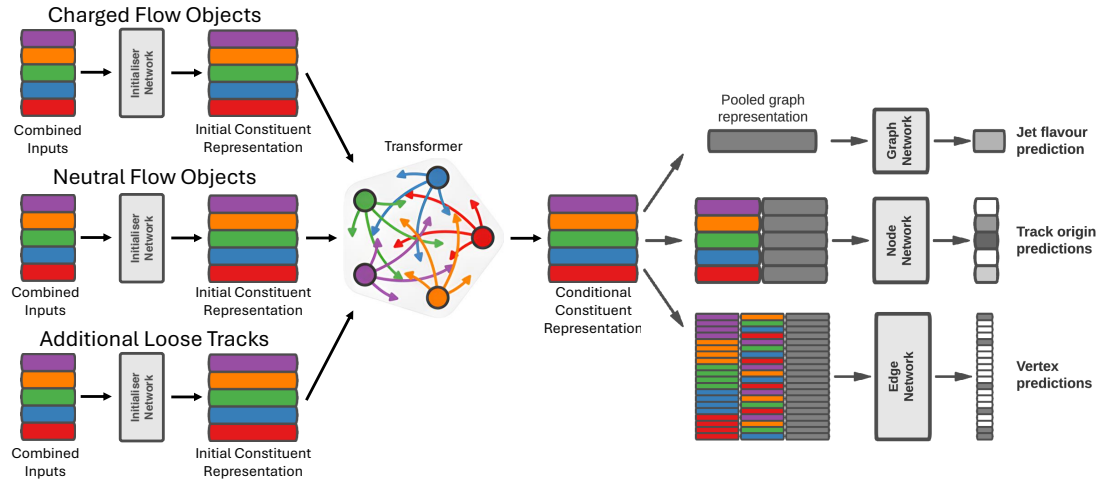
2022: DL1d

2023: GN1

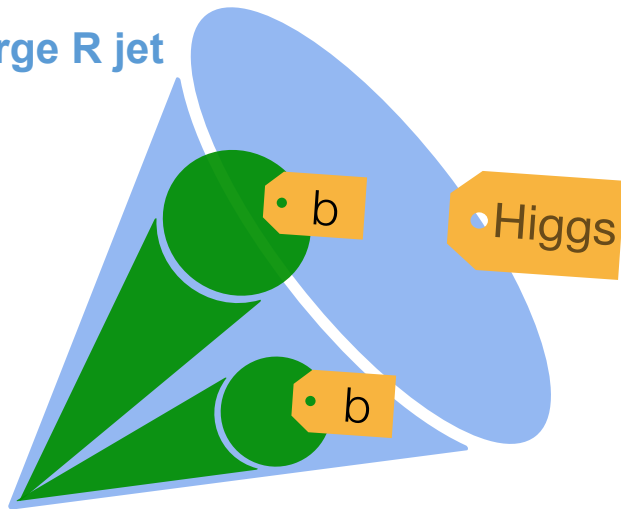
2024: GN2



Versatile: Xbb tagging

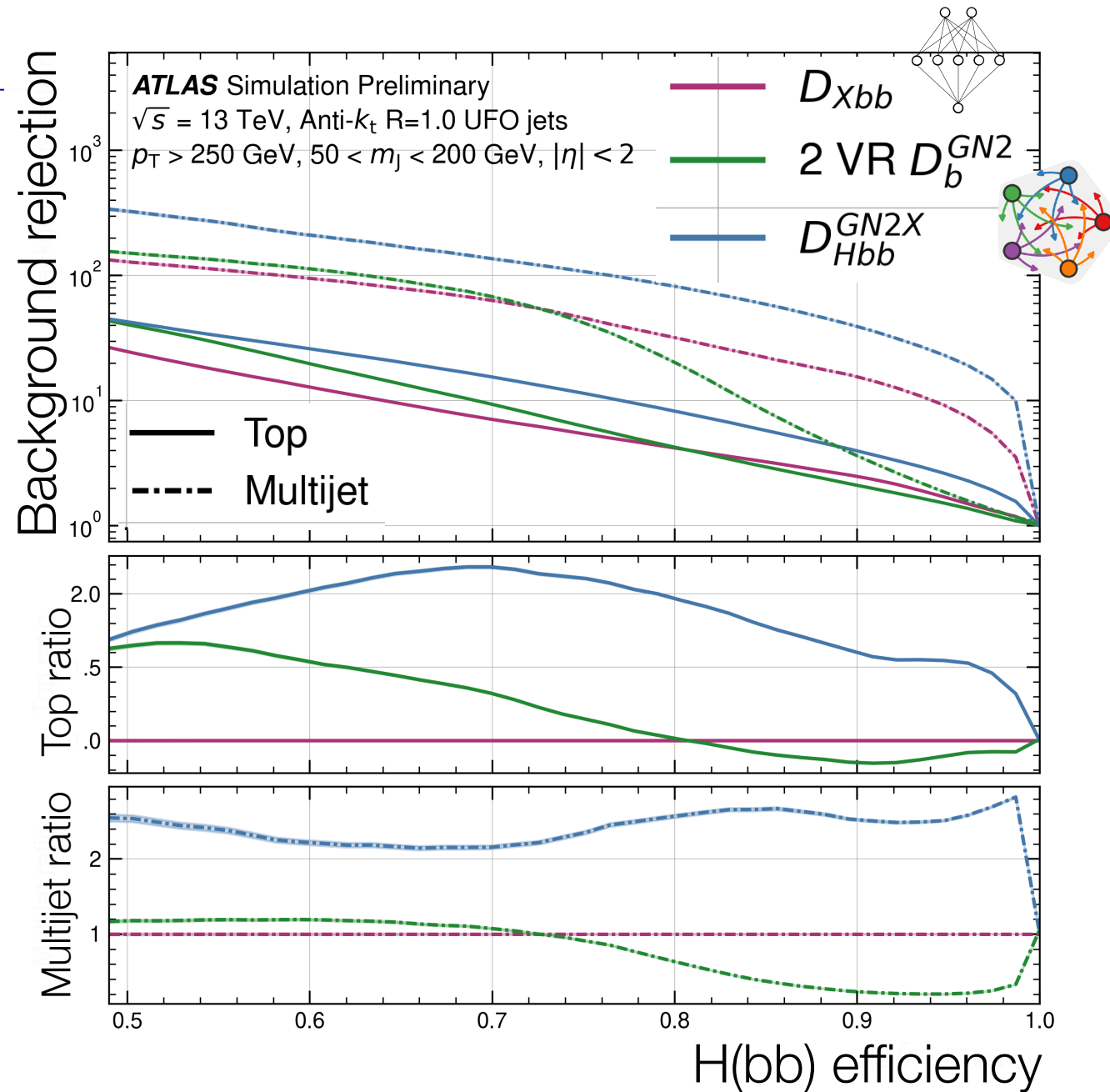


Large R jet

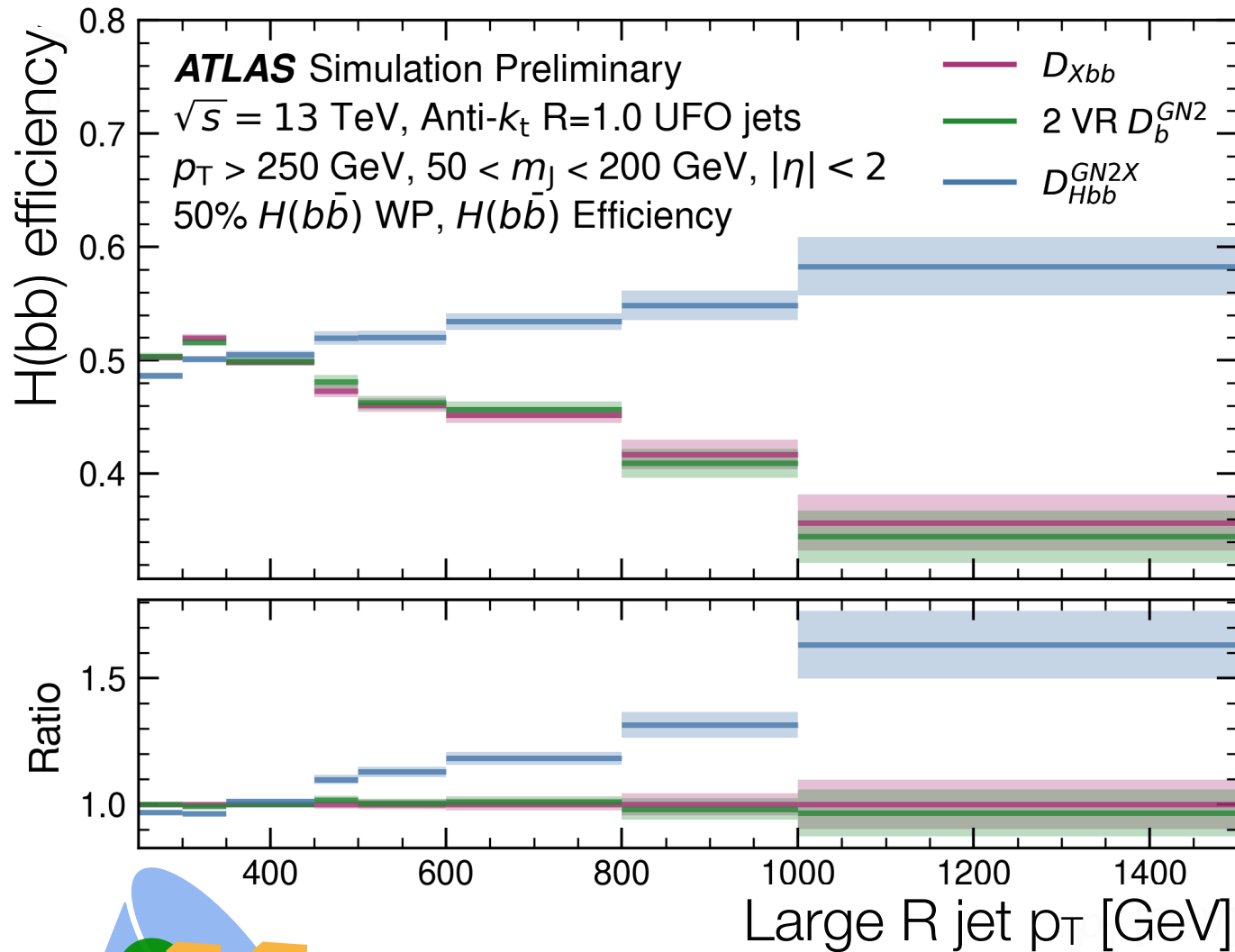


Variable radius
track jets

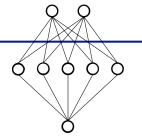
50



Versatile: Xbb tagging



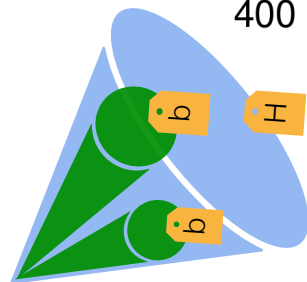
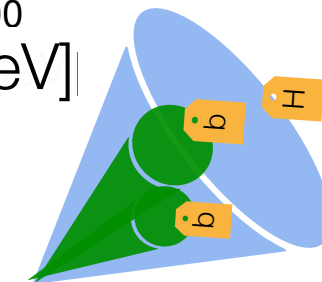
baseline NN



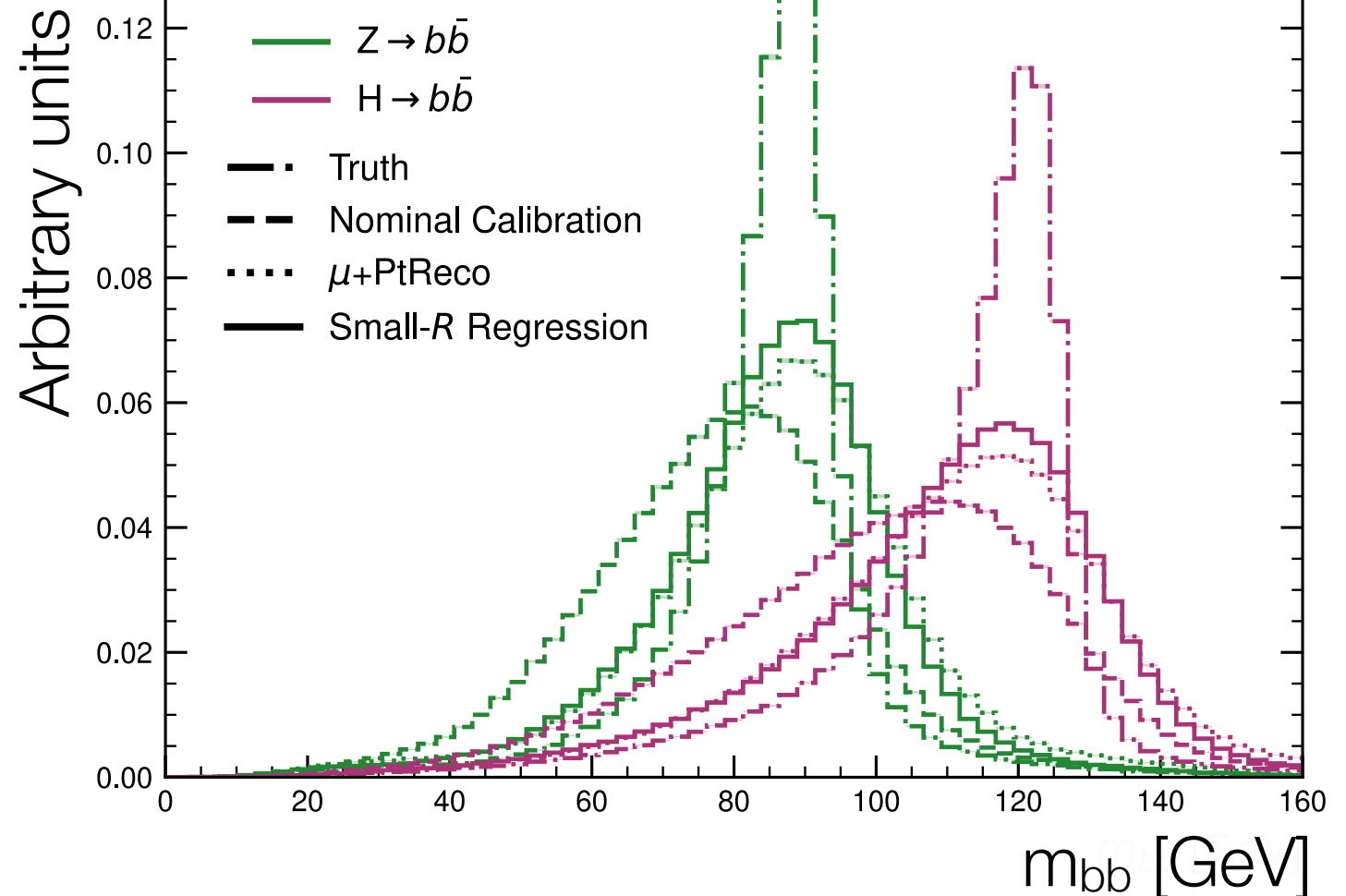
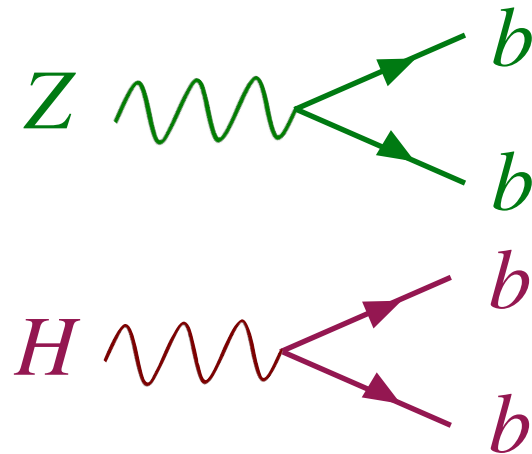
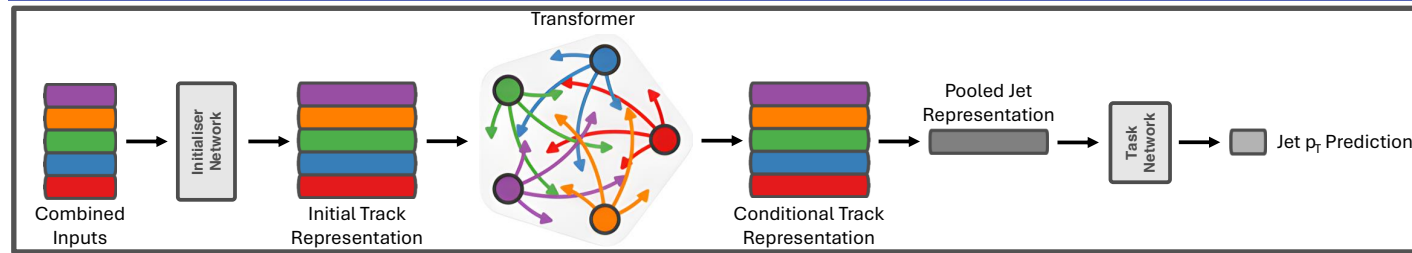
tagging subjects

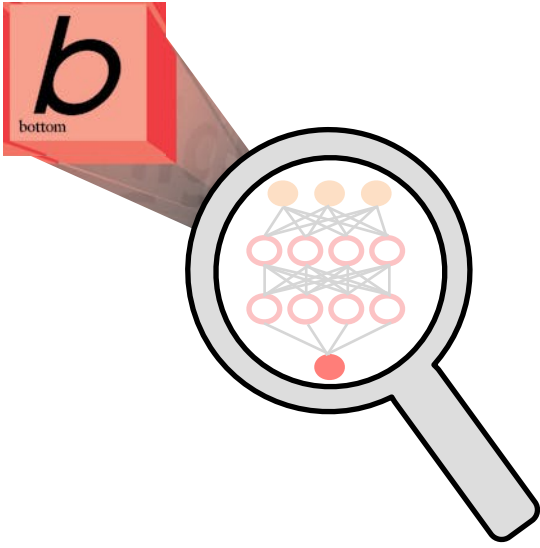


large-R jet transformer

Increasing Higgs p_T 

Versatile: regression





Future Outlook

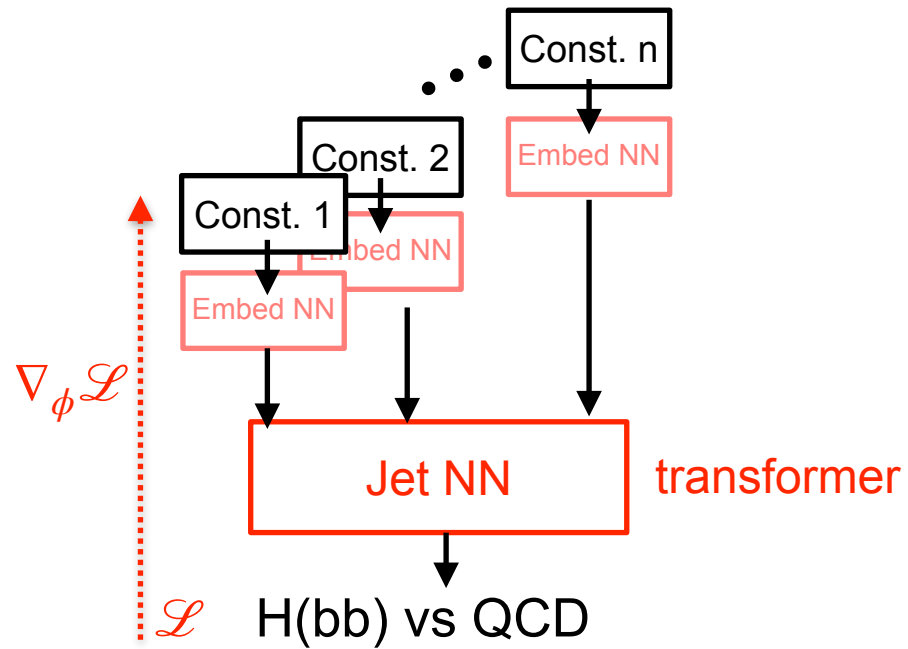
- End-to-end optimization
- Use with novel calibrations

Transformer-era

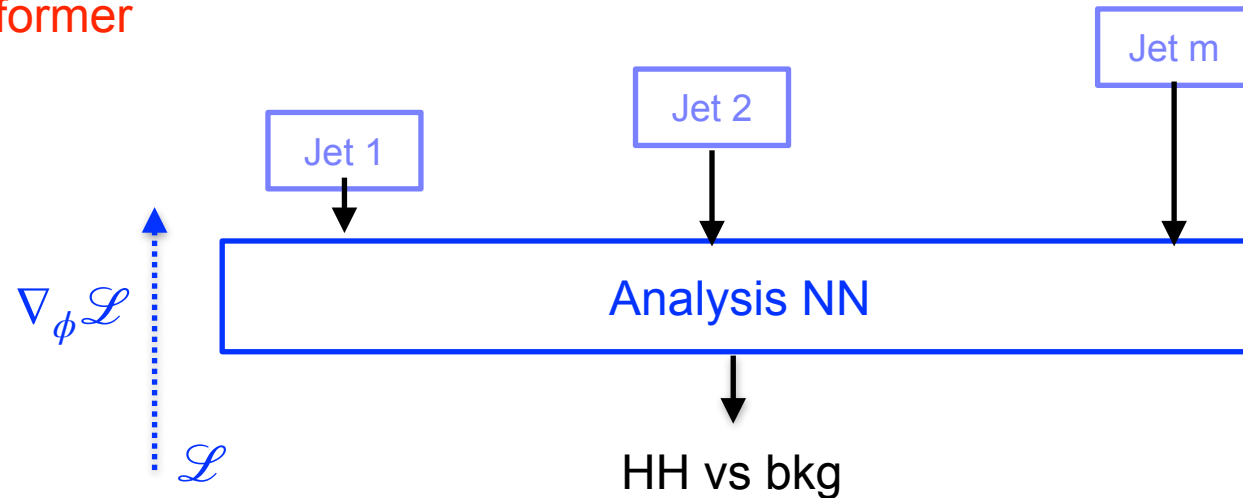
Deep Learning in FTAG

Traditional Analysis

Step 1: Train the tagger

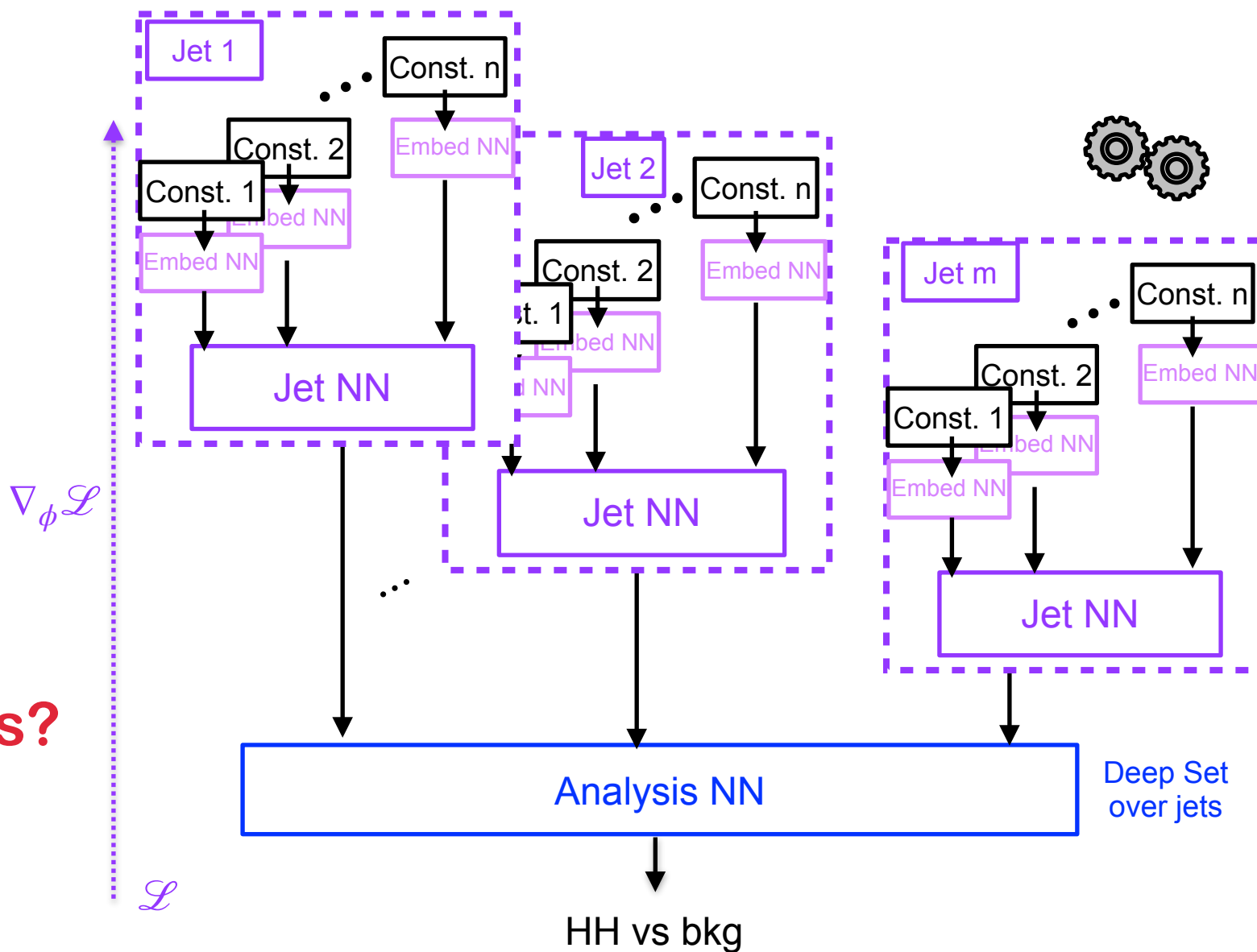


Step 2: Optimize the analysis



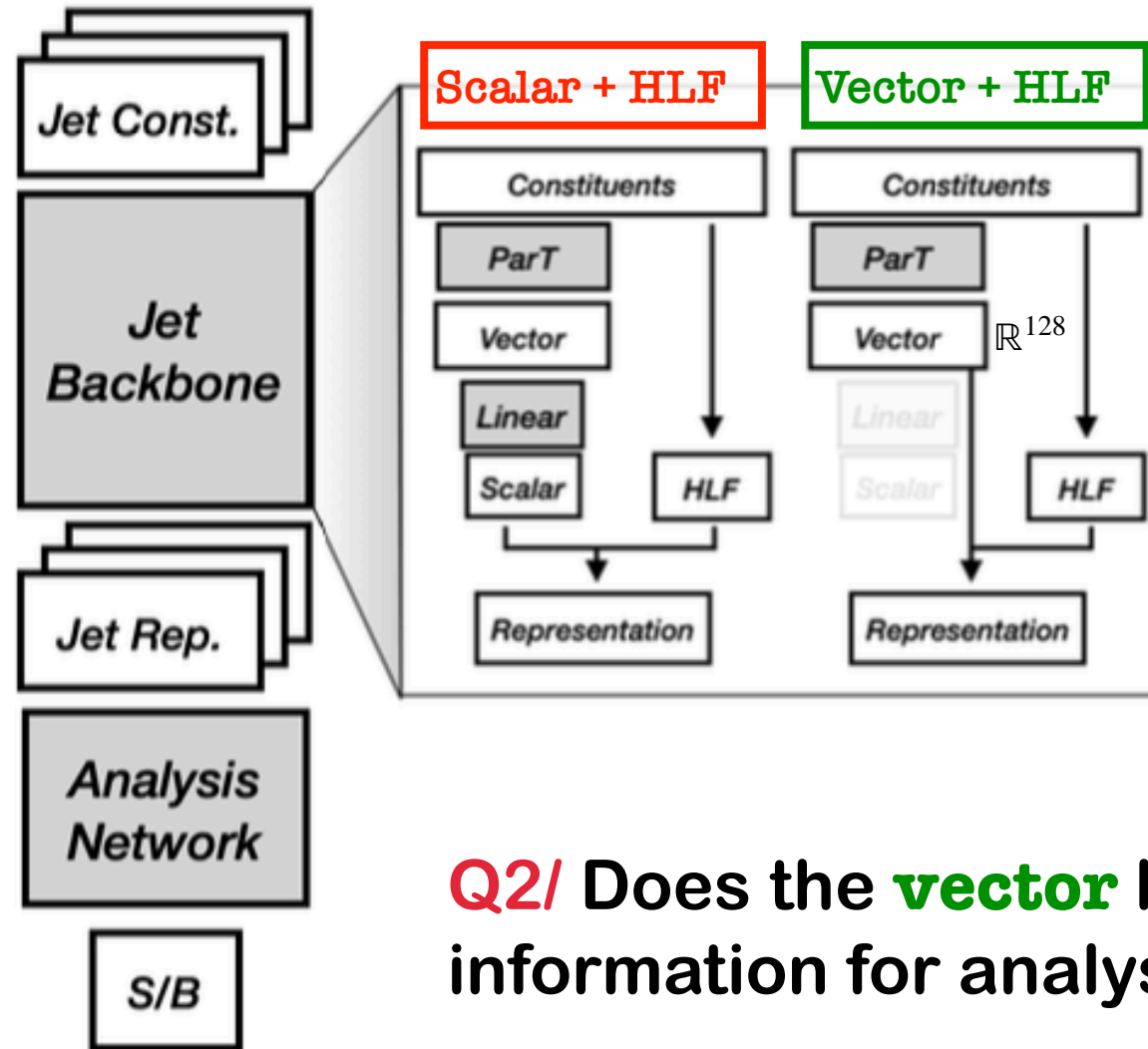
New paradigm

End to end



Q1/ Does this help us?

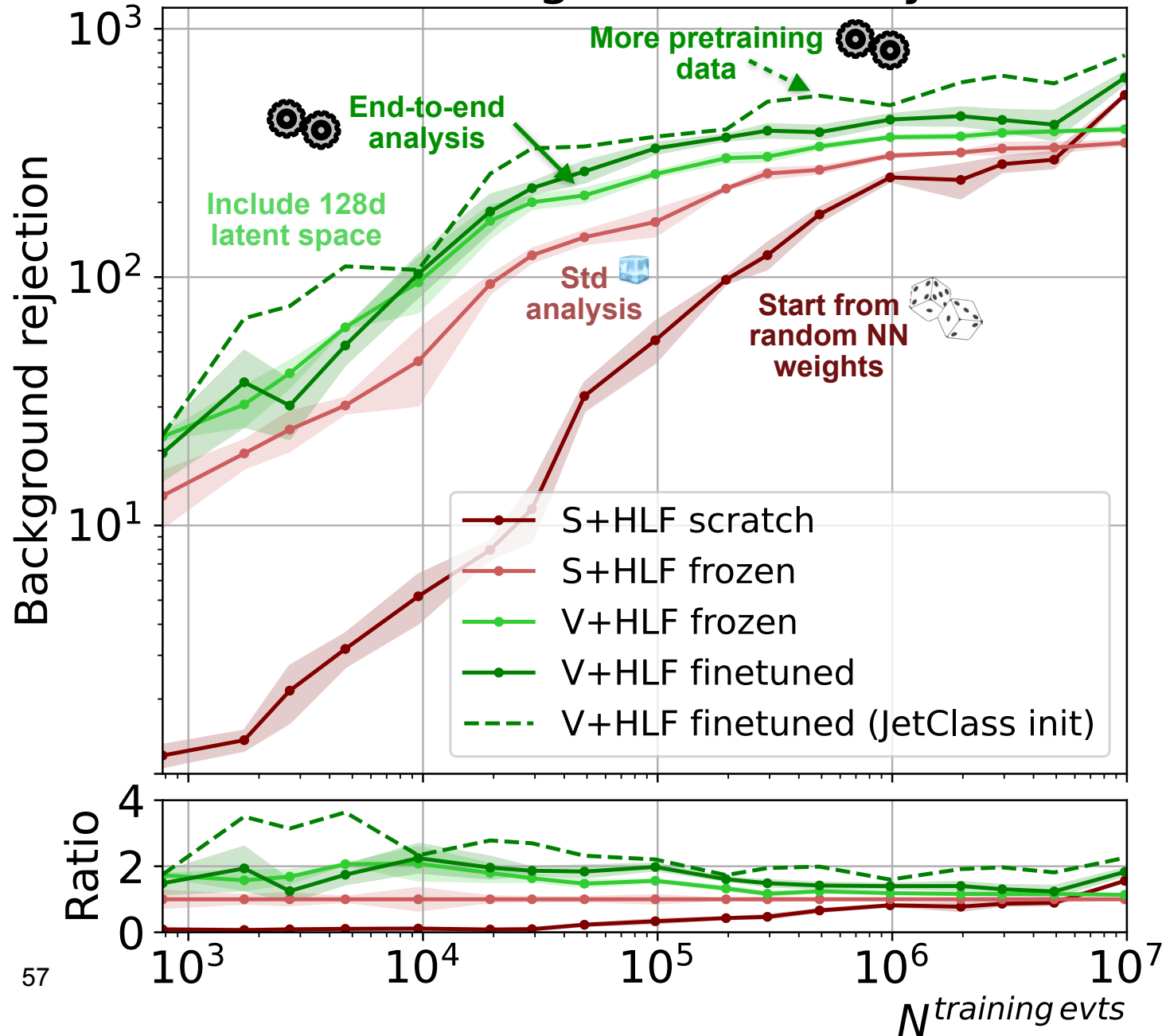
Backbone architectures



High Level Features (HLF):
 $\rho, \tau, \eta, \varphi, m$, soft drop m

Q2/ Does the **vector** latent dimensions hold extra information for analyses?

90% signal efficiency



A better Higgs tagger helps analysis performance

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



Xbb + 4-vec

Standard ML HEP

<

vector + 4-vec

More info from latent

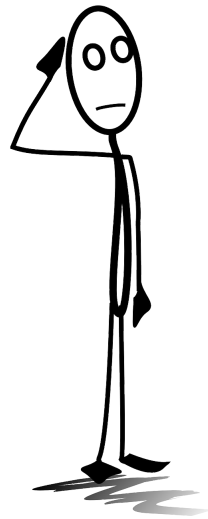
→ Need a continuous, multi-dim calibration

<

Xbb + 4-vec

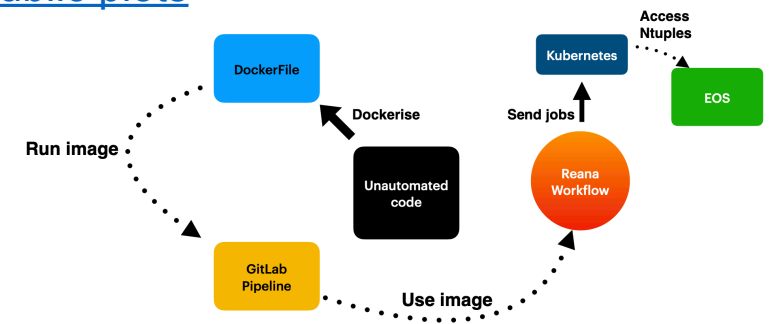
Custom “Xbb” (scalar) for each analysis

→ Need custom calibration for each analysis

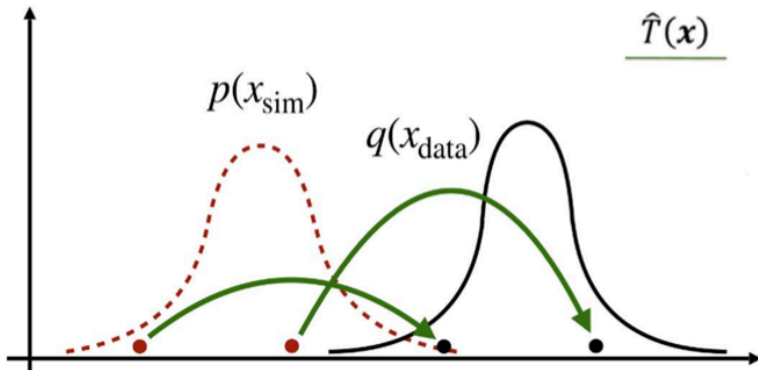


Recent work

Will need automated calibrations
[public plots](#)

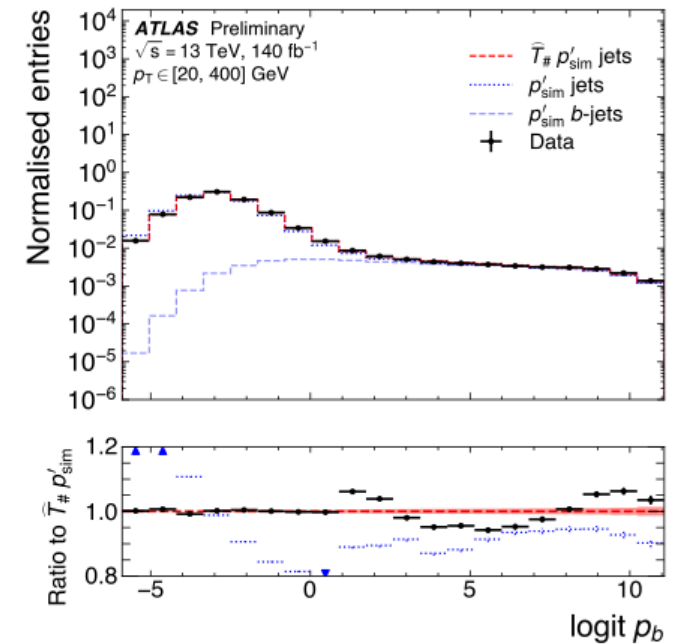
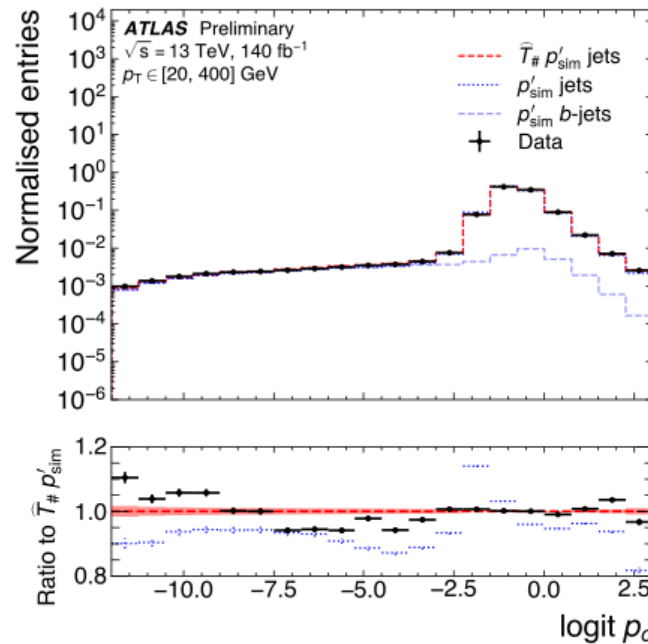
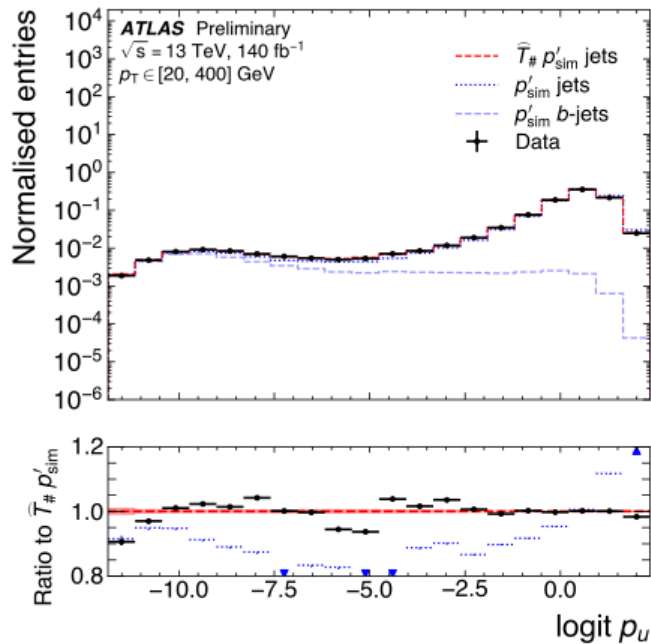


What about calibration ...

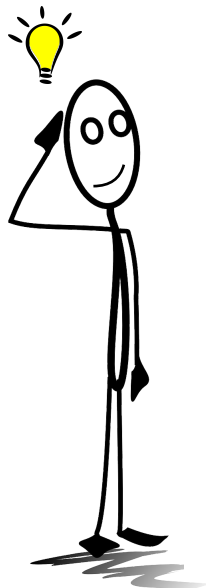


Idea: How do you have a mapping from $p_b^{\text{MC}} \rightarrow p_b^{\text{data}}$

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

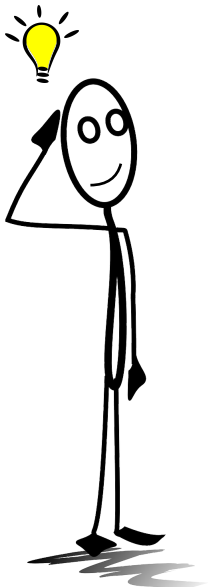
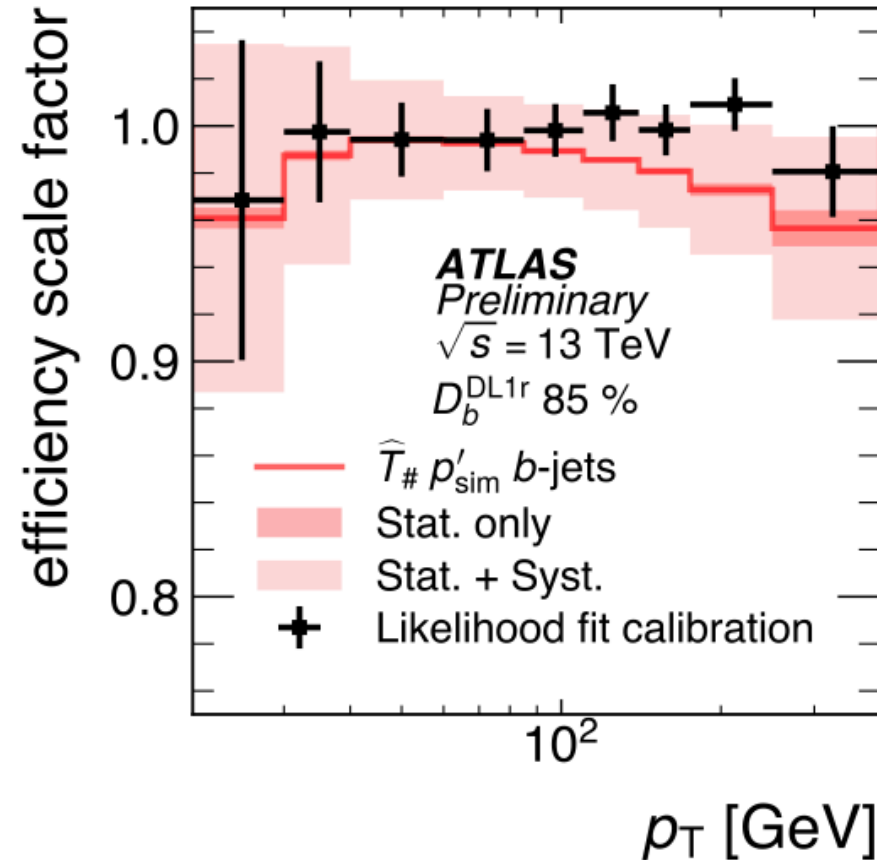
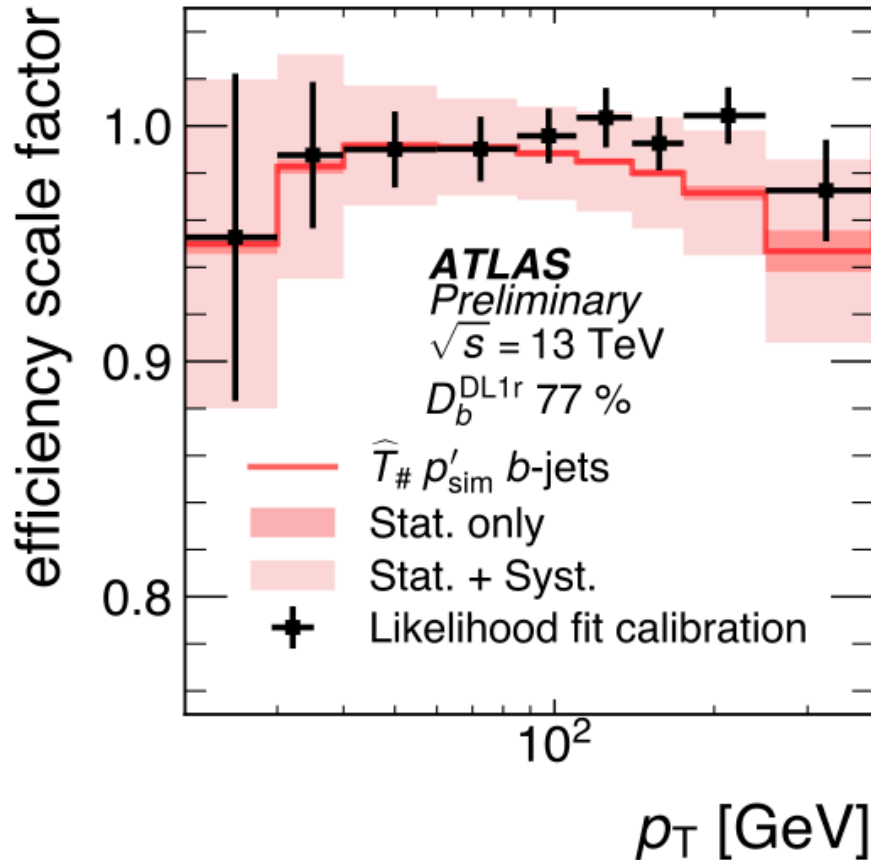


And result consistent with the standard calibration backup



What about calibration ...

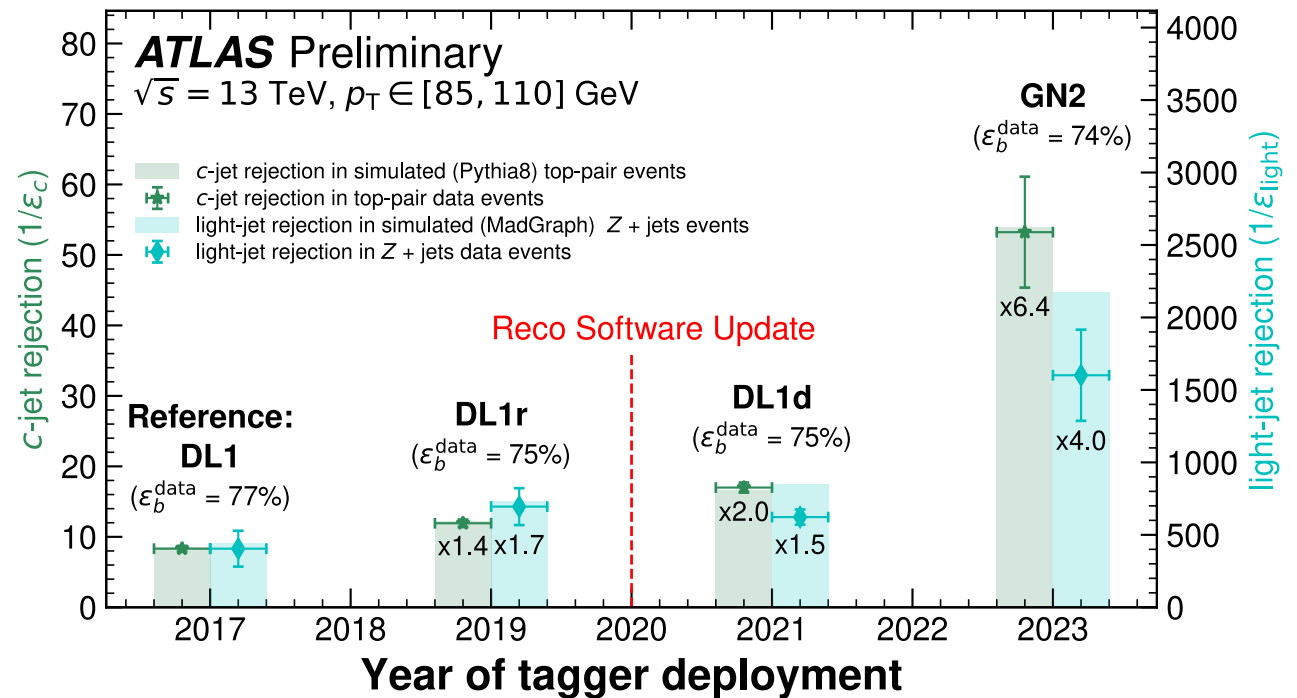
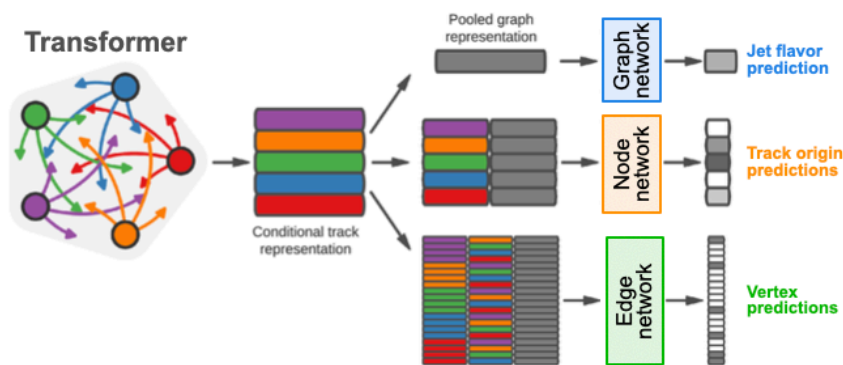
And matches the “standard calibration” with the same WPs



In Σ -mary

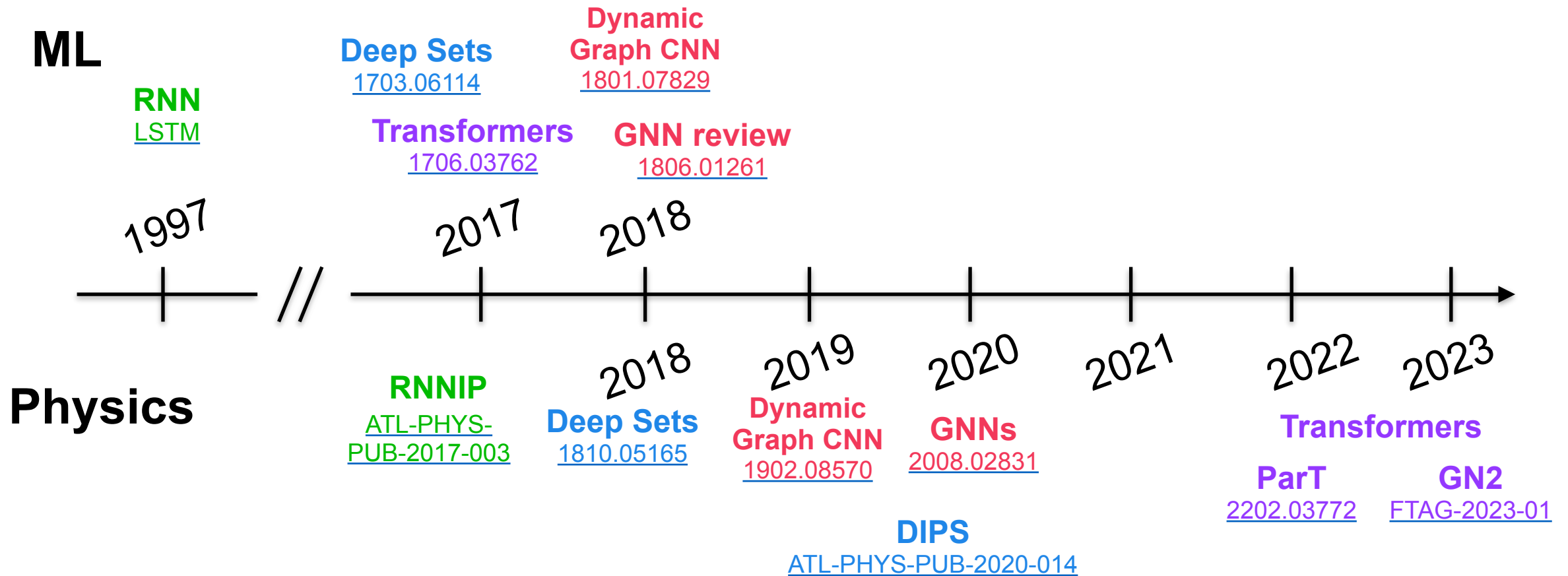
Deep Learning has transformed ATLAS FTAG our physics program in the past 7 years

- **RNNs**: model the jet as a SEQUENCE
- **Deep Sets (DIPS)**: model the jet as a SET
- **Transformers (GN2)**: Monolithic all-in-one architecture
 - Calibrated, ready for physics
 - Now collecting HH events in Run 3 triggers

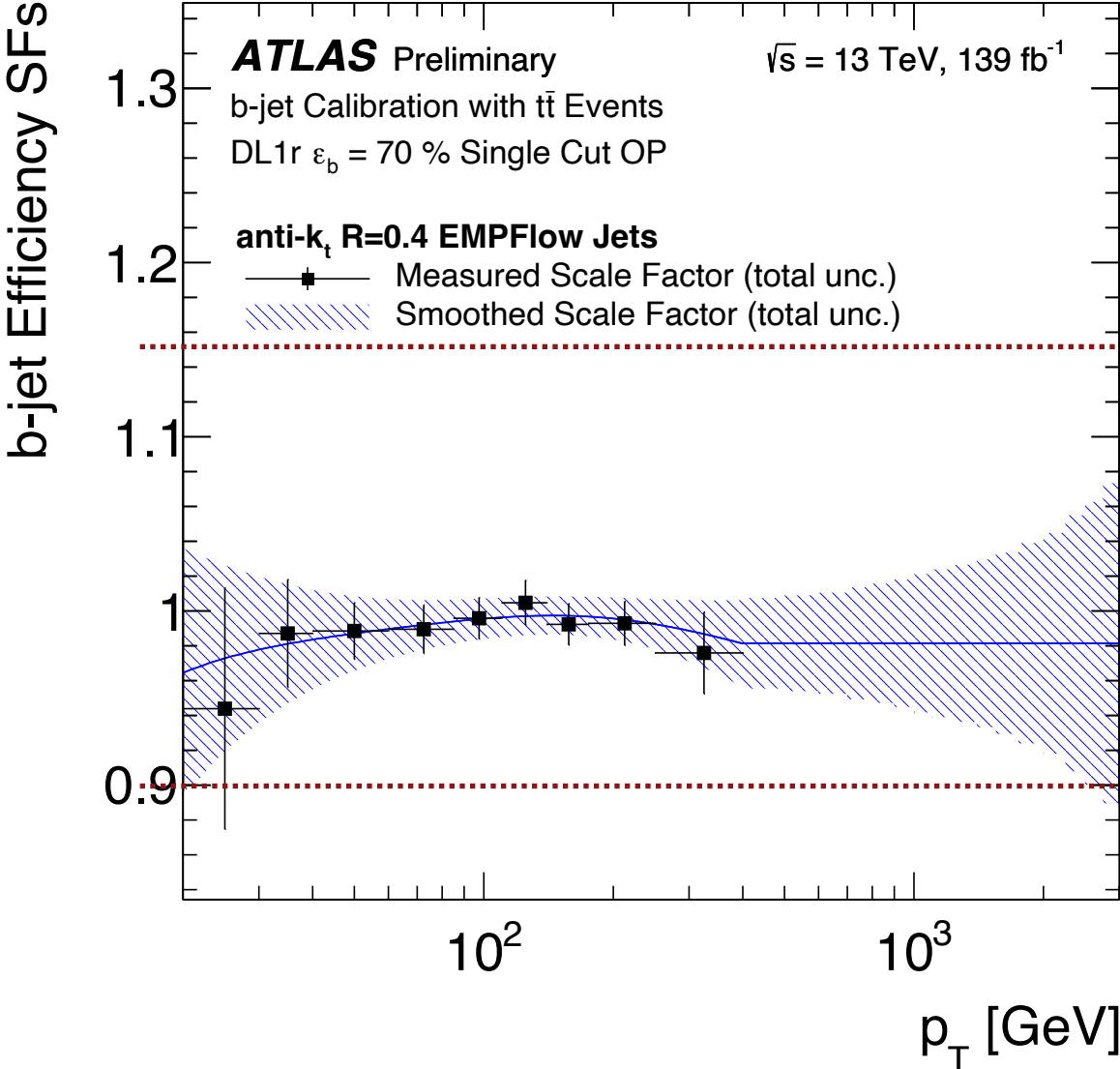


Backup

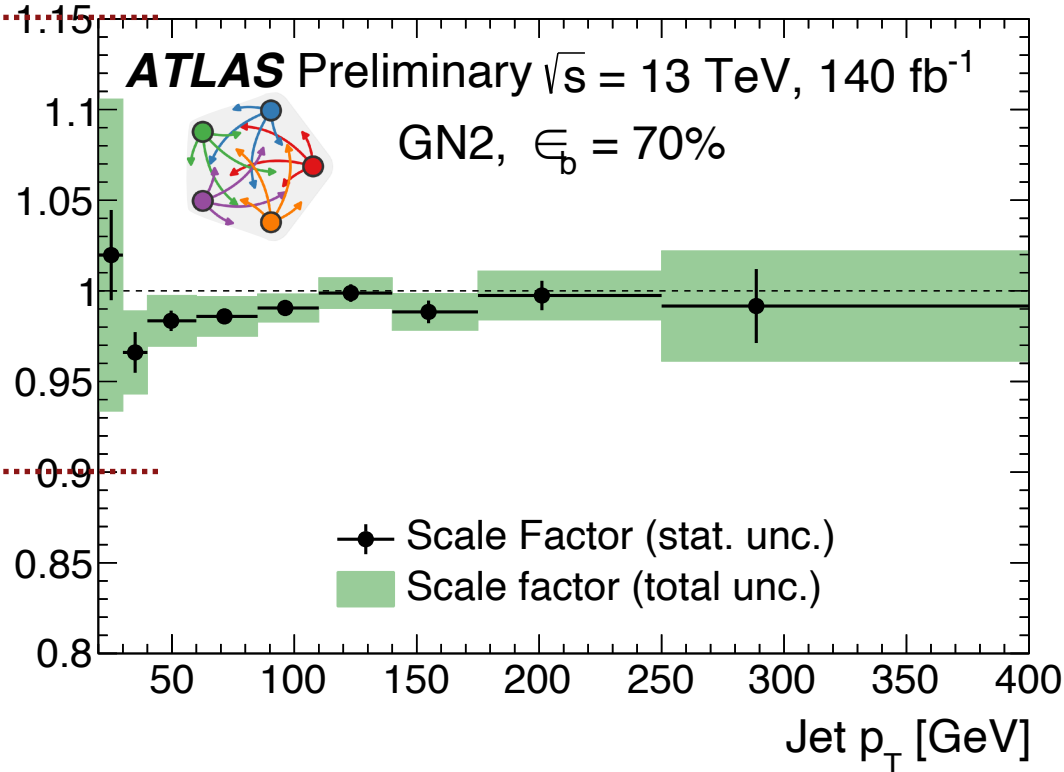
Timeline



Calibration

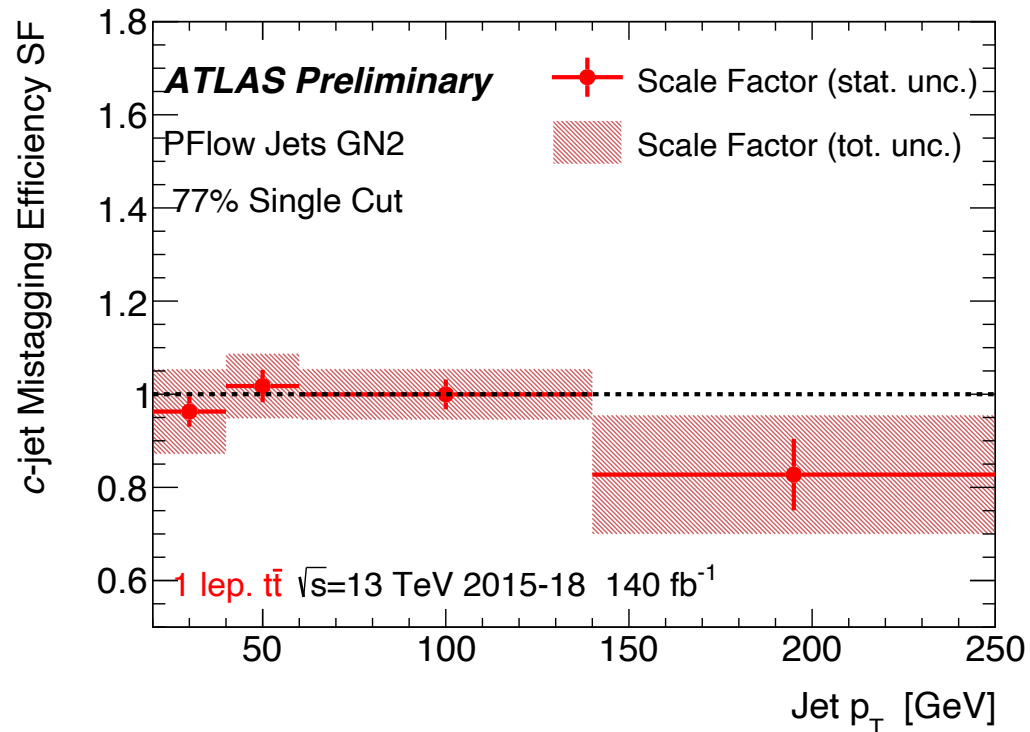
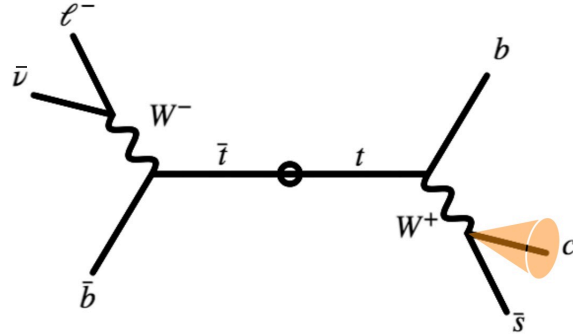


And no worse than the previous DL1r tagger 🥳



Calibration

c-jet mistag



light-jet mistag

