# TRAINING & OPTIMISATION OF LARGE TRANSFORMERS
# ATLAS CASE STUDY ON KUBEFLOW

**University of Oxford**
**MAXENCE DRAGUET**
**DANIELA BORTOLETTO**

## JET FLAVOUR TAGGING

### CLASSIFICATION TASK
**Labels**
$b, c, light, \tau$

**Inputs**
Tracks + jet information

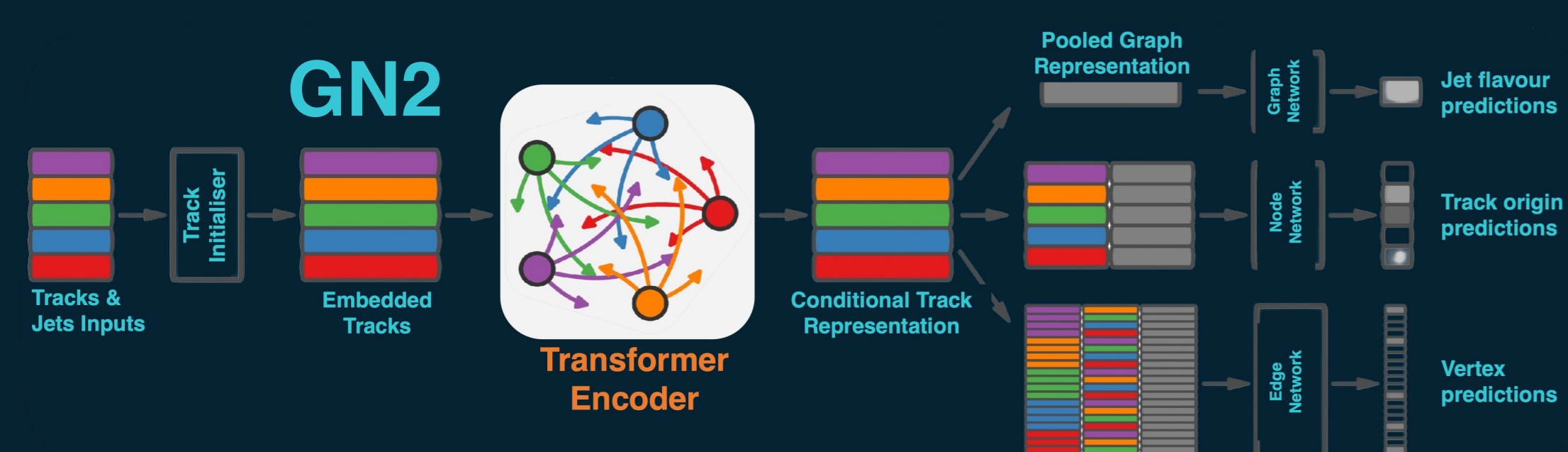**Outputs**
Per-flavour probability + Discriminant

**Used in ATLAS:**
$H\left(\to b\bar{b}\mid c\bar{c}\right)$, di-Higgs, ...

ATLAS event display of a Higgs boson decaying to two b-quarks
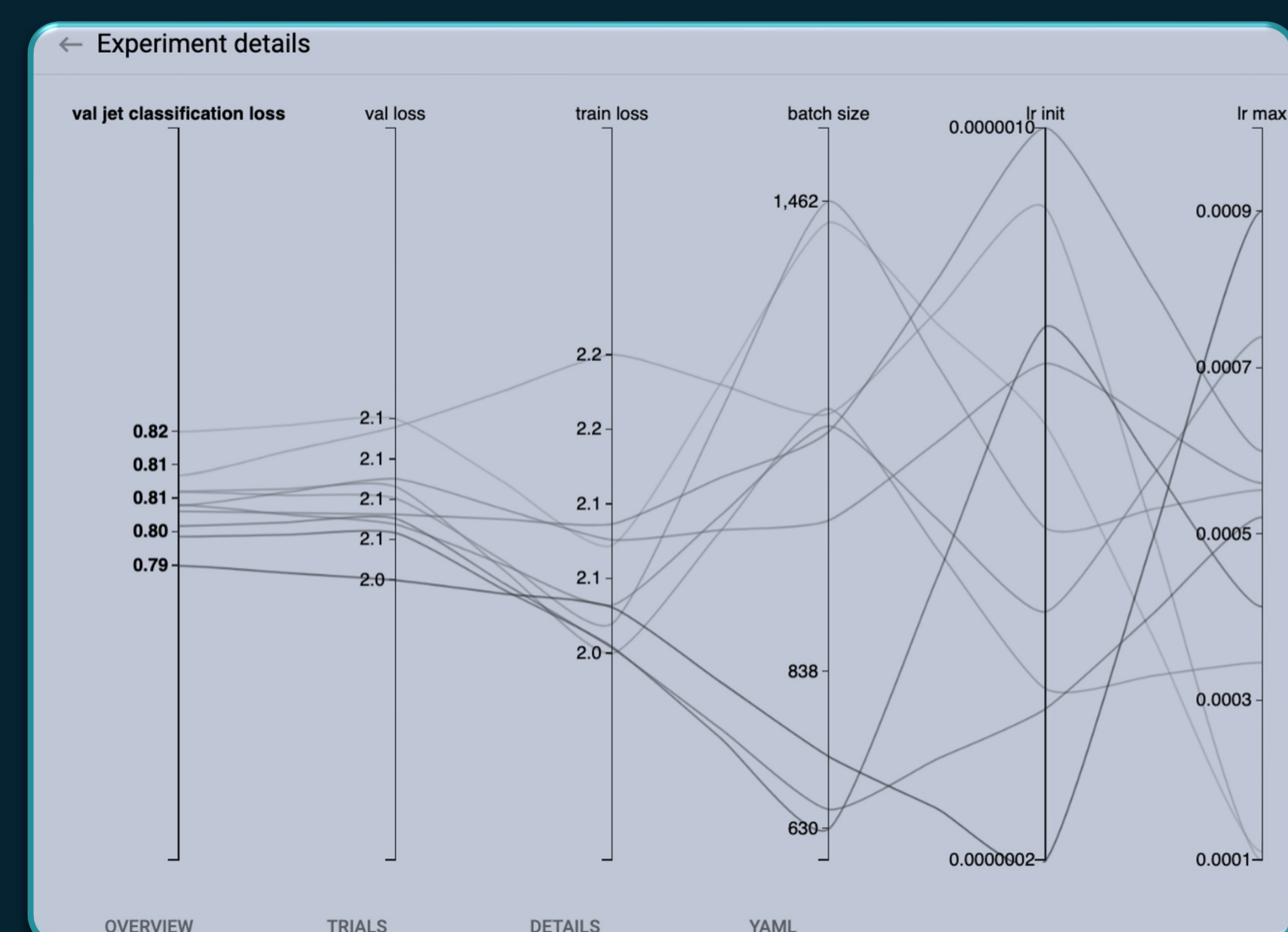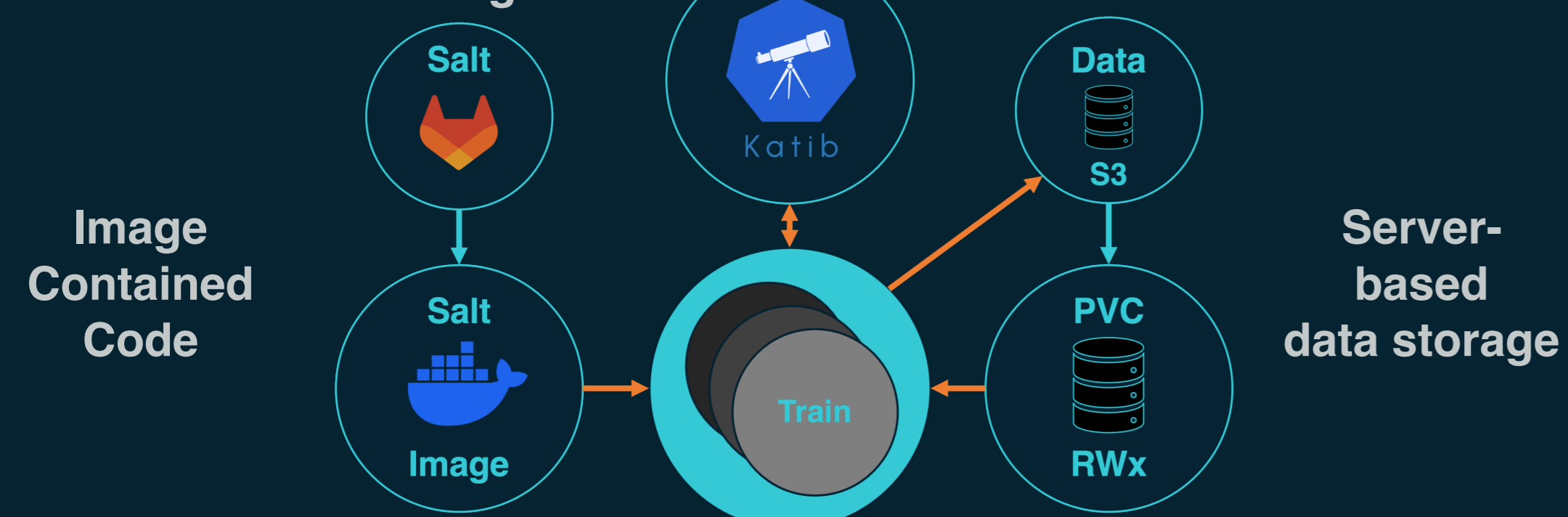
### SOTA: Multimodal Multitask Transformer Model

**GN2**

Tracks & Jets Inputs → Track Initialiser → Embedded Tracks → Transformer Encoder → Conditional Track Representation → Pooled Graph Representation → Graph Network → Jet flavour predictions; Node Network → Track origin predictions; Edge Network → Vertex predictions

### Goal Hyperparameter Optimisation (HPO)

### Kubeflow

**Built on Kubernetes**
Open-source container orchestration engine

**Designed for MLOps**
Training, Inference, Katib HPO + AutoML

Salt → Salt → Image → Train
Image Contained Code
Katib
Data S3 → PVC → RWx
Server-based data storage

← Experiment details

### EFFICIENT DESIGN
1. Multi-platform
2. Resource Usage
3. Monitoring
4. AutoML Algorithms

## PARAMETRIZATION MATTERS!

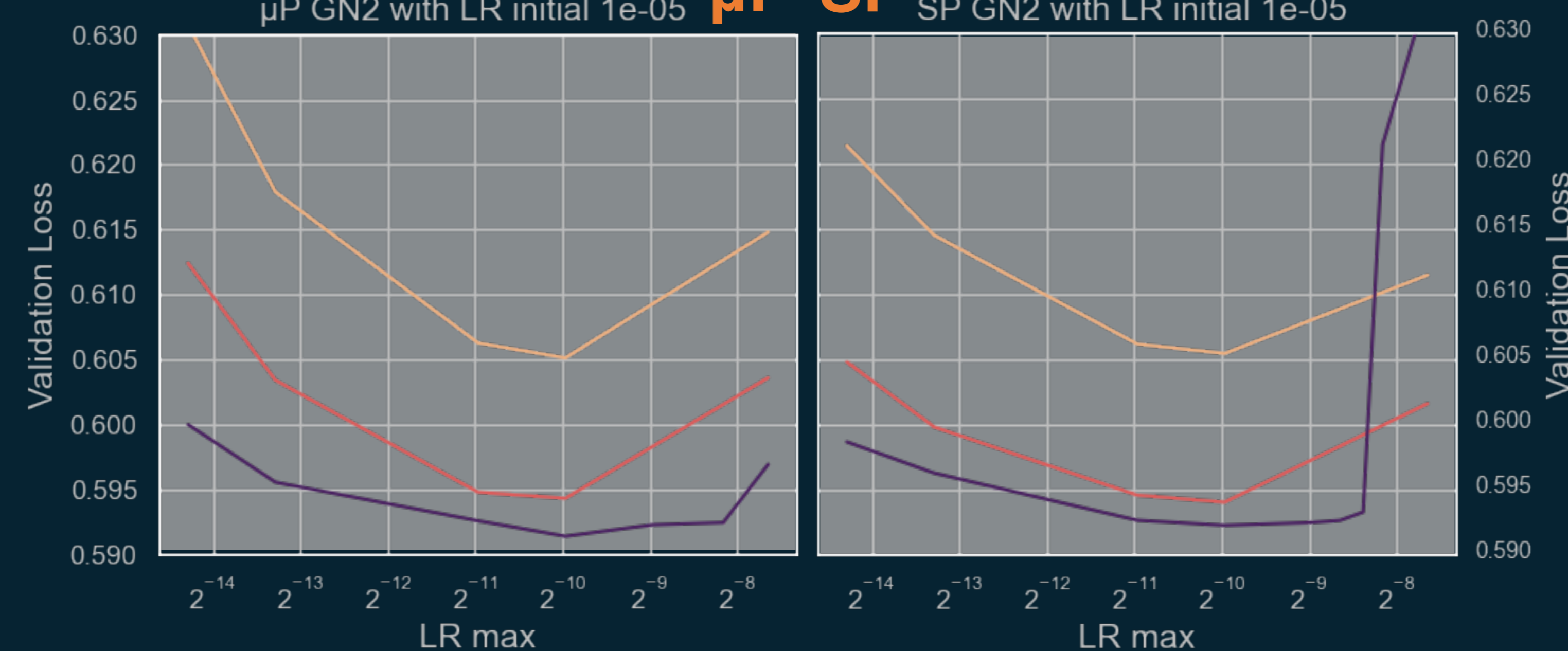➤ **Standard Parametrization (SP)**

LeCun et al; 1998

**Initialisation:** $W^{L_{in}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{in}}^{in}}\right), W^{L_{hid}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{hid}}^{in}}\right), W^{L_{out}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{out}}^{in}}\right), b^{L\cdots} = 0$

**SGD & Adam LR:** $\forall$ weights $\eta$

➤ **Maximal Update Parametrization (μP)**

Yang et Hu; 2022

**Initialisation:** $W^{L_{in}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{in}}^{in}}\right), W^{L_{hid}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{hid}}^{in}}\right), W^{L_{out}} \sim \mathcal{N}\left(0, \frac{1}{d_{L_{out}}^{in} \times d_{out}^{in}}\right), b^{L\cdots} = 0$

**SGD LR:** $\eta_{W_{L_{in}}} = \eta_{b_{L\cdots}} = \eta \times d_{L_{in}\mid L\cdots}^{out}$ ; $\eta_{W_{L_{hid}}} = \eta$ ; $\eta_{W_{L_{out}}} = \eta / d_{L_{out}}^{in}$

**Adam LR:** $\eta_{W_{L_{in}}} = \eta_{b_{L\cdots}} = \eta$ ; $\eta_{W_{L_{hid}}} = \eta / d_{L_{hid}}^{in}$ ; $\eta_{W_{L_{out}}} = \eta / d_{L_{out}}^{in}$

**ATLAS** Simulation Preliminary, μP GN2 with LR initial 1e-05 — **μP**

**SP** — **ATLAS** Simulation Preliminary, SP GN2 with LR initial 1e-05

(Validation Loss vs LR max)

### μTransfer Algorithm
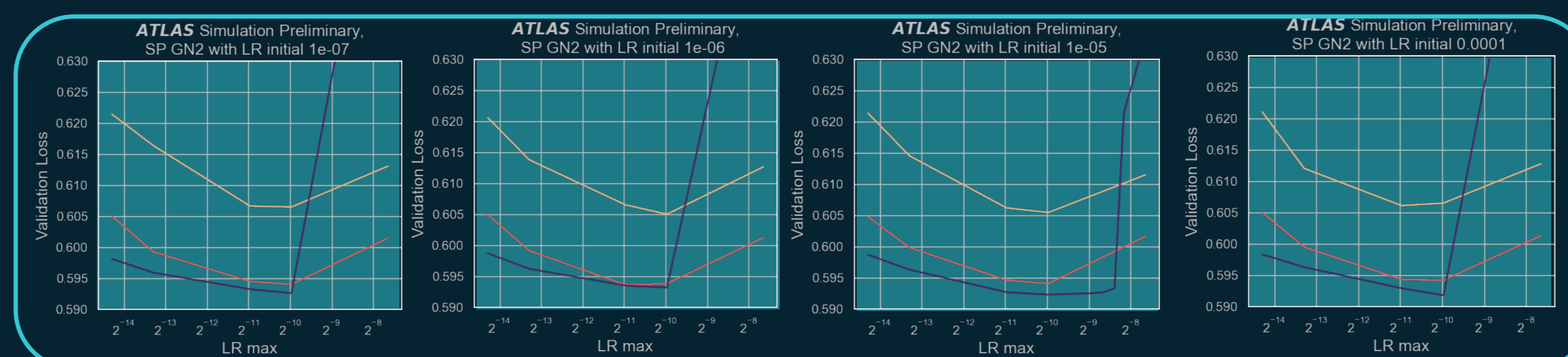
Yang et al; 2022

1. Similar performance hierarchy between small and large models in μP
2. Hyperparameter optimisation on a smaller model
3. Zero-shot transfer: best low-complexity = best high-complexity model
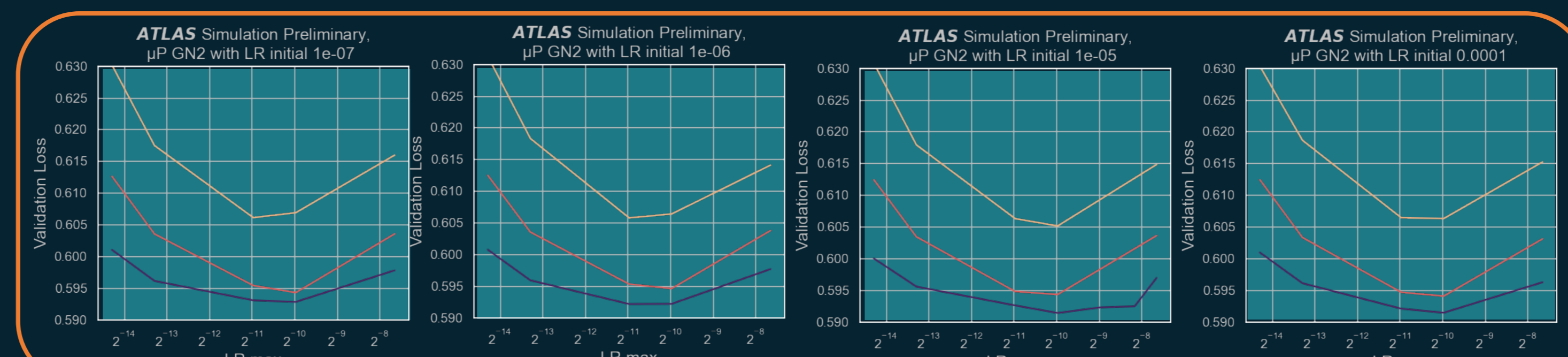4. Wider is always better

**Embedding Width**
64
128
256

### Real Deployment: LR Scheduler Optimisation (initialisation & maximal values)

**SP**

ATLAS Simulation Preliminary, SP GN2 with LR initial 1e-07 | SP GN2 with LR initial 1e-06 | SP GN2 with LR initial 1e-05 | SP GN2 with LR initial 0.0001

**μP**

ATLAS Simulation Preliminary, μP GN2 with LR initial 1e-07 | μP GN2 with LR initial 1e-06 | μP GN2 with LR initial 1e-05 | μP GN2 with LR initial 0.0001

## VERIFICATION — A look at pre-activation weights

**ATLAS** Simulation Preliminary, SP GN2-like, lr=0.01, nseeds=5 — **SP**
Initialisation | 1st training step | 2nd training step
$L_1(m) = \sum |w_i^{(m)}|$

**ATLAS** Simulation Preliminary, μP GN2-like, lr=0.01, nseeds=5 — **μP**
Initialisation | 1st training step | 2nd training step

**Stable for μP, blows up for SP!**

## ADVANTAGES

### Kubeflow
1. Multi-platform & flexible
2. Hardware agnostic MLOps & optimised resource usage
3. Advanced AutoML algorithms for improved HPO
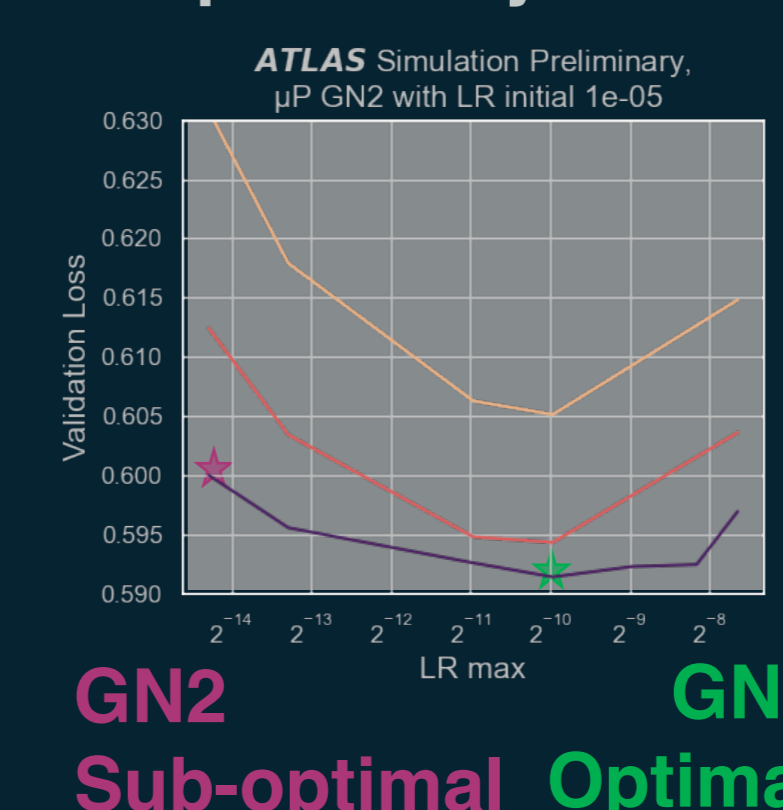4. Powerful visualisation
5. CERN-wide access

**C O M B I N E**

### Maximal Update Parametrization (μP)
1. Wider is always better → simple neural architecture search
2. μTransfer HP zero-shot small → large models
3. Reduced computing requirement per HP set test
   ➤ Width **256** (2.30 M params): 2 GPUs → 39 min / epoch
   ➤ Width **64** (0.29 M params): 1 GPU → 20 min / epoch
4. Better coverage of the HP search space
   ➤ With **μP**, 4 small-model tests ≈ 1 full-model test

## HPO MATTERS

**Significant performance dependency on HP**

**ATLAS** Simulation Preliminary, μP GN2 with LR initial 1e-05

**GN2 Sub-optimal** | **GN2 Optimal**

$$D_b = \frac{p_b}{f_c p_c + (1-f_c) p_{light}}$$

**ATLAS** Simulation Internal
$\sqrt{s}$ = 13 TeV, ttbar sample, $f_c$ = 0.1
GN2 sub-optimal | GN2 optimal
Light-jets | c-jets
(Background rejection vs b-jet efficiency; Light-jet ratio; c-jet ratio)