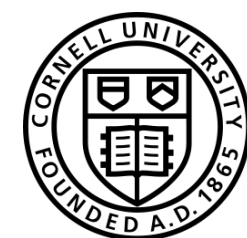# Improving Tracking Algorithms with ML:
# A case for Line Segment Tracking at the HL-LHC

## On behalf of the CMS Collaboration

October 12th, 2023

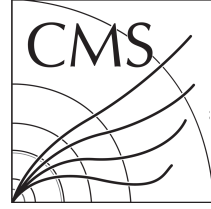# Challenge: HL-LHC Tracking
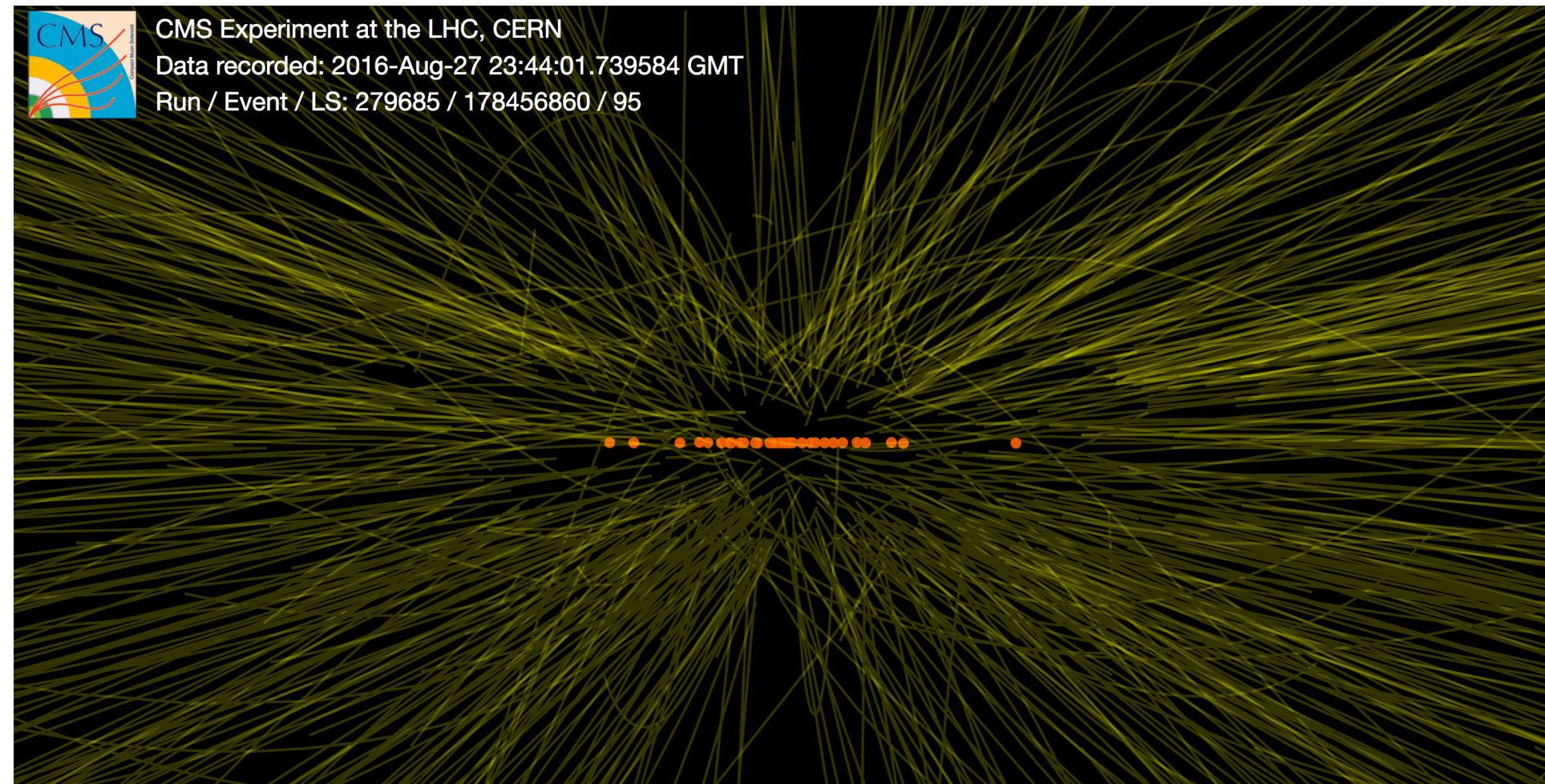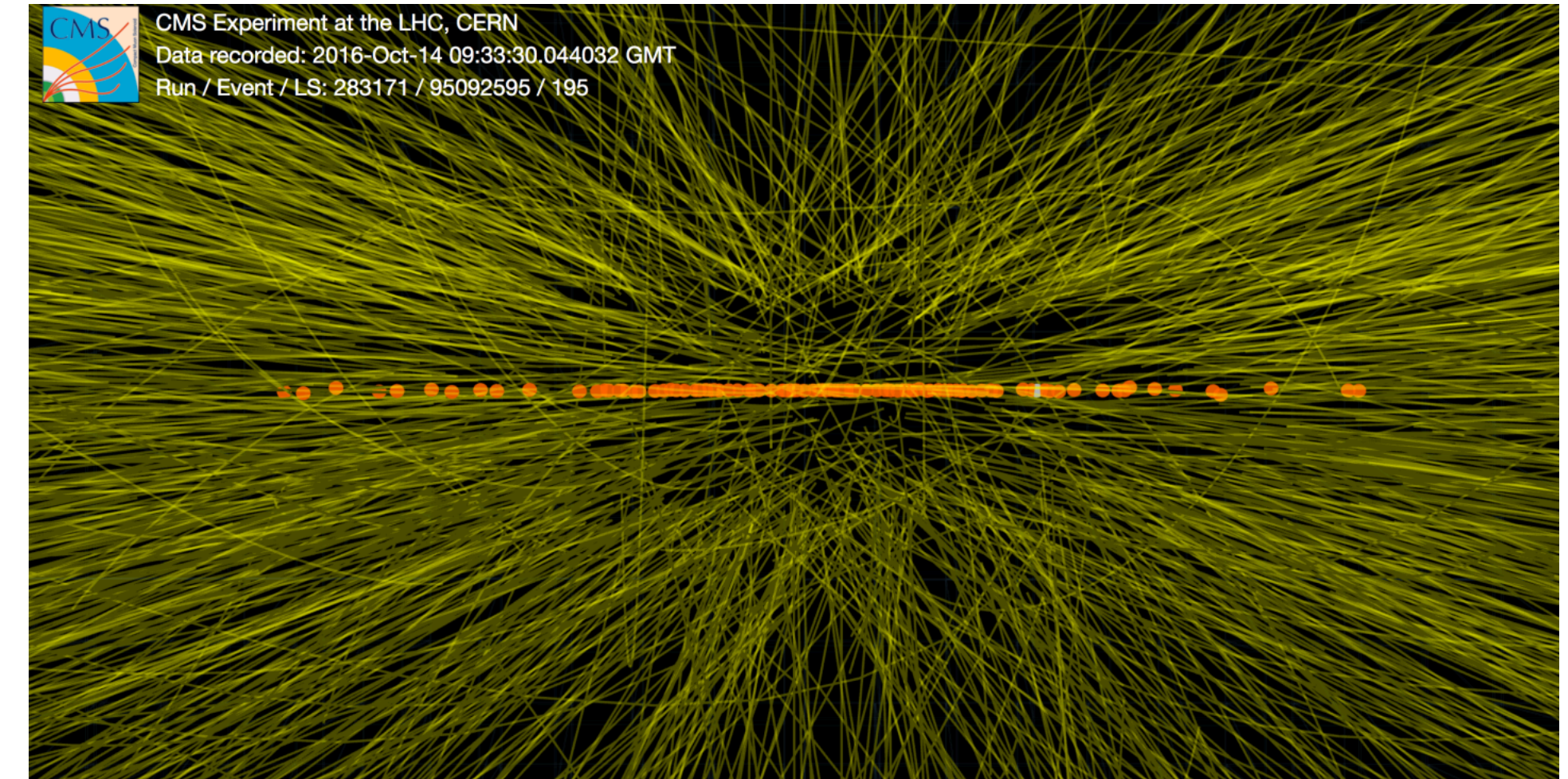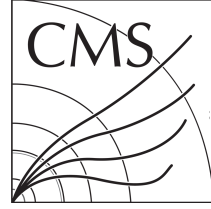## "High Luminosity" LHC (HL-LHC) planned for 2030s



Nominal Run 2 event (PU 30)

HL-like event (PU 130)

O(10x) concurrent collisions ("pile up") = O(10x) tracks
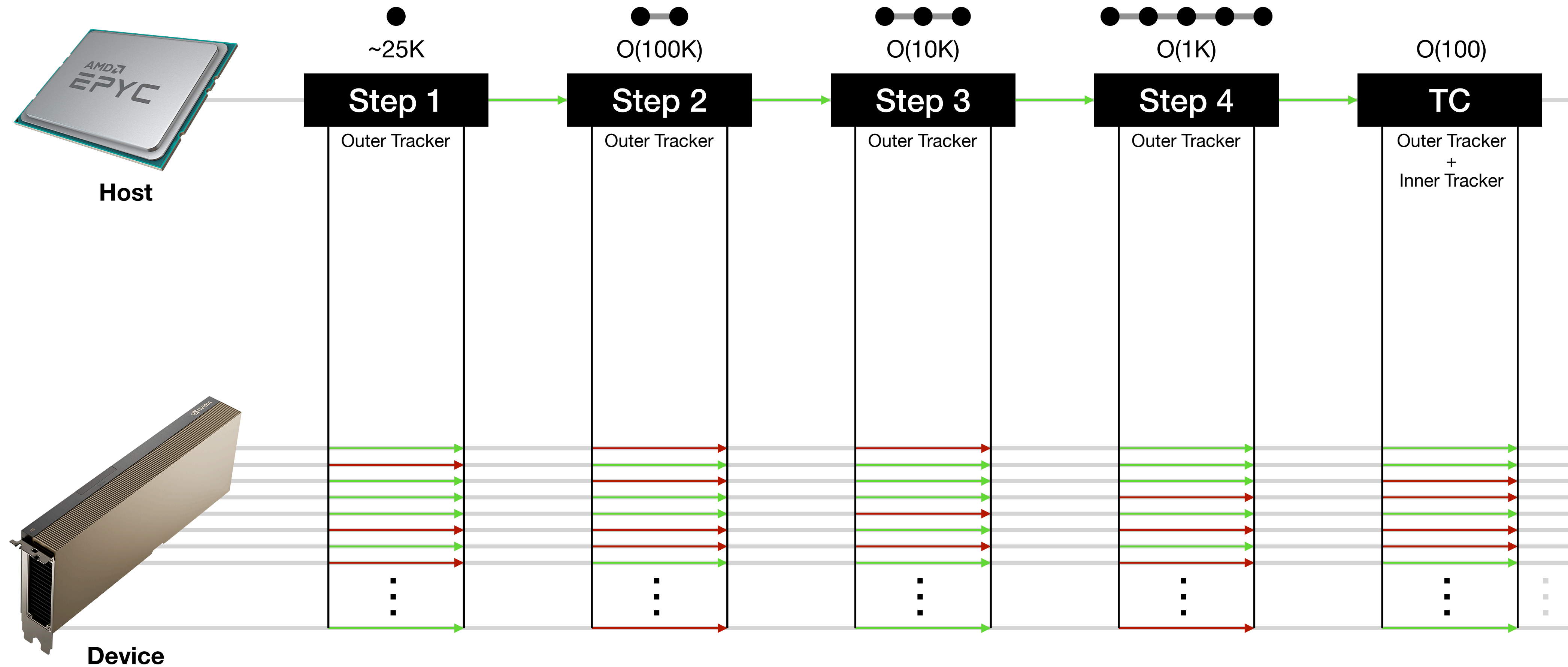⇒ **need a fast tracking algorithm** to keep up

**Current tracking algorithm is <span style="color:red">inherently sequential</span> ⇒ poor scaling**

We propose **LST: a <span style="color:red">highly parallelizable</span> tracking algorithm**
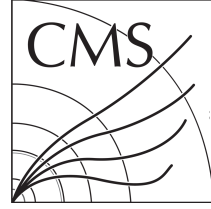
UC San Diego

# Solution: Line Segment Tracking (LST)

## At each step, one thread per object: deciding **keep** or **discard**



Each step is designed such that objects can be assessed independently
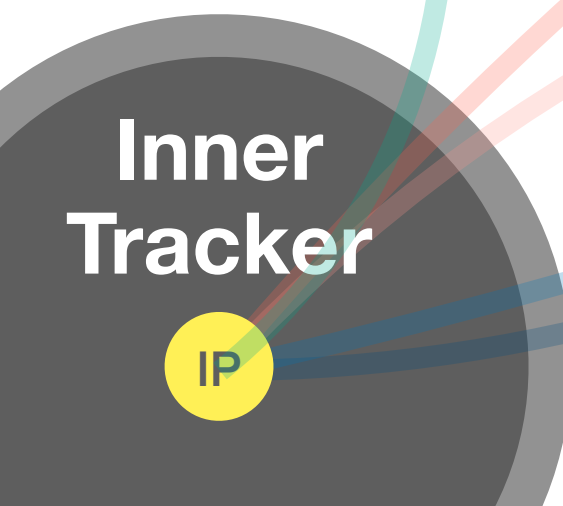⇒ **massively parallelizable!**

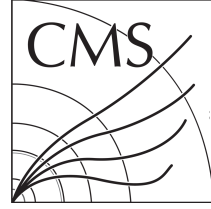# Solution: Line Segment Tracking (LST)

We will provide a basic description
of each step of the LST algorithm here

We show the steps here only to **introduce
the LST terminology/context**

**More info on LST can be found here:**
CTD 2022, CHEP 2023, CMS DP Note

**Inner
Tracker**
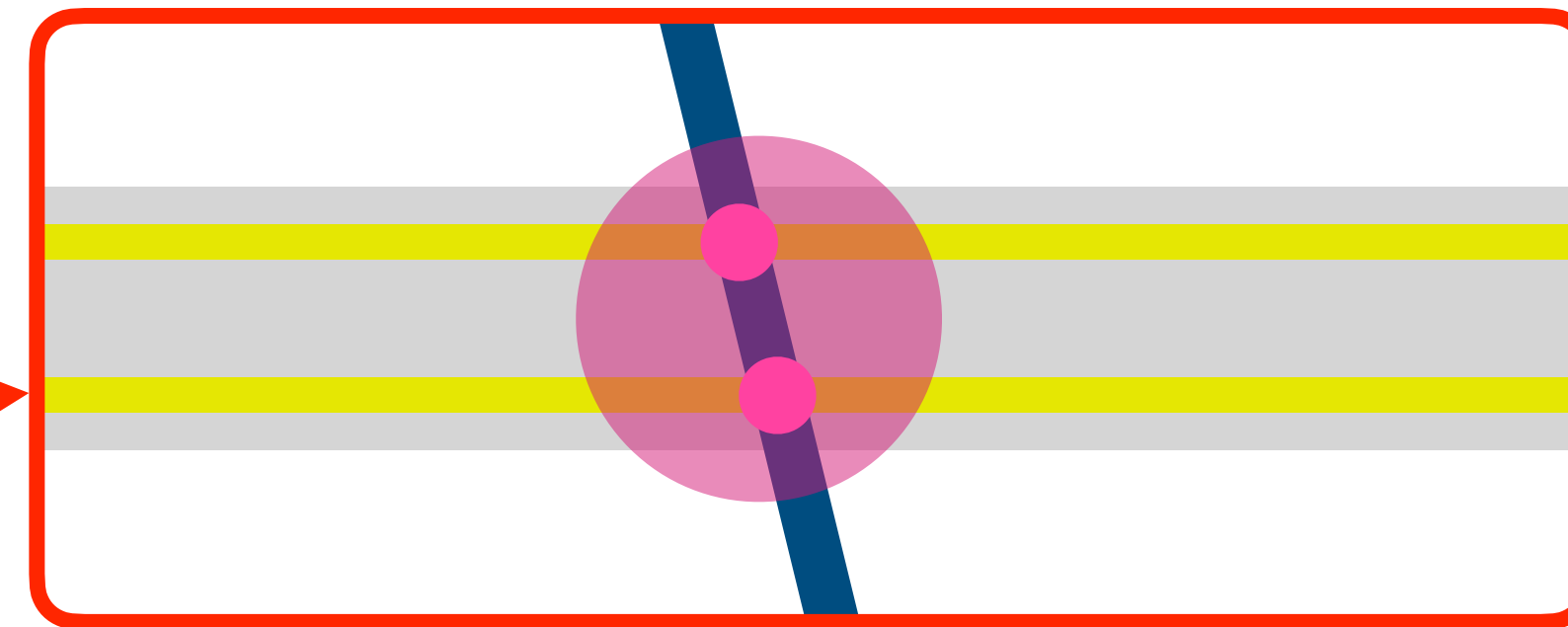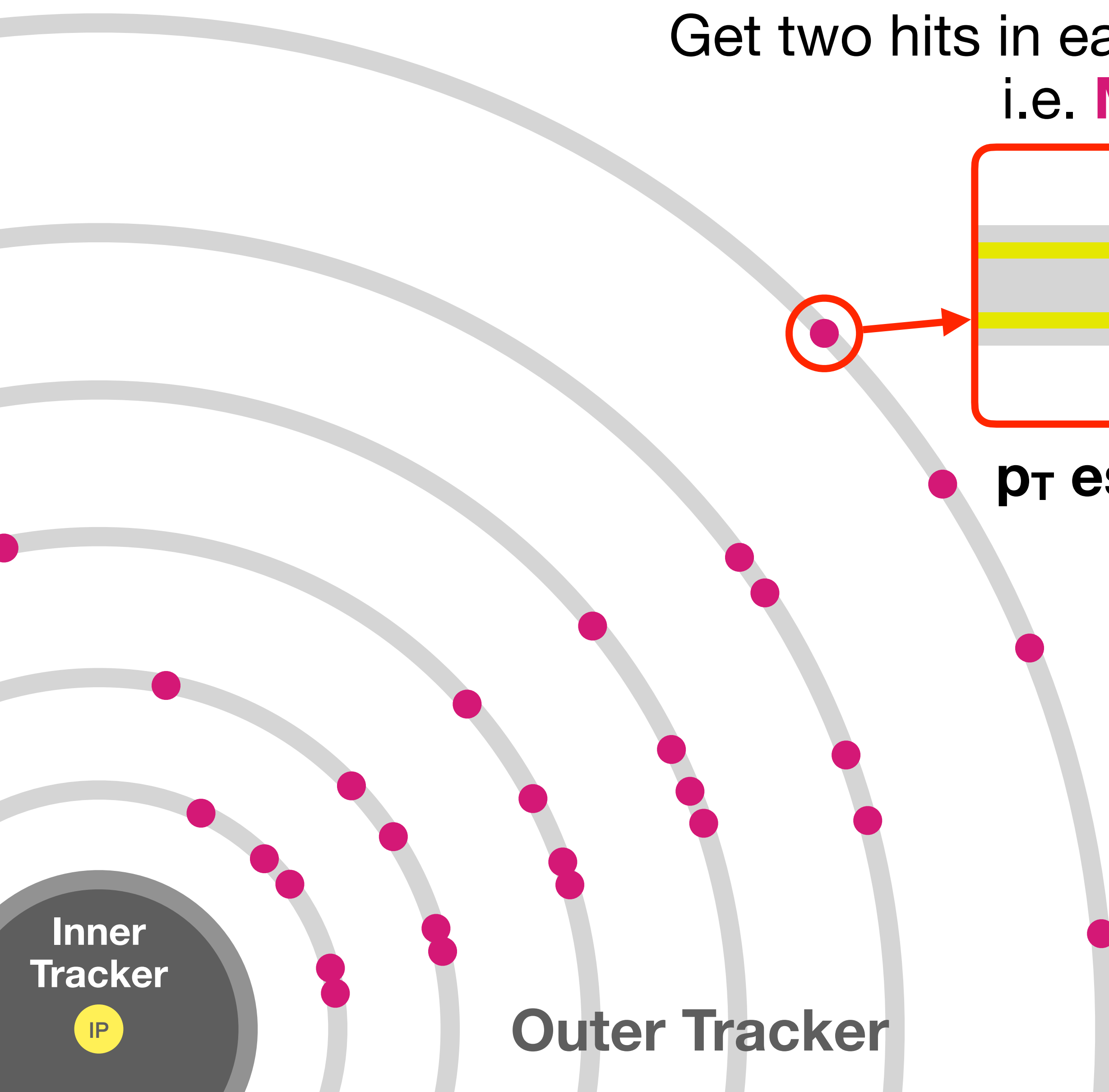
IP

**Outer Tracker**

UC San Diego

# LST in a Nutshell: Mini-Doublets

Get two hits in each layer in Phase 2 Outer Tracker:
i.e. **Mini-Doublets (MDs)**

"$p_T$ module"

$p_T$ **estimate for each MD**

**1:** 2S silicon sensor
**2:** Al-CF spacer
**3:** Front-end hybrid
**4:** Service hybrid
**5:** CFRP support
**6:** High voltage tab
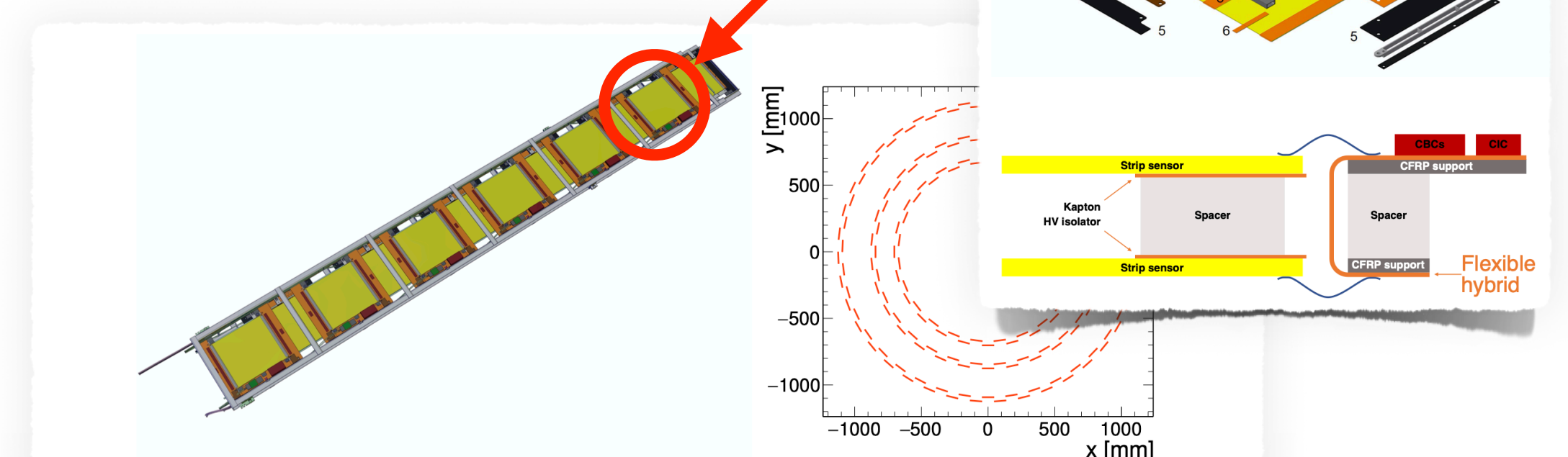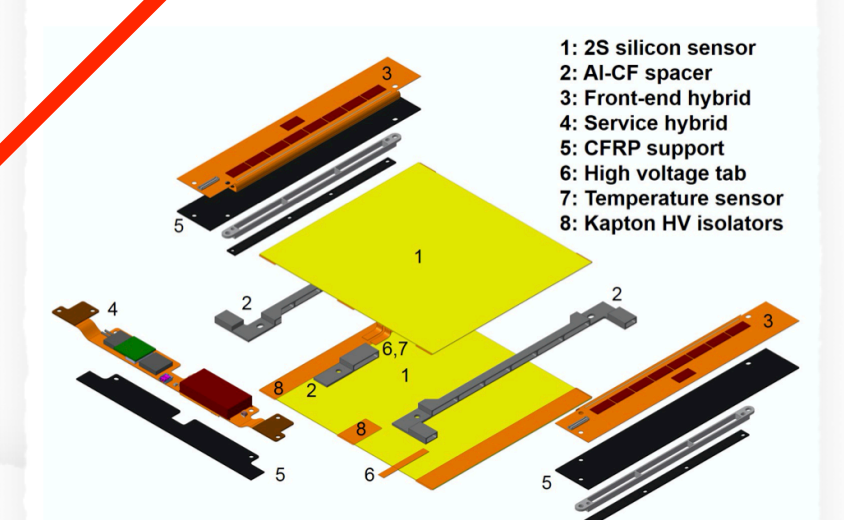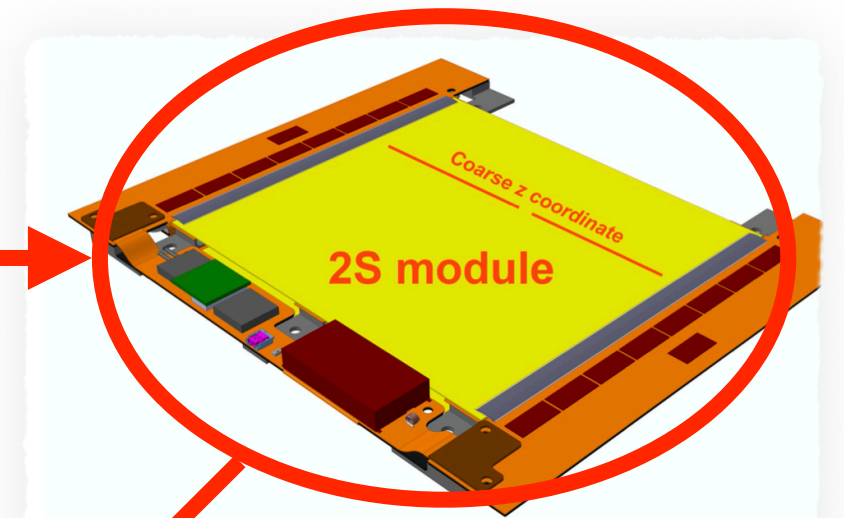**7:** Temperature sensor
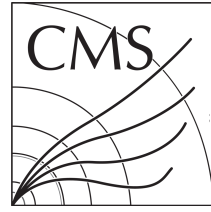**8:** Kapton HV isolators

Figure 3.3: Left: model of a TB2S ladder, housing twelve 2S modules. Right: x-y view of the TB2S, showing the staggering of neighbouring ladders.
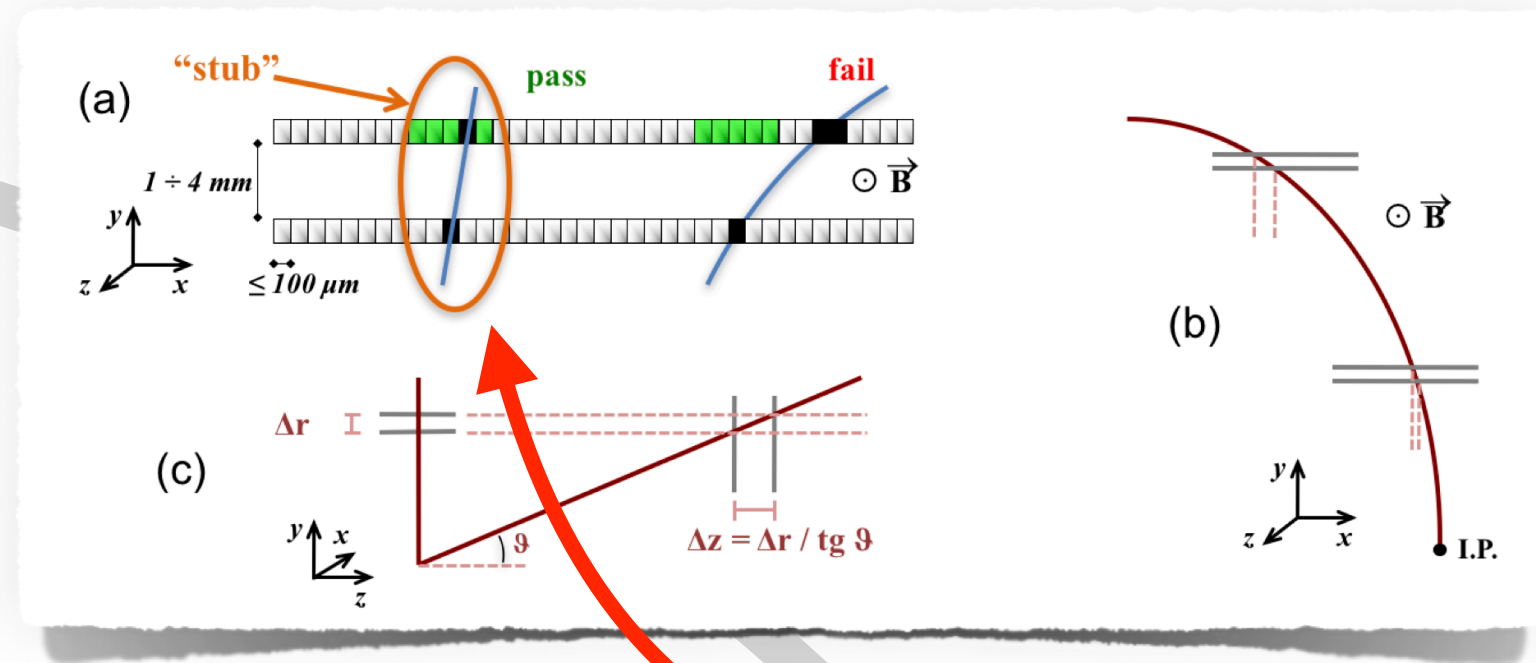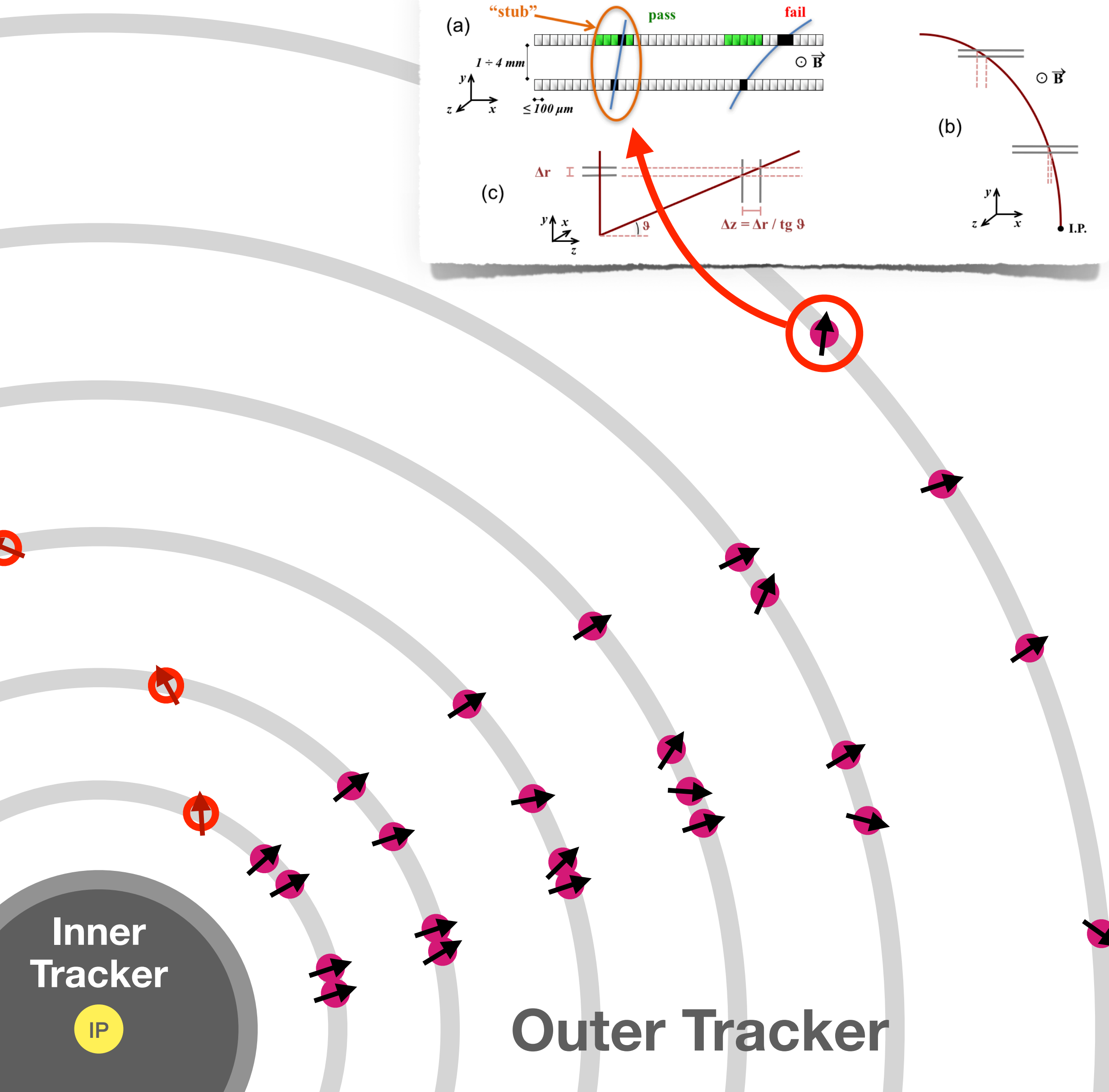
**Inner Tracker**

IP

**Outer Tracker**

5

# LST in a Nutshell: Mini-Doublets

Build all **good MDs**

good = $p_T > 0.8$

**One thread per MD**

Host

Device

MD

Inner Tracker

IP

Outer Tracker

UC San Diego

# LST in a Nutshell: Line Segments

Build all **valid** connections of two MDs:
i.e. **Line Segments (LSs)**

Derived a "**module map**" that
pre-determines valid LSs

**Not allowed to skip layers**

**Unnatural LS also not allowed**

**One thread per LS**

MD → LS

Host

Device

**Inner Tracker**

IP

**Outer Tracker**

UC San Diego

# LST in a Nutshell: Line Segments

Keep **good LSs**

good = **consistency between MD $p_T$**

**One thread per LS**

MD → LS

Host

Device

**Inner Tracker**

IP

**Outer Tracker**

UC San Diego

# LST in a Nutshell: Triplets

Keep **good** pairs of LSs that share a MD:
i.e. **Triplets (T3s)**

good = **p$_T$ consistency + other constraints**

**One thread per T3**

**Inner Tracker**

IP

**Outer Tracker**

Host

Device

| MD | LS | T3 |
|----|----|----|

UC San Diego

# LST in a Nutshell: Quintuplets

Keep **good** pairs of T3s that share a MD:
i.e. **Quintuplets (T5s)**

good = **p_T consistency + circle fit quality**

**One thread per T5**



**Inner Tracker**

IP

**Outer Tracker**

# LST in a Nutshell: Track Candidates



Take all **T5s** and match to **pixel seeds (pLS)**:
i.e. **pT5s**

Keep **good pT5s** as **Track Candidates (TCs)**

Take unmatched **T5s**, **good pT3s** (pLS + T3),
and unmatched **pLS** also as **TCs**

good = **p$_T$ consistency + circle fit quality**

pT5

pT3

**Inner Tracker**

IP

**Outer Tracker**

UC San Diego

# Improving LST with ML

- Line Segment Tracking (LST) is already highly performant and parallelizable

- Central question: **where can Machine Learning (ML) *realistically* be used to improve LST?**
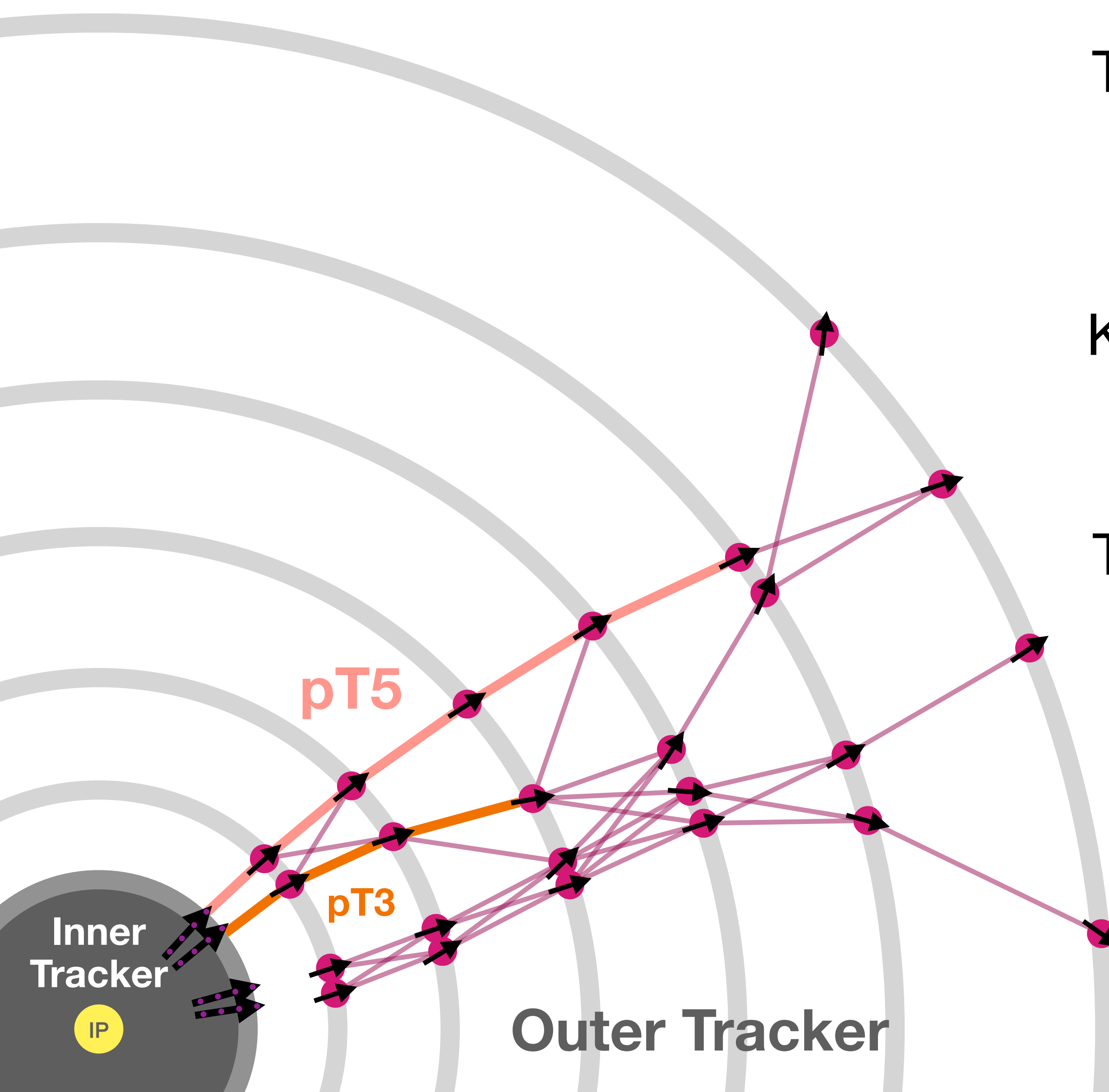
- In this talk we will…

  - Outline a suitable step in LST to try a **simple ML solution**

  - Show significant **improvements to LST!**

  - Present a prospectus for **more ambitious ML solutions/algorithms**

# ML Opportunity: Quintuplets



**pT5s (pLS + T5) + T5s give most of the TC efficiency
T5s have a high fake rate**

# ML Opportunity: Quintuplets



pT5s + T5s give most of the TC eff.
T5s have a high fake rate
⇒ **there is room for improvement**

**Can ML do better without heavily
impacting the total LST runtime?**

**Next:** we train/deploy a small neural
network for classifying real vs. fake **T5s**

UC San Diego

# T5 DNN Training Data

**Objective**
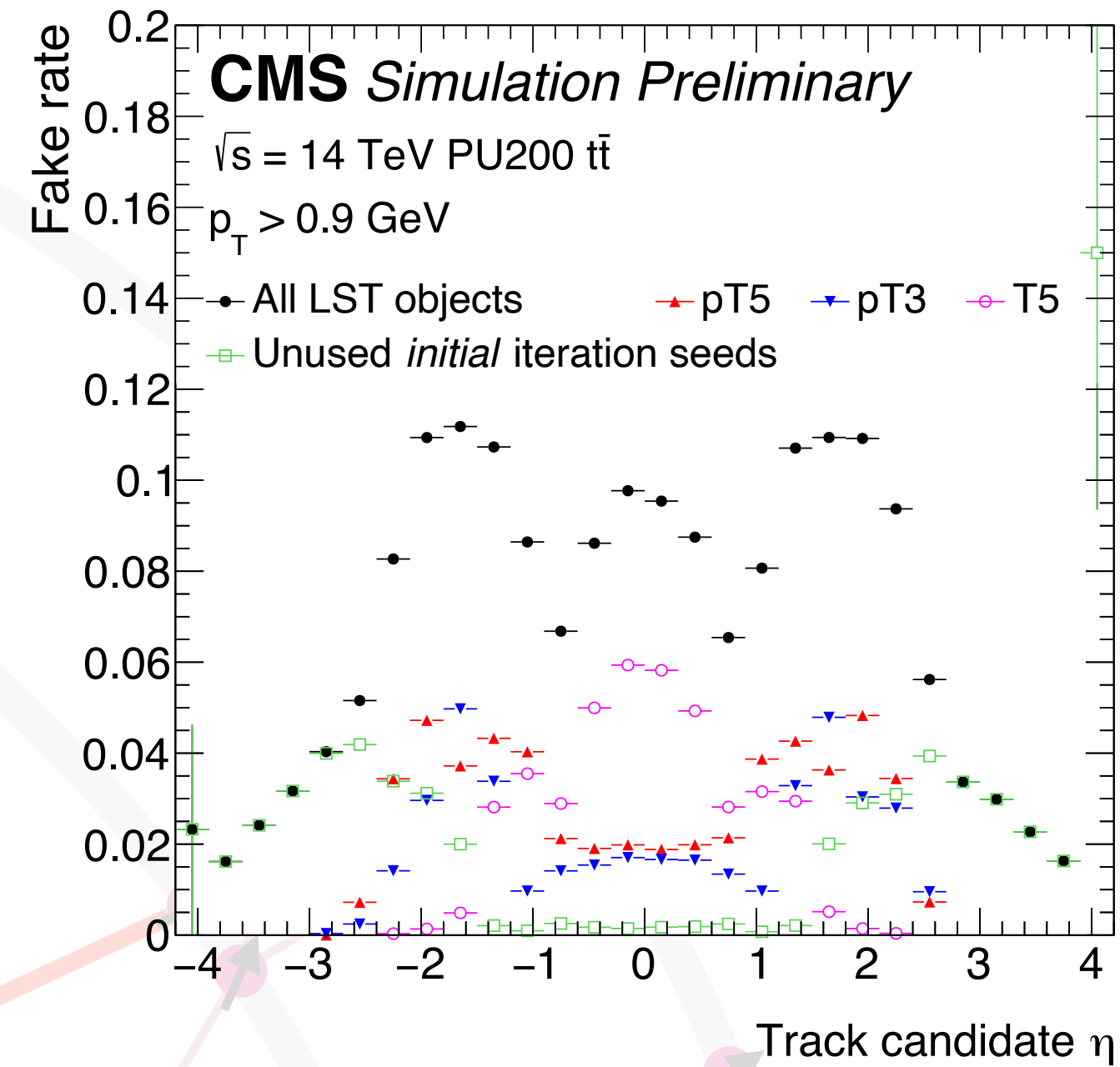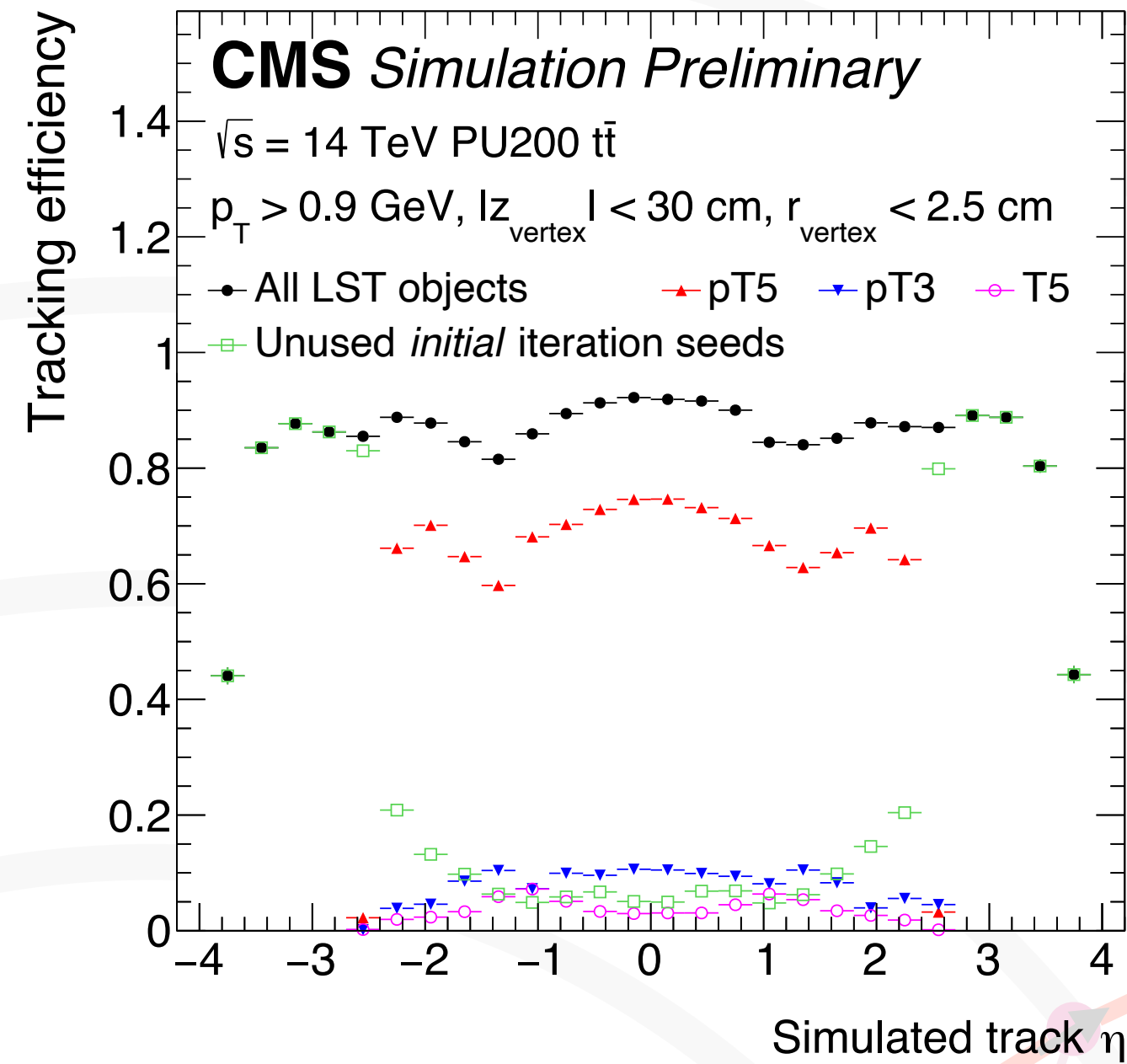Train DNN to classify
"real" vs. "fake" **T5s**

**Real:** > 75% of hits are from the same sim track

**Fake:** not "real"

**Baseline LST T5s**
Pass basic quality cuts
Pass r-z $\chi^2$ cut $\Big]$
Pass r-φ $\chi^2$ cut $\Big]$ Quality of circle fit

Remove r-φ $\chi^2$ cut

T5 circle fit

Uses variables that the
DNN might use better

**DNN Training T5s**
Pass basic quality cuts
Pass r-z $\chi^2$ cut
Pass r-φ $\chi^2$ cut

**~2.1 million T5s**
**40% real**

Inner
Tracker

Outer Tracker

UC San Diego

**Objective**
Train DNN to classify
"real" vs. "fake" **T5s**

**Real:** > 75% of hits are from the same sim track
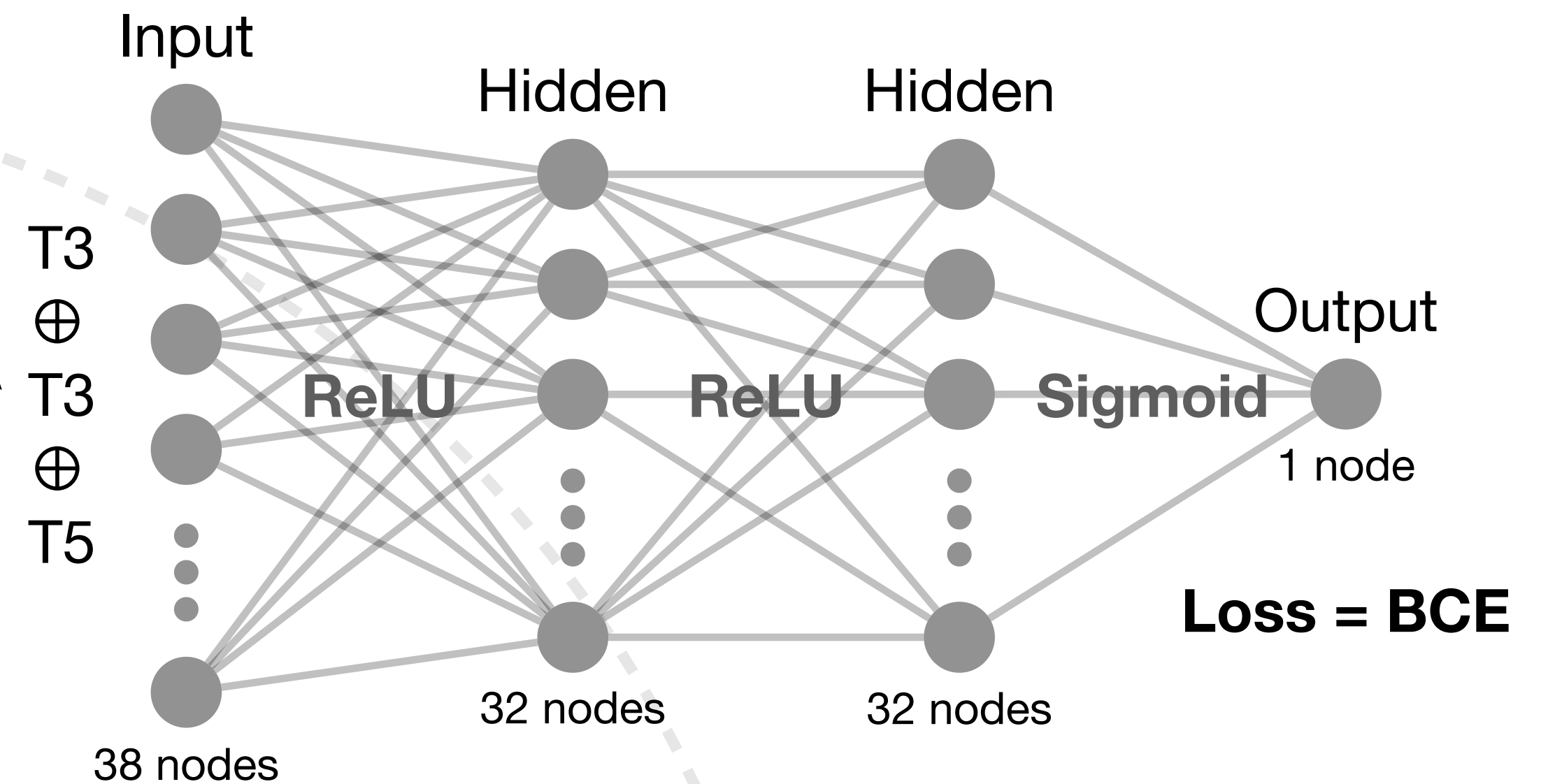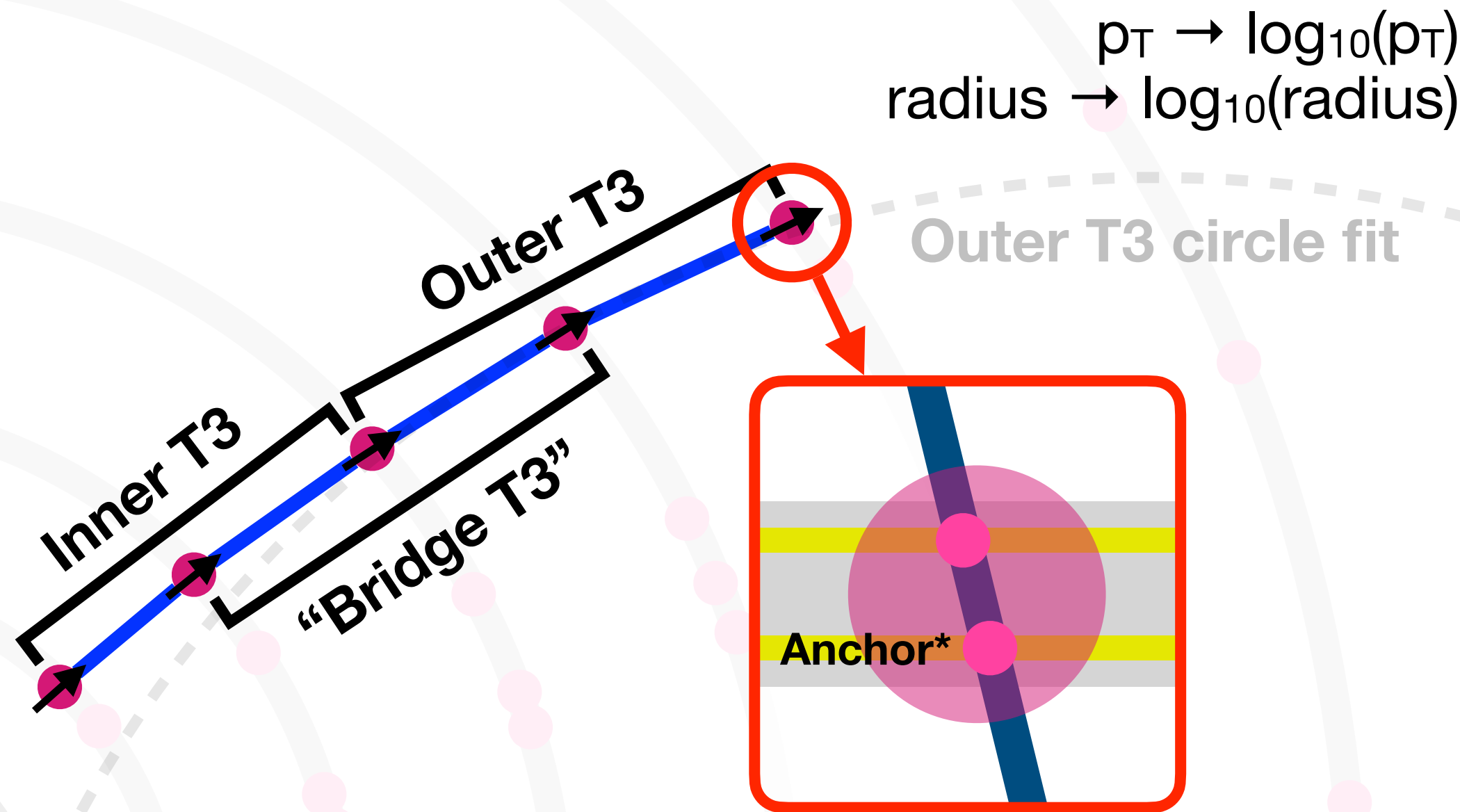
**Fake:** not "real"

| Object | Feature |
|---|---|
| T3 (x2) | $p_T$ |
| | Inner anchor hit r, z, φ, η, layer |
| | Middle anchor hit r, z, φ, η, layer |
| | Outer anchor hit r, z, φ, η, layer |
| | Radius of circle fit |
| **T5** candidate | $p_T$, η, φ |
| | Radius of circle fit for Bridge "T3" |

$p_T \rightarrow \log_{10}(p_T)$
$radius \rightarrow \log_{10}(radius)$

Outer T3
Inner T3
"Bridge T3"
Outer T3 circle fit
Anchor*

Input
Hidden
Hidden
T3 ⊕ T3 ⊕ T5
ReLU
ReLU
Sigmoid
Output
1 node
38 nodes
32 nodes
32 nodes
**Loss = BCE**

Inner Tracker

*Not always the inner hit: for PS modules, it is always the pixel (P) hit, which is not necessarily the hit that is closest to the beamline

Outer Tracker

16

UC San Diego

# T5 DNN Training and Selection



**Choice of epoch is arbitrary after 400 (we choose 500)**

UC San Diego

# Aside: Other T5 DNN Architectures
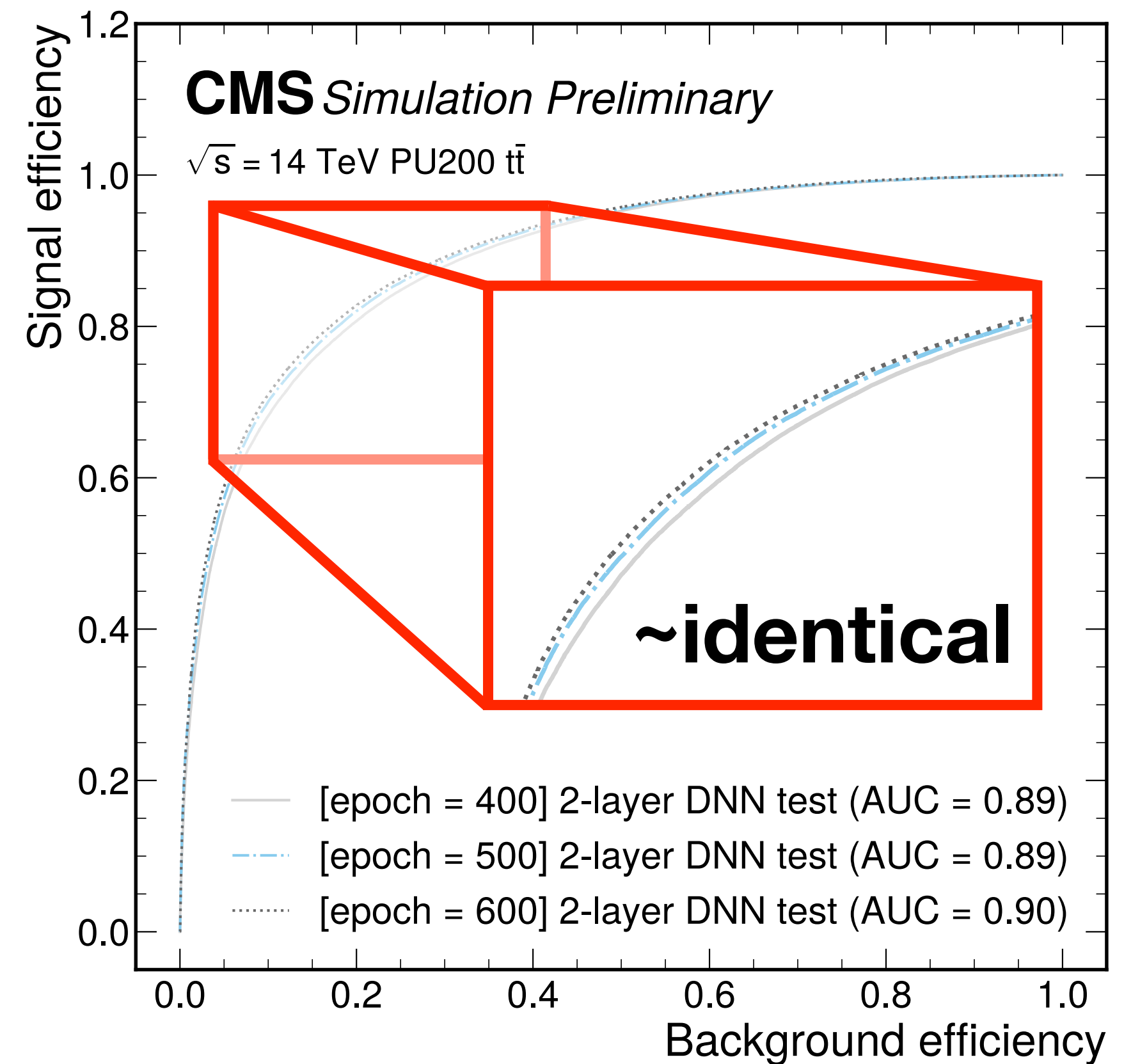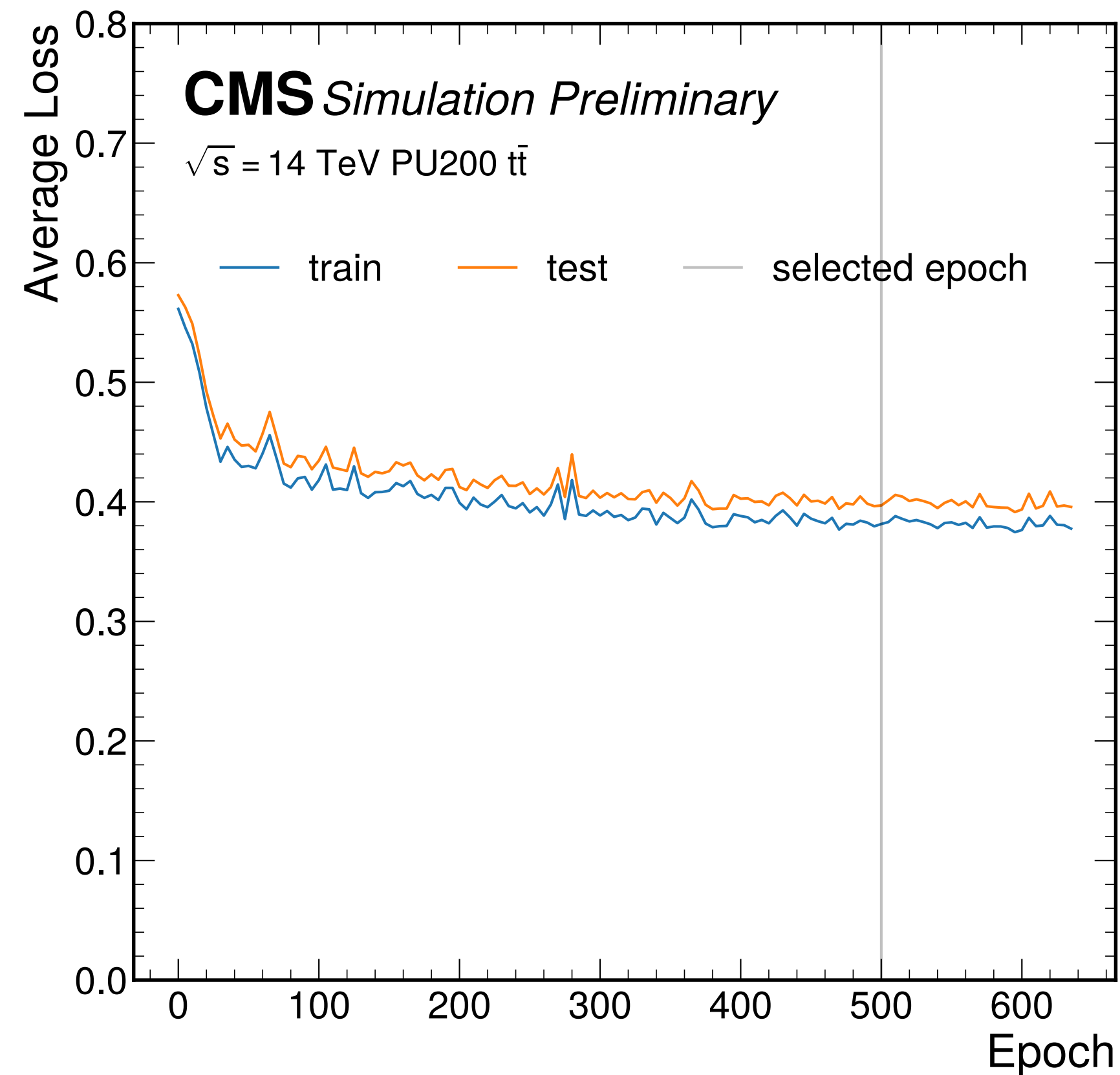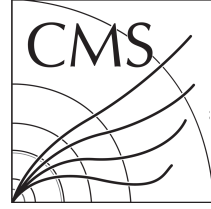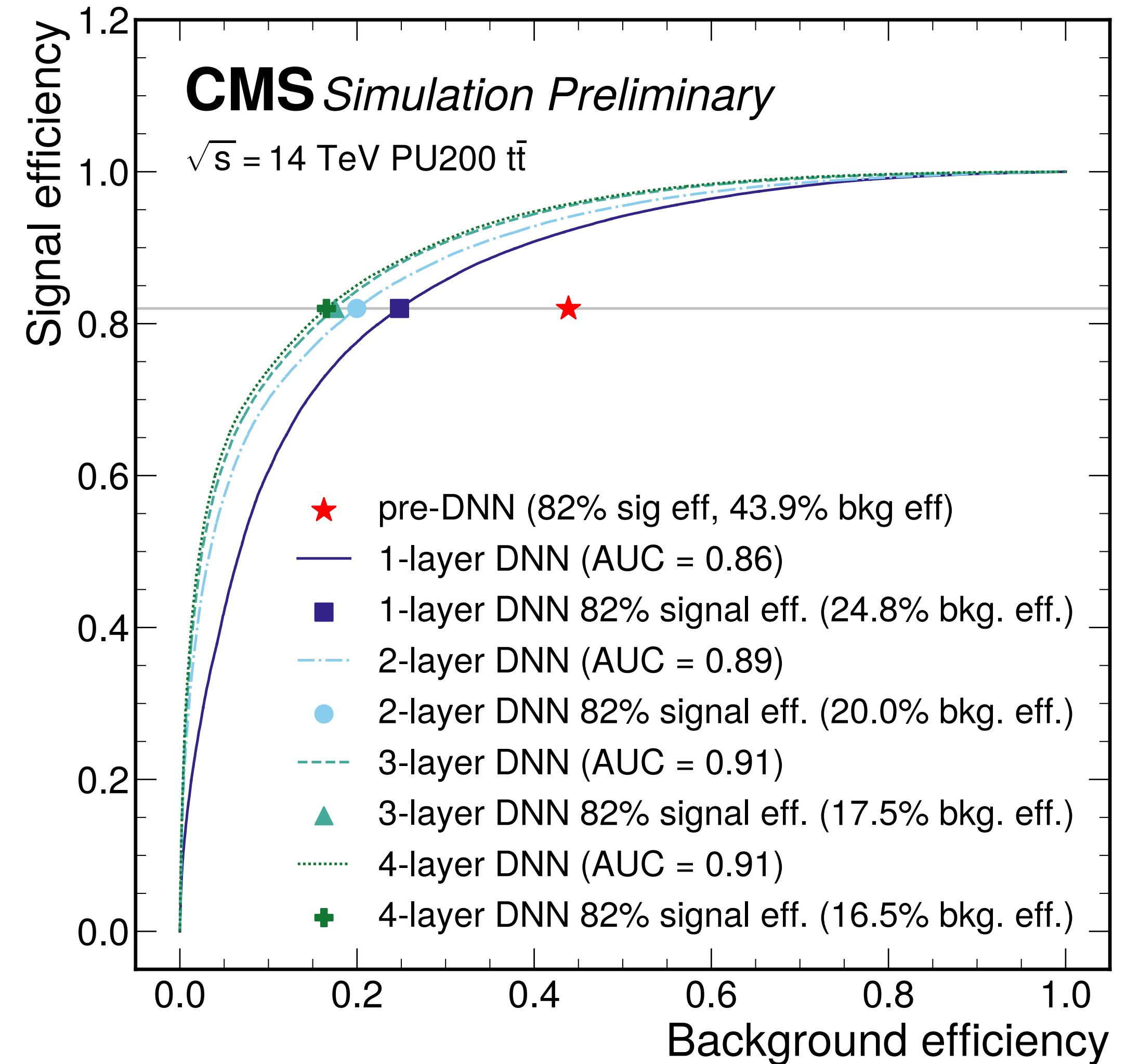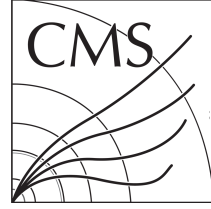
- Additional layers (up to 4) gives an extra ~3% decrease in bkg. efficiency (i.e. fake rate)

- Selected 2-layer DNN as a balance between performance and computational complexity

- A larger DNN could be used, but seems to be diminishing returns



CMS *Simulation Preliminary*

$\sqrt{s}$ = 14 TeV PU200 t$\bar{t}$

★ pre-DNN (82% sig eff, 43.9% bkg eff)

— 1-layer DNN (AUC = 0.86)

■ 1-layer DNN 82% signal eff. (24.8% bkg. eff.)

—·— 2-layer DNN (AUC = 0.89)

● 2-layer DNN 82% signal eff. (20.0% bkg. eff.)

— — 3-layer DNN (AUC = 0.91)

▲ 3-layer DNN 82% signal eff. (17.5% bkg. eff.)

····· 4-layer DNN (AUC = 0.91)

✚ 4-layer DNN 82% signal eff. (16.5% bkg. eff.)

UC San Diego

# T5 DNN Performance

- Performance gain over pre-DNN baseline (★)

  - Background efficiency = Fake rate

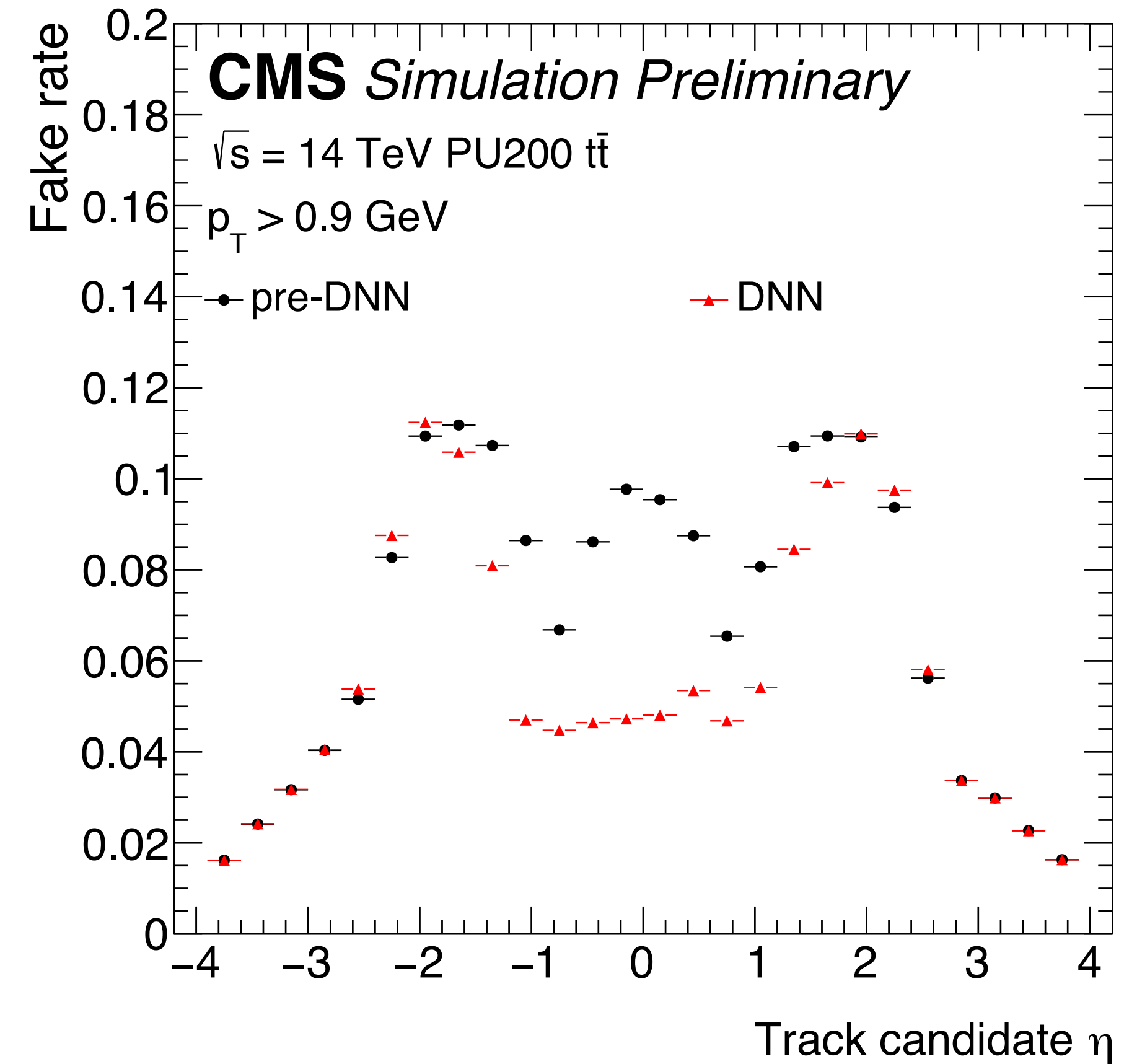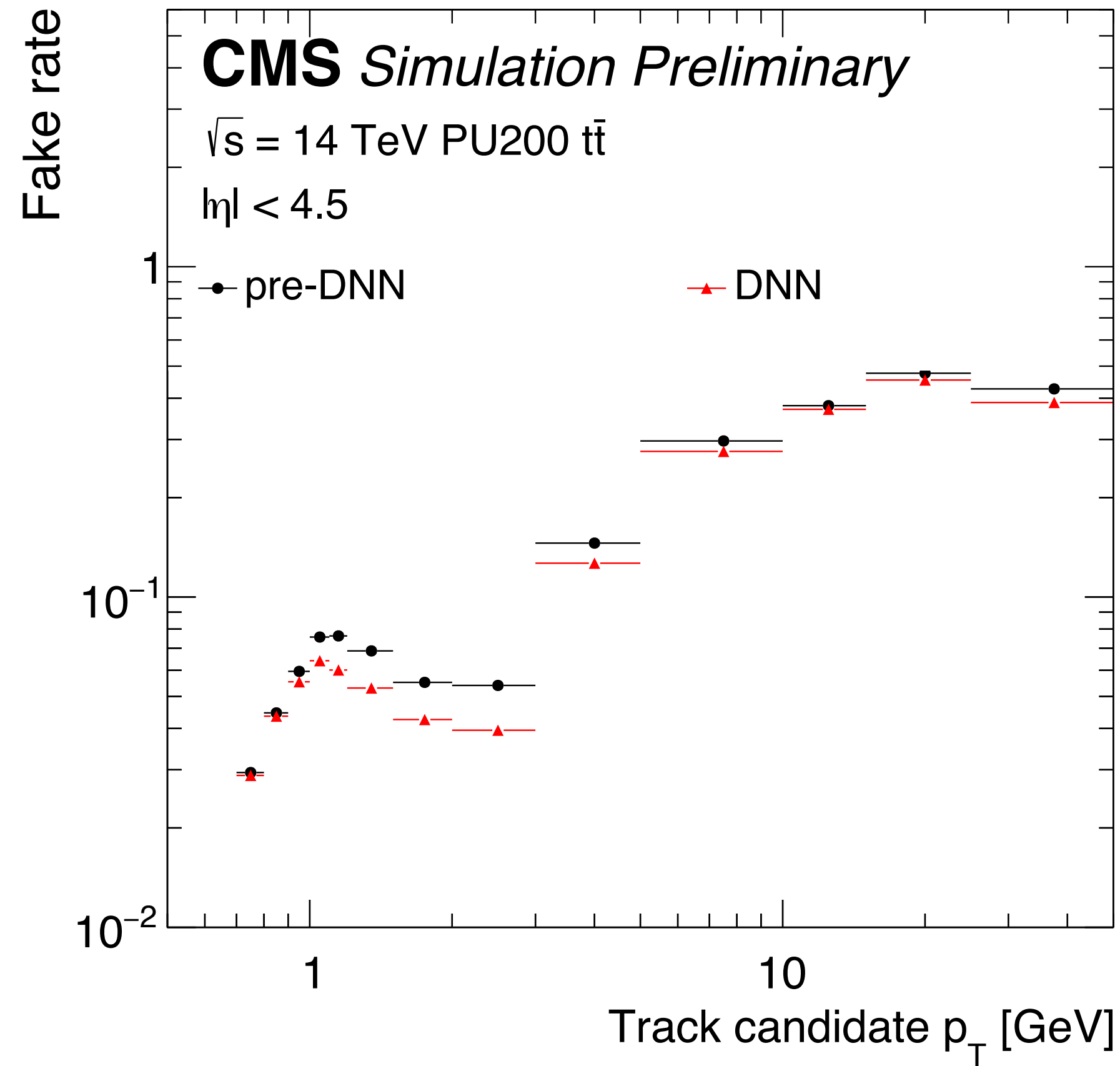- **Cannot get full picture of DNN in LST from this ROC curve alone:**

  - Only looking at T5s, but the TCs are more diverse

  - Effects of duplicate removal, TC selection, etc. are non-trivial

- **Next:** implement DNN inference in LST and compare to pre-DNN performance
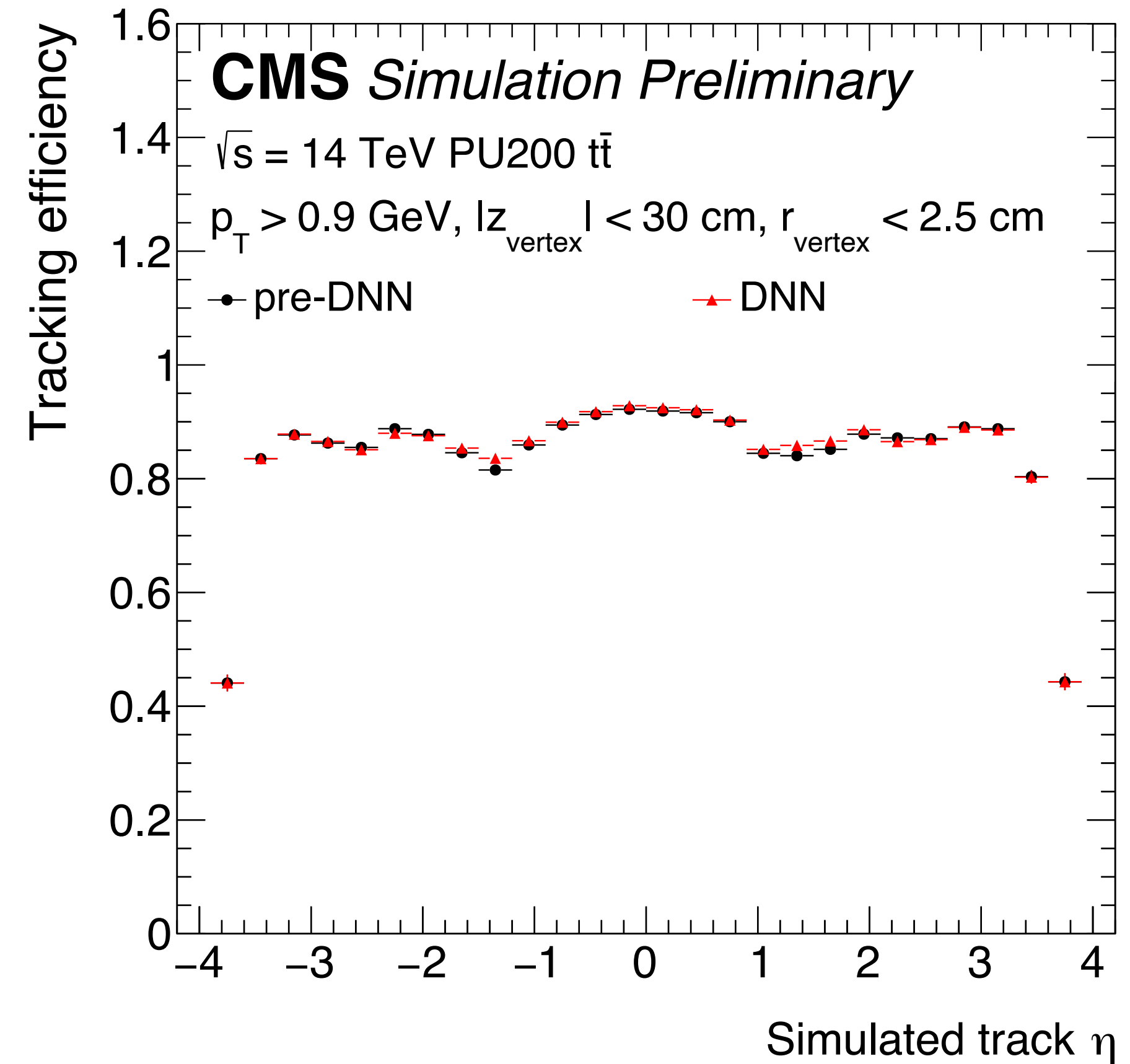
- Details of implementation are in the backup



**CMS** *Simulation Preliminary*

$\sqrt{s}$ = 14 TeV PU200 $t\bar{t}$

Significant performance gain!

★ pre-DNN (82% signal eff, 43.9% bkg. eff)

2-layer DNN train (AUC = 0.90)

2-layer DNN test (AUC = 0.89)

● 2-layer DNN 82% signal eff. (20.0% Bkg. eff.)

Signal efficiency
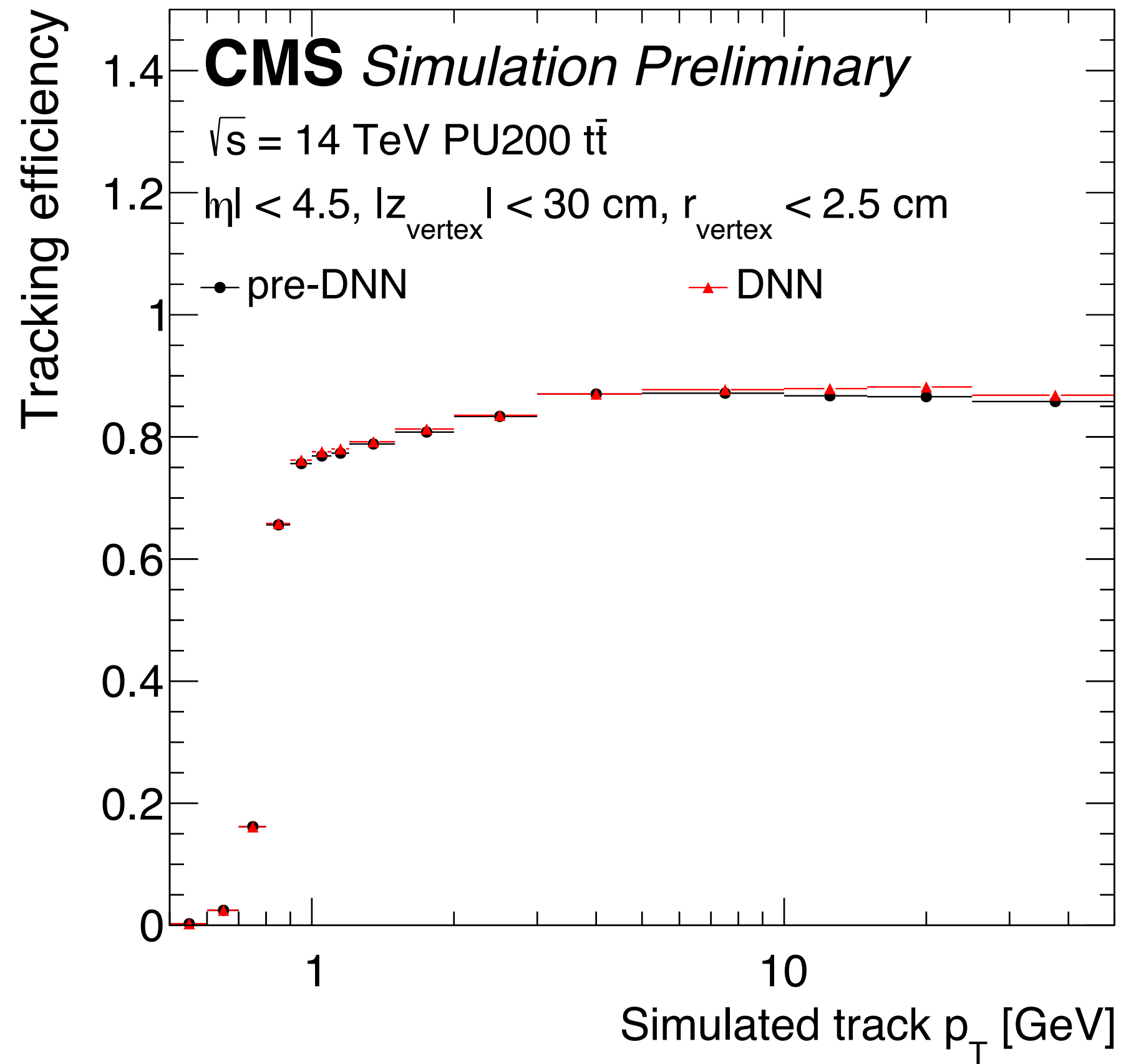
Background efficiency

UC San Diego
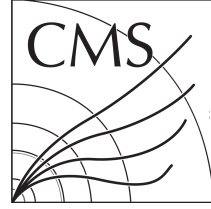
# Fake Rate Comparison



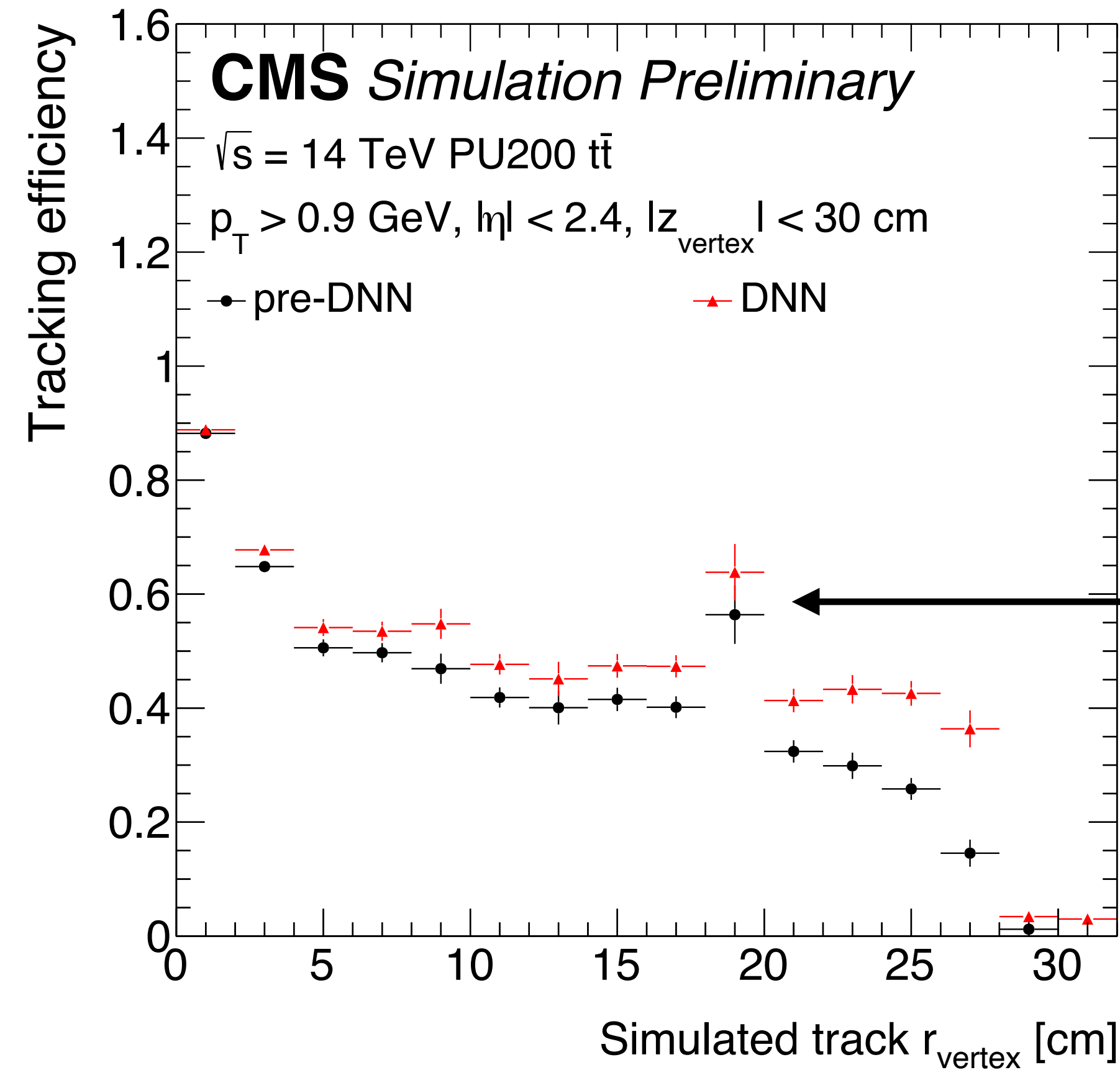**DNN gives ~40% reduction of fake rate in the barrel**

UC San Diego

# Efficiency Comparison



**DNN gives no loss in efficiency**
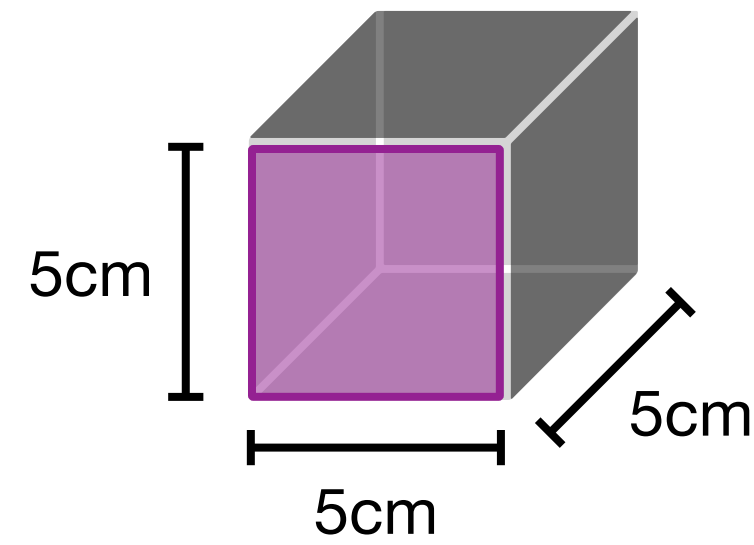**Using DNN WP that matches LST signal efficiency**

UC San Diego
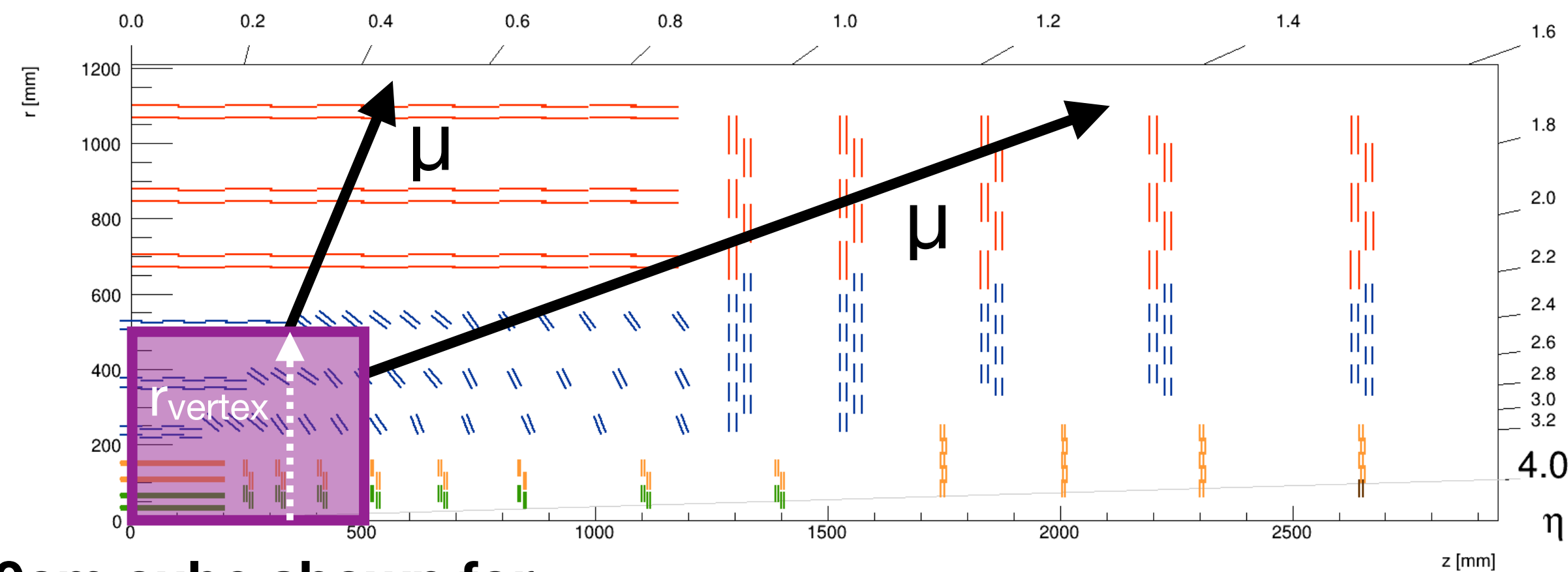
# Efficiency vs. r$_{vertex}$ Comparison



Contributions from material interactions are lower than neighboring bins due to the detector geometry

**Significant gain in efficiency for displaced tracks**
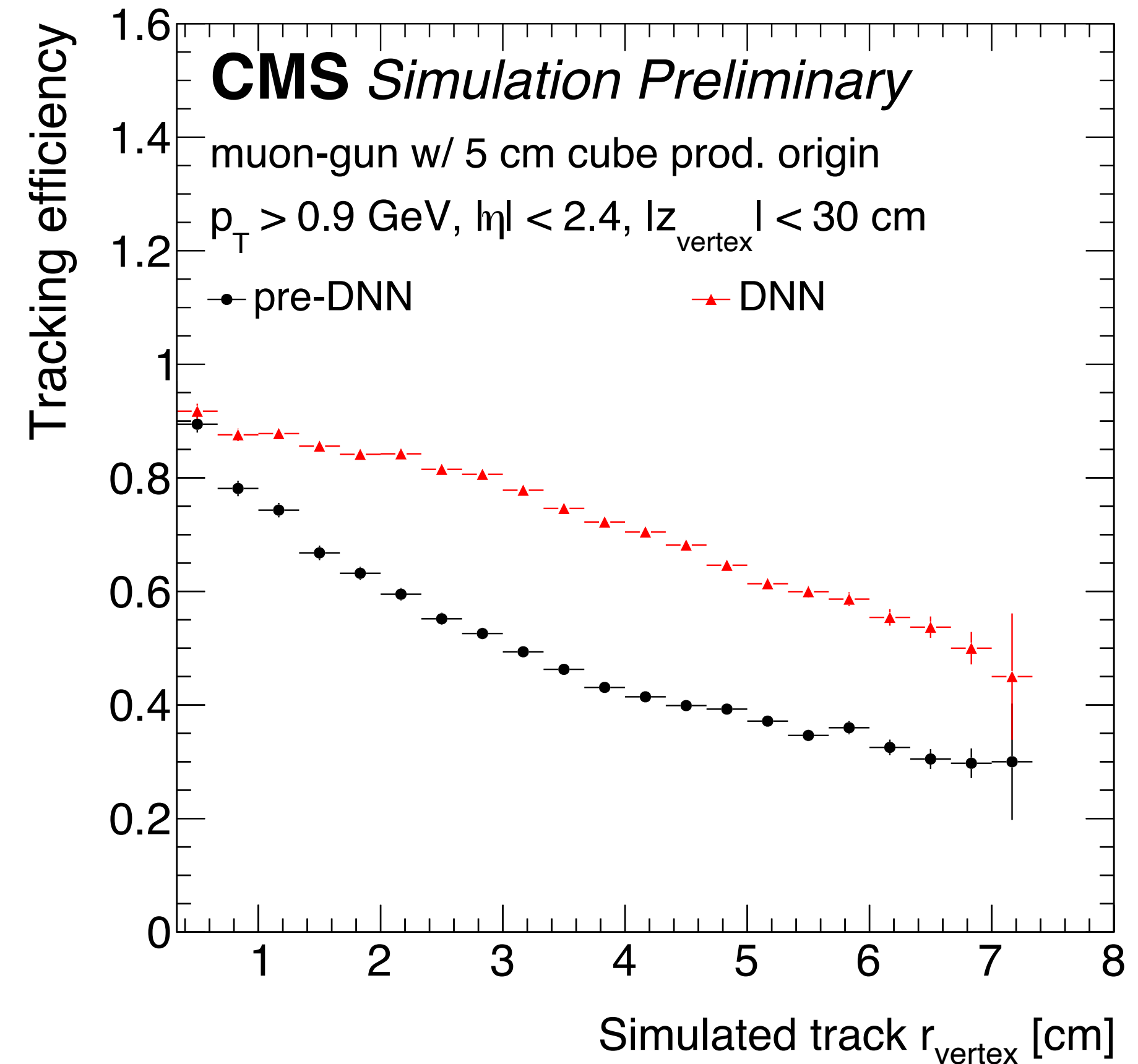
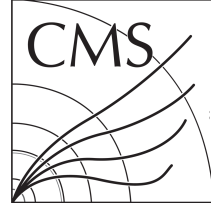UC San Diego

# Efficiency Comparison: Muon Cube



**Muon cube sample:**
Production vertex uniformly distributed across 5cm cube (no pileup), **i.e. displaced tracks**

**50cm cube shown for visualization purposes**

**CMS** *Simulation Preliminary*

muon-gun w/ 5 cm cube prod. origin

$p_T > 0.9$ GeV, $|\eta| < 2.4$, $|z_{vertex}| < 30$ cm

- pre-DNN
- DNN

Tracking efficiency

Simulated track $r_{vertex}$ [cm]

**Significant recovery of efficiency for displaced tracks**

UC San Diego

# T5 DNN Timing Impact

**Units of milliseconds (ms)**

|  | T5 | 1/Throughput | N streams |
|---|---|---|---|
| pre-DNN | **3.37 ± 0.13** | **28.4 ± 1.5** | 1 |
| DNN | **3.39 ± 0.07** | **28.7 ± 1.1** | 1 |

Error is 1 standard deviation for 10 trials

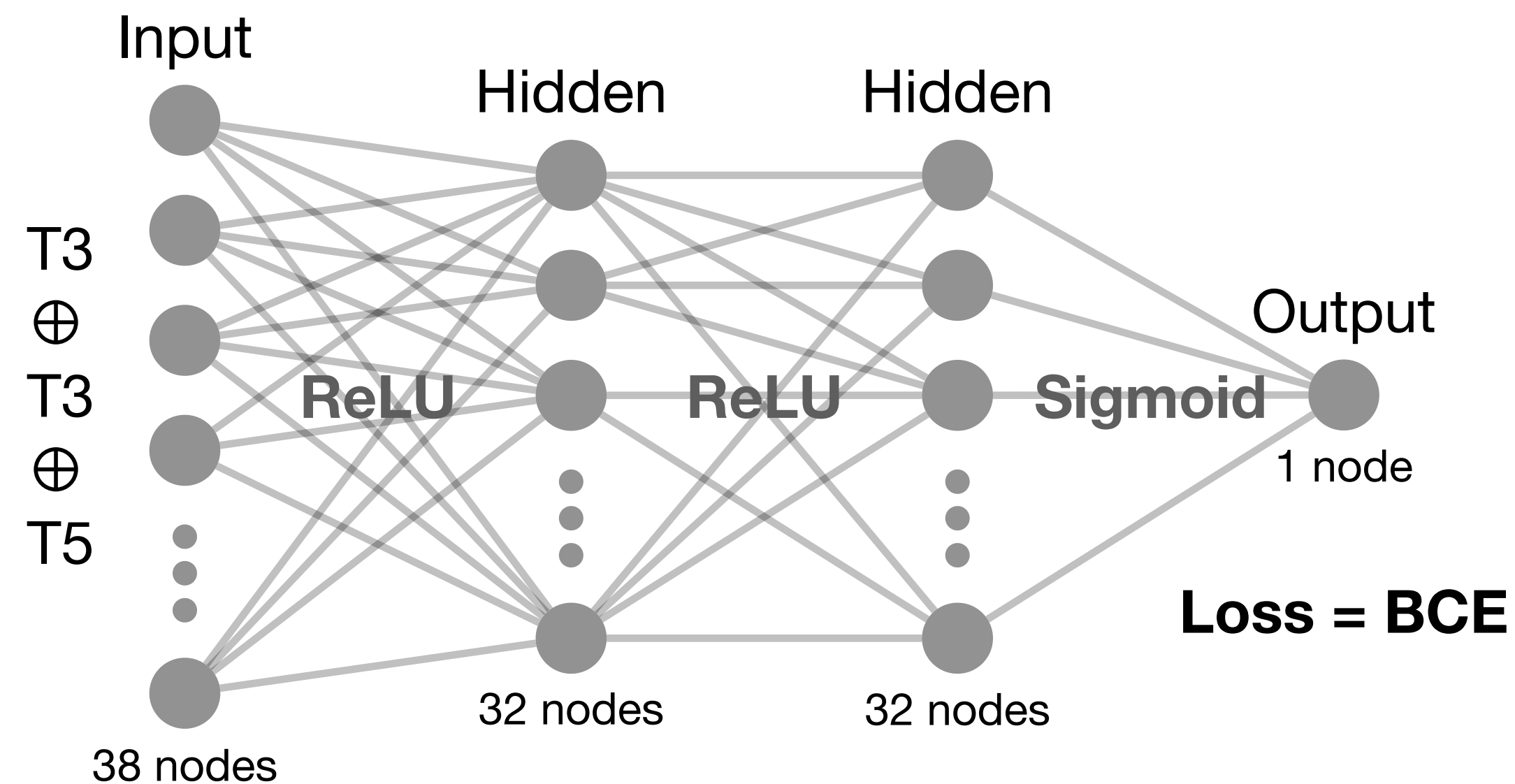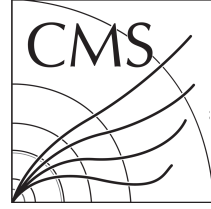**CMS** *Simulation Preliminary*



**DNN has no measurable impact on LST runtime**
(Measured on an NVidia A30 GPU)

UC San Diego

# T5 DNN Summary

- We have trained a lightweight DNN to classify real vs. fake LST T5s

- We have shown that **the T5 DNN improves LST w/ no impact on runtime**

- We have **established a pipeline for training ML algorithms on LST data**

  - Includes a full simulation of the CMS detector

# GNN Prospectus

- Interaction Networks (DeZoort, Thais, Duarte et al.)
  https://doi.org/10.1007/s41781-021-00073-z

- Exa.TrkX (Ju, Murnane, Calafiura et al.)
  https://doi.org/10.1140/epjc/s10052-021-09675-8

- GNNs (review by DeZoort, Battaglia, Biscarat et al.)
  https://doi.org/10.1038/s42254-023-00569-0

- And much more (see CTD 2023 agenda)!

**Promising results from TrackML GNNs**
No single best graph-building algorithm
No results using full CMS simulation

Use **LST to build input graph**,
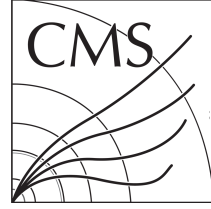most obvious: **MD nodes**, **LS edges**

**If GNN can build good TCs: LST becomes
a *very fast* graph building algorithm!**
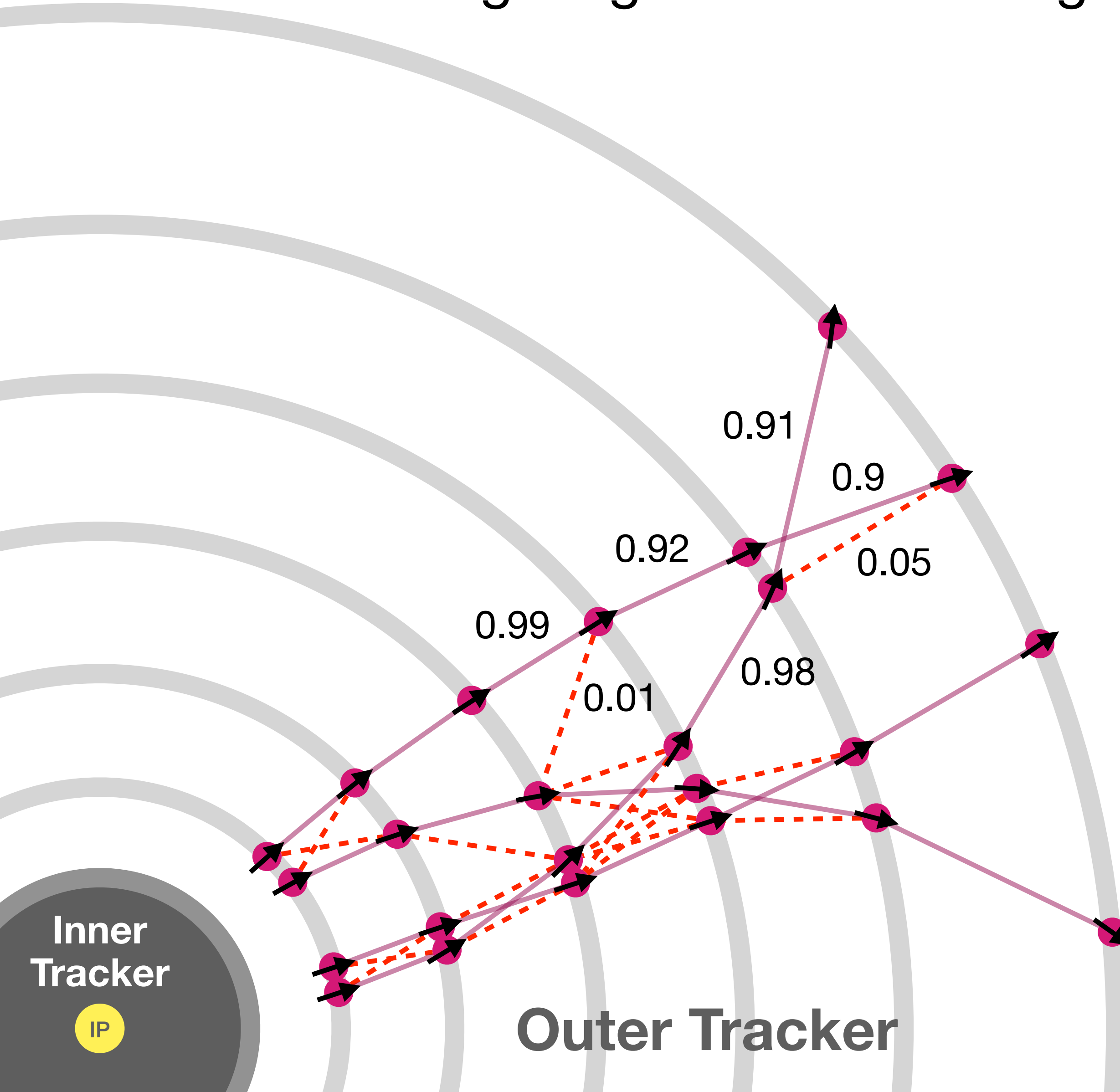O(ms/event)

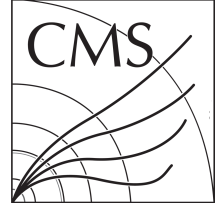**Inner Tracker**
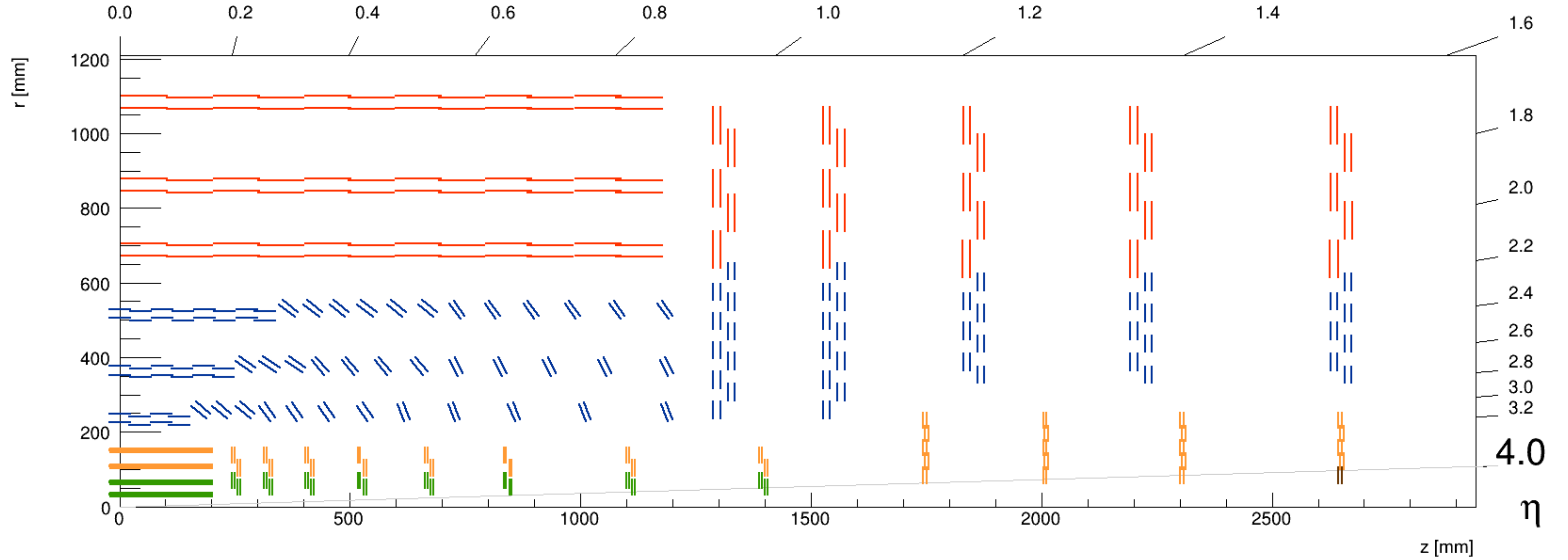
IP

**Outer Tracker**

UC San Diego

# Conclusion

- LST is a highly performant and parallelizable tracking algorithm

- We are investigating how LST could be improved with ML

- We have **improved LST with a lightweight DNN**

  - The DNN has **no impact on the runtime**

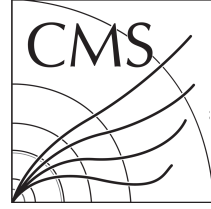- We are working towards training a LST GNN→track candidate pipeline
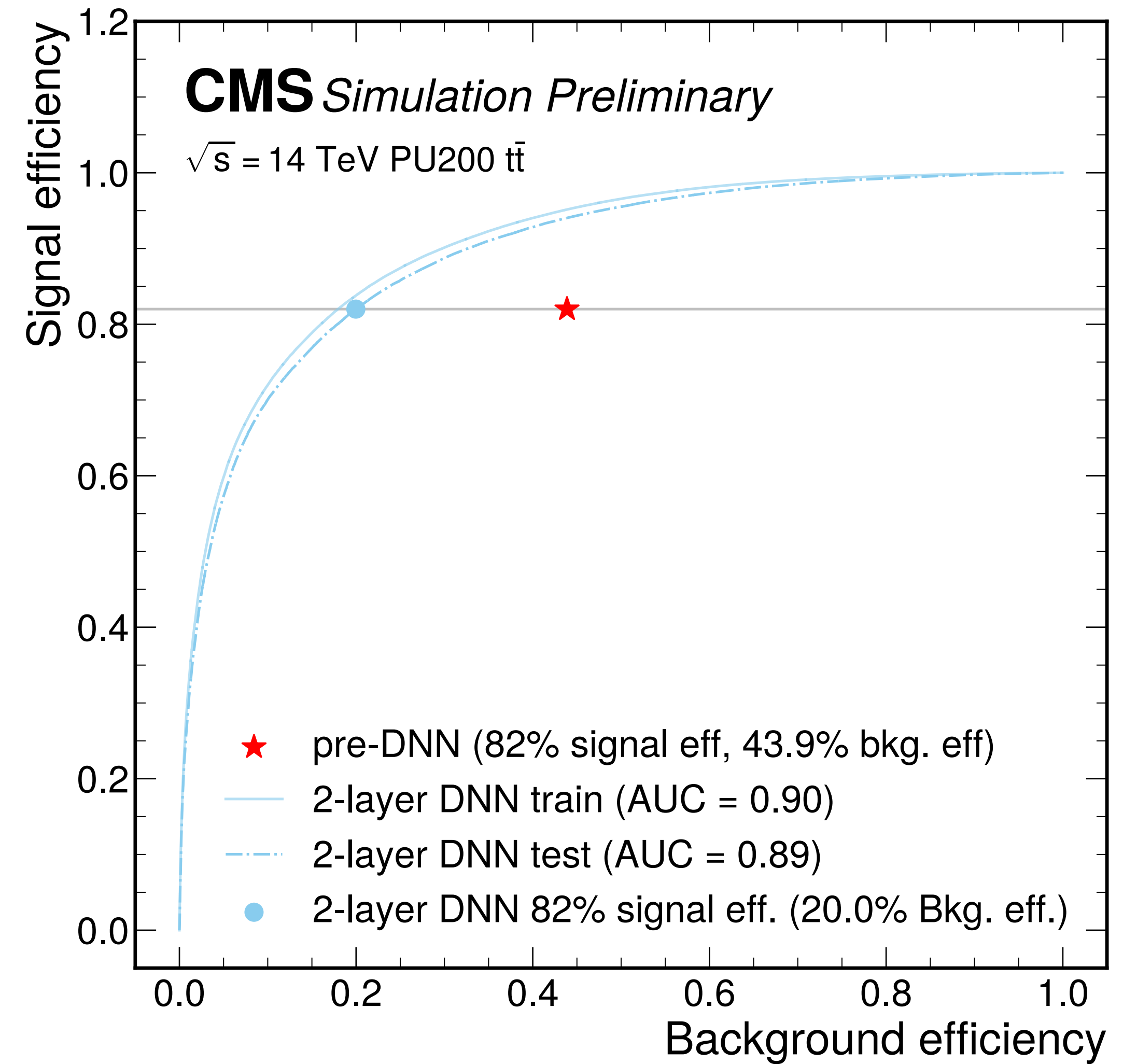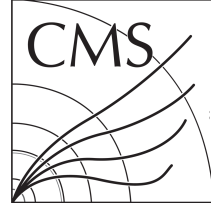

**Thank you!**

# Backup

# Phase 2 Tracker Layout

- Background efficiency = FPR = FP/N

  - i.e. Fake rate

- Signal efficiency = TPR = TP/P

- For DNN ROC curve:

  - TP|FP = # real|fake T5s after DNN cut

  - P|N = # real|fake T5s before DNN cut

- For pre-DNN performance (★):

  - TP|FP = # real|fake T5s after r-φ $\chi^2$ cut

  - P|N = # real|fake T5s before r-φ $\chi^2$ cut



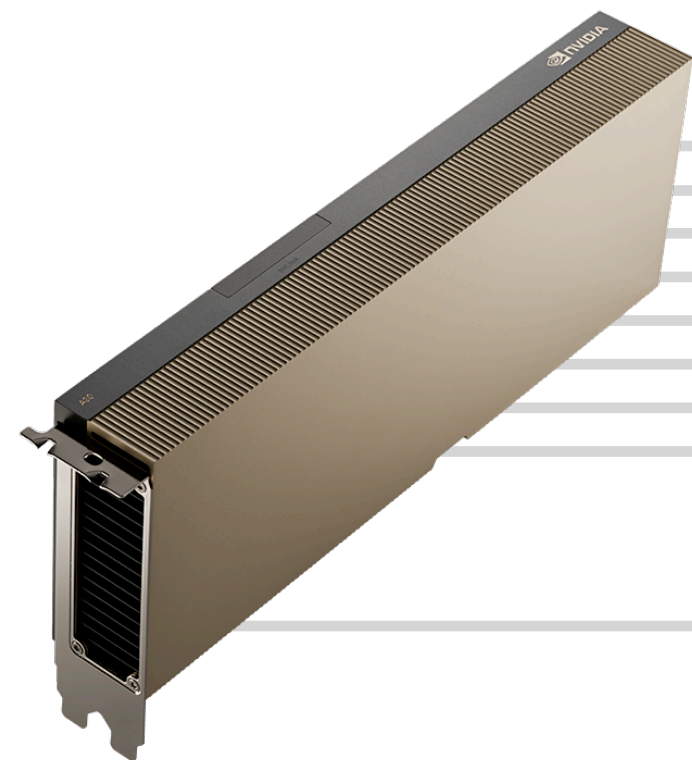CMS *Simulation Preliminary*

$\sqrt{s}$ = 14 TeV PU200 t$\bar{\text{t}}$

★ pre-DNN (82% signal eff, 43.9% bkg. eff)
— 2-layer DNN train (AUC = 0.90)
-·- 2-layer DNN test (AUC = 0.89)
● 2-layer DNN 82% signal eff. (20.0% Bkg. eff.)

Signal efficiency

Background efficiency

UC San Diego

# T5 DNN in LST

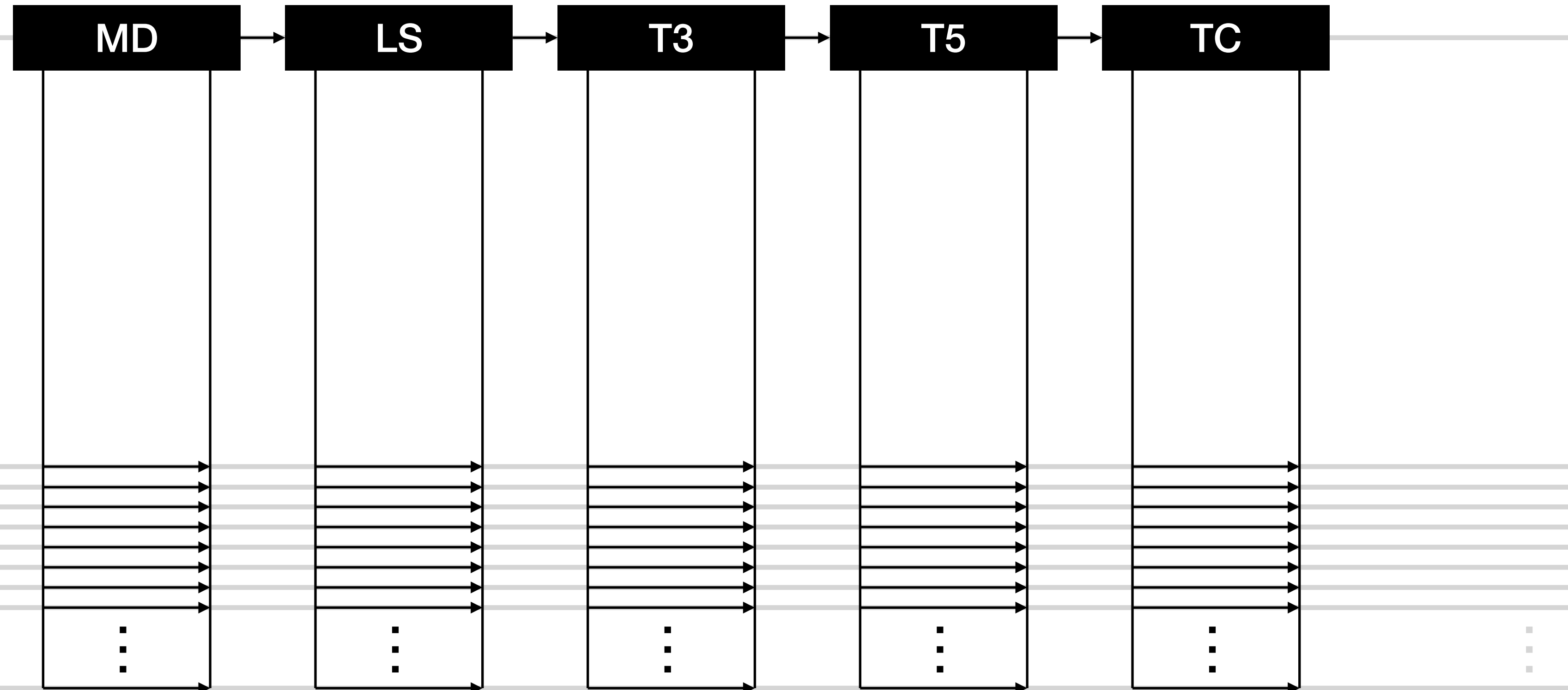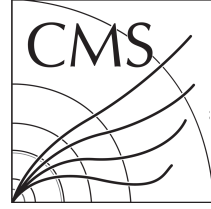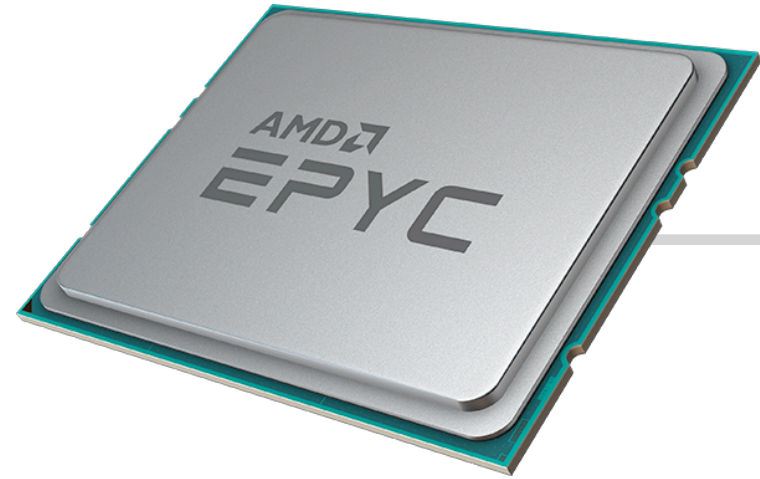## For each object, one thread per candidate



e.g. for T5 kernel, each thread determines whether
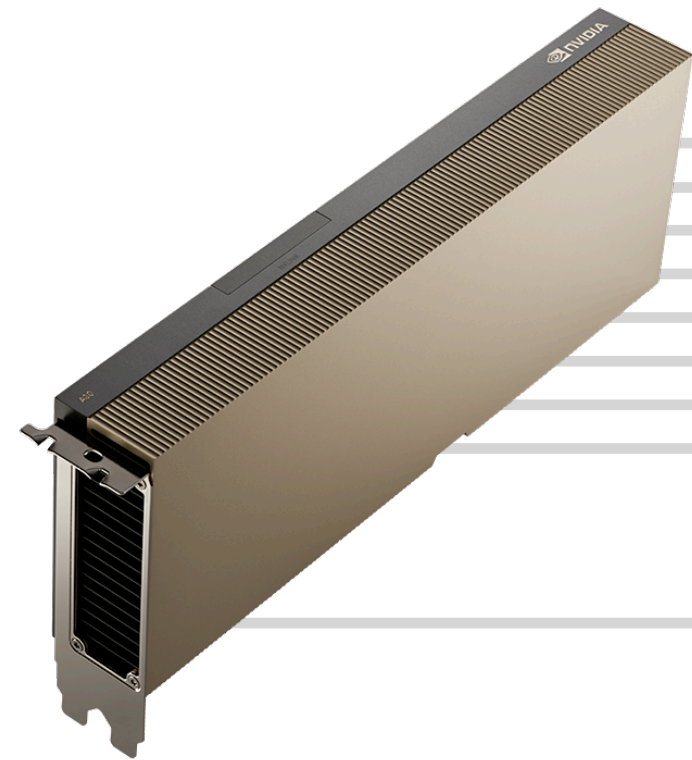or not to pass a given T5 candidate to the next step

# T5 DNN in LST

## Each thread (one per T5) runs the ML inference!

Host

Device

```
def passT5QualityCutsPseudoCode(...)
{
    // Original LST cuts
    if (!passBasicT5QualityCuts(...)) { return false; }
    if (!passRZChi2T5Cut(...)) { return false; }

    // Build DNN input features vector
    float x[38] = {
        log10(innerT3pT),
        innerT3eta,
        ...
    };

    // Input -> first hidden layer
    float hidden0[32];
    for (int col = 0; col < 32; ++col) {
        hidden0[col] = 0.f;
        for (int inner = 0; inner < 38; ++inner) {
            hidden0[col] += x[inner]*wgts0[inner][col];
        }
    }
    hidden0 = leakyReLU(hidden0);
    ... // and so on…

    // Last hidden layer -> output
    float inference = 0.f;
    for (int i = 0; i < 32; ++i) {
        inference += hidden1[i]*wgts4[i][0];
    }
    inference = sigmoid(inference);

    if (inference < LSTDNN::WP95) { return false; }

    return true;
}
```
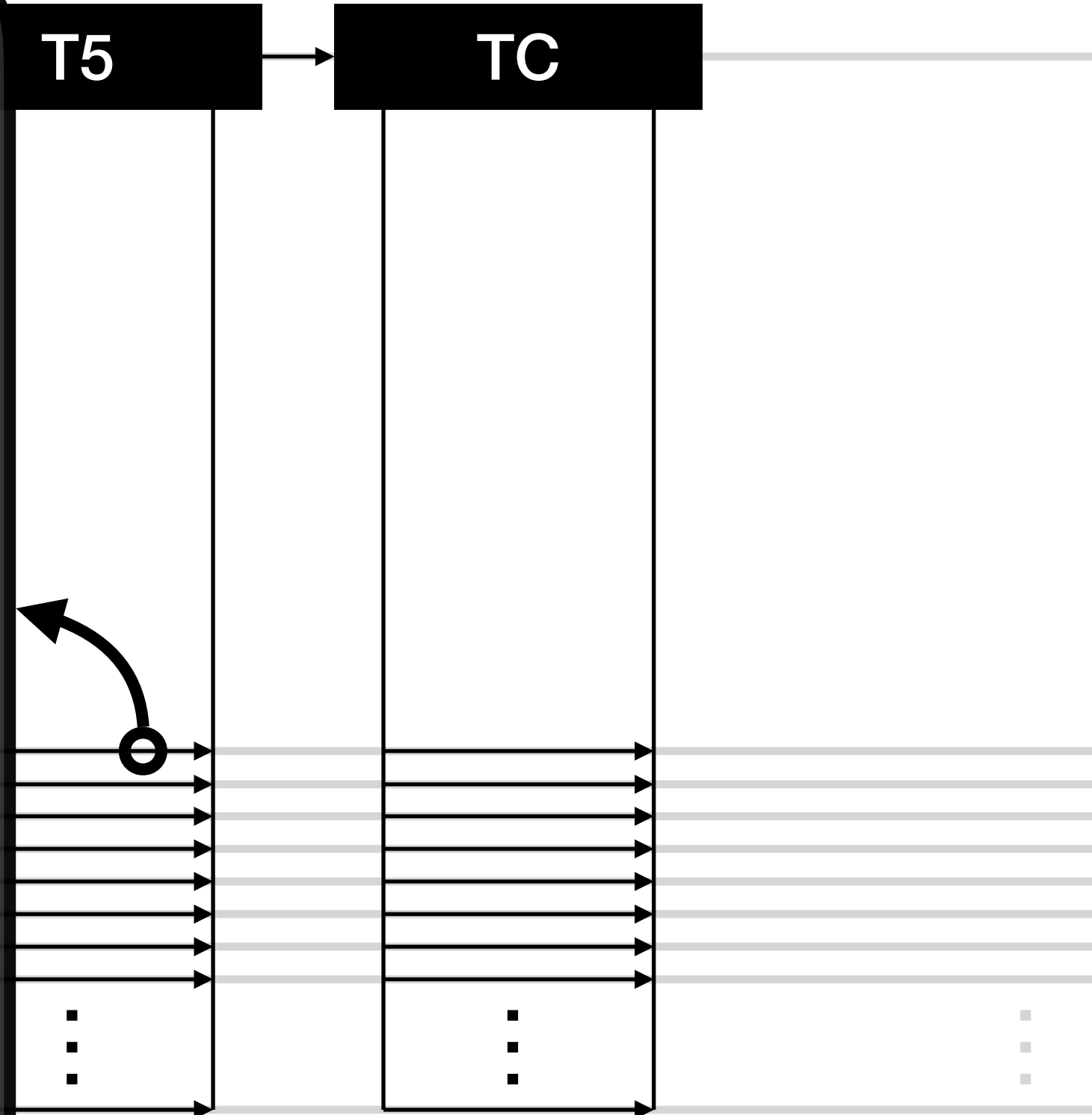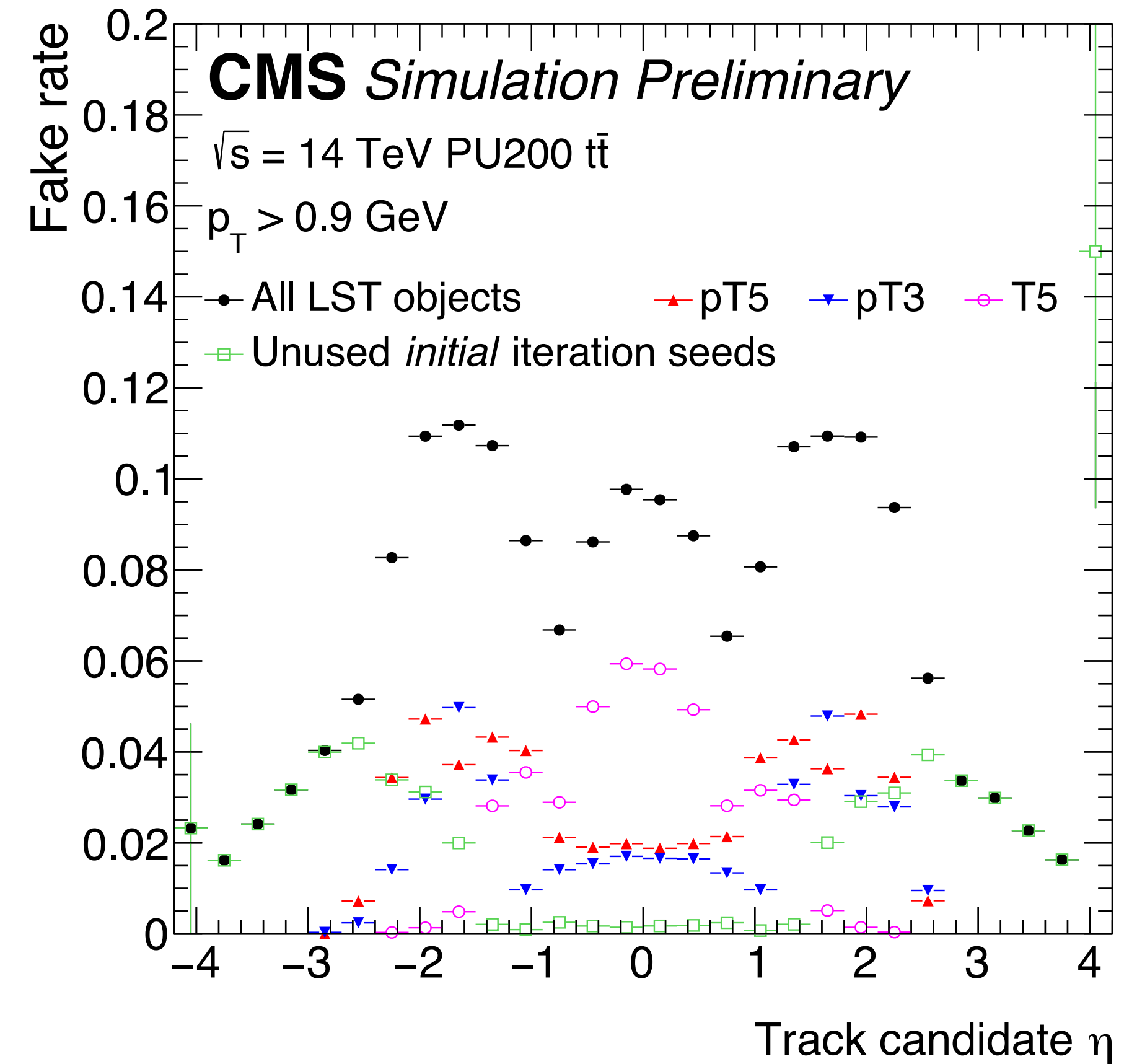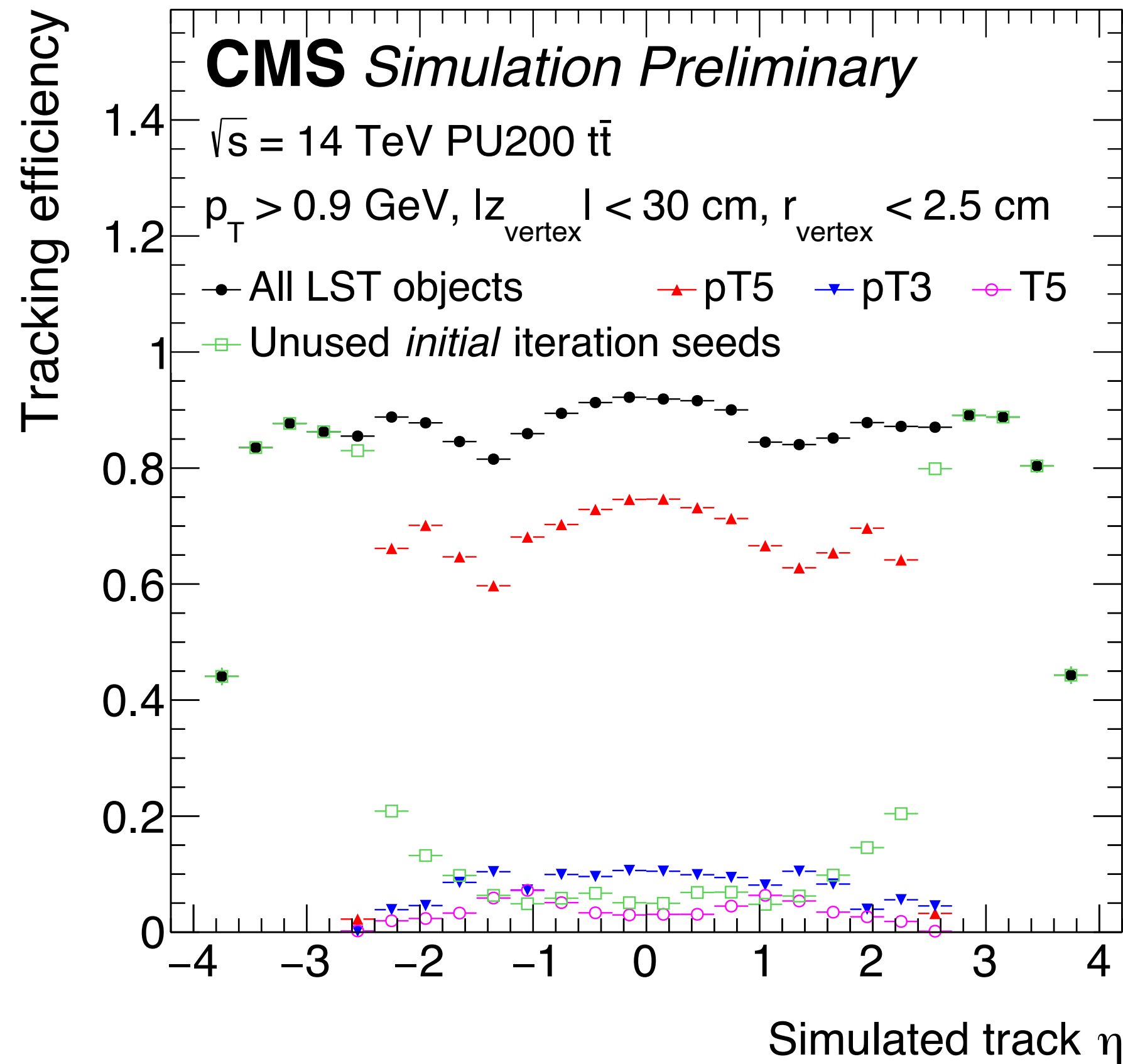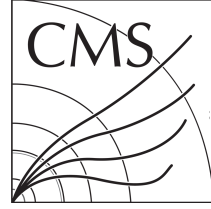
T5    TC

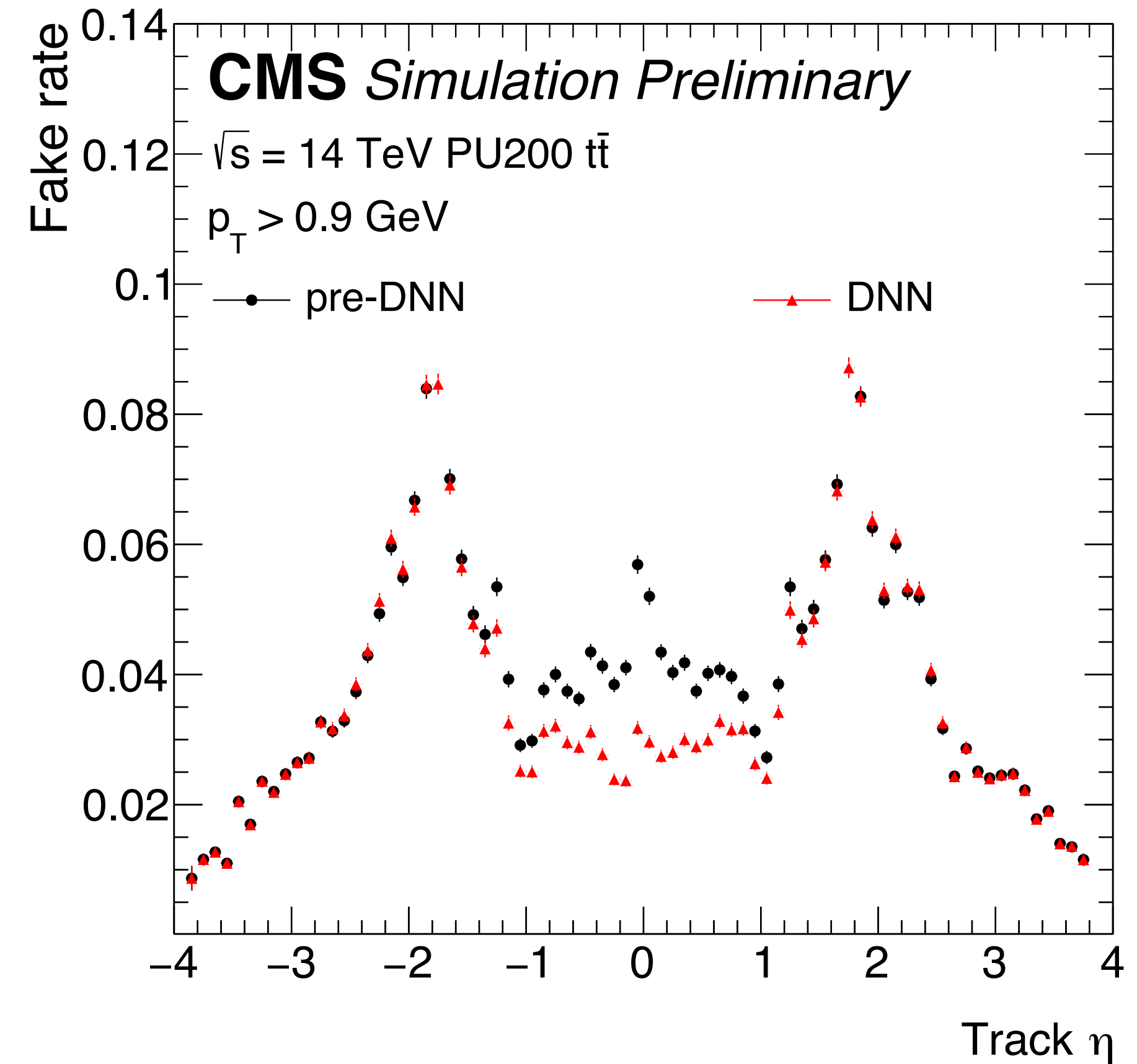**Very simple implementation (certainly not optimal)**
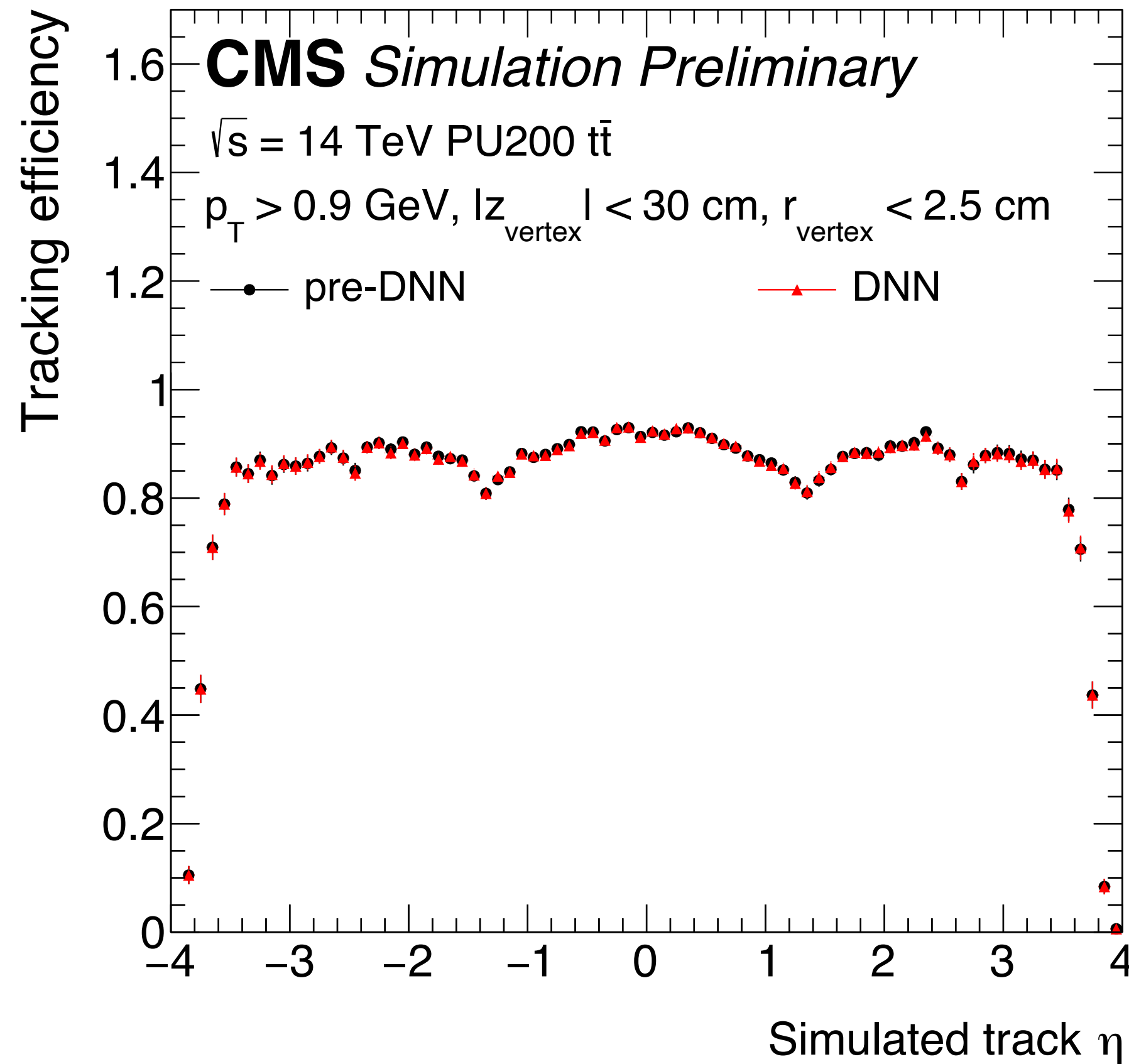
UC San Diego

# Pre-DNN Performance



**T5s account for much of the LST efficiency**
**T5s have a high fake rate**

UC San Diego

# Performance Comparison



**Significant reduction of fake rate in the barrel**
**No loss in efficiency**

UC San Diego