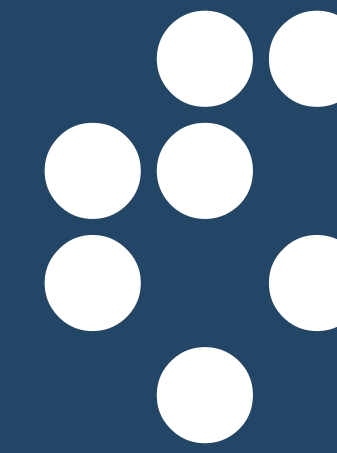




SMASH
machine learning for science and humanities postdoctoral program



Jožef Stefan
Institute

SIMULATION OF PARTICLE PHYSICS SILICON DETECTORS WITH TRANSFORMERS

1st SMASHING WORKSHOP
October 7, 2024

I FEEL
SLOVENIA

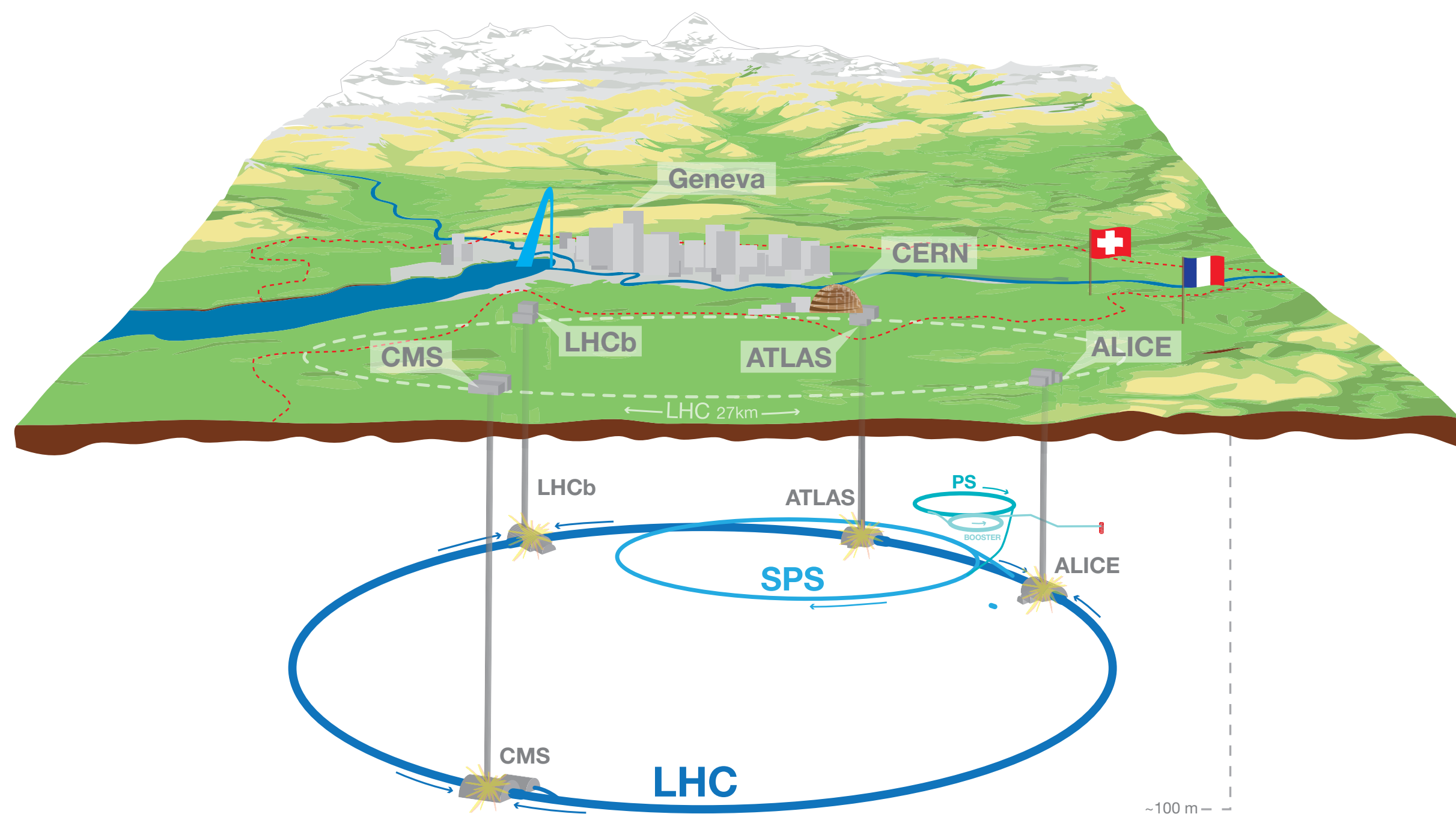


Co-funded by
the European Union

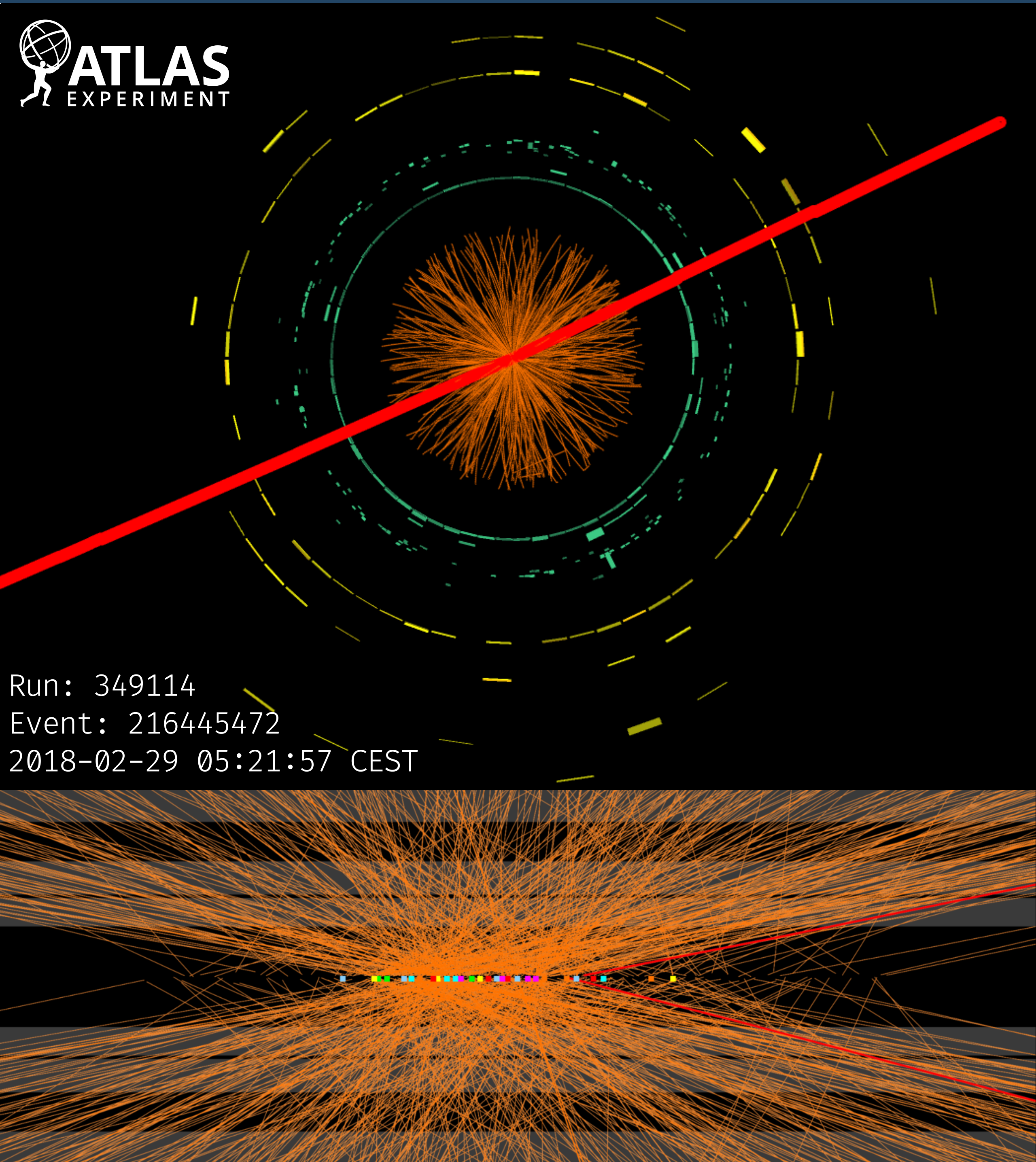
Tadej Novak
Jožef Stefan Institute



- Largest particle collider — circumference of 27 km:
 - up to 40 million proton-proton collisions per second
- HL-LHC upgrade targeting 2030.
 - data rate 7-10 times greater
 - average number of collisions per bunch crossing rising to as much as 200, from 30-60 currently



Source: CERN



- Largest particle collider — circumference of 27 km:
 - up to 40 million proton-proton collisions per second
- HL-LHC upgrade targeting 2030.
 - data rate 7-10 times greater
 - average number of collisions per bunch crossing rising to as much as 200, from 30-60 currently
- ATLAS detector a general purpose experiment.
 - Need to measure particle momentum and energy.

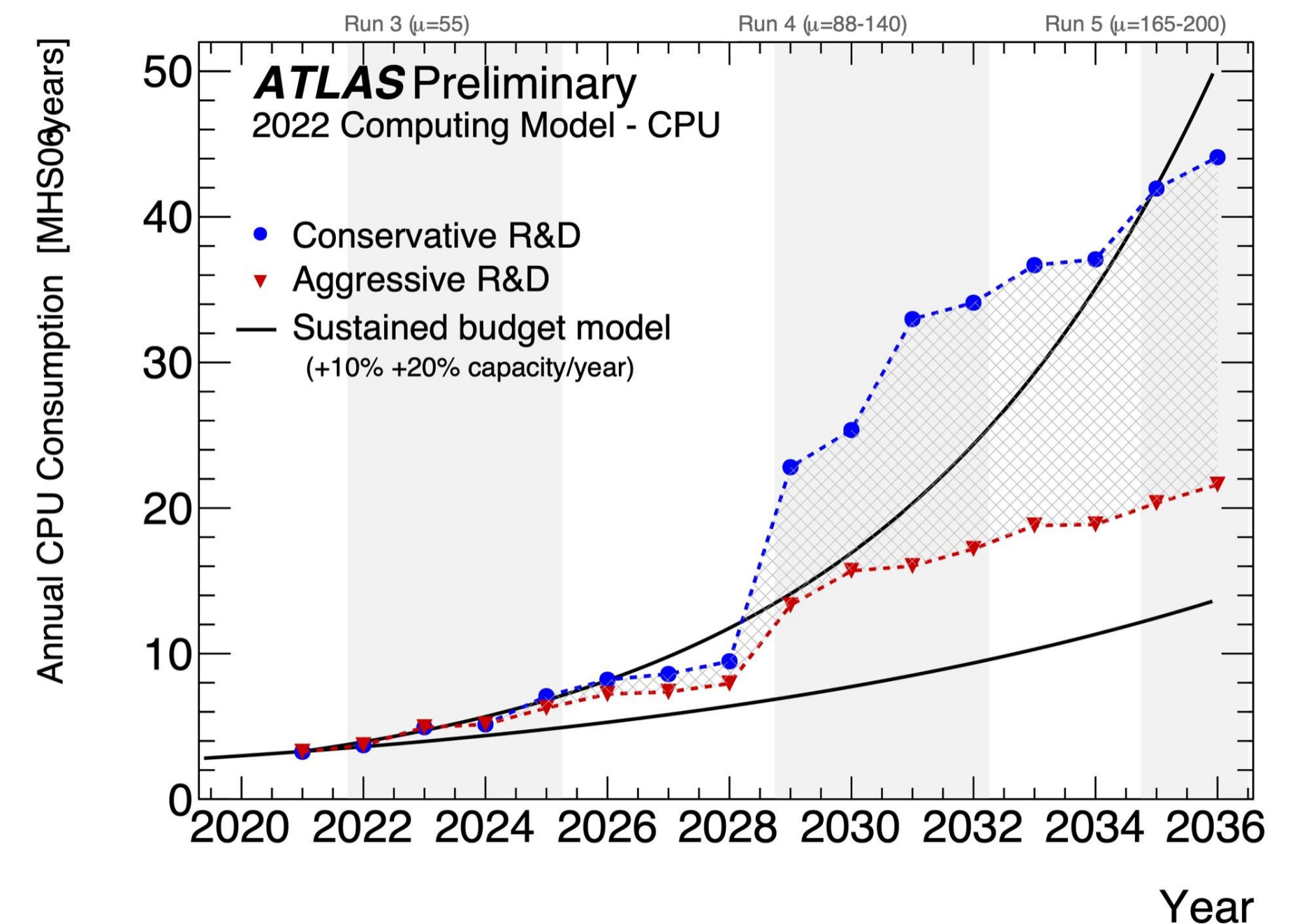


- A large part of the LHC physics programme relies on **accurate Monte Carlo simulation of collision events**.
 - every single particle needs to be simulated
 - detailed (full) detector response simulation most intensive
- Producing simulated samples → majority of experiments' CPU requirements
 - CMS used 85% CPU for Monte Carlo production during 2009-2016
 - half spent detector simulation



- A large part of the LHC physics programme relies on **accurate Monte Carlo simulation of collision events**.
 - every single particle needs to be simulated
 - detailed (full) detector response simulation most intensive
- Producing simulated samples → majority of experiments' CPU requirements
 - CMS used 85% CPU for Monte Carlo production during 2009-2016
 - half spent detector simulation

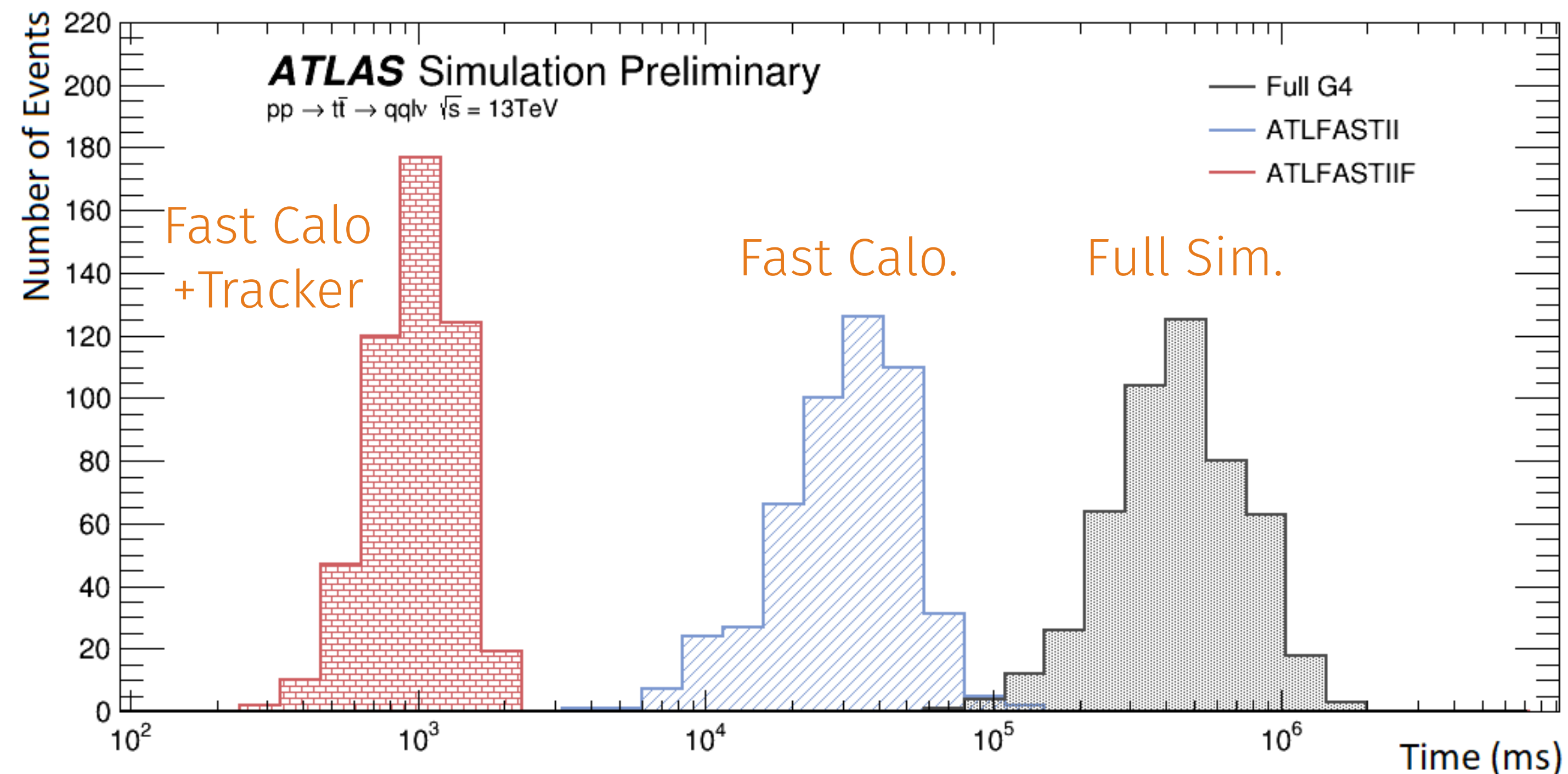
Source: ATLAS Software and Computing HL-LHC Roadmap



- Current methods do not scale with HL-LHC data rates and **more aggressive R&D is needed**.

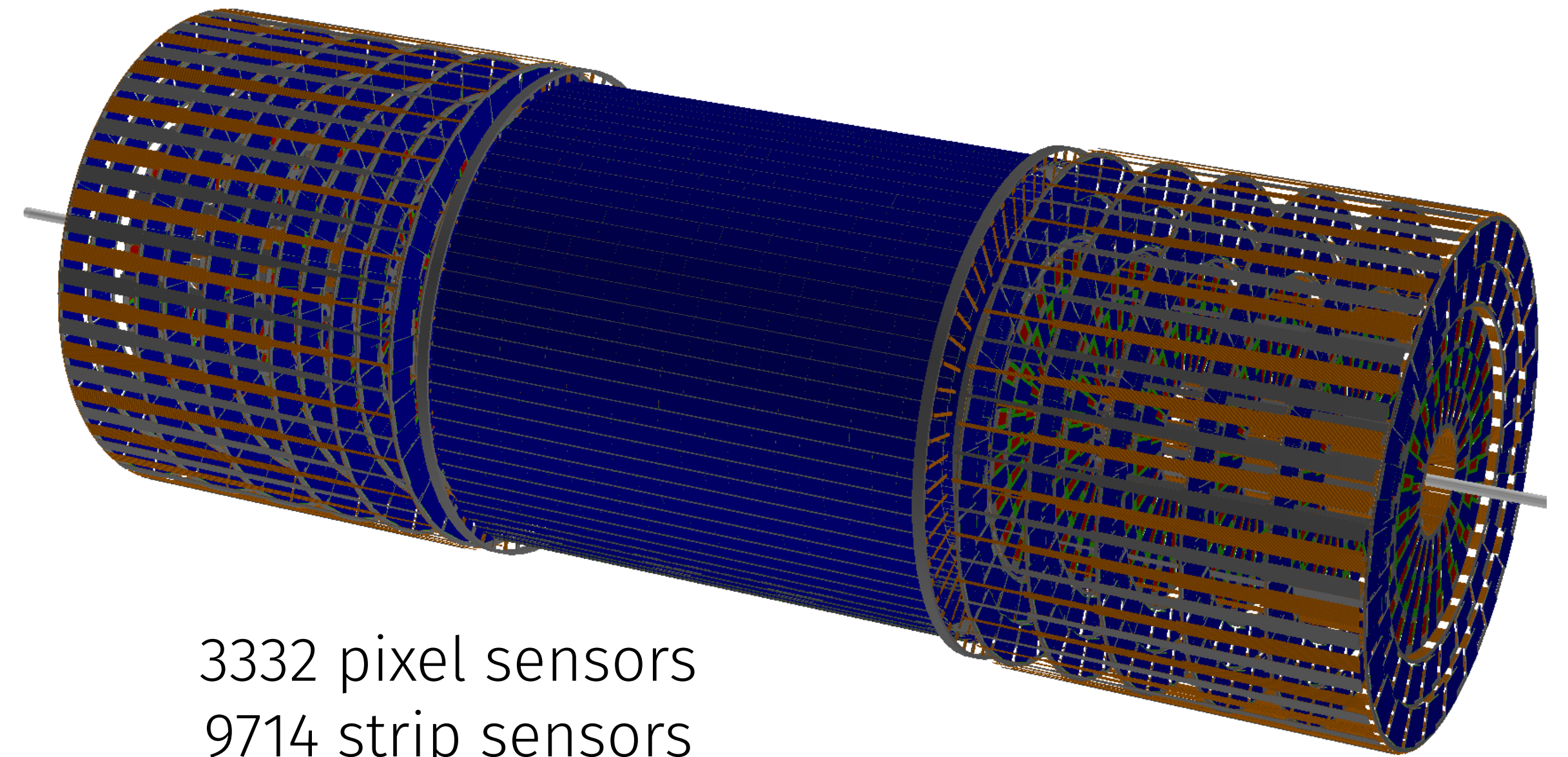


- Large efforts to speed-up simulation — **fast simulation**.
 - Detector response to a particle is **parameterised**.
- Fast simulation for particle physics successfully applied at calorimeter level.
 - Derived for different particles at different energies and parts of calorimeter.
 - Generative neural networks also used.
 - Order of magnitude speed-up achieved.
- Tracking detectors fast simulation not production-ready yet.
 - **Machine learning target of this project.**



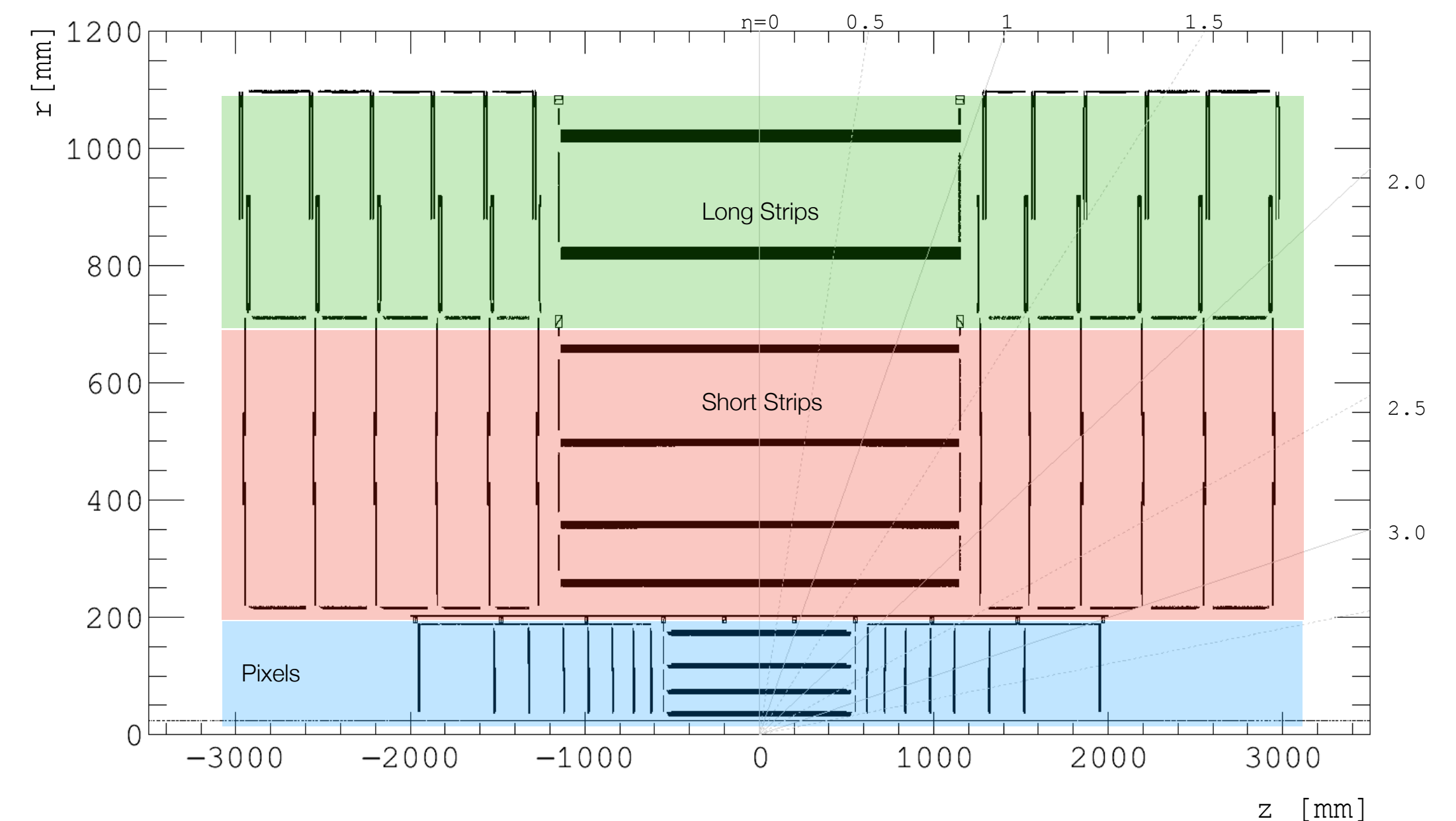


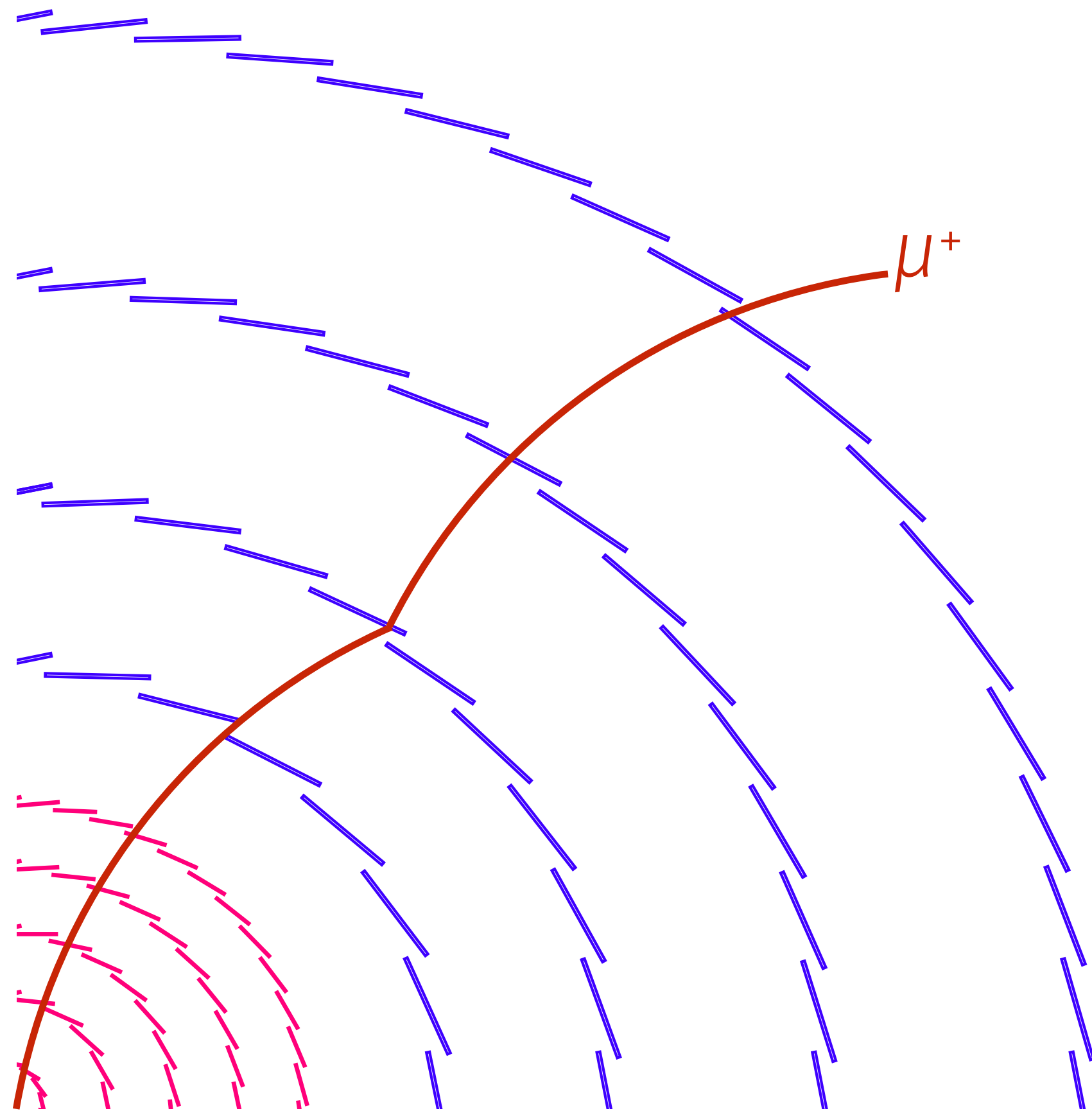
- A generic, HL-LHC style tracking detector.
- Each sensor split into multiple readout channels.
 - Can be described as a 2D surface.
- Goal to be reasonably close to a real-world detector.
 - Loosely modelled after the ATLAS ITk (58700 sensors, ~5 billion electronic channels).
- Ensures the ability to generalise R&D projects for silicon tracking detectors.



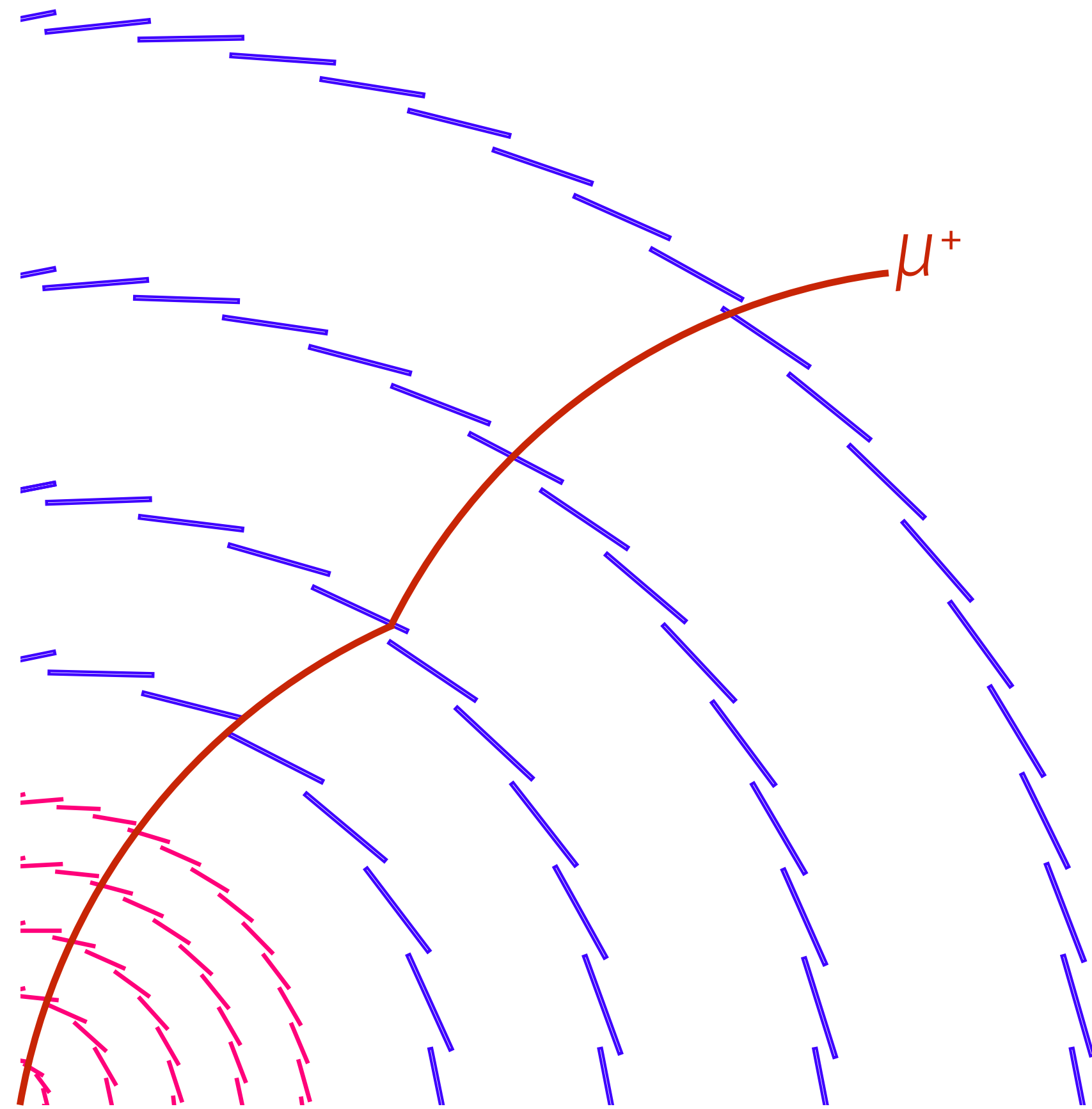
3332 pixel sensors
9714 strip sensors

Source: The Open Data Detector Tracking System

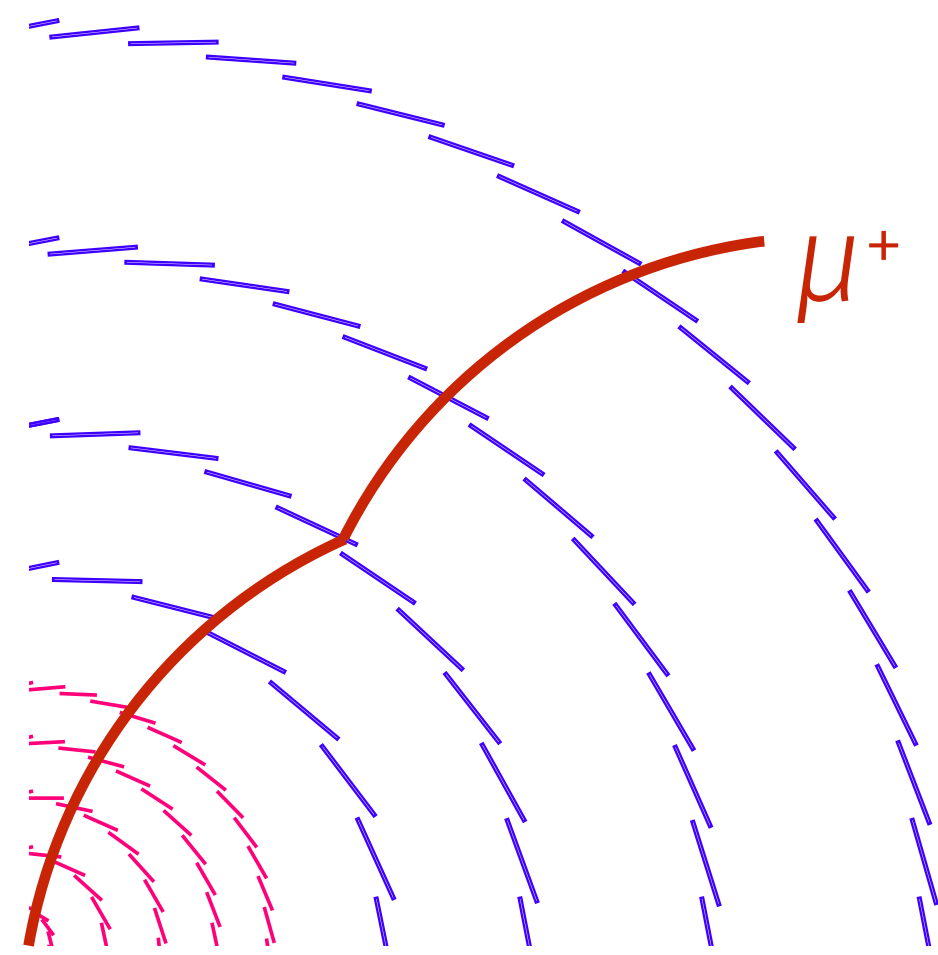




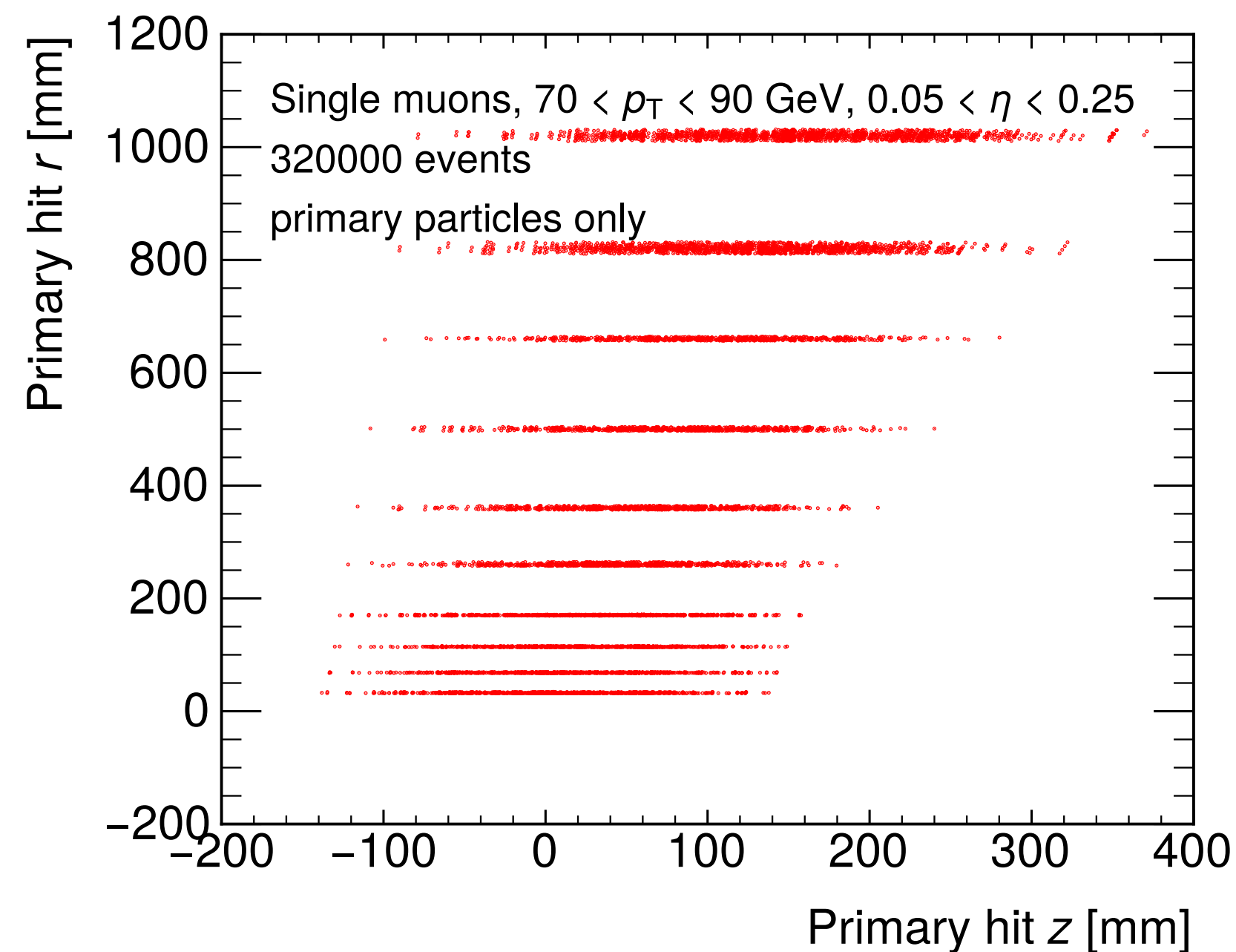
- Sensitive to **charged particles** (electrons, muons, charged hadrons).
 - curved track (helix) due to magnetic field
 - simulating induced electric charge in sensitive detector elements



- Sensitive to **charged particles** (electrons, muons, charged hadrons).
 - curved track (helix) due to magnetic field
 - simulating induced electric charge in sensitive detector elements
- Multiple processes need to be simulated but for now limited to:
 - single muon events
 - multiple-scattering

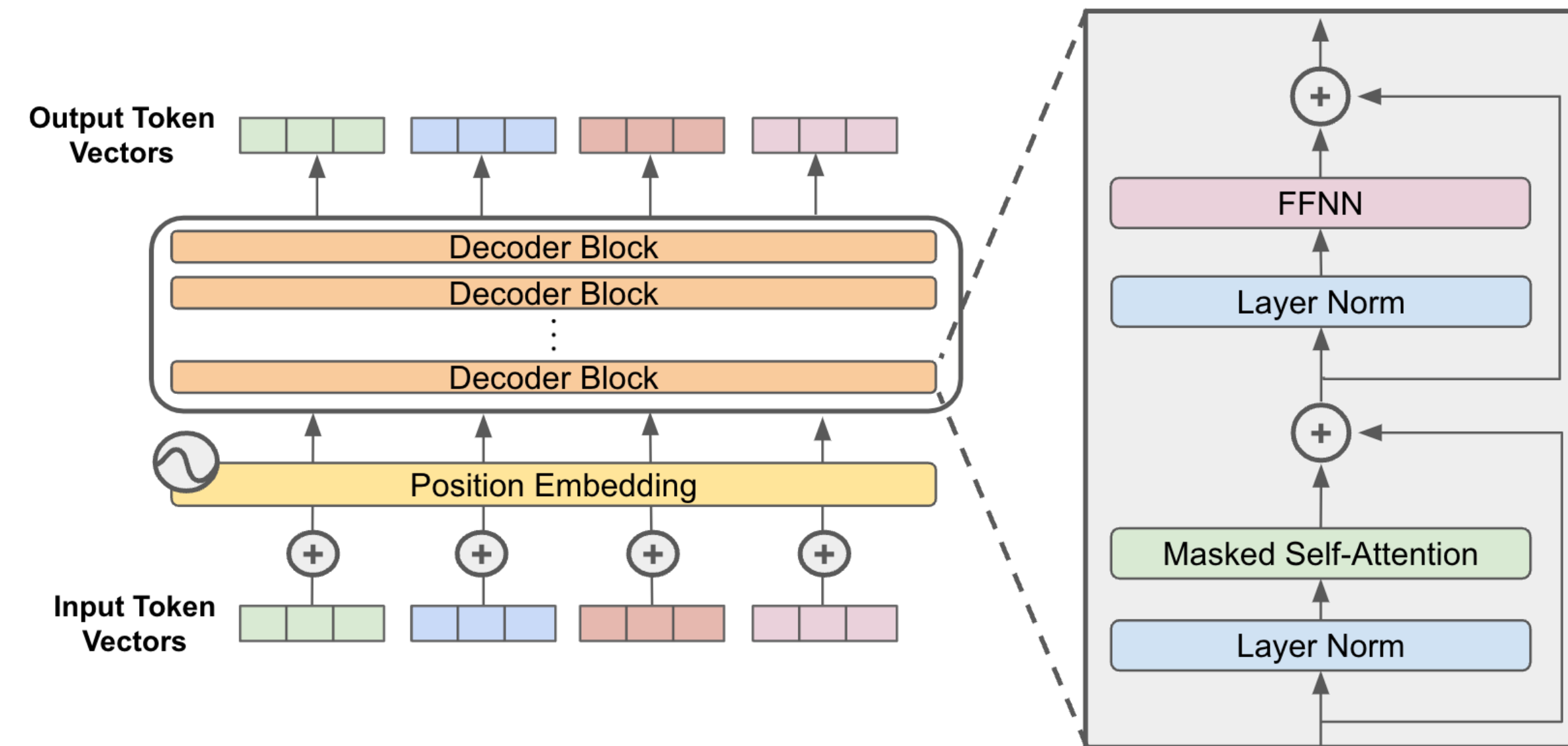


- A sequence of detector hits.
 - With additional start and end “virtual hit” to describe input and output state with the same data structure.
- 7 features per hit:
 - particle ID + geometry ID
 - particle momentum (after the hit)
 - hit position on the sensitive detector (local)
- Each hit is an element of a sequence, each particle has its own sequence.
- Local coordinates taken to **constrain hits on the sensitive parts** and prevent them happening in the vacuum.





- Transformers popular nowadays to deal with sequential data (most commonly LLMs), see [1706.03762](#).
- Using **decoder-only** architecture.
 - Input/output data are the same.
 - Target to predict the next element of the sequence.
 - The famous example are the GPT family of models.
- Specialised on discrete sequences which are **tokenised** (sequential integers).
 - Can be anything e.g. words, detector modules, ...
- For this application all **continuous data is discretised** (rounded to two decimal points) and **each feature is tokenised separately**.



Source: Cameron R. Wolfe

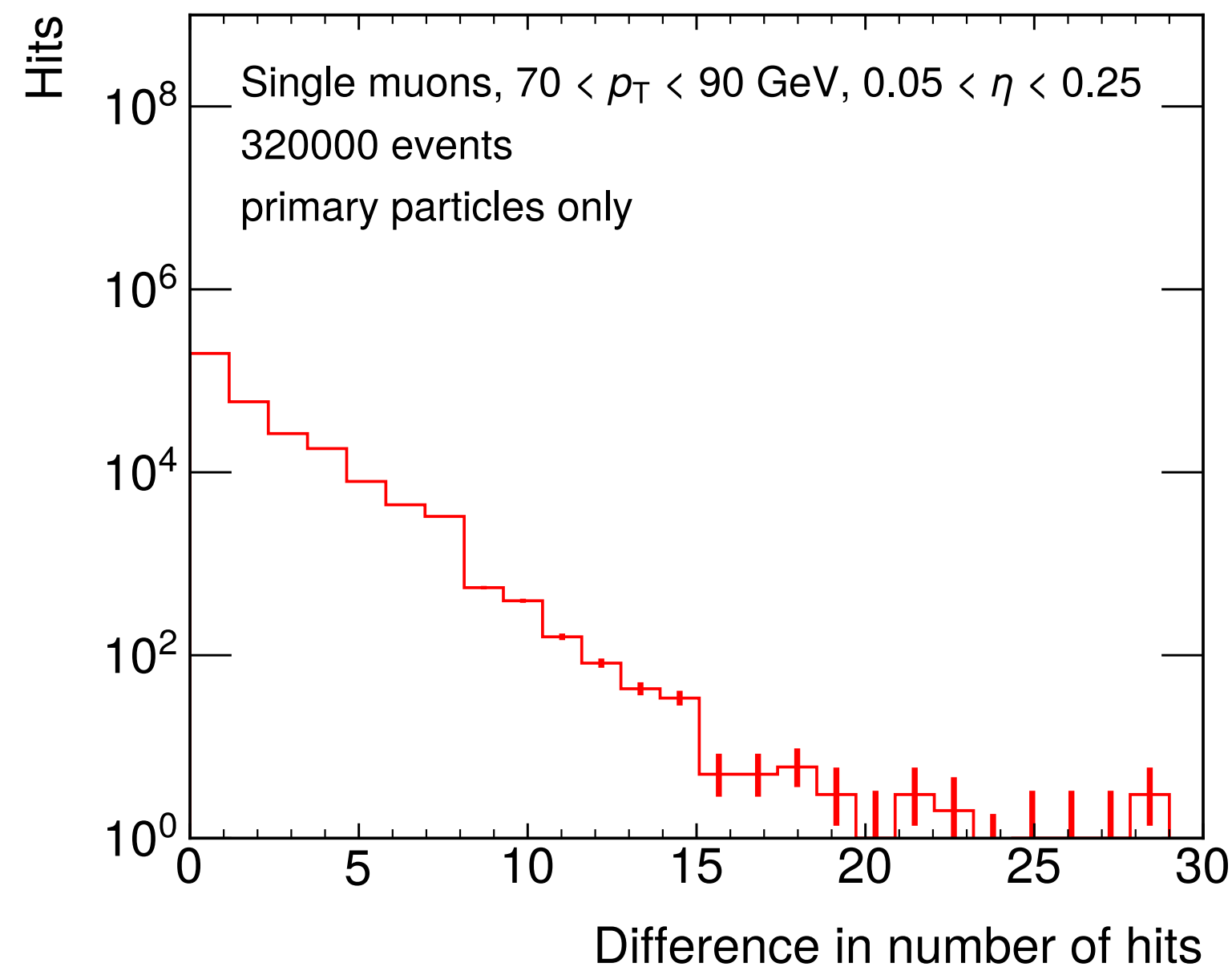
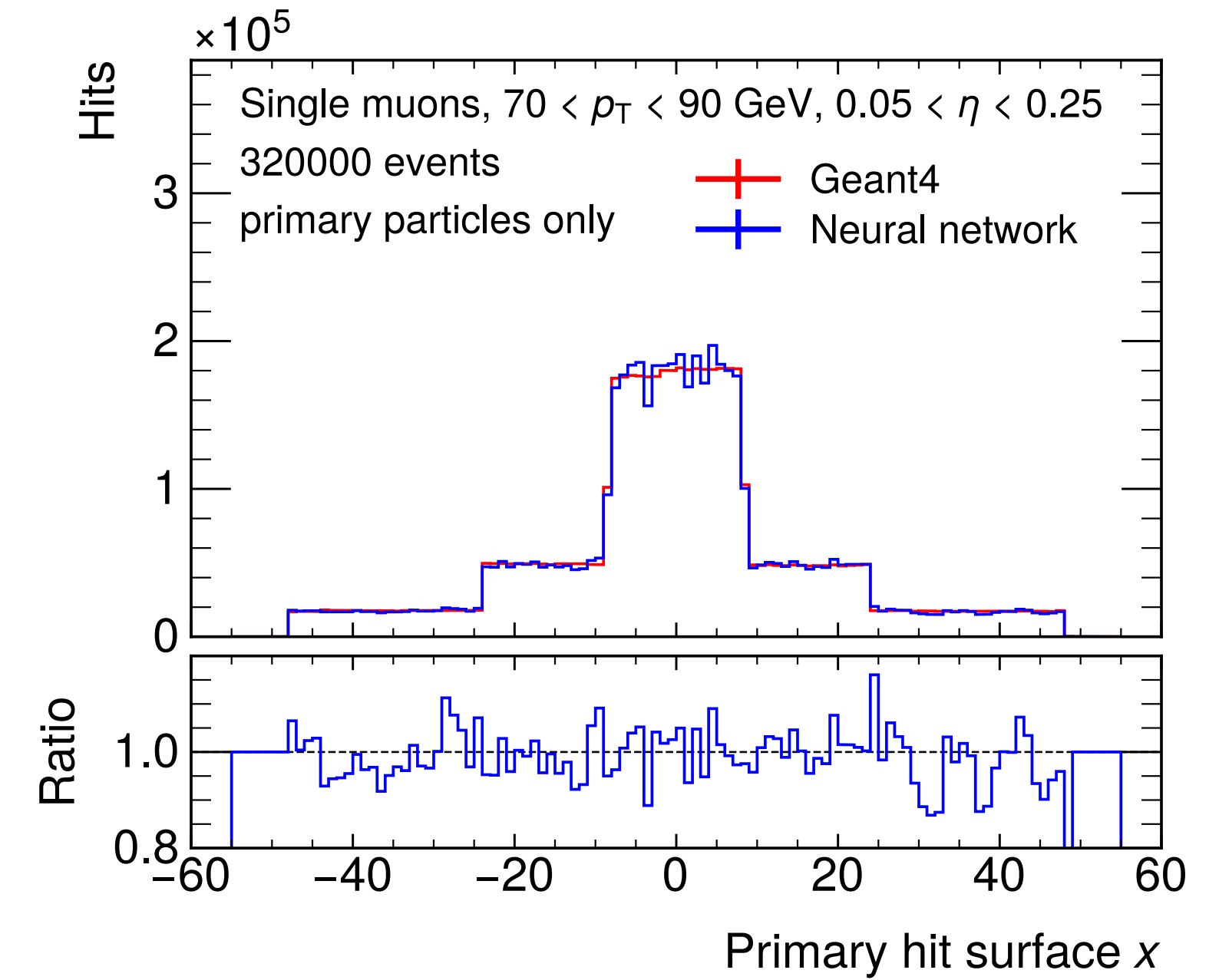
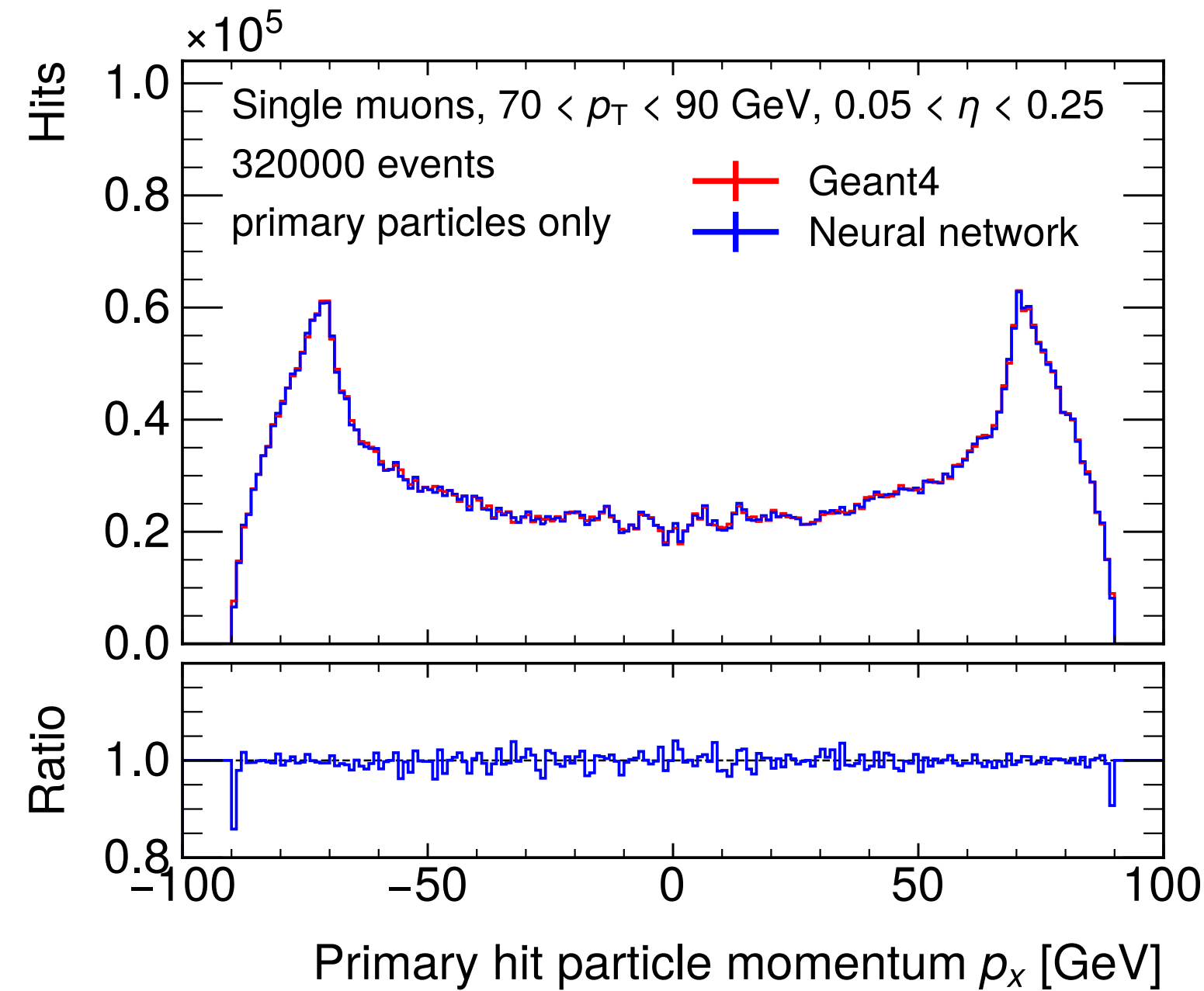
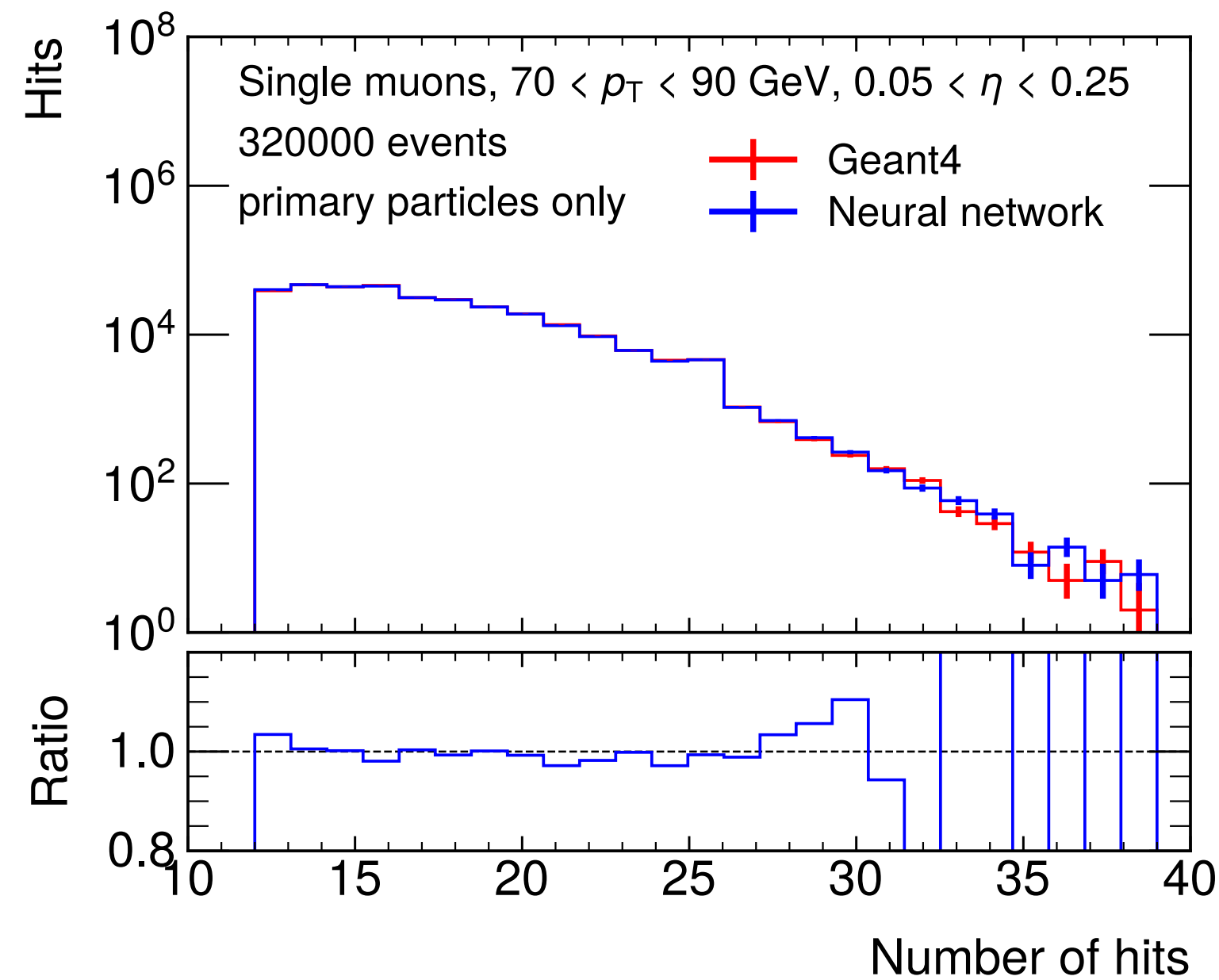


- **Sample details:**
 - single muons, $70 < p_T < 90$ GeV, $0.05 < \eta < 0.25$
 - 320000 events
 - training : validation : test = 2 : 1 : 1
 - augmented with random numbers between 1 and 10000
- Training performed on a workstation with a **NVIDIA GeForce RTX 4090 GPU**.
 - Duration ~1 week.
 - Vega also used but not for these specific results.
- Learning rate variation using cosine annealing with warm restarts with a period of one epoch and fixed amplitude.
- **Inference:** most probable next sequence element

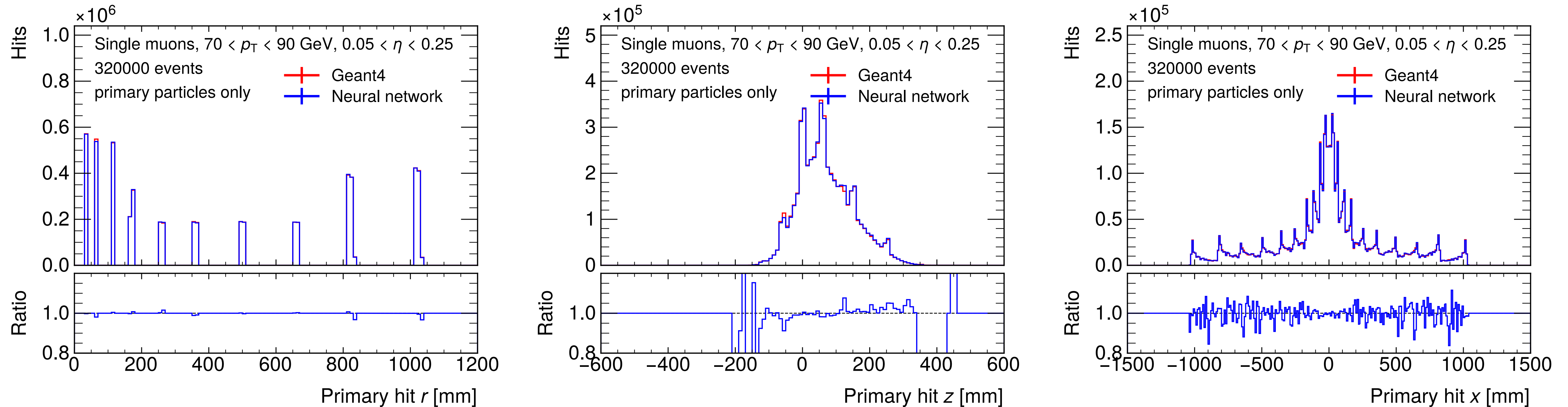
Model Parameter	Value
input dimension	128
layers	3
heads	4
feedforward dim.	512
activation	GELU
dropout	0.1

Training Parameter	Value
epochs	6000
optimizer	AdamW
learning rate	0.001
weight decay	0.01
gradient clipping	5.0
batch size	512

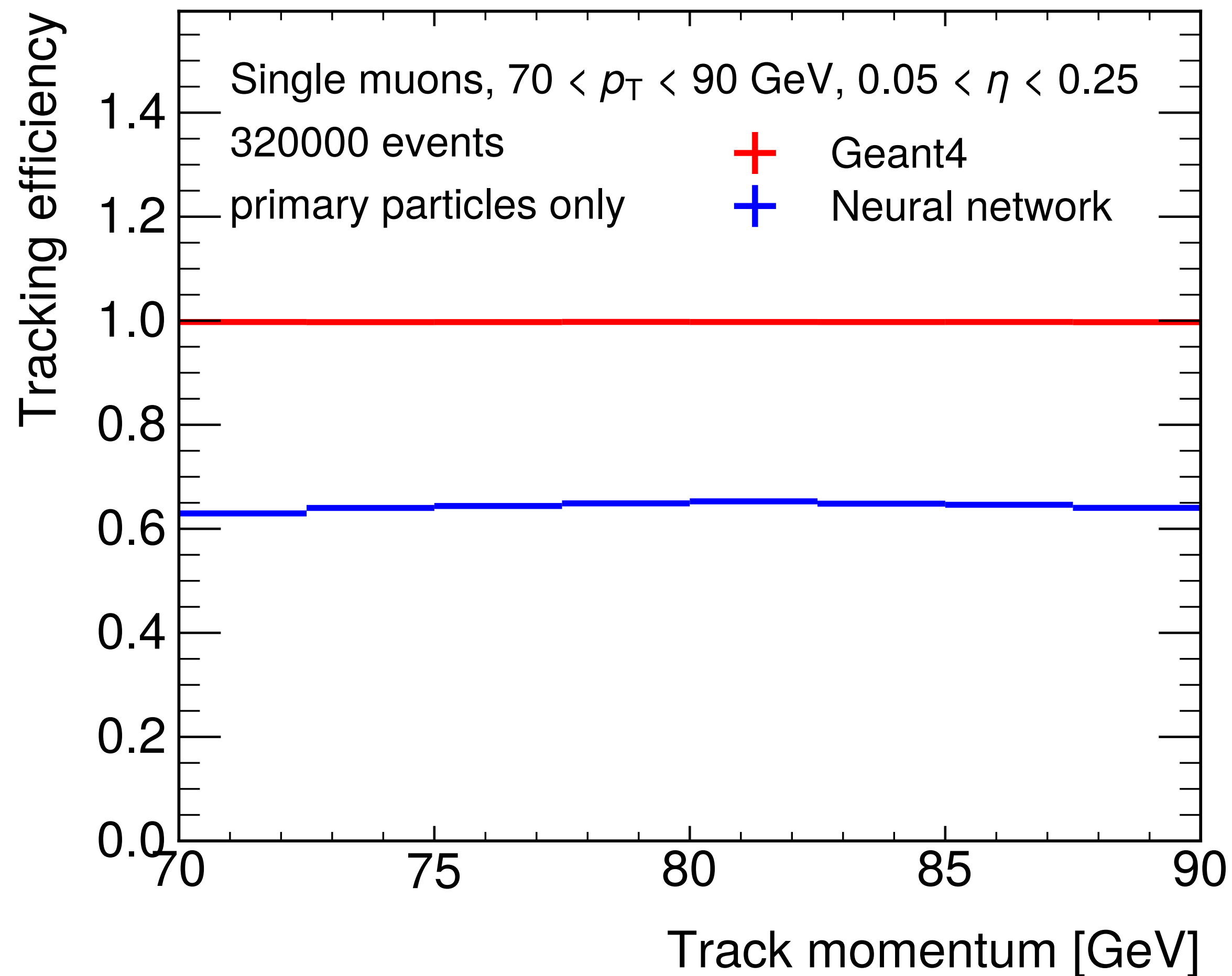




- Good agreement with full simulation.
 - Coordinates fluctuating a bit up to ± 10 %.
- Number of hits accurately reproduced.
 - Some difference seen but random numbers included in inference.



- Global coordinates show good agreement describing complex detector structure.
- Larger deviations in tails of the z -coordinate due to lower statistics.
- Inference also performed on the same GPU: ~ 4 s / 10k particles



- Evaluating performance using the **ACTS** (A Common Tracking Software) framework.
 - Default test setup for the Open Data Detector.
- Seeding efficiency only ~65 % compared to 99 % for full simulation.
 - Hit displacement from the estimated helix is too large.
 - Rounding has no significant effect on the reference sample.
 - Detailed investigations ongoing.



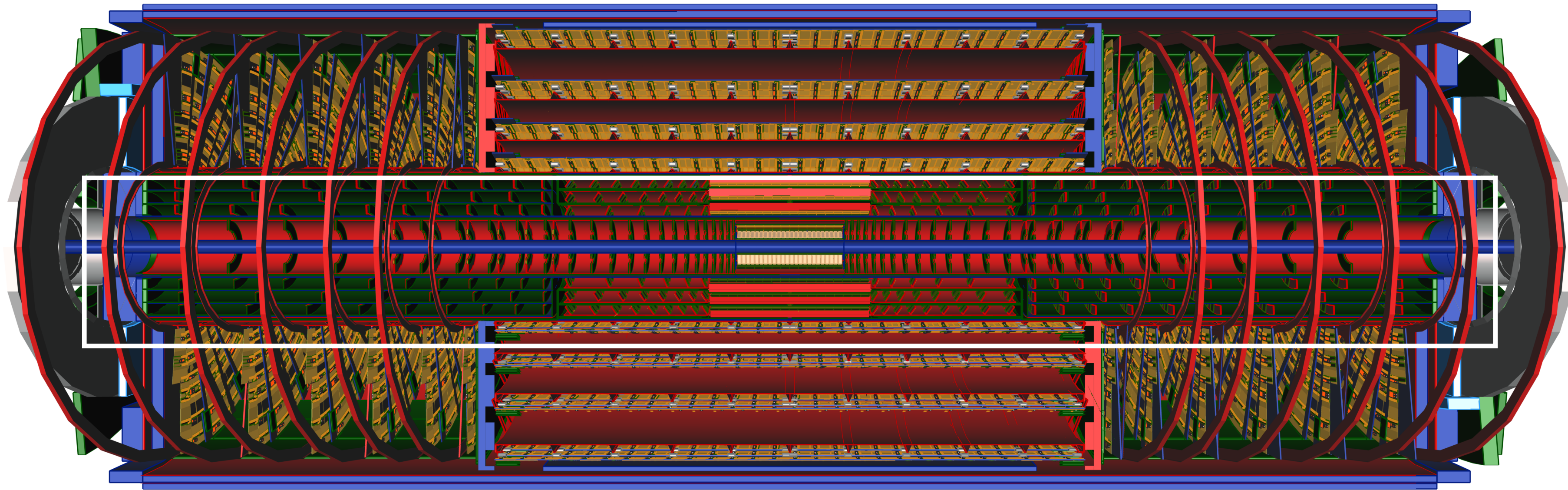
- Transformers can describe a sequence of physics data very well.
 - Physics performance not sufficient yet, needs optimisation.
- Training relatively long, but inference is fast.
- Future plans:
 - Optimise the current setup for better tracking performance.
 - Describe continuous features with floating point numbers.
 - Try proper generative sampling of a transformer.
- I want to thank for ideas and tips from my ATLAS collaboration colleagues and the Visual Cognitive Systems Laboratory at Faculty of Computer and Information Science.



**Co-funded by
the European Union**

This project has received funding from the European Union's Horizon Europe research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 101081355.

The operation (SMASH project) is co-funded by the Republic of Slovenia and the European Union from the European Regional Development Fund.



Source: [ATL-PHYS-PUB-2021-024](#)

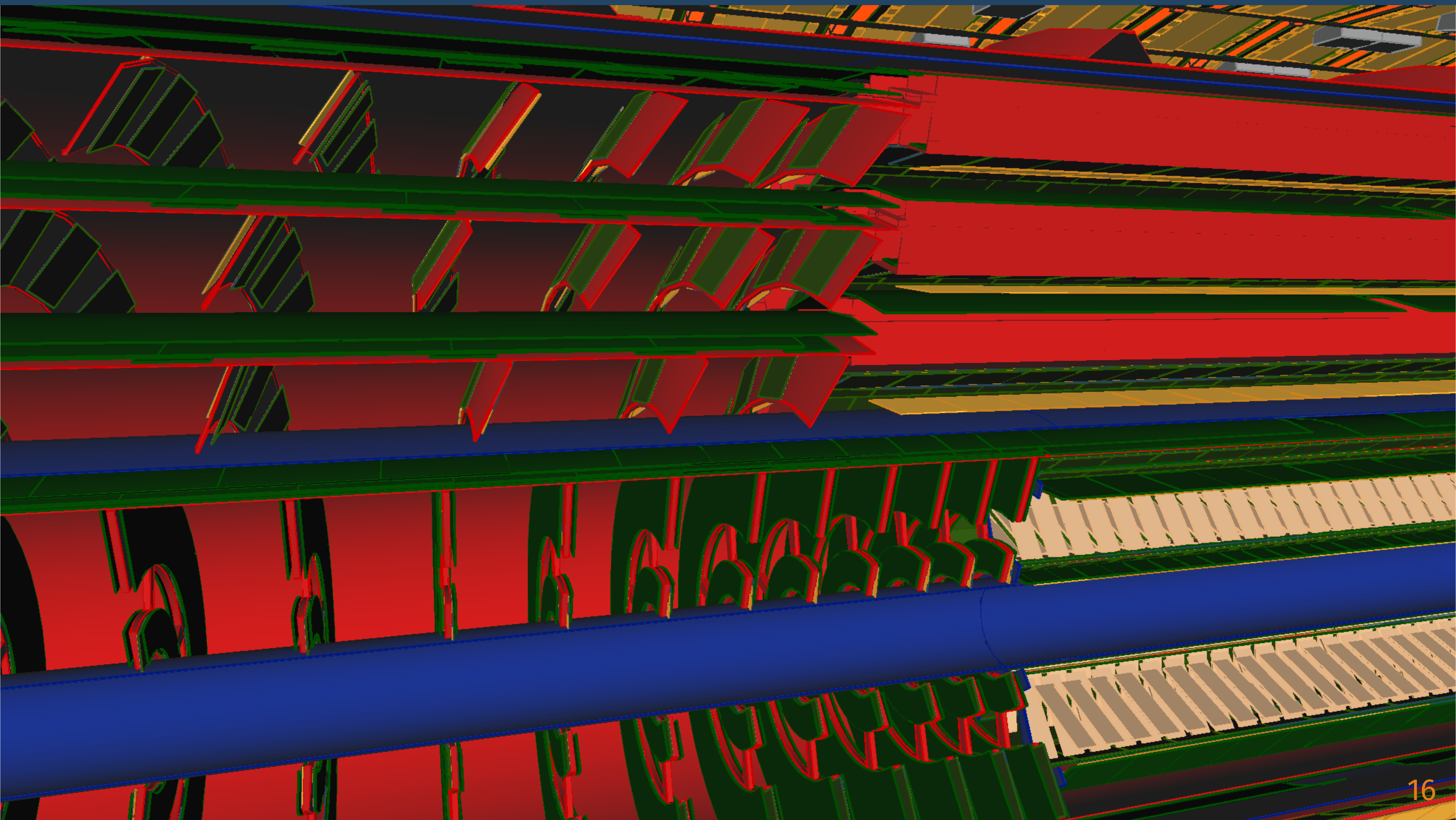
Pixel detectors

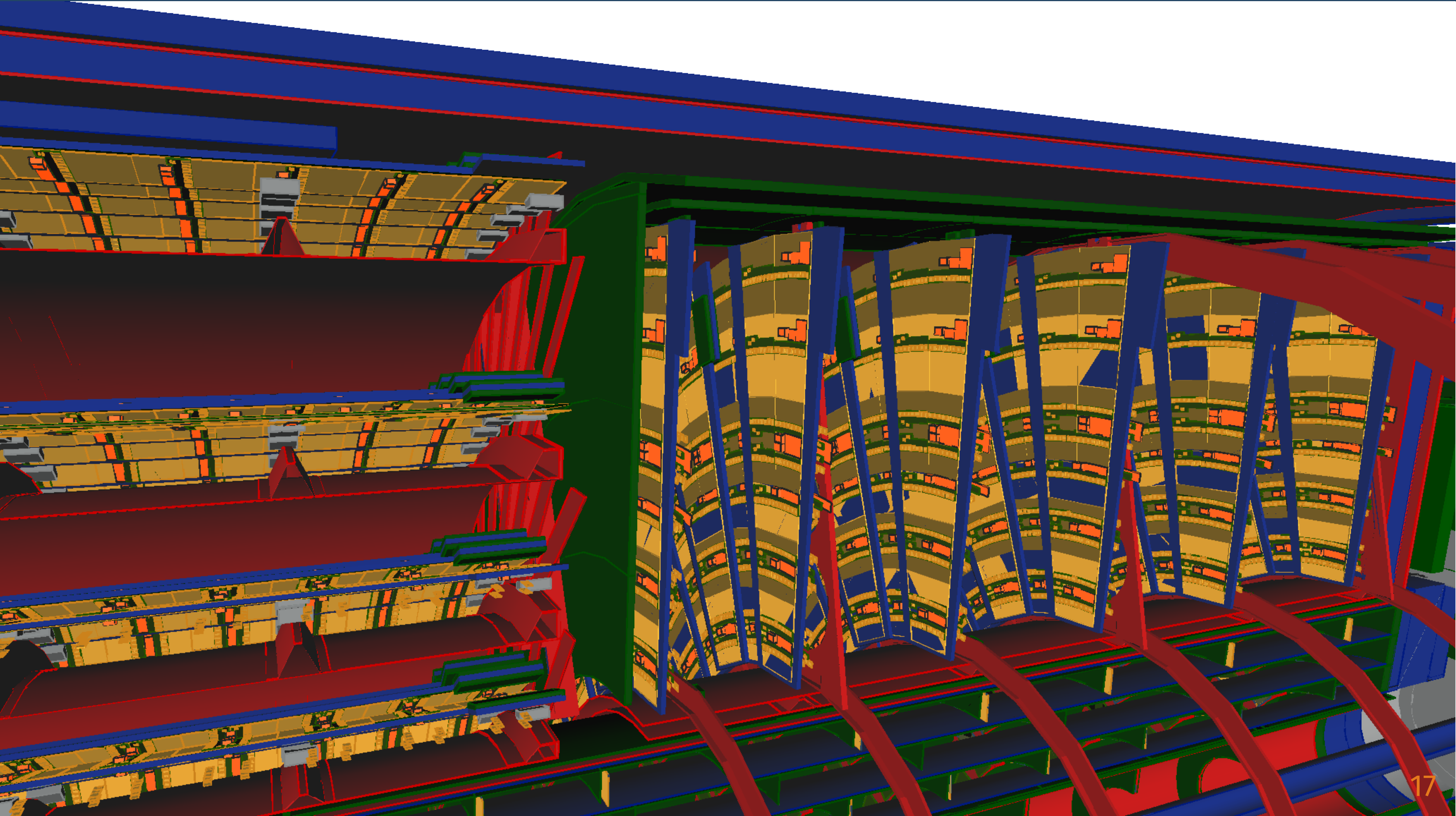
- 2D silicon detectors
- 5 barrel, 9 endcap layers
- 9164 modules
- up to 614400 readout channels per module

Strip detectors

- 1D silicon detectors
 - double-modules with 90° rotation to gain 2D detection
- 4 barrel, 6 endcap layers
- 49536 modules
- up to 1536 readout channels per module

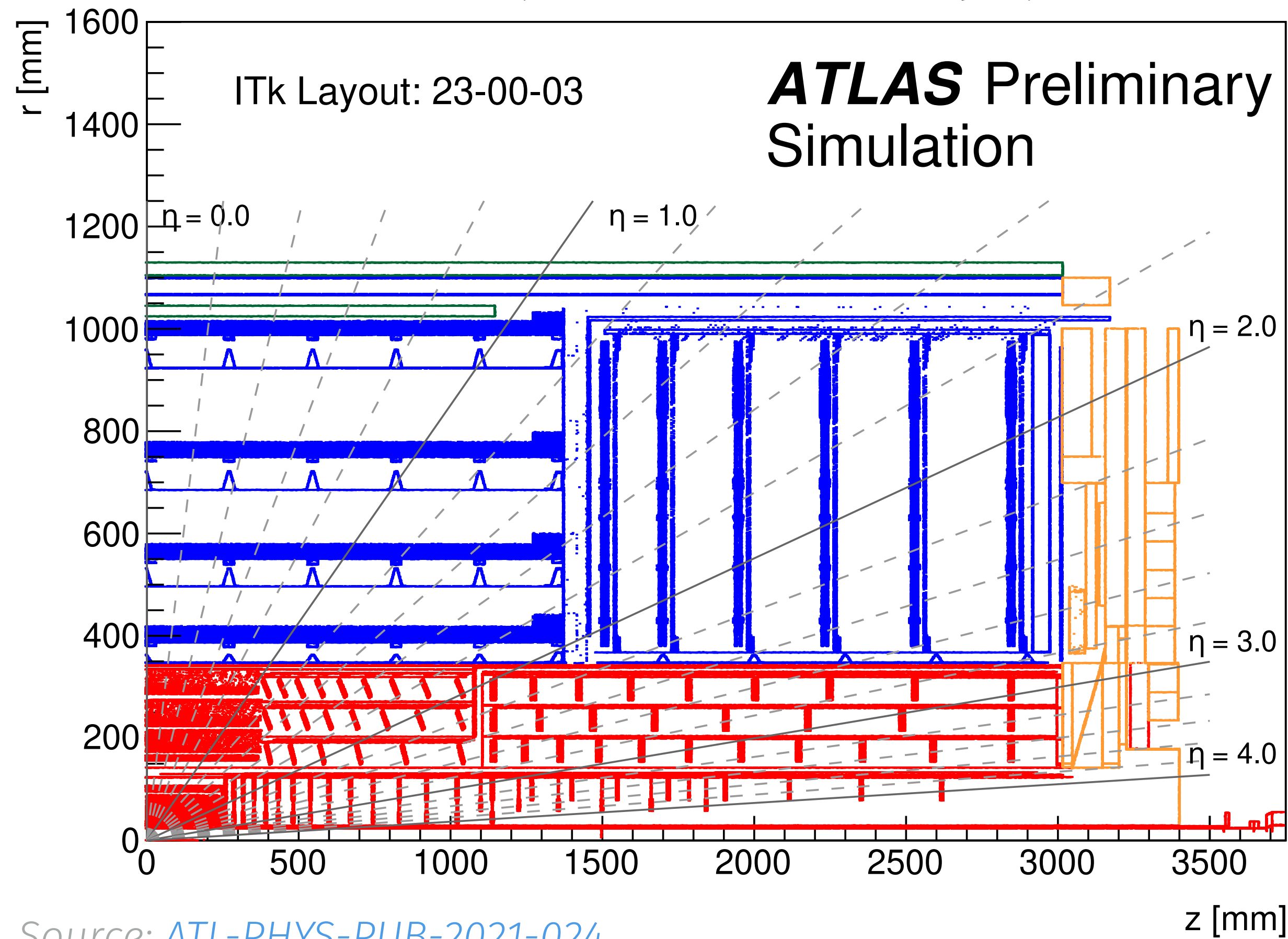
ITK PIXEL MODULES DETAIL







red: ITk Pixel System, blue: ITk Strip System

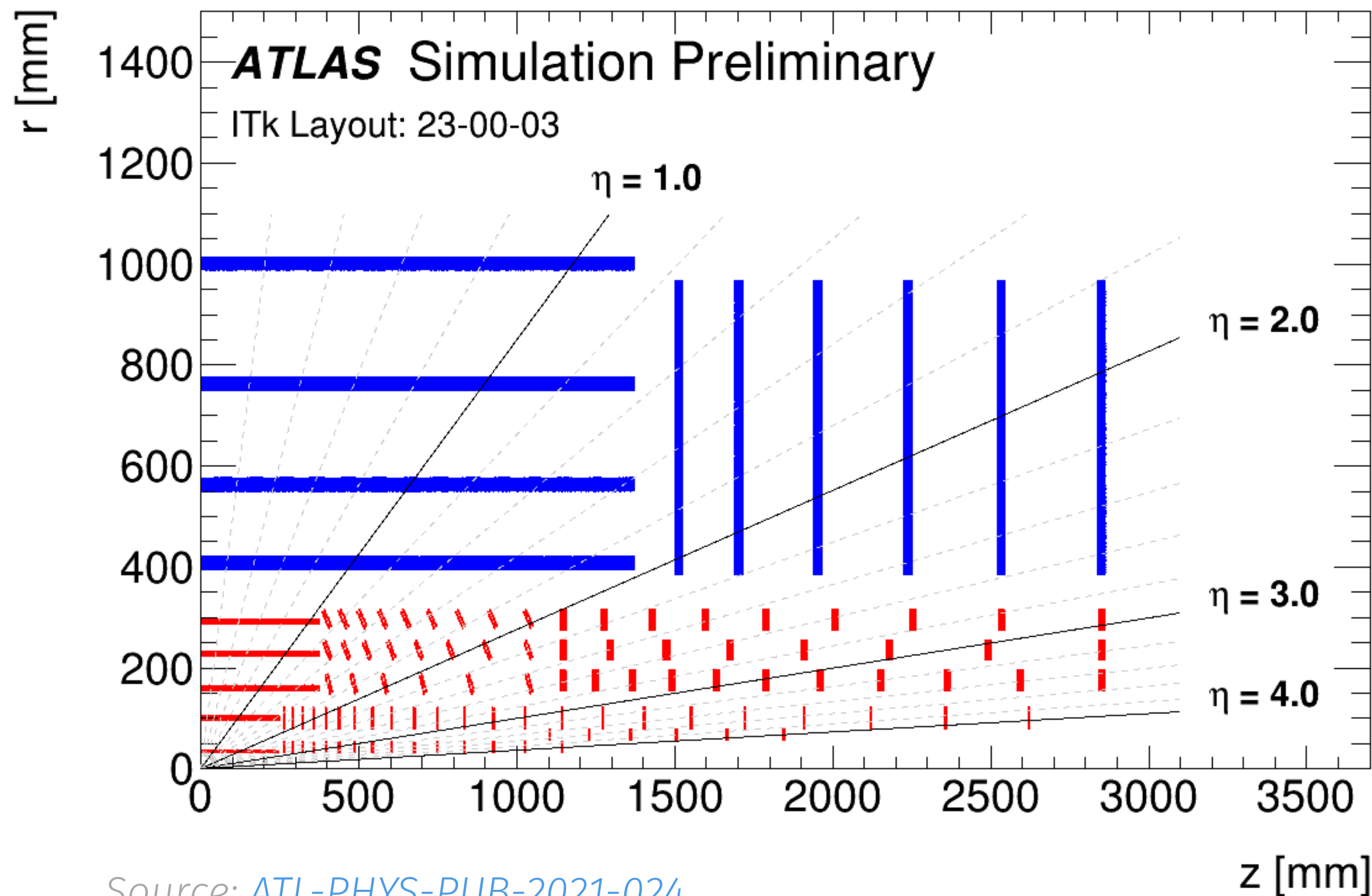


Source: [ATL-PHYS-PUB-2021-024](https://arxiv.org/abs/2102.024)

- Particles interact with all the material of the detector.



red: ITk Pixel System, blue: ITk Strip System

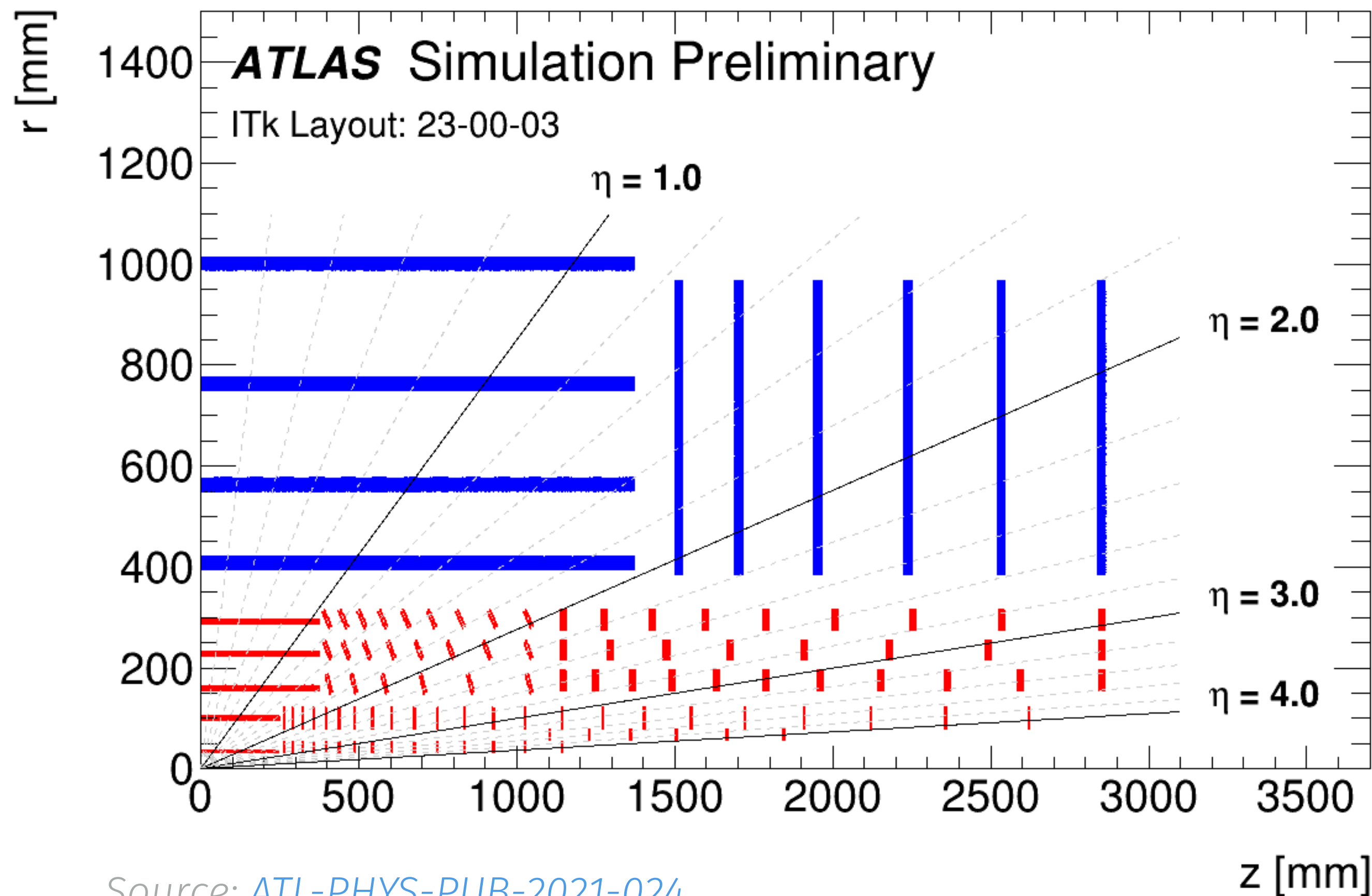


Source: [ATL-PHYS-PUB-2021-024](https://arxiv.org/abs/2102.02441)

- Particles interact with all the material of the detector.
- Only the **sensitive detector** part can actually detect charged particles.
 - Measuring the energy loss via induced electric charge.



red: ITk Pixel System, blue: ITk Strip System

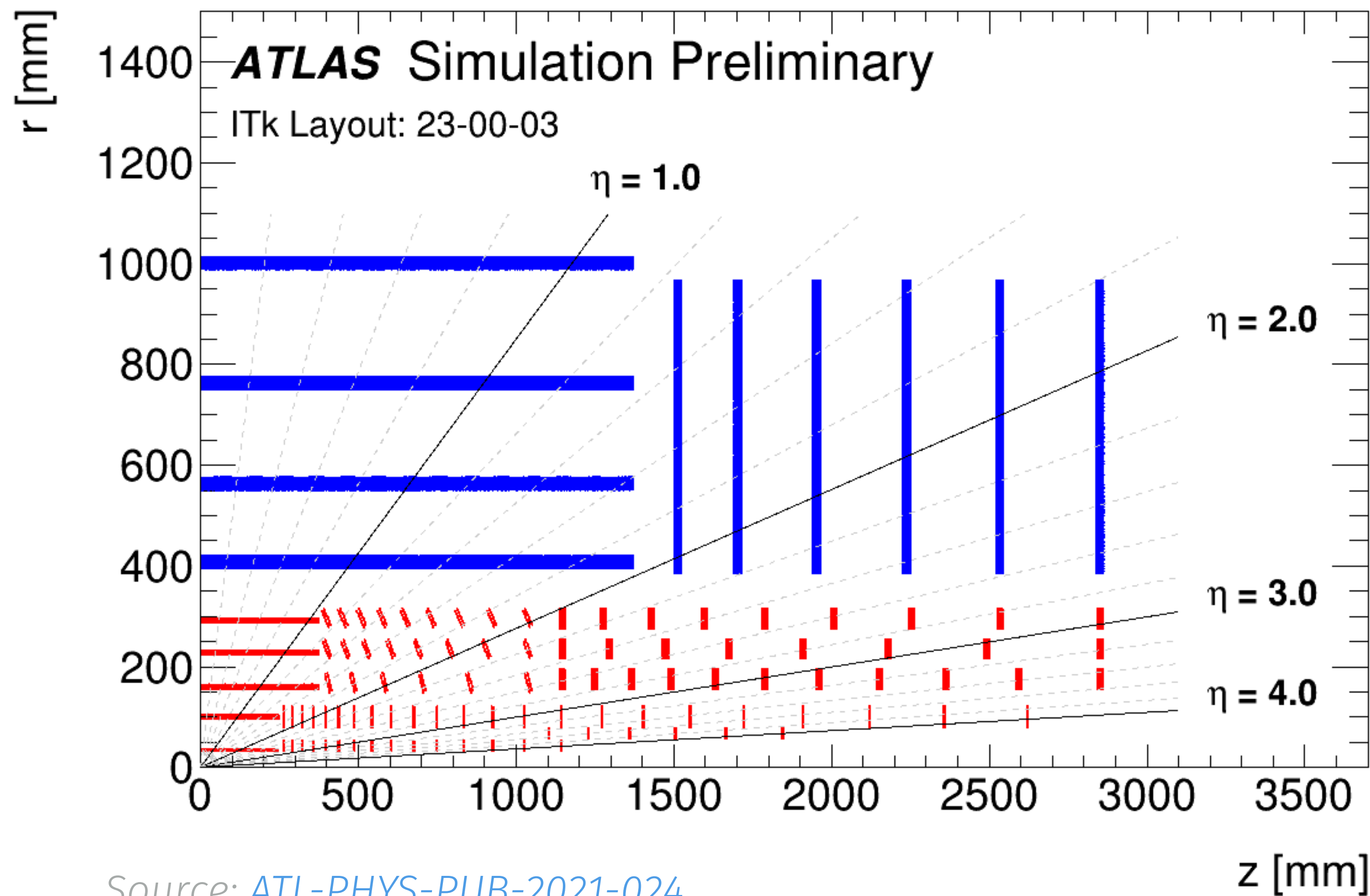


Source: [ATL-PHYS-PUB-2021-024](https://arxiv.org/abs/2102.024)

- Particles interact with all the material of the detector.
- Only the **sensitive detector** part can actually detect charged particles.
 - Measuring the energy loss via induced electric charge.
- Energy loss independent of the readout — simulation performed **per-module**.
 - Much lower number of elements to simulate.
- Each particle makes a “hit”:
 - coordinates
 - particle energy lost

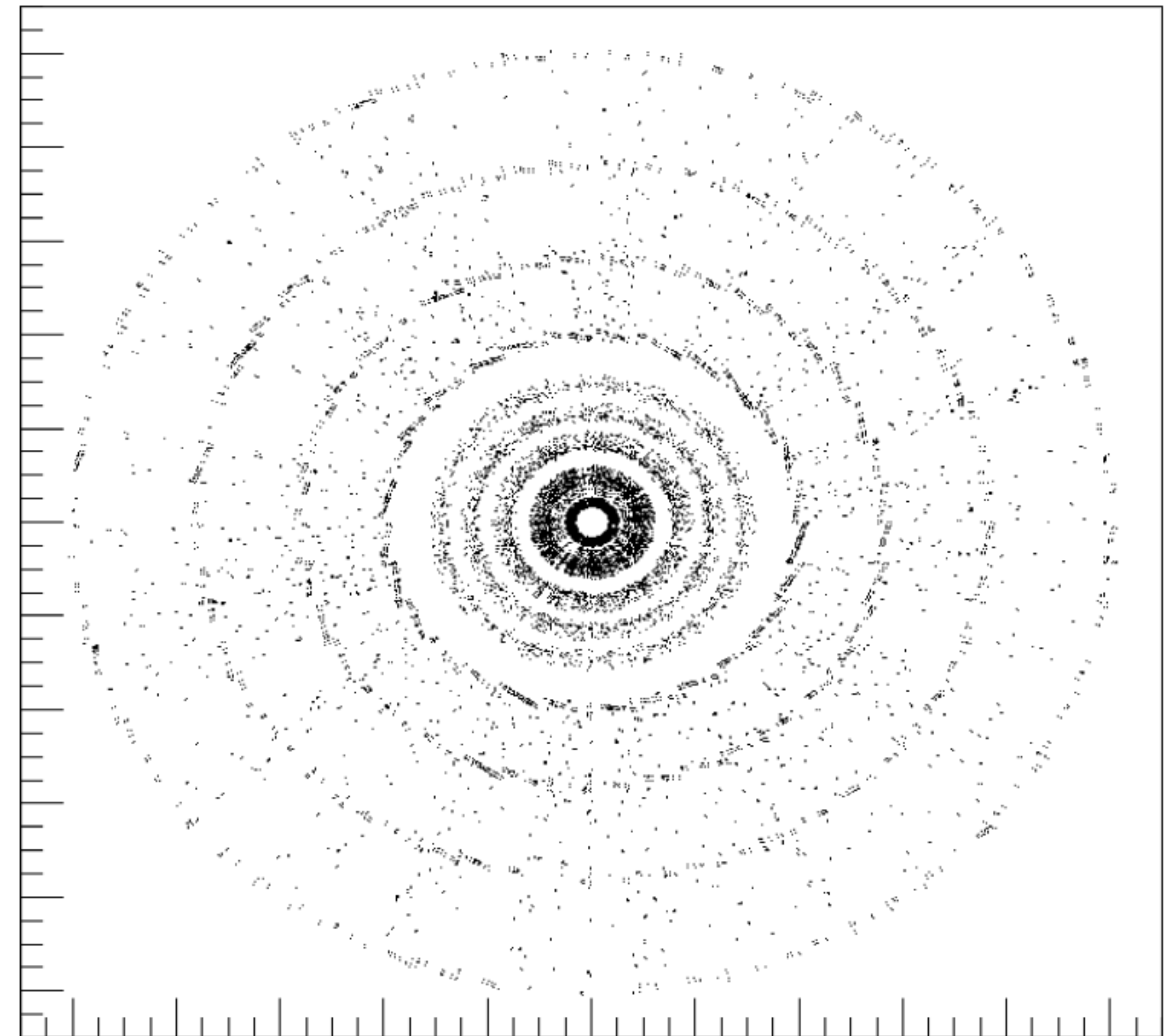


red: ITk Pixel System, blue: ITk Strip System



Source: [ATL-PHYS-PUB-2021-024](https://arxiv.org/abs/2102.02441)

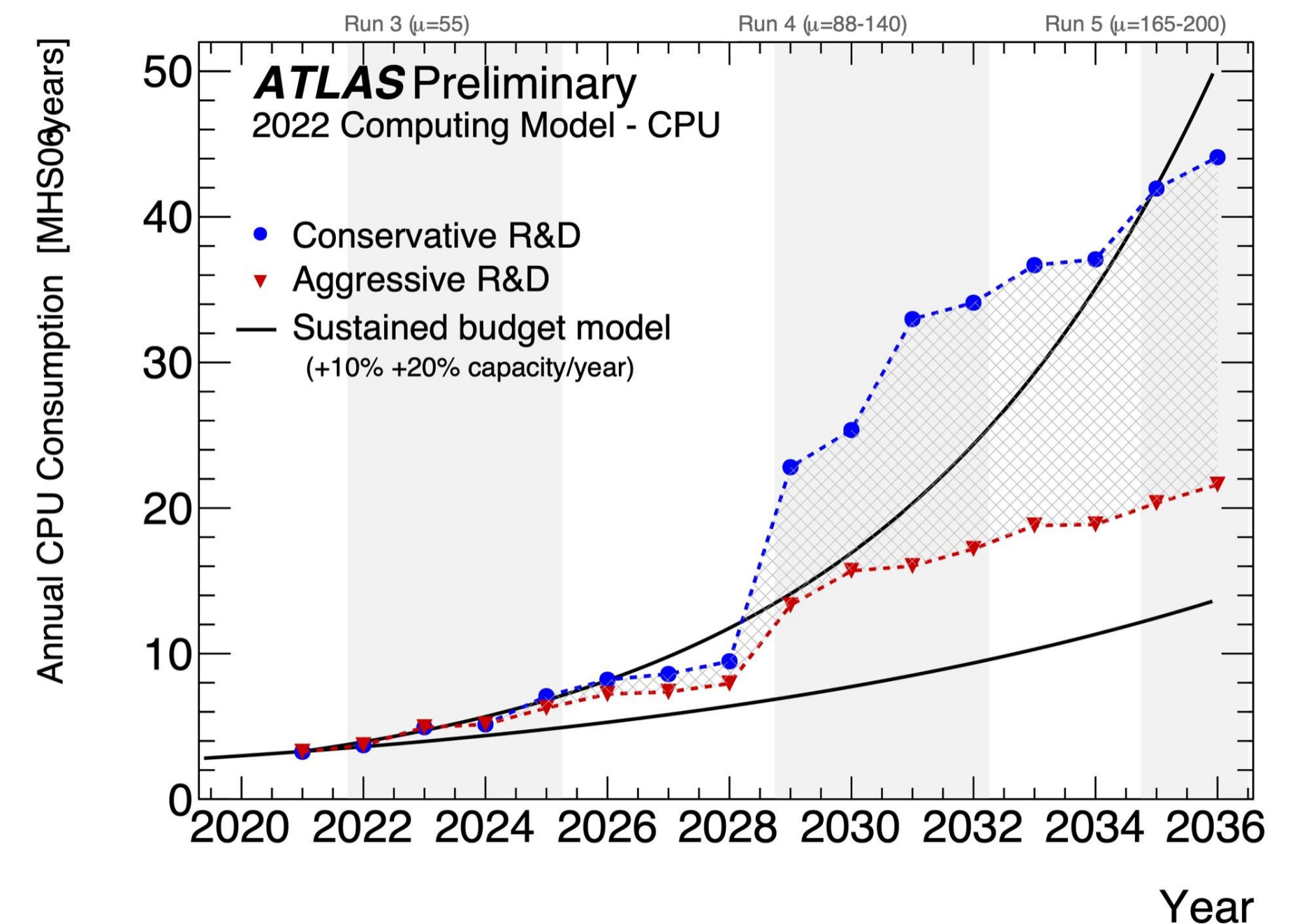
- Particles interact with all the material of the detector.





- A large part of the LHC physics programme relies on **accurate Monte Carlo simulation of collision events**.
 - generation of physics events and their immediate decays
 - simulation of the detector and physics interactions
 - digitisation of the detector response (readout)
 - (simulated events are reconstructed the same way as collected data)
- Producing simulated samples → majority of experiments' CPU requirements (CMS used 85% CPU for sim. during 2009-2016, half was spent detector simulation).

Source: ATLAS Software and Computing HL-LHC Roadmap

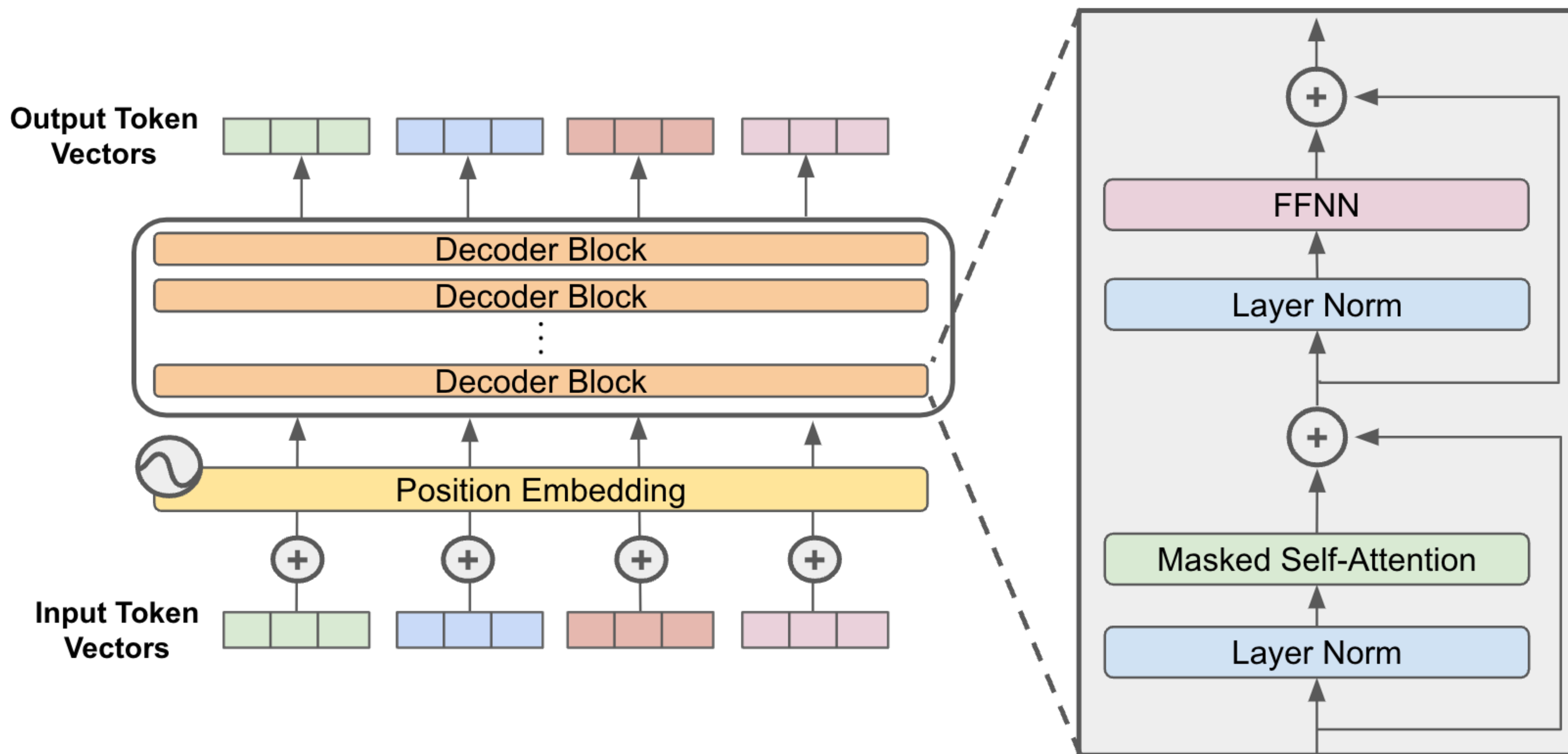


- Current methods do not scale with HL-LHC data rates and more aggressive R&D is needed.

DECODER-ONLY TRANSFORMER



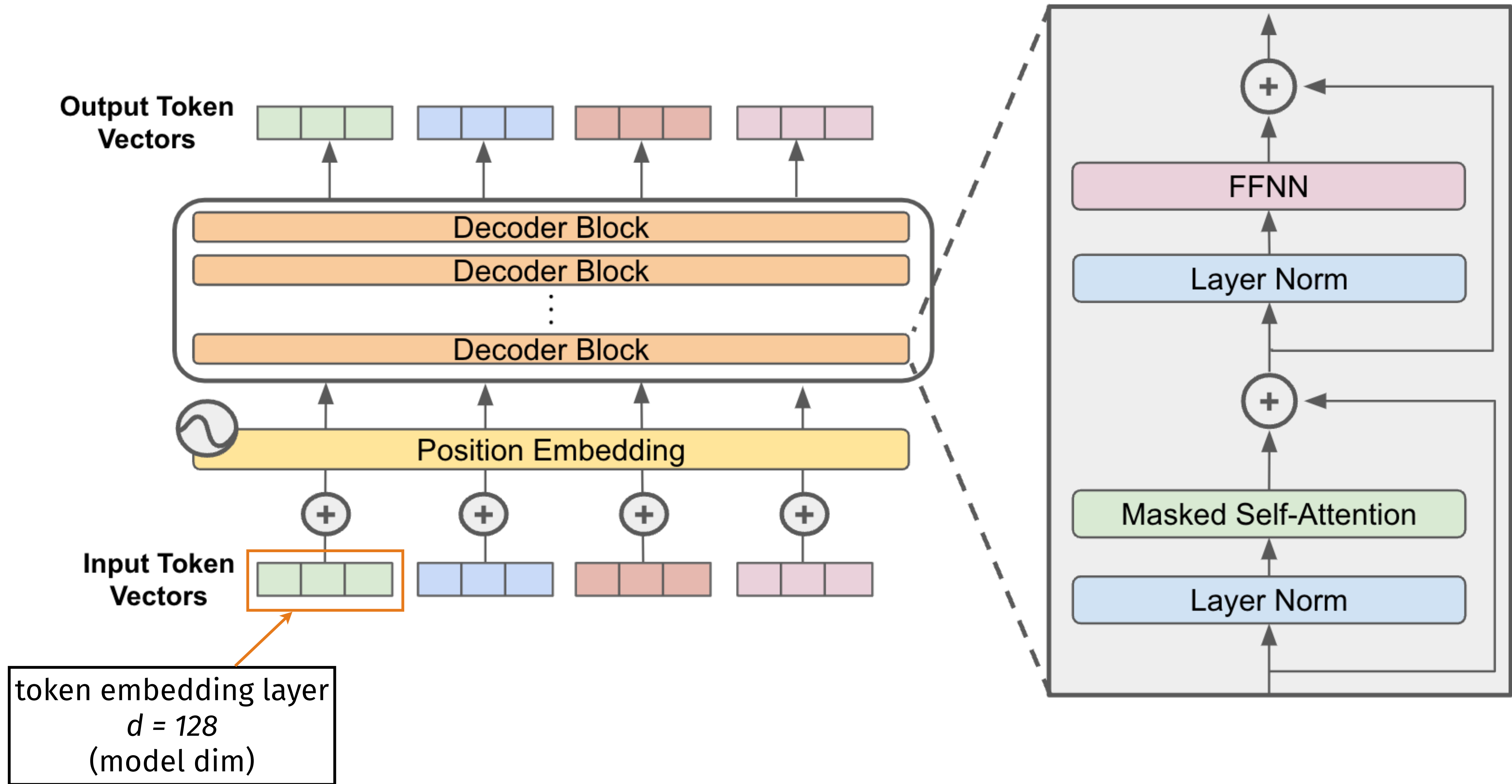
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



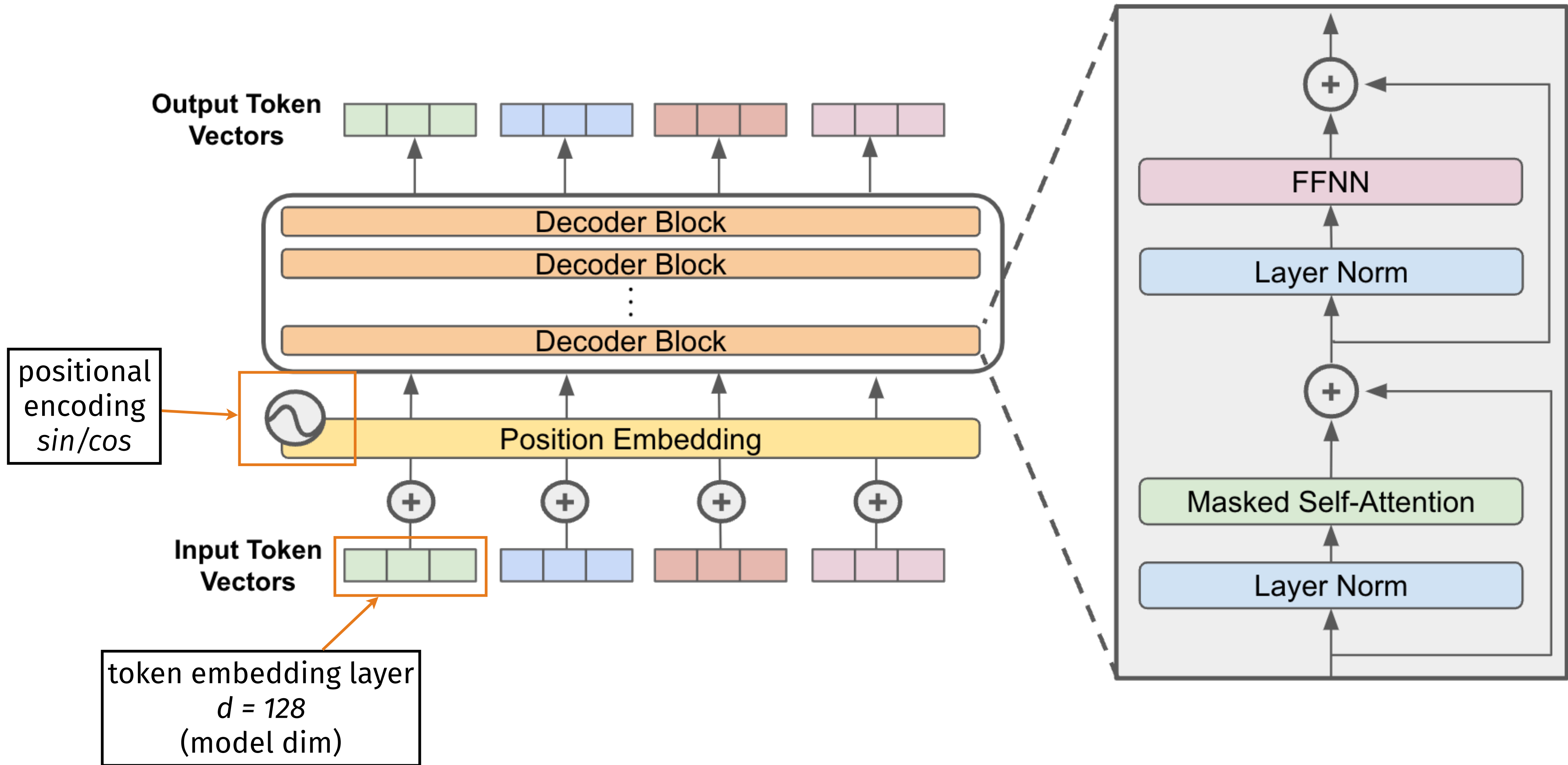
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



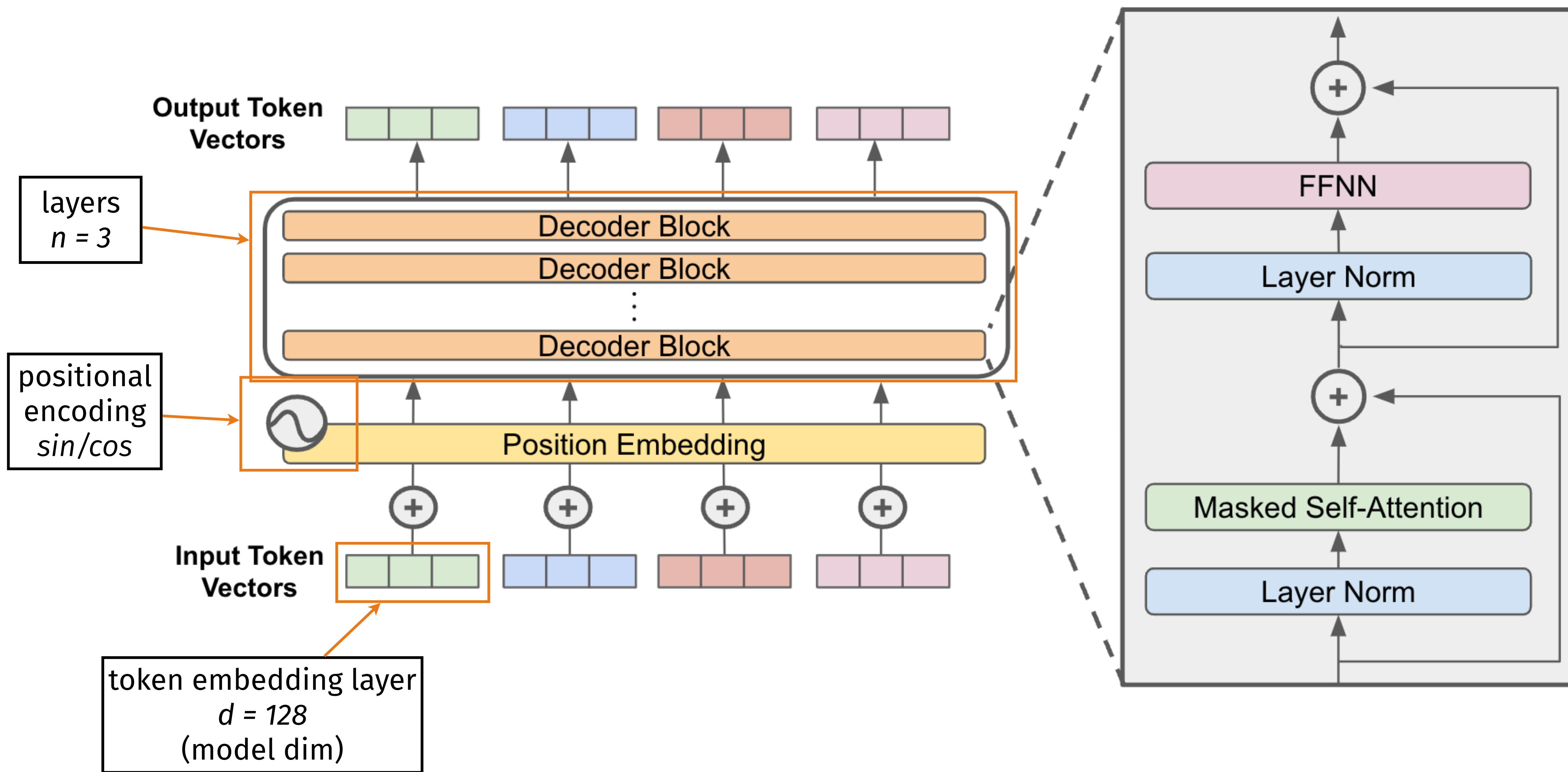
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



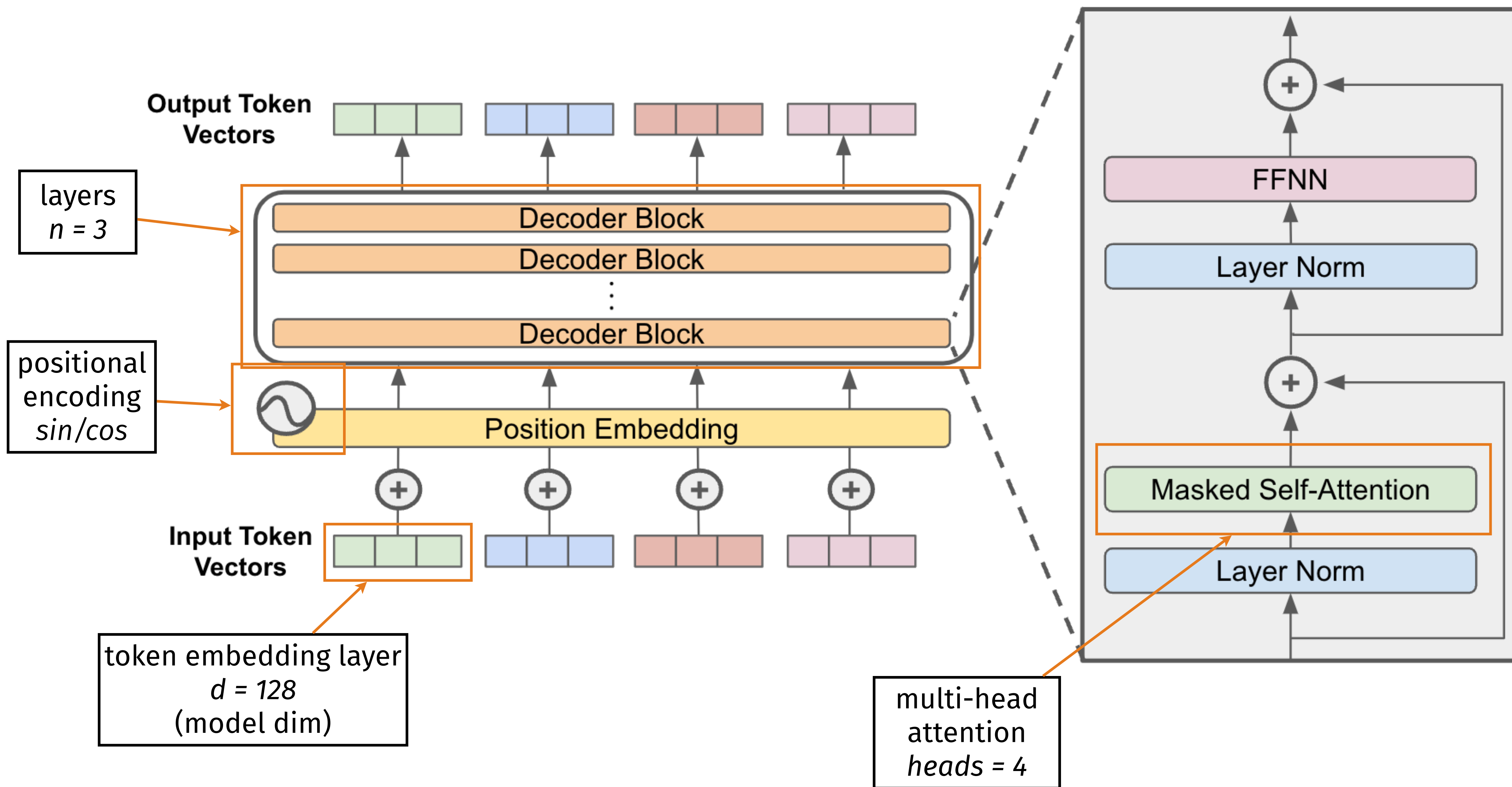
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



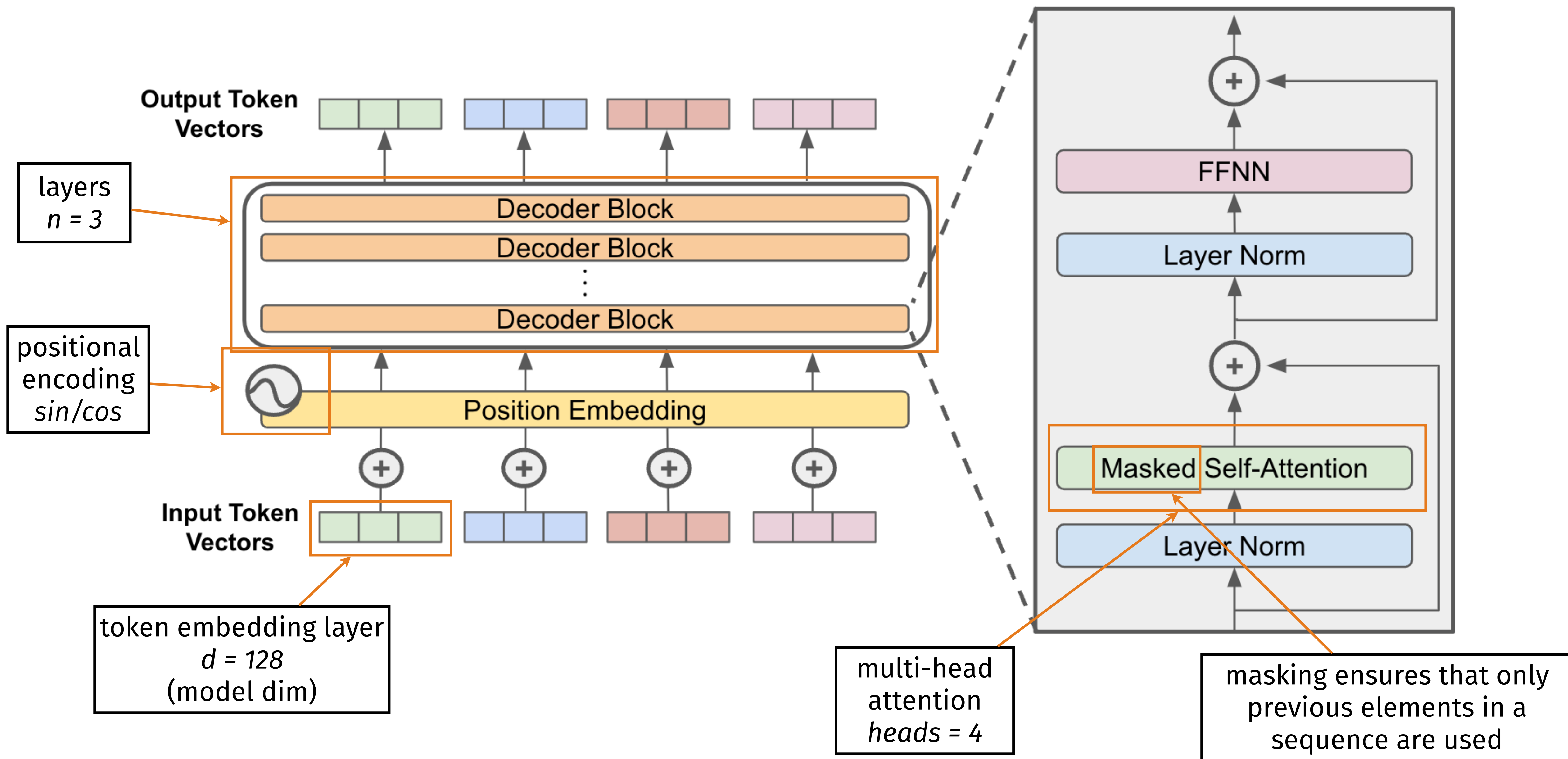
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



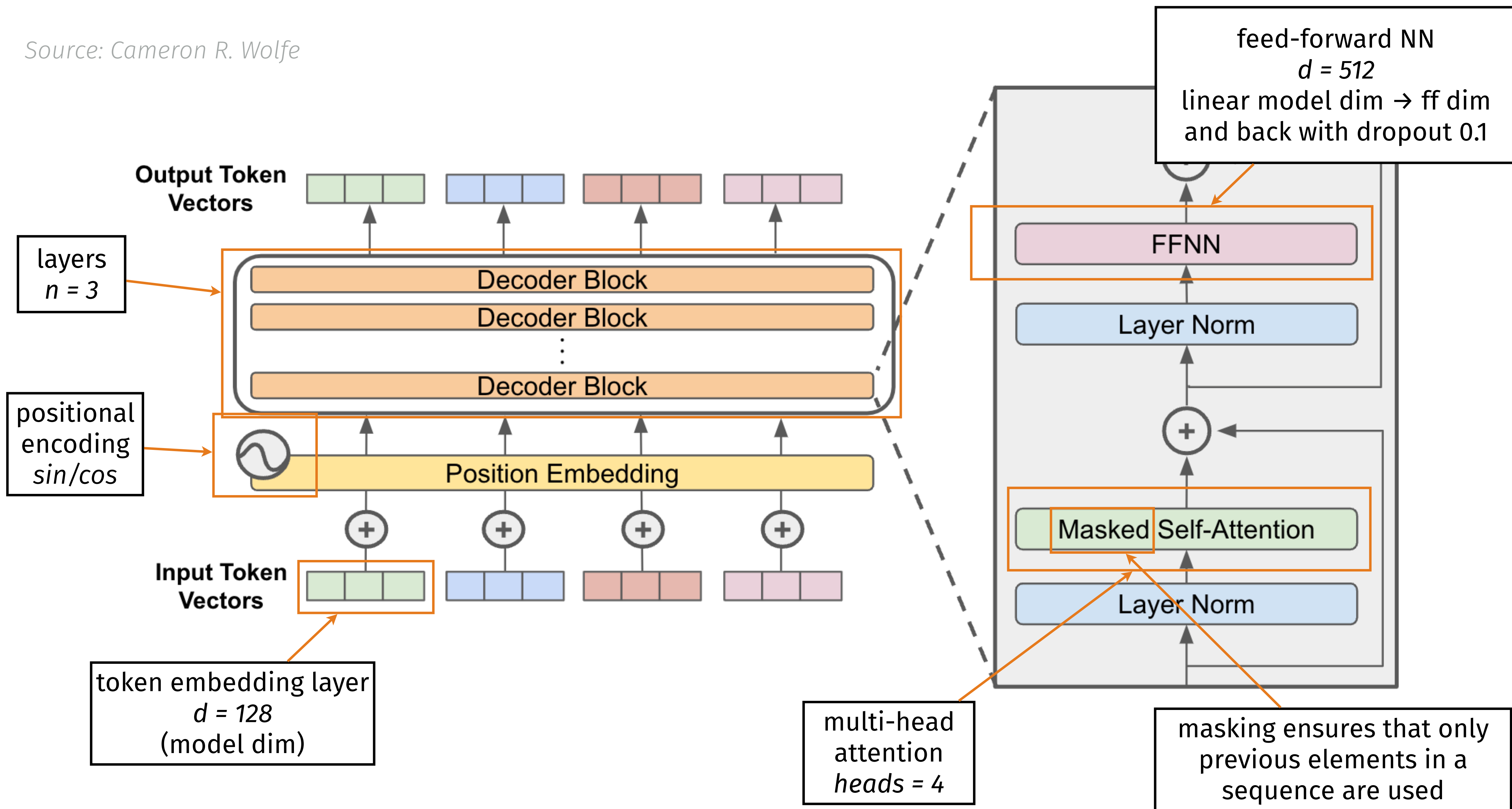
Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



Source: Cameron R. Wolfe



DECODER-ONLY TRANSFORMER



Source: Cameron R. Wolfe

