



Machine Learning @ CMS

Abhijith Gandrakota
for CMS MLG group

MLFP school, Georgia Tech

*Many slides from Javier Duarte

Machine Learning

Machine learning is the use of data to make predictions or decisions without explicit programming

Why? When did we use it in HEP ?


Machine Learning

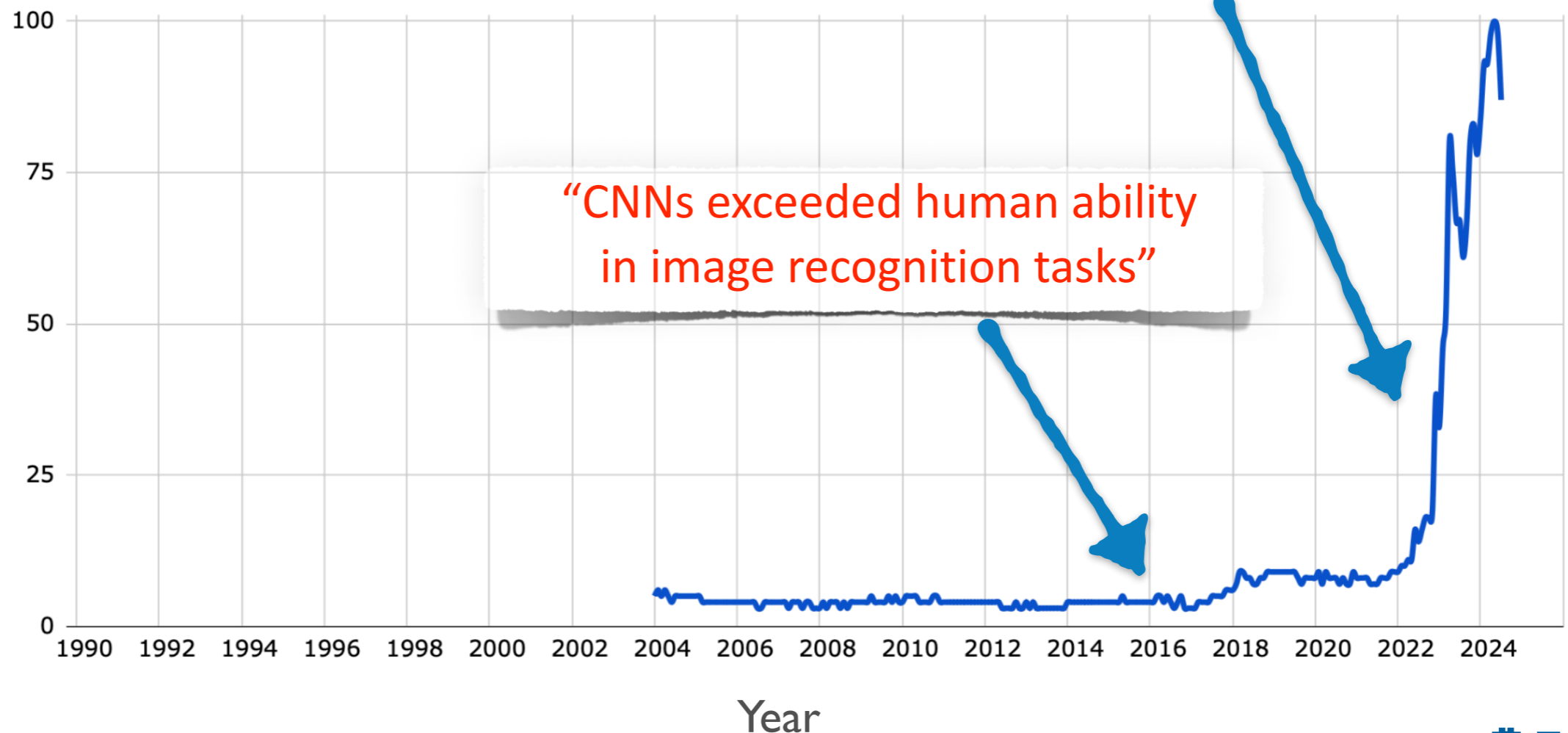
Machine learning is the use of data to make predictions or decisions without explicit programming

Why? When did we use it in HEP ?

Searches for “Artificial Intelligence” in google over the years

AI /ML is more main stream than ever now!

 OpenAI
 Gemini  Claude



Machine Learning

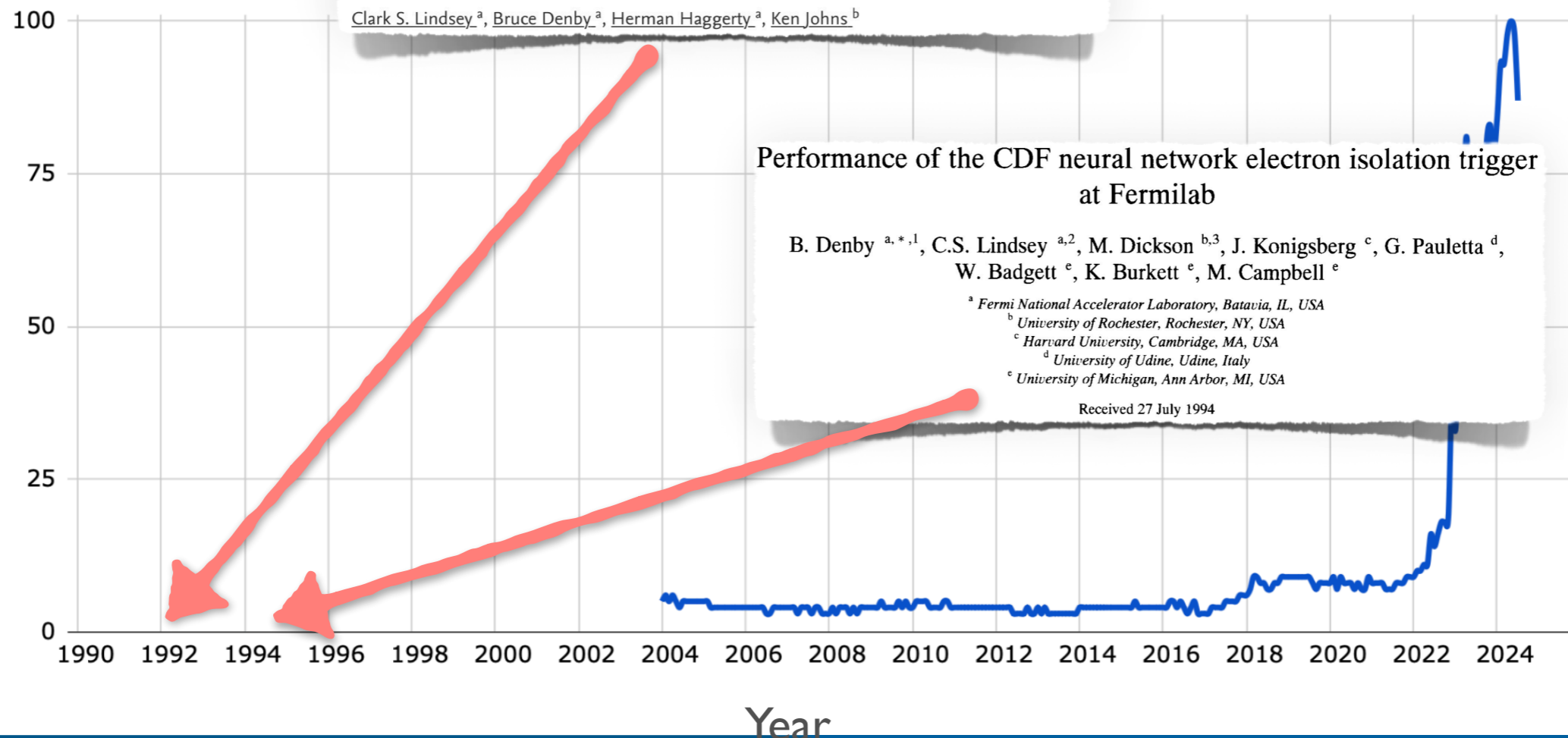
Machine learning is the use of data to make predictions or decisions without explicit programming

Why? When did we use it in HEP ?



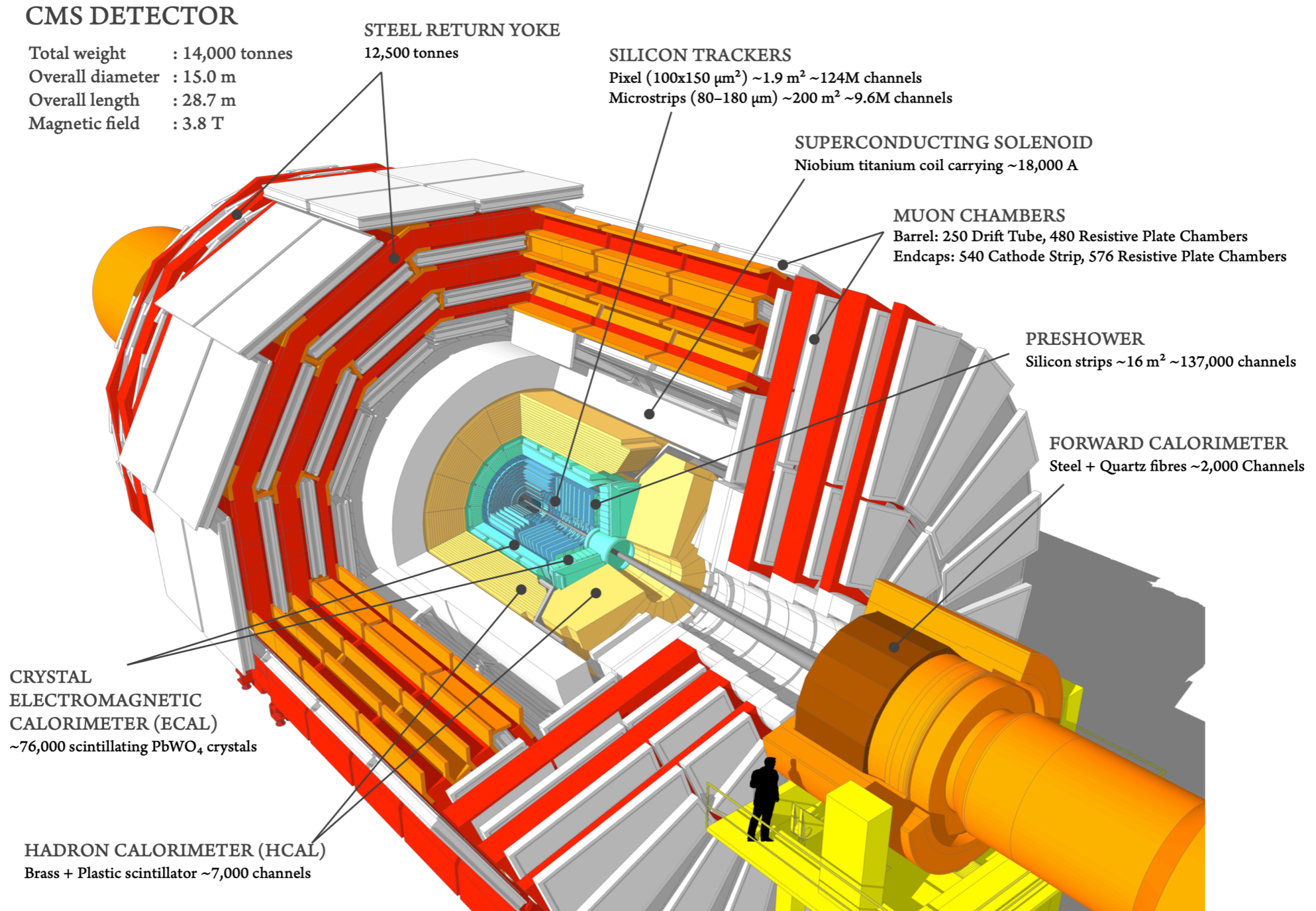
Real time track finding in a drift chamber with
a VLSI neural network

Clark S. Lindsey^a, Bruce Denby^a, Herman Haggerty^a, Ken Johns^b

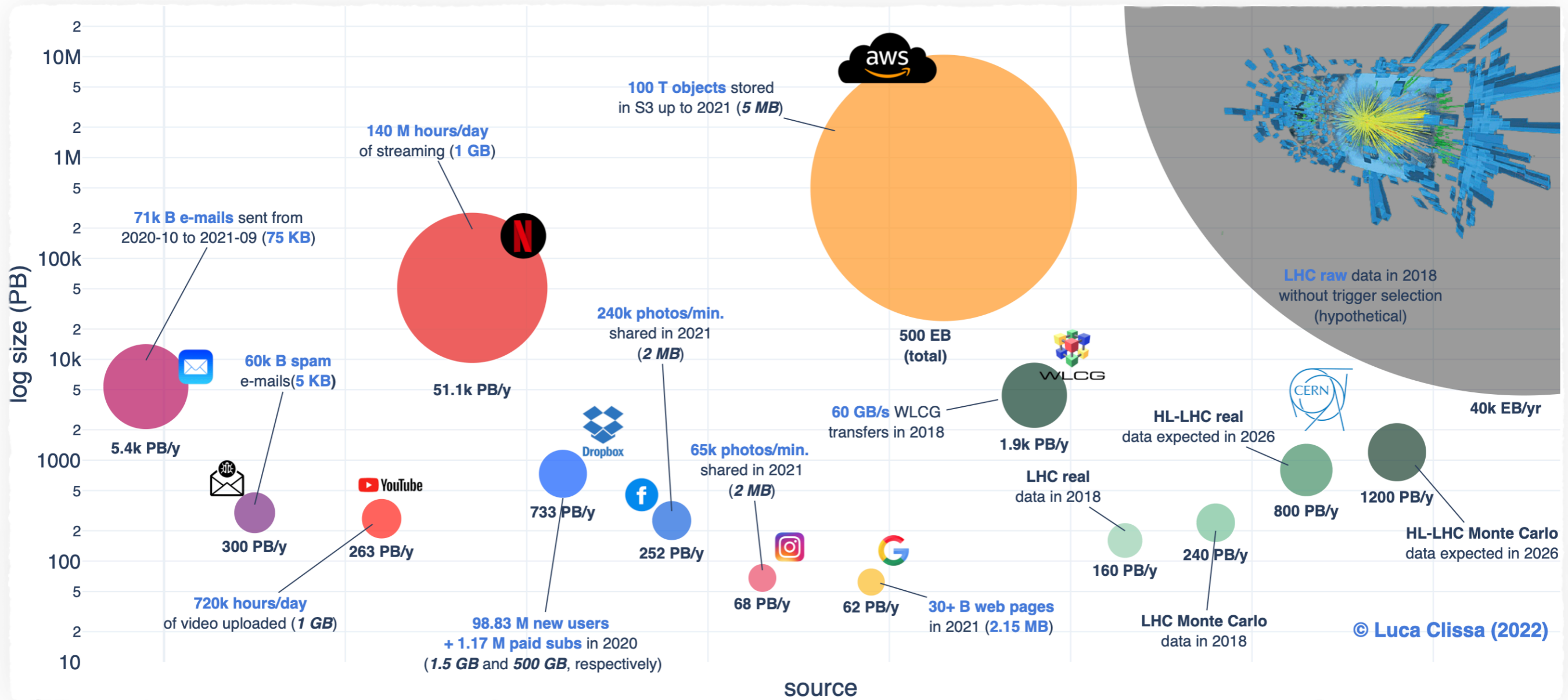


Compact Muon Solenoid

- Lots of readout channels from the detector
 - Data rates ~ 10x more than North American internet traffic!



Data from the LHC

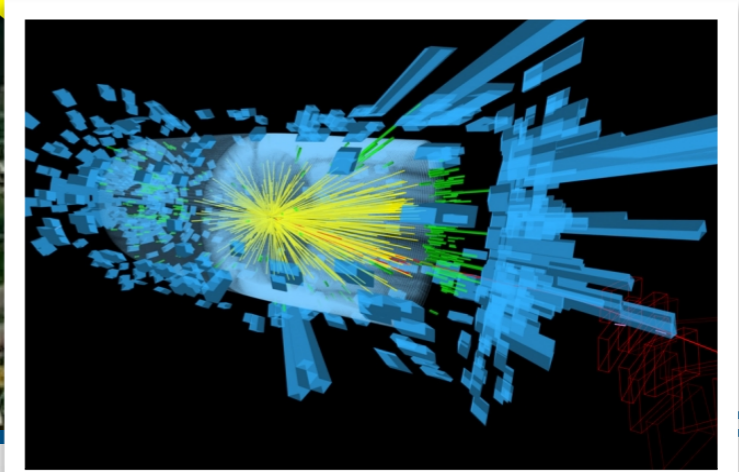
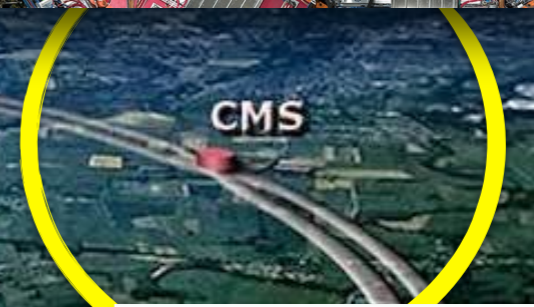
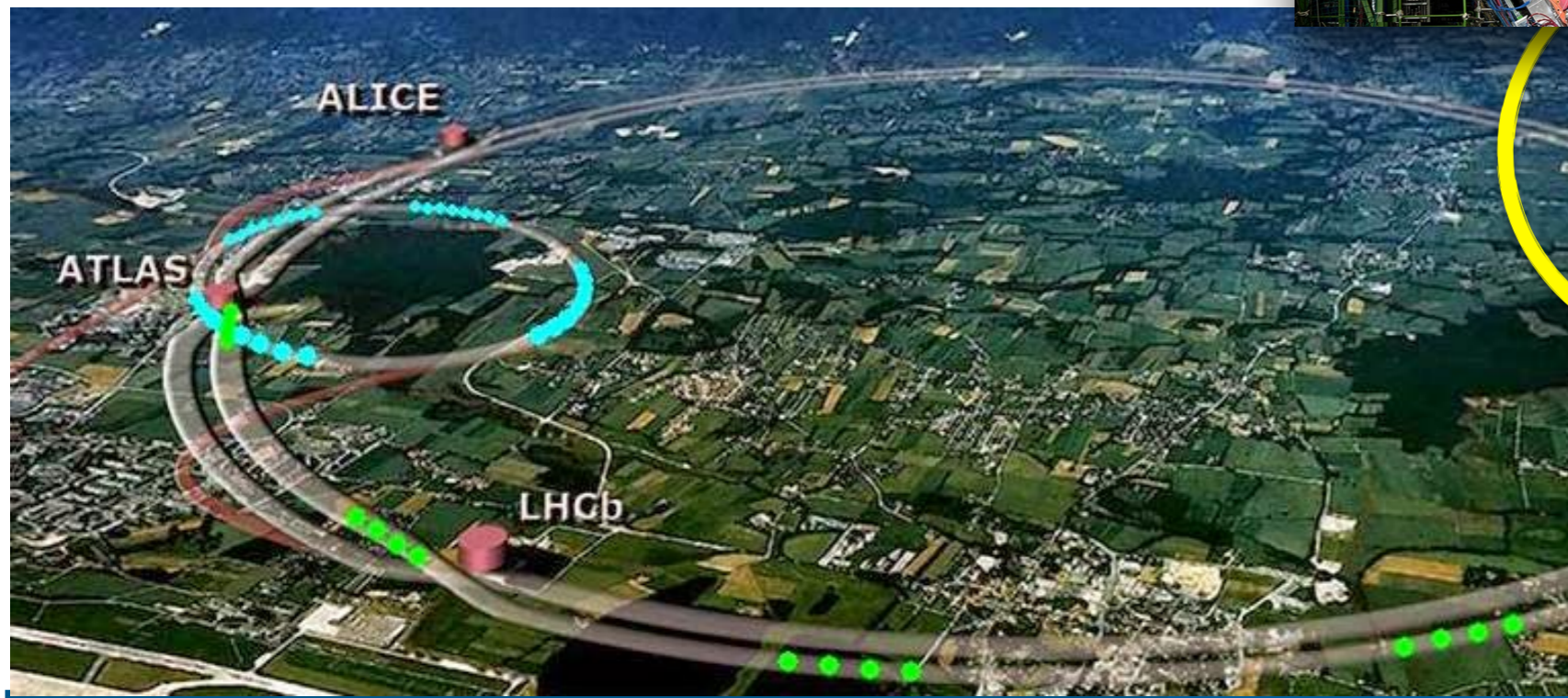
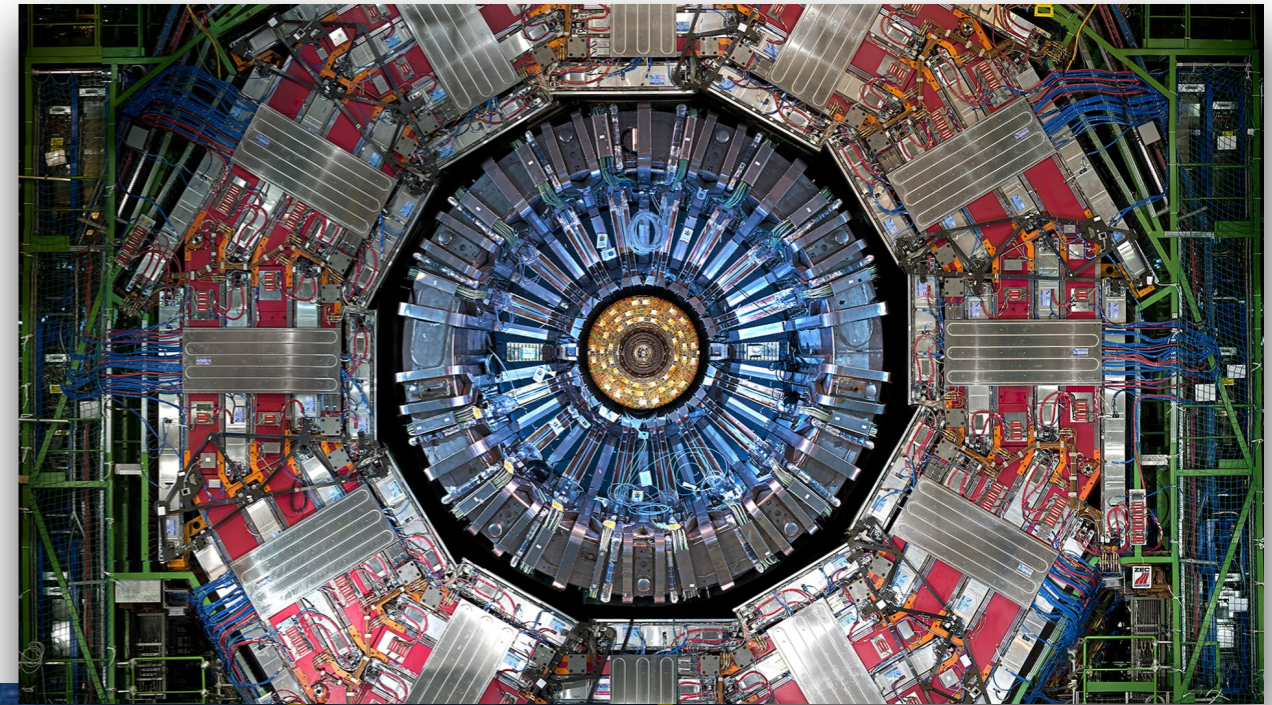


<https://arxiv.org/abs/2202.07659>

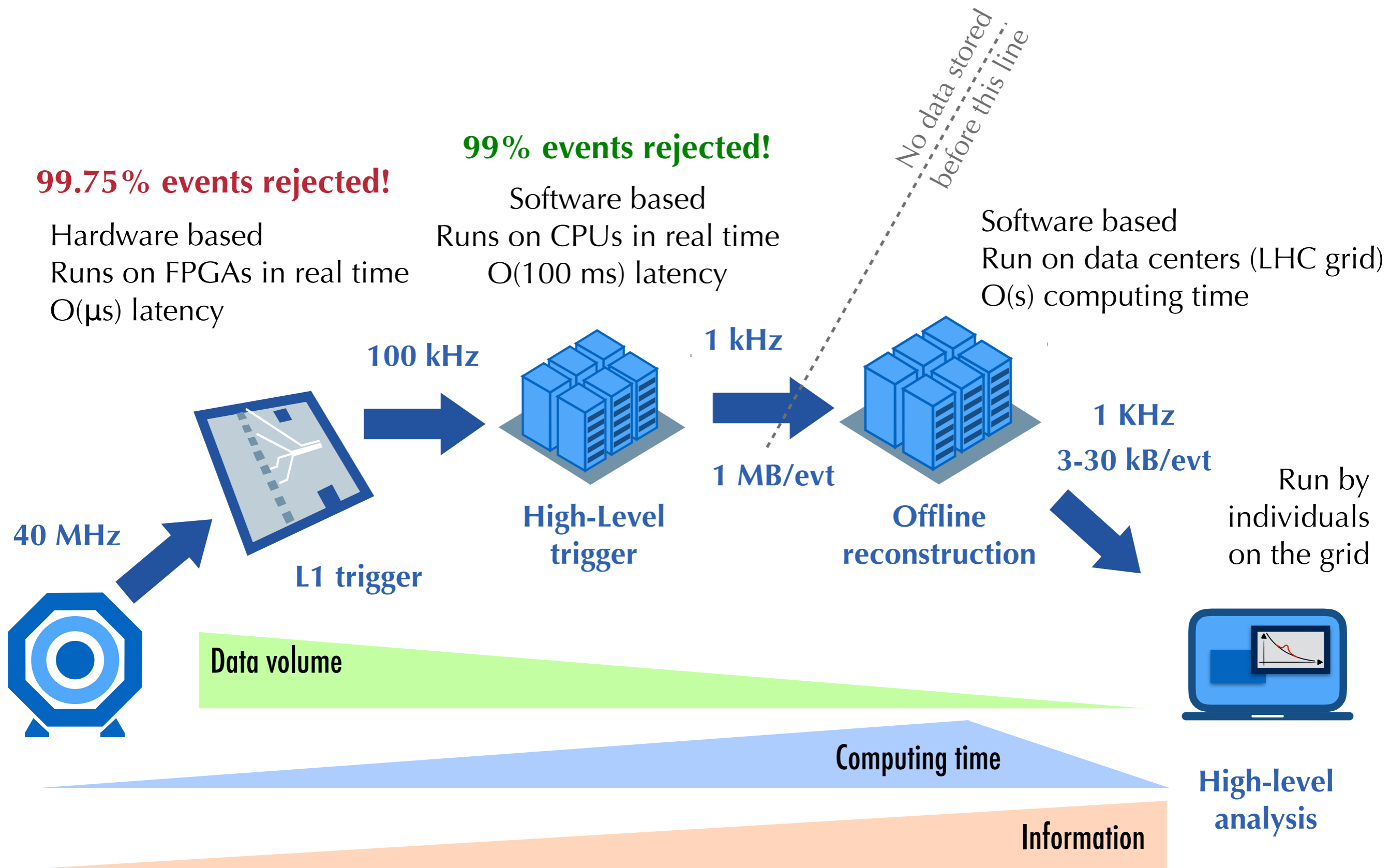
Big data @ LHC

40 MHz proton-proton collision frequency
 $O(10^3)$ particles per collision
 $O(10^8)$ sensors
→ **$O(100)$ TB/s data rates!**

The Compact Muon Solenoid



Data reduction workflow @ LHC



Make physics discoveries with 0,0025% of the events! (the rest is lost...)



No data before the

Software based
Run on data centers (LHC grid)
O(s) computing time



1 KHz
3-30 kB/evt

Run by
individuals
grid

40



Data volume



The High-Luminosity LHC challenge

instantaneous luminosity

x 0.75

x 1-2

NOW

x 2.5

2029

x 5-7

LHC Run 1
 $\sqrt{s} = 7-8 \text{ TeV}$
30/fb

Long
Shutdown 1

LHC Run 2
 $\sqrt{s} = 13 \text{ TeV}$
150/fb

Long
Shutdown 2

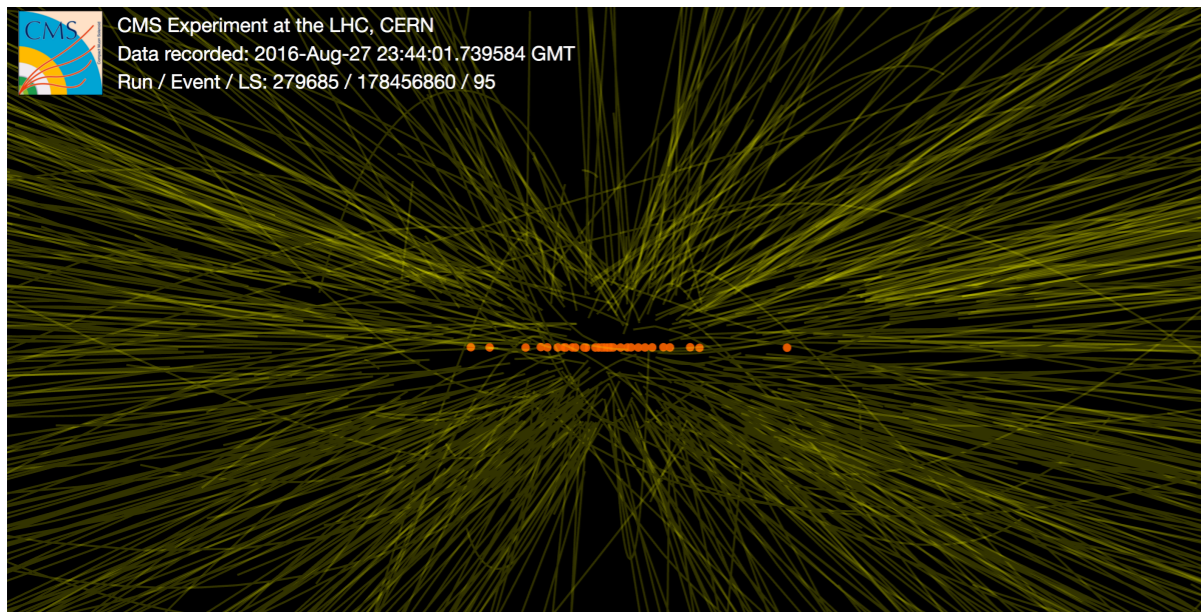
LHC Run 3
 $\sqrt{s} = 14 \text{ TeV}$
300/fb

Long
Shutdown 3

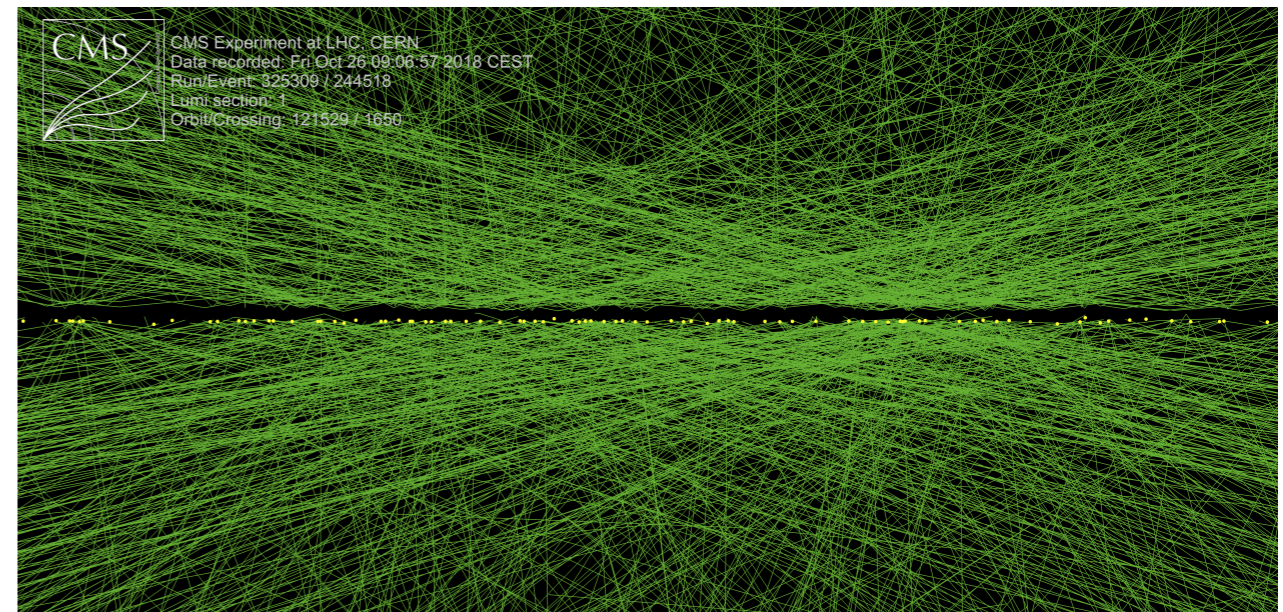
Run 4: HL-LHC
 $\sqrt{s} = 14 \text{ TeV}$
3000/fb

LHC TODAY

HL-LHC



40 simultaneous collisions
per bunch crossing



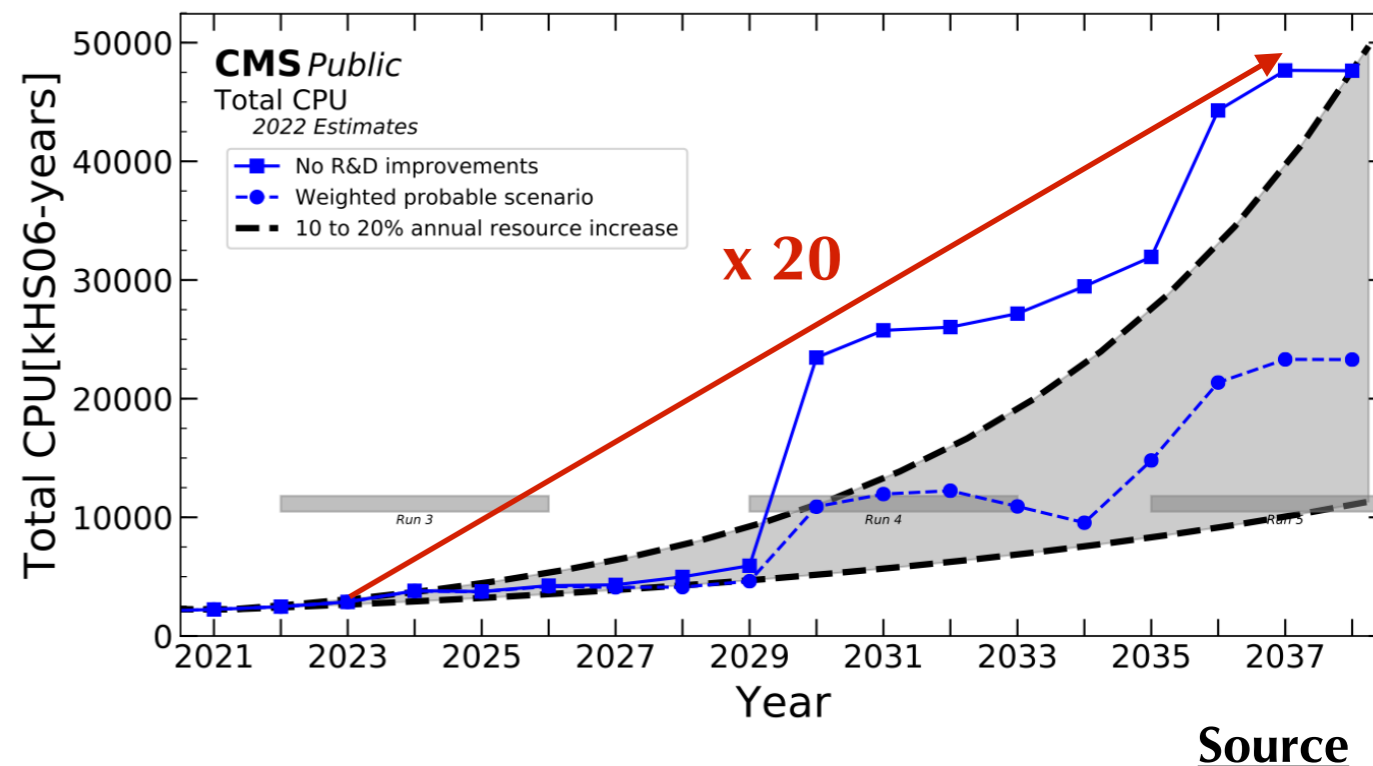
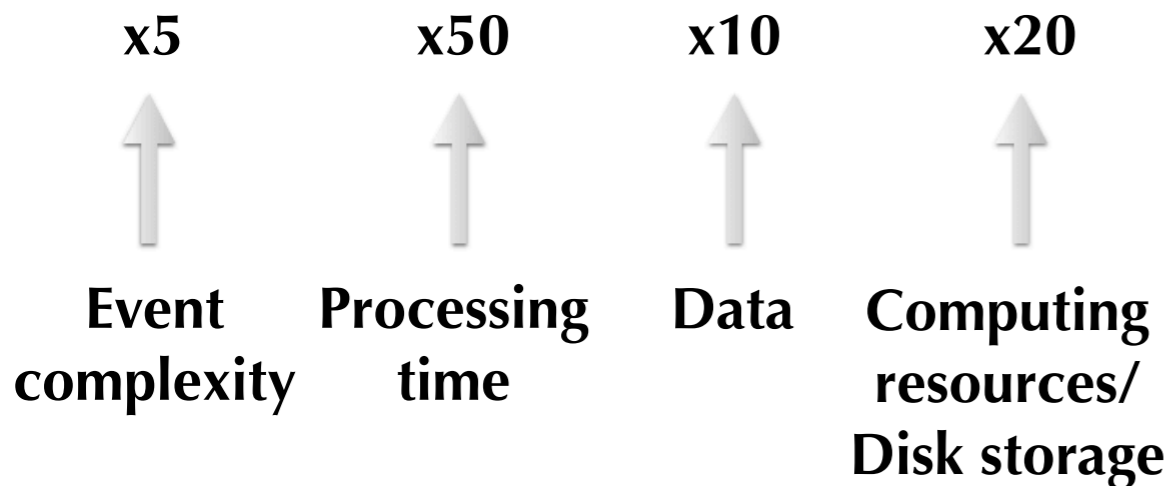
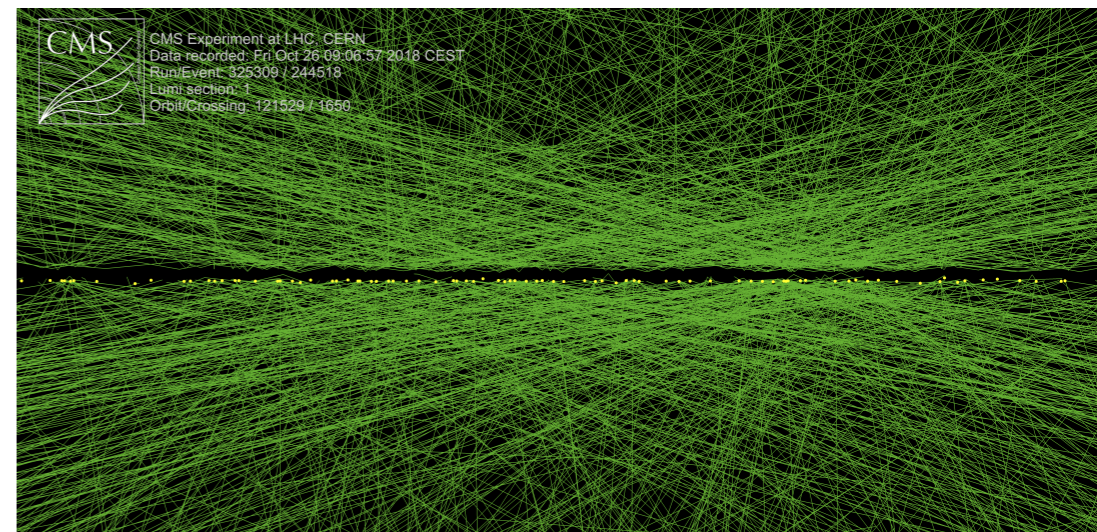
200 simultaneous collisions
per bunch crossing
+
more granular detector!

The HL-LHC challenge

With more particles per collision and more readout channels to combine, the reconstruction to become even more computing intensive

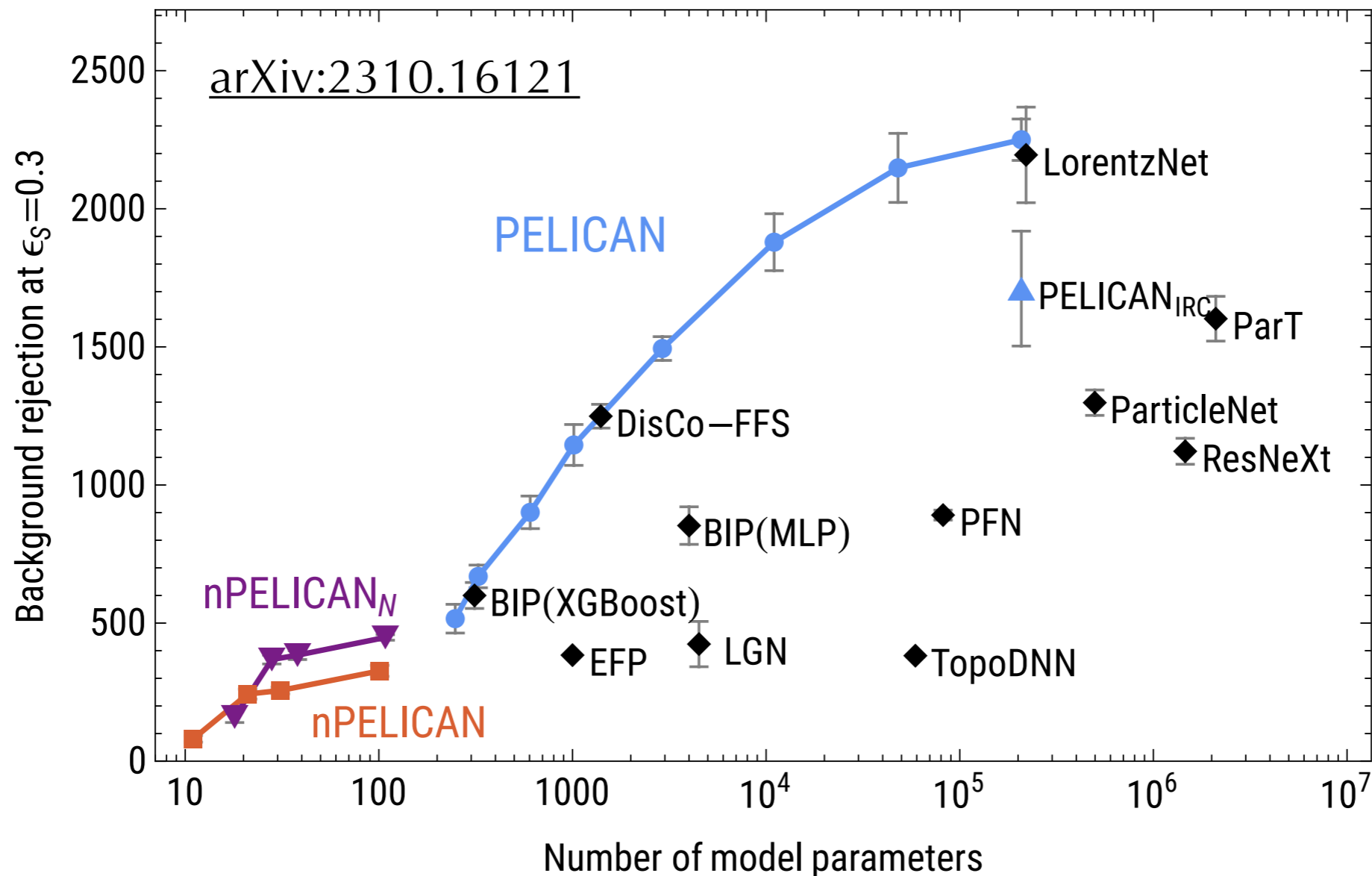
We cannot throw away more data
We must increase the throughput
with at most flat budget for computing resources

**We must do more with less
to preserve the physics!**



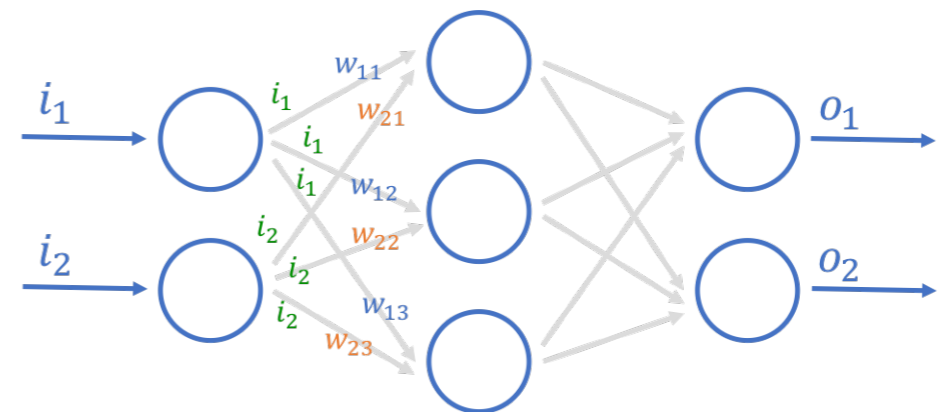
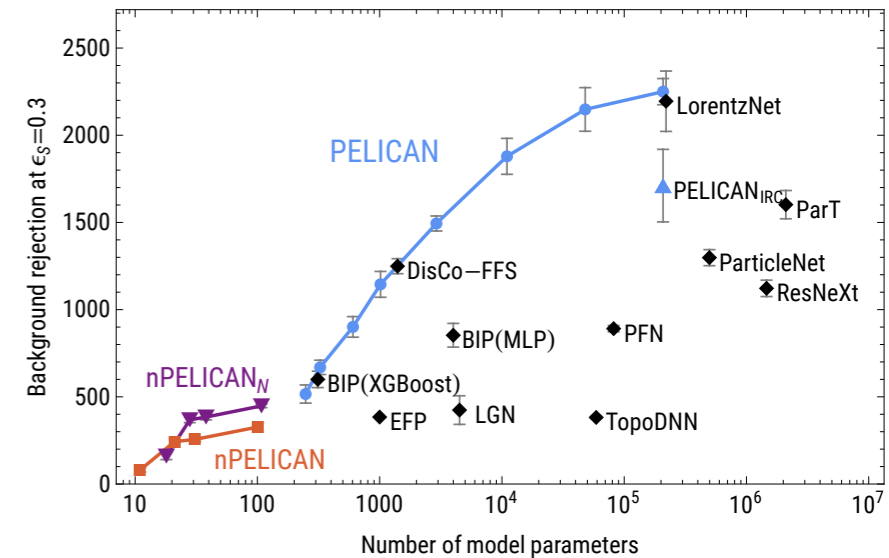
ML Advantages

- High accuracy
 - improved physics performance:
expressive models can extract most useful info from complex datasets



ML Advantages

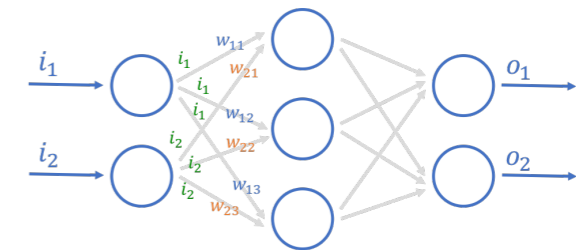
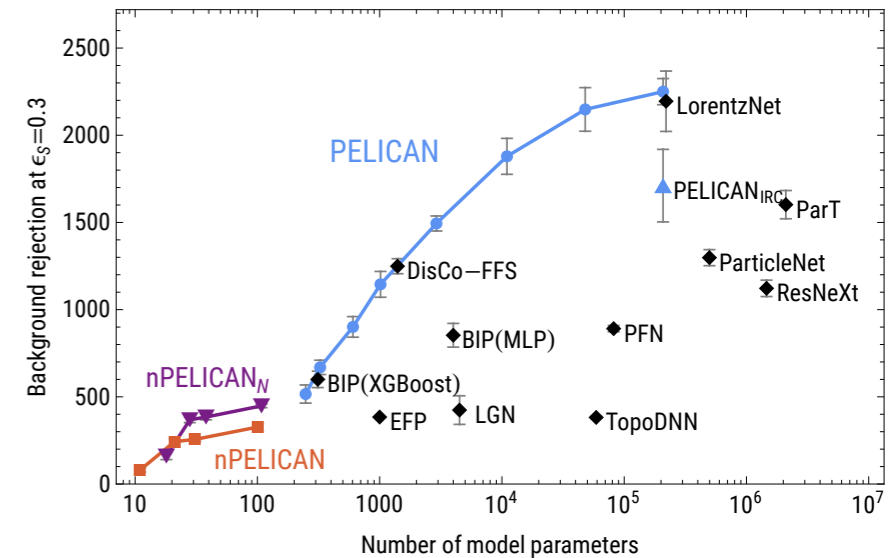
- **High accuracy**
 - improved physics performance: capacity to extract most useful info from complex datasets
- **Fast speed**
 - matrix multiplication can be massively parallelized
 - Leveraging GPUs/FPGAs, and other modern AI chips
 - take advantage of reduced precision
 - reduced development time: adaptive, generalizable



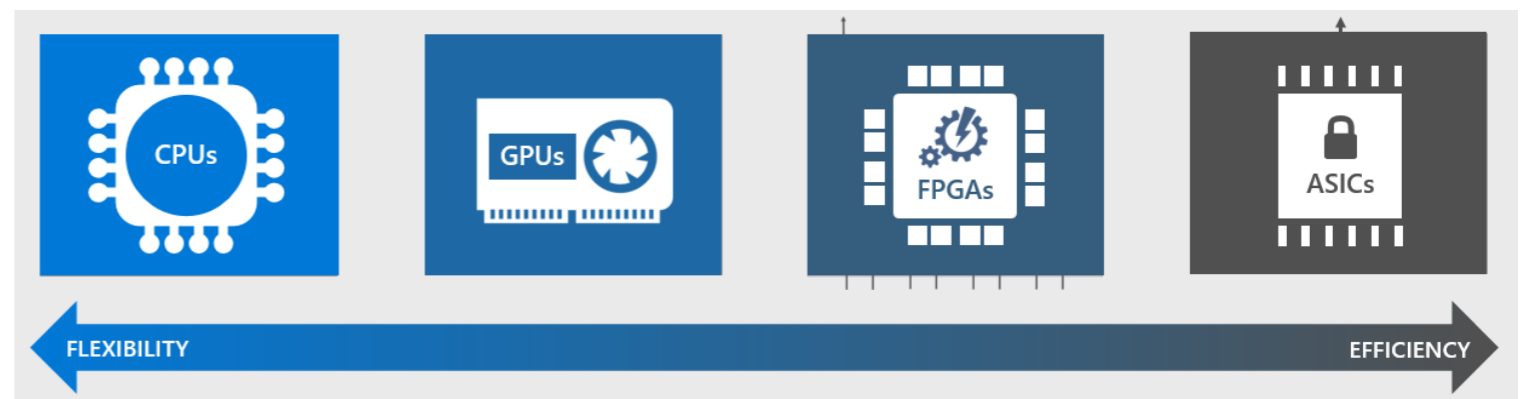
$$\begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \\ w_{13} & w_{23} \end{bmatrix} \cdot \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} = \begin{bmatrix} (w_{11} \times i_1) + (w_{21} \times i_2) \\ (w_{12} \times i_1) + (w_{22} \times i_2) \\ (w_{13} \times i_1) + (w_{23} \times i_2) \end{bmatrix}$$

ML Advantages

- **High accuracy**
 - improved physics performance: capacity to extract most useful info from complex datasets
 - better scaling with data complexity (HL-LHC)
- **Fast speed**
 - reduced development time
 - matrix multiplication can be massively parallelized
 - take advantage of reduced precision
- **Better portability**
 - many dedicated processors: GPUs, FPGAs, TPU, IPU, ...
 - large investment in tools to compile and optimize ML models for the hardware



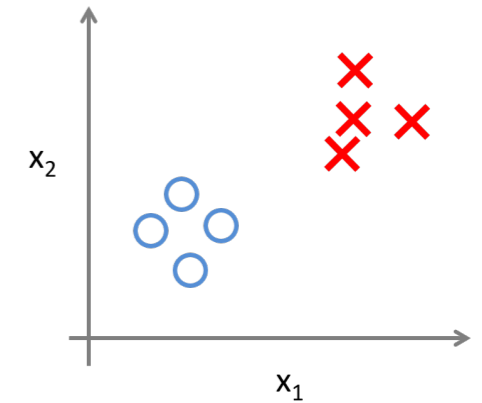
$$\begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \\ w_{13} & w_{23} \end{bmatrix} \cdot \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} = \begin{bmatrix} (w_{11} \times i_1) + (w_{21} \times i_2) \\ (w_{12} \times i_1) + (w_{22} \times i_2) \\ (w_{13} \times i_1) + (w_{23} \times i_2) \end{bmatrix}$$



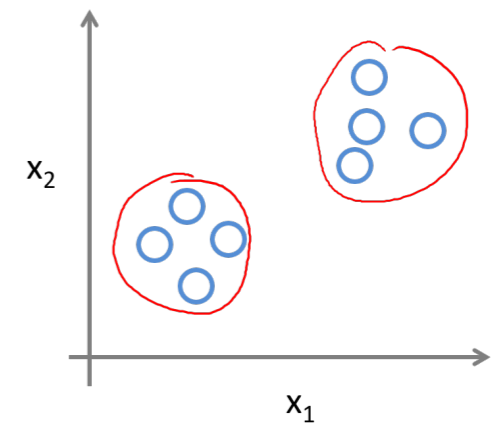
Types of learning

- **Supervised learning:**
 - attempt to predict/infer some target (truth label)
 - classification or regression target
 - train on data with known label (often MC simulation)
- **Unsupervised learning:**
 - target is not given
 - learn the data underlying distribution directly from data
 - can be used for sampling, clustering, anomaly detection, ...
- **Weakly/semi-supervised learning:**
 - Only noisy or partial targets are known

Supervised Learning

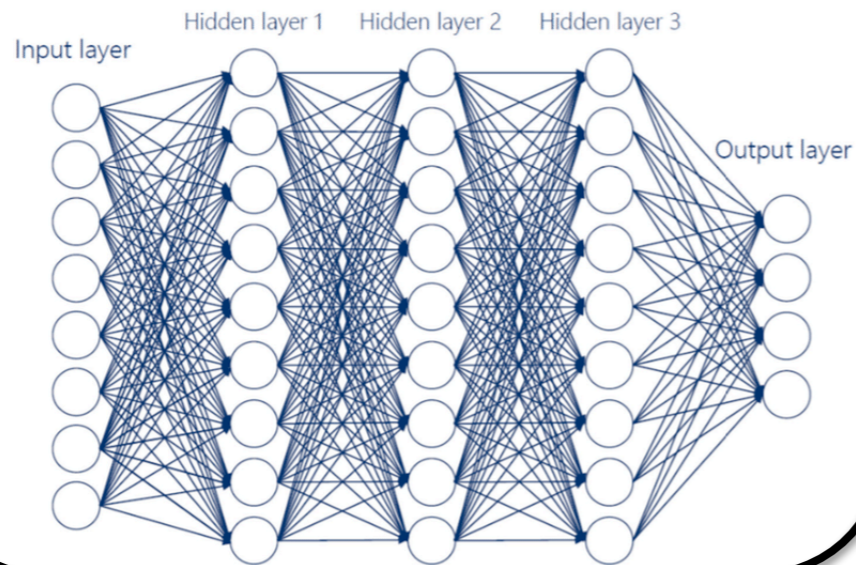


Unsupervised Learning

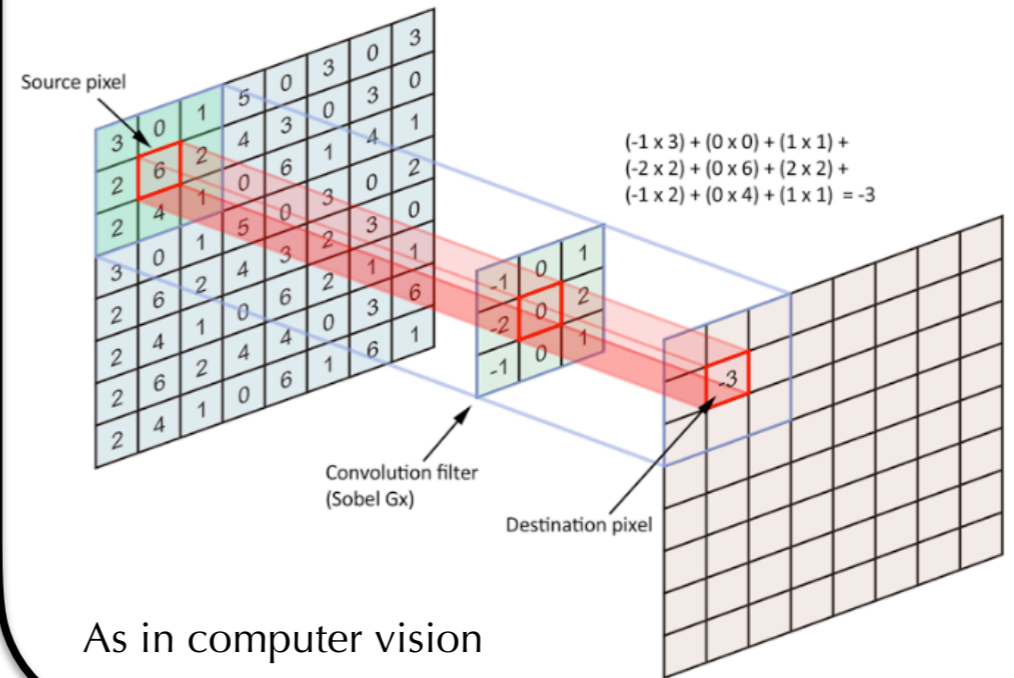


Data representations

**High level/tabular:
fully connected NN, BDTs**

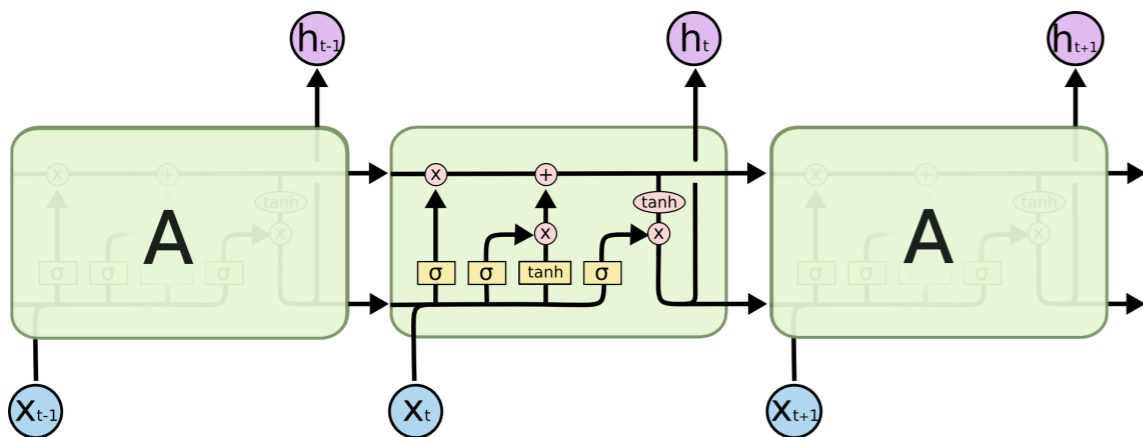


Regular grid: convolutional NN



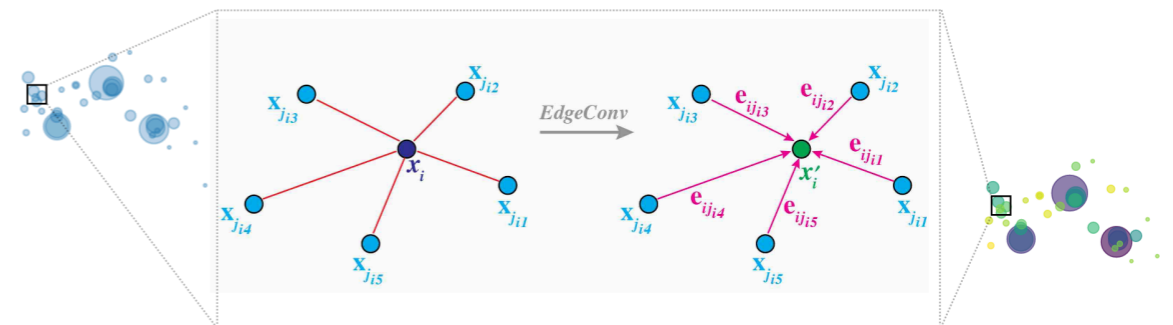
As in computer vision

**Ordered sequence/time series:
recurrent NN, transformers**



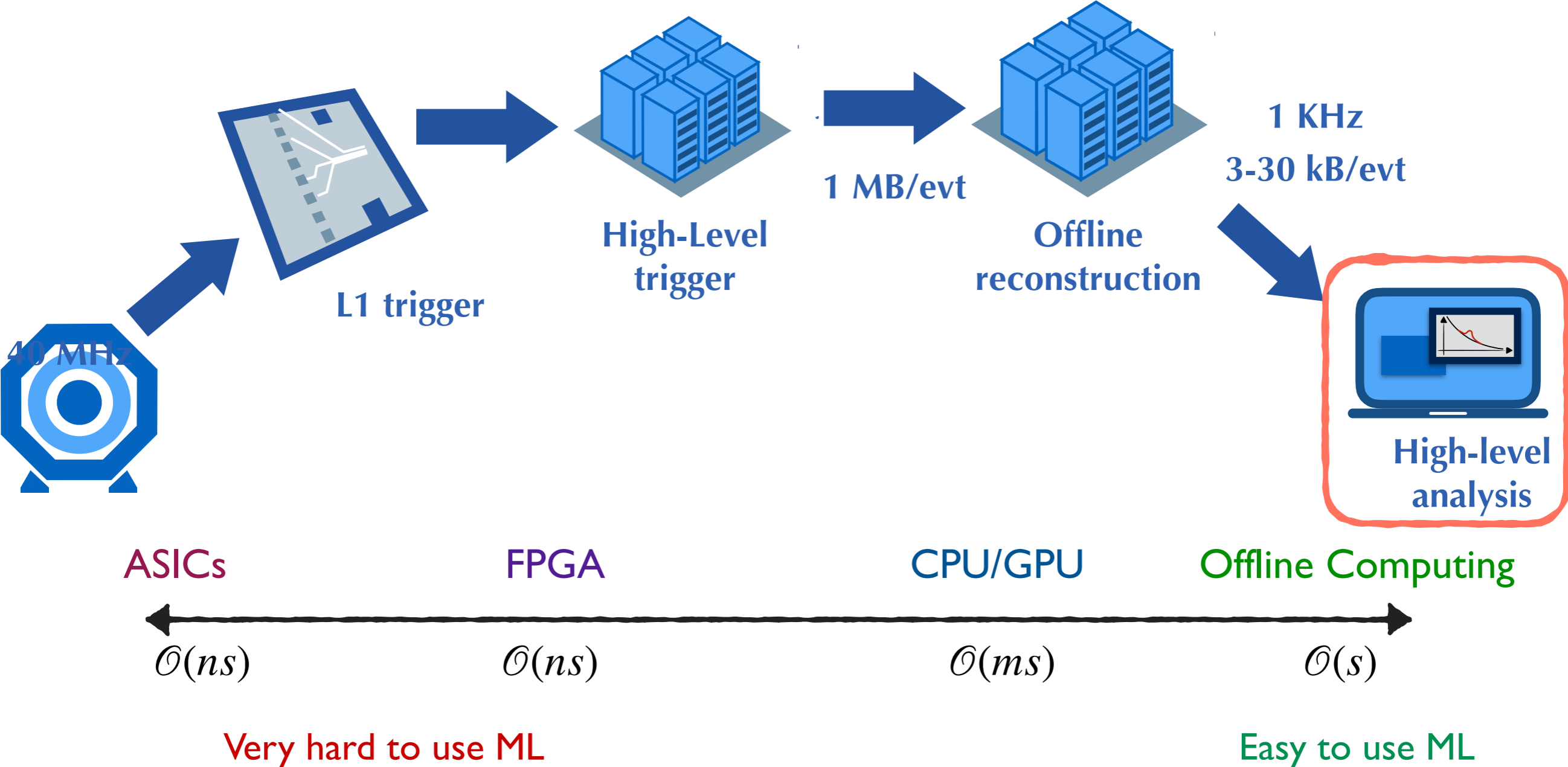
As in natural language processing

**Point cloud:
Deep sets, graph NN, transformers**



As in social media analysis

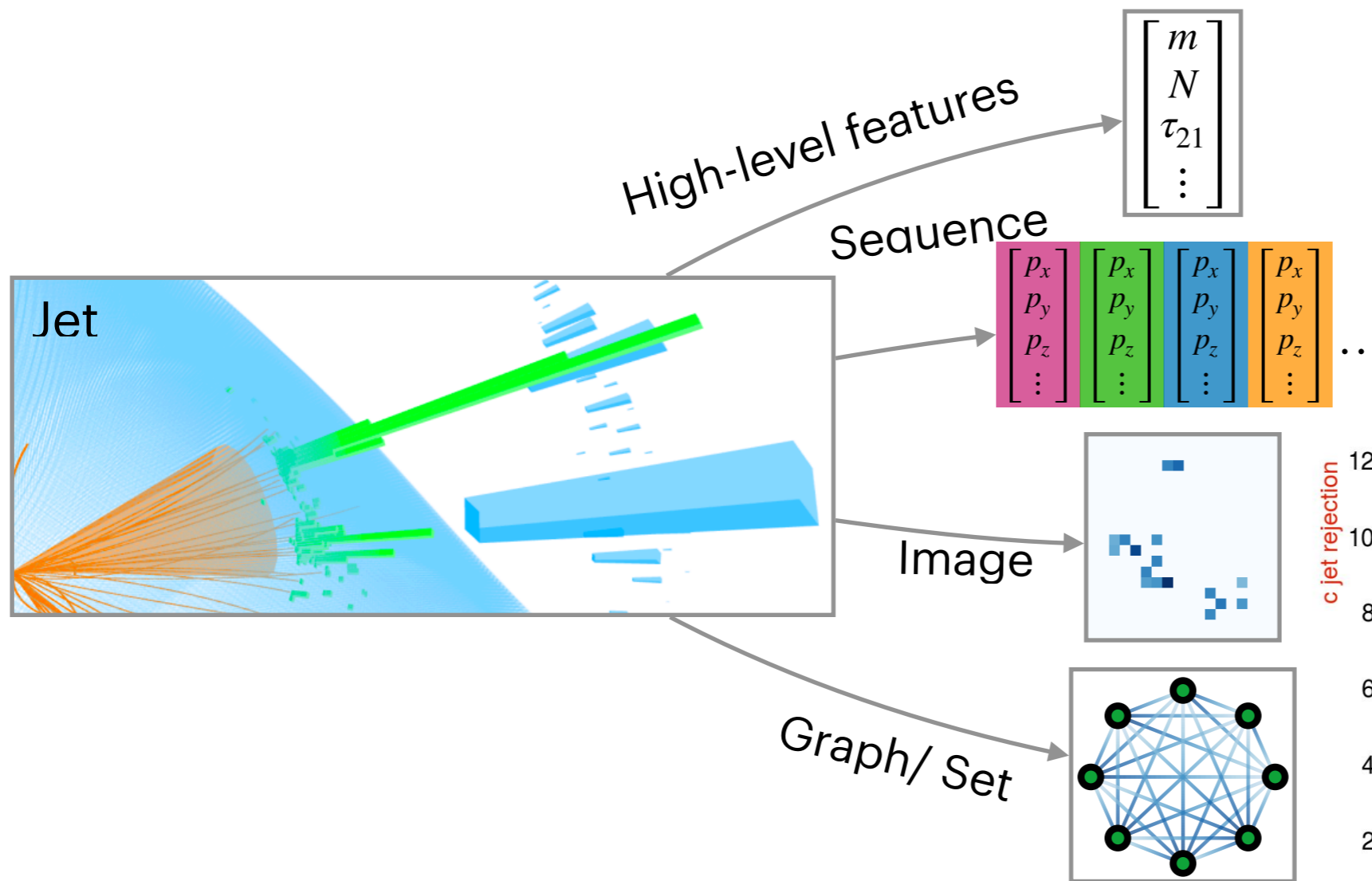
Data reduction workflow @ LHC



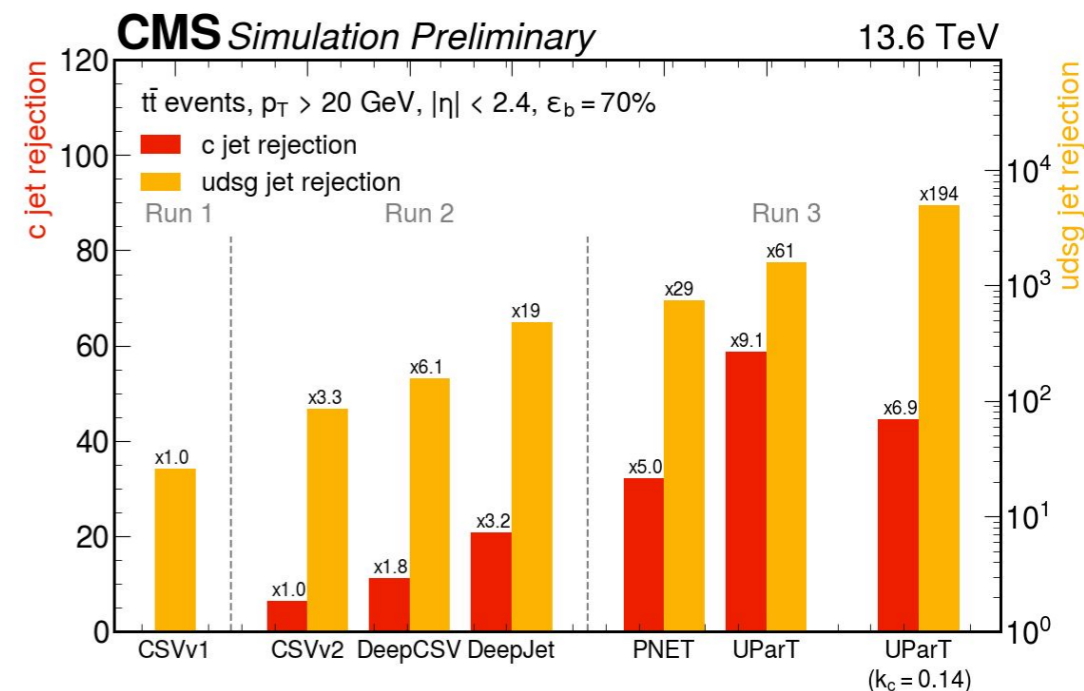
Taggers: From BDTs to Transformers

- ▶ Great strides made to leverage **rich low-level information** with **graph neural networks** and **transformers** for a variety of tasks including **jet classification**

[CMS-DP-2024-066](#)

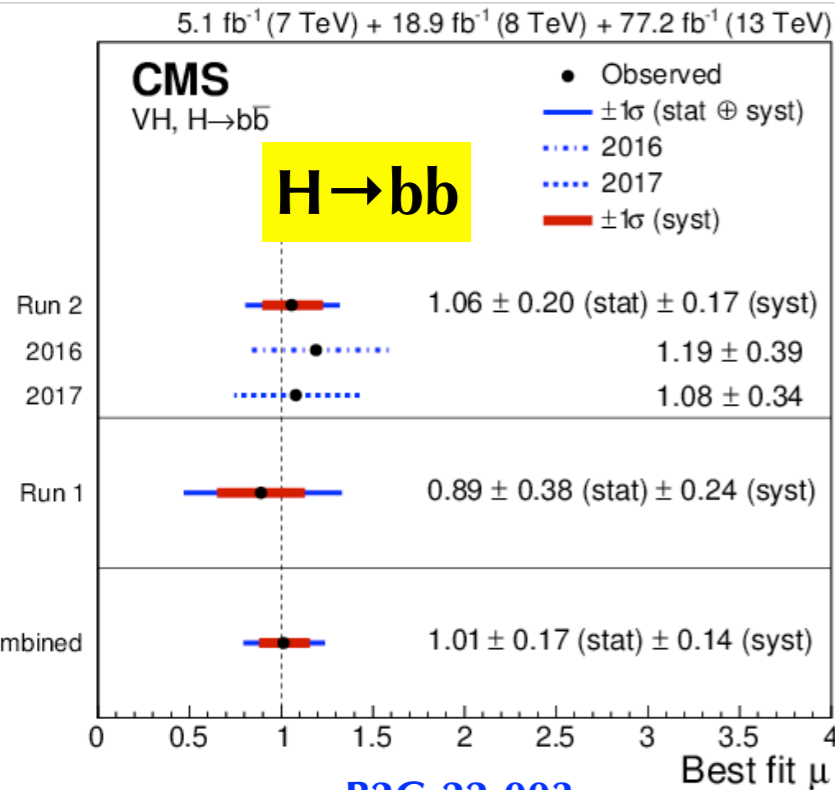


- ▶ Major improvements demonstrated in CMS (ParticleNet/Transformer)

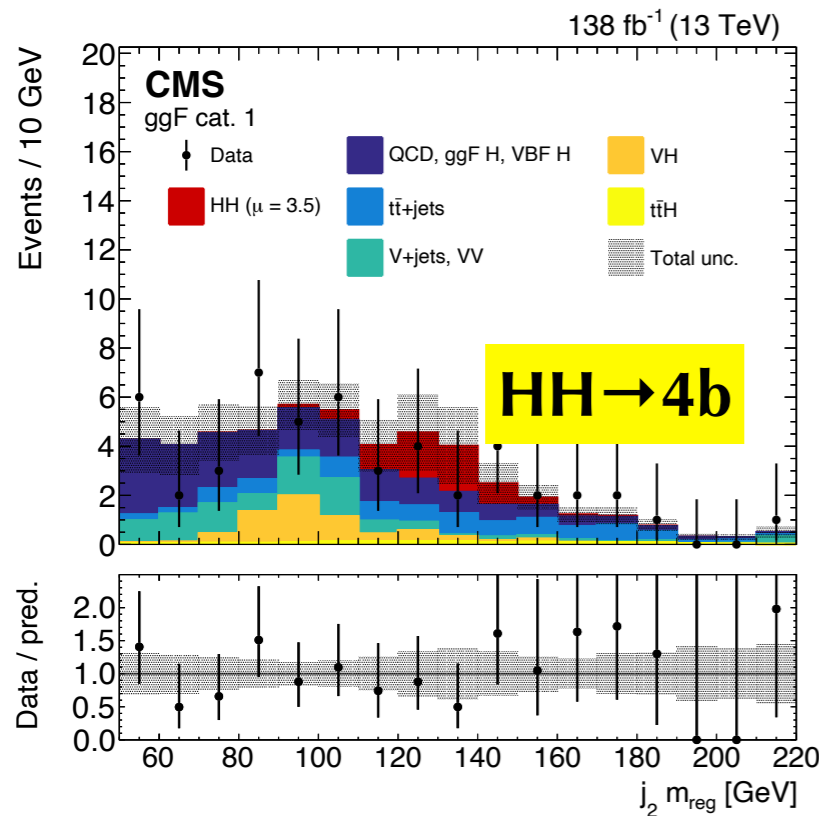


Impact: Better searches and measurements

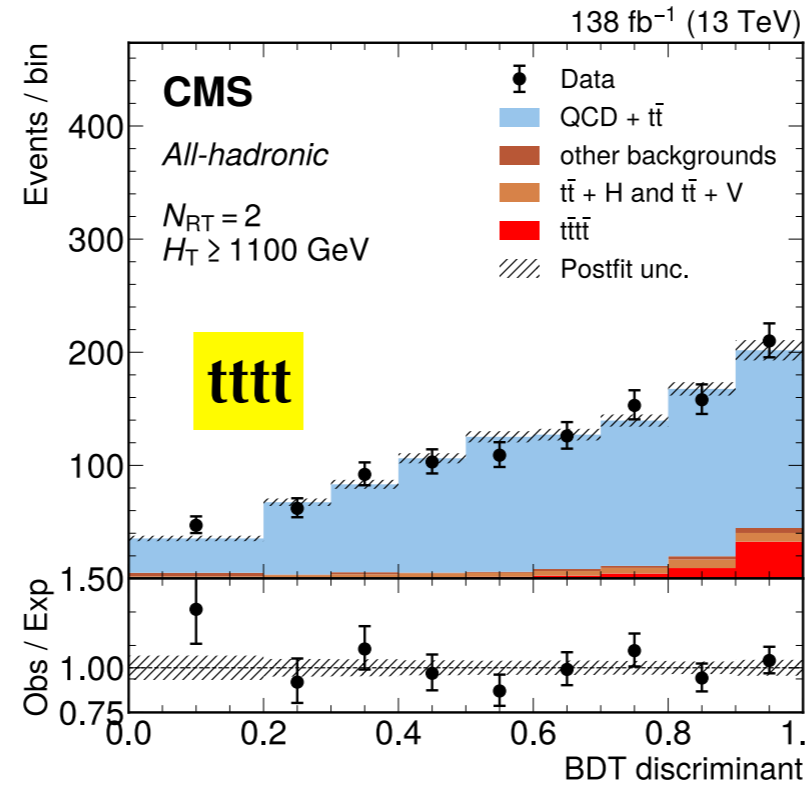
[HIG-18-016](#)



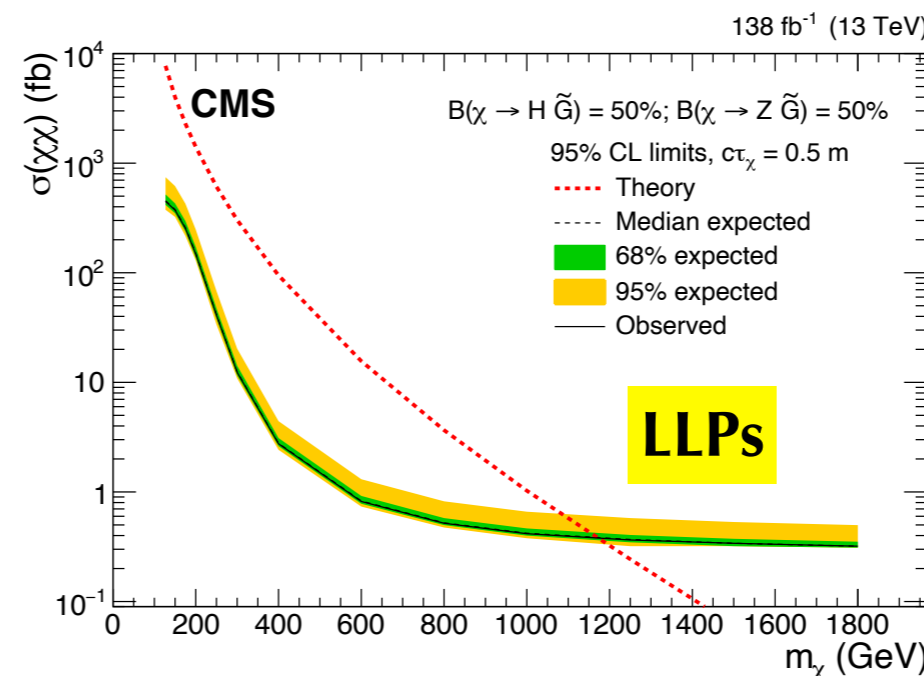
[B2G-22-003](#)



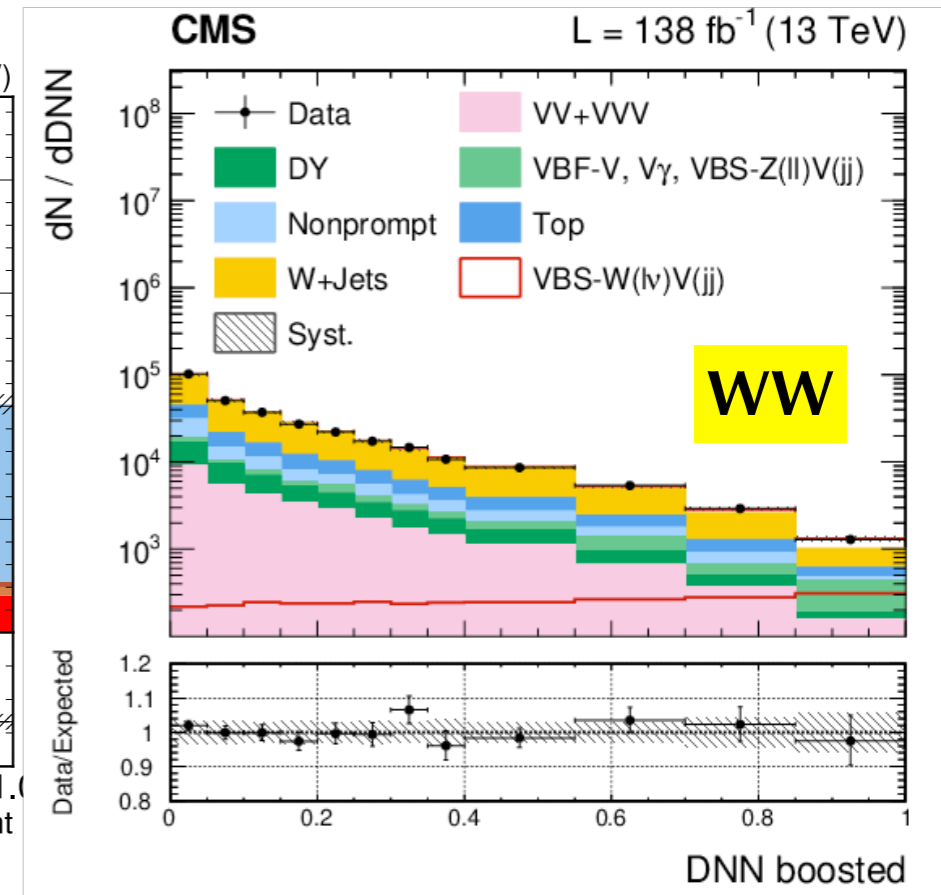
[TOP-21-005](#)



[EXO-21-014](#)



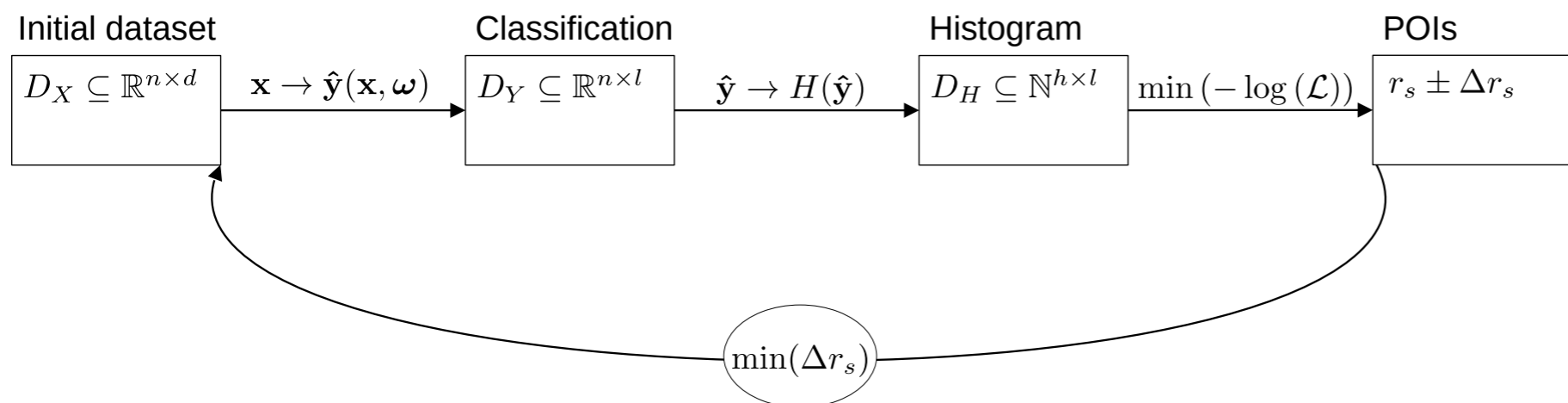
[SMP-20-013](#)



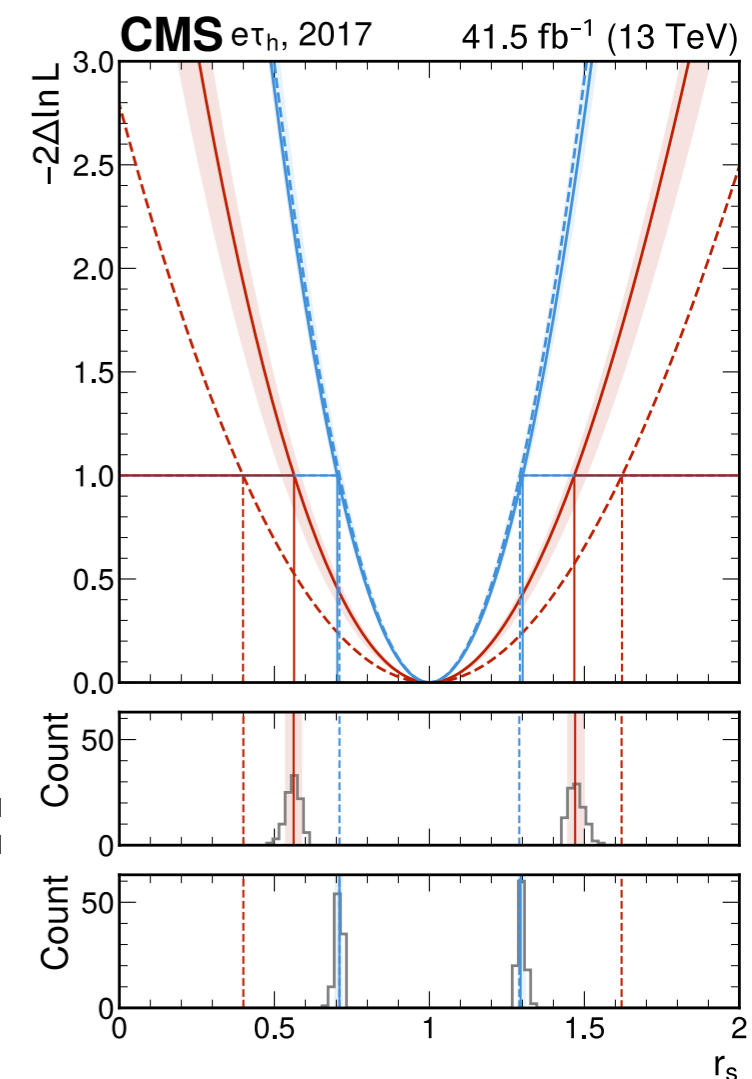
- Increased sensitivity with ML used by each PAG (mostly supervised):
 - improved event classification and object reconstruction

Systematic-aware learning

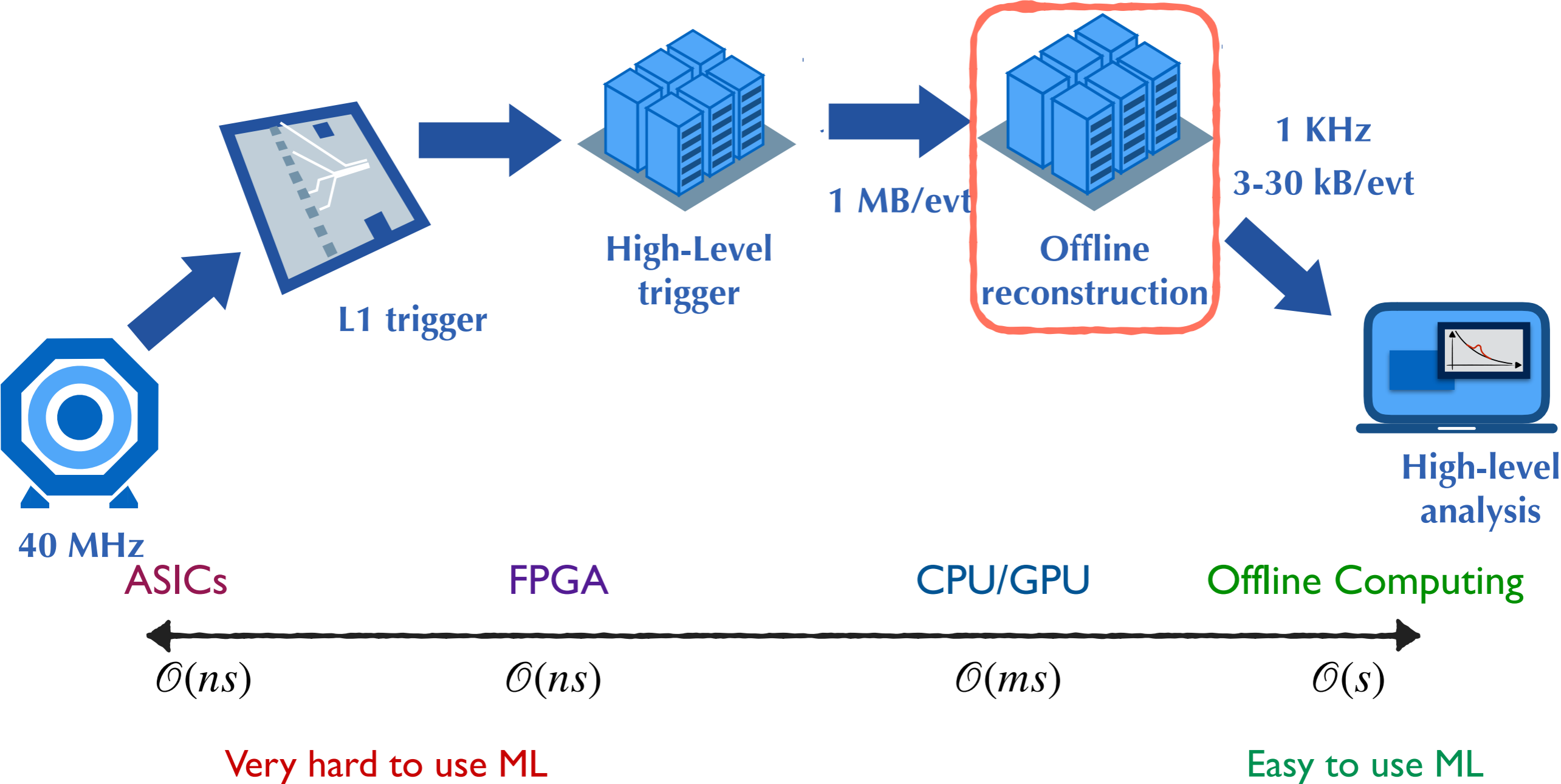
- ▶ Standard classifier training (CENNT) optimizes for signal vs. background discrimination without considering systematics and other effects that affect the ultimate figure: uncertainty Δr_s on a physics parameter r_s
- ▶ By implementing the analysis chain (including systematics) in a **differentiable way**, we can directly optimize for min. Δr_s in the neural network training!



SANNT		
—	$\Delta r_s = -0.44 \quad +0.47$	68% CI
—	$\Delta r_s^{\text{stat}} = -0.30 \quad +0.30$	68% CI
CENNT		
- - -	$\Delta r_s = -0.60 \quad +0.62$	
- - -	$\Delta r_s^{\text{stat}} = -0.29 \quad +0.29$	



Data reduction workflow @ LHC

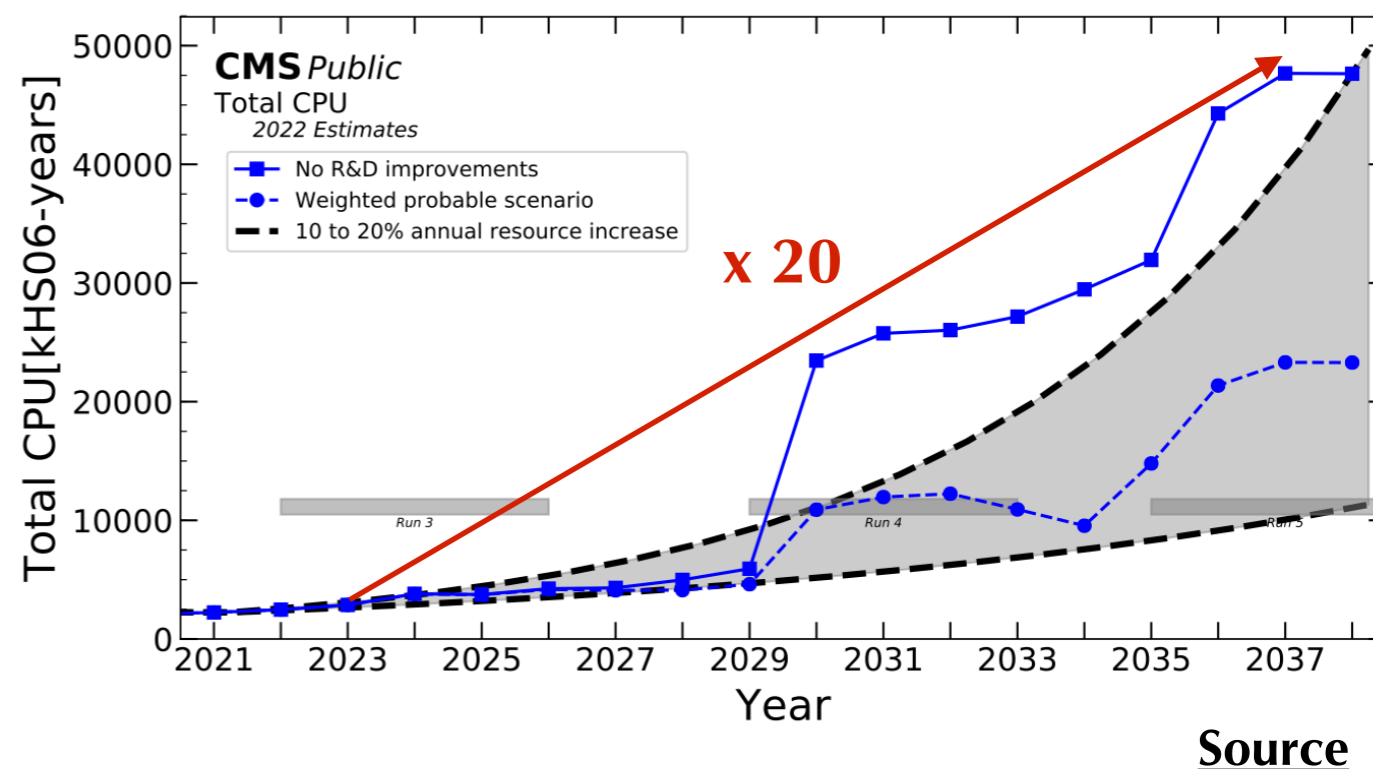
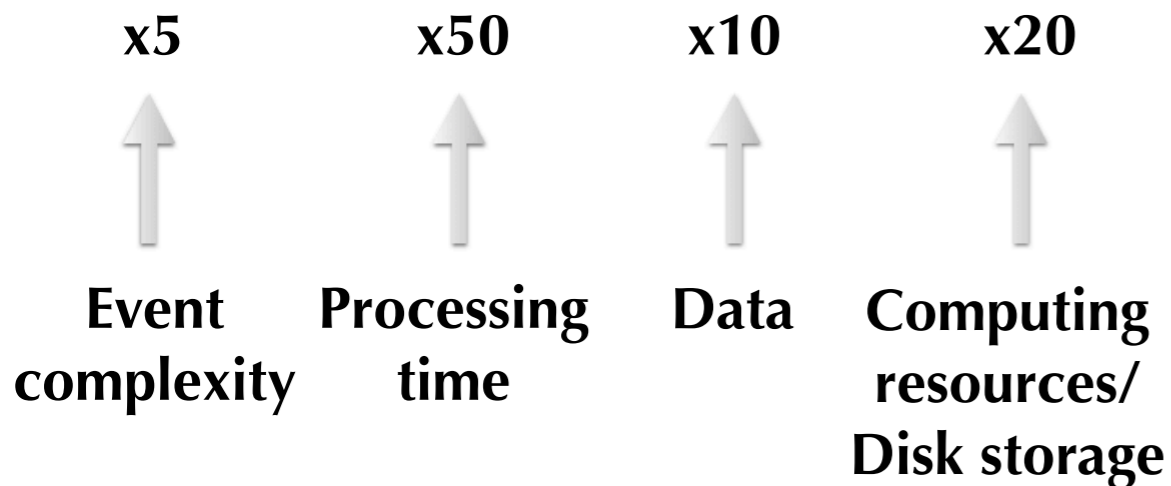
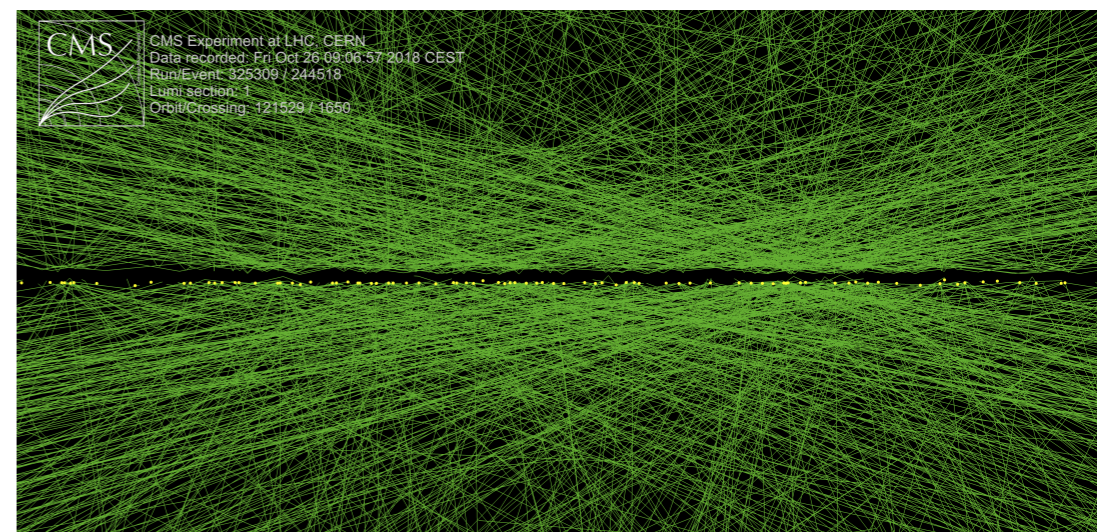


The HL-LHC challenge

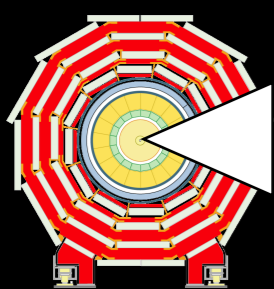
With more particles per collision and more readout channels to combine, the reconstruction to become even more computing intensive

We cannot throw away more data
We must increase the throughput
with at most flat budget for computing resources

**We must do more with less
to preserve the physics!**

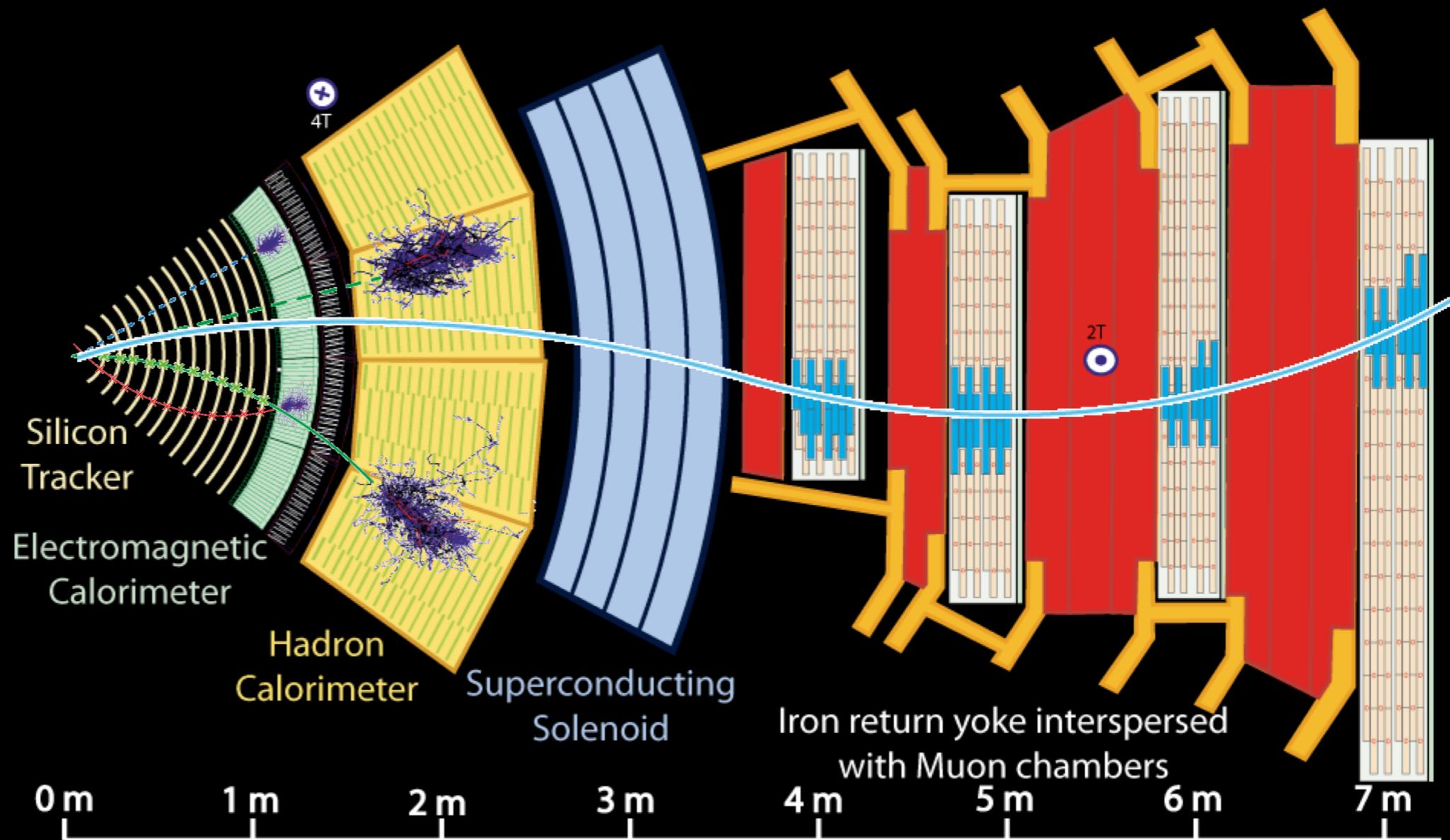


Multilayered Detectors, e.g. CMS



Current and future multilayered detectors...

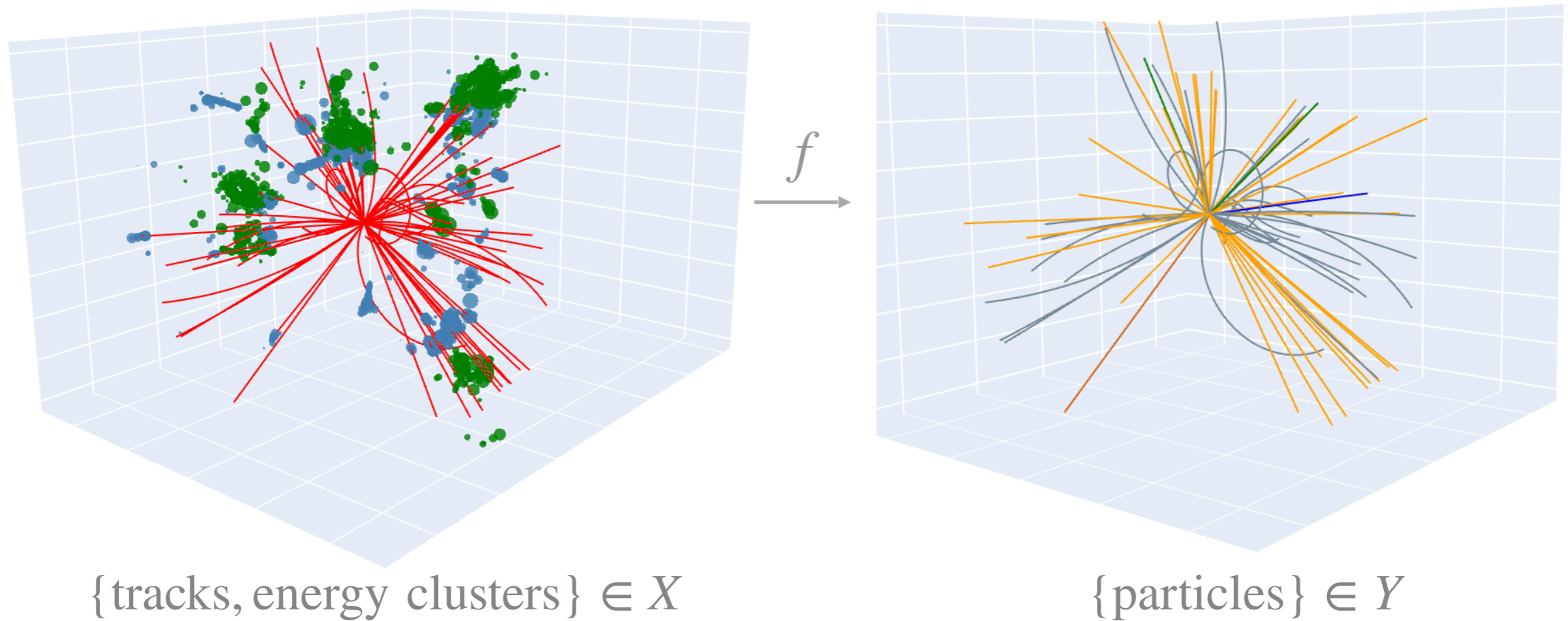
Require complex pattern recognition



- Key:
- Muon
 - Electron
 - Charged Hadron (e.g. Pion)
 - Neutral Hadron (e.g. Neutron)
 - Photon

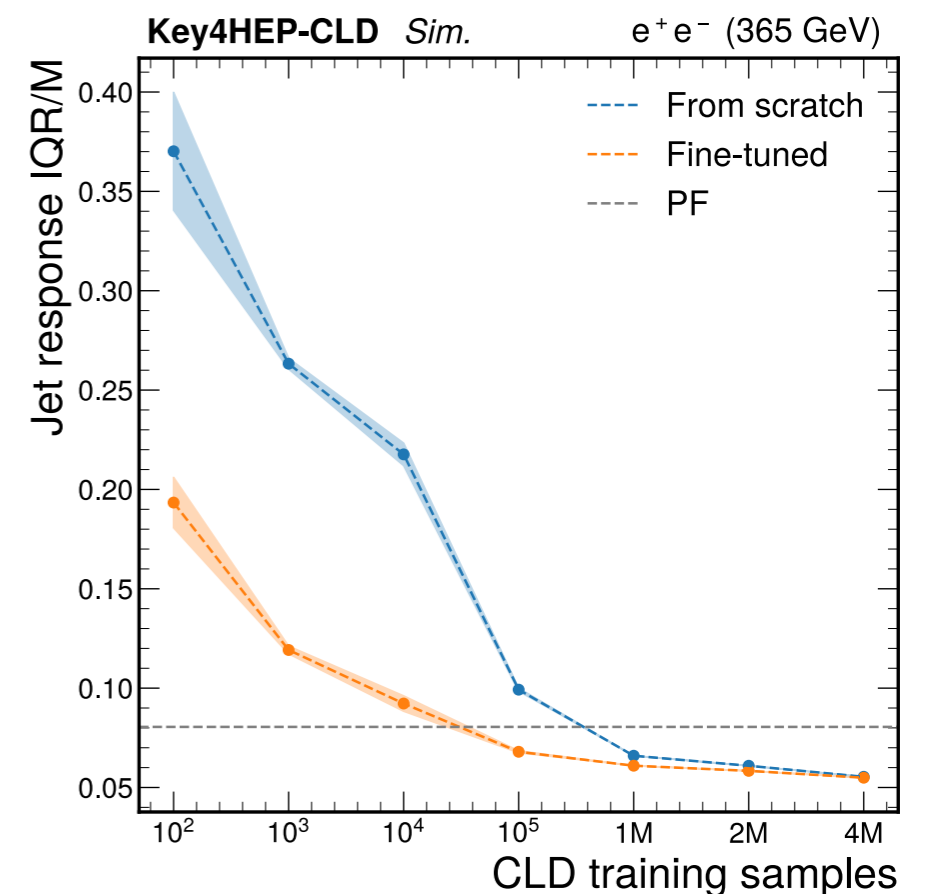
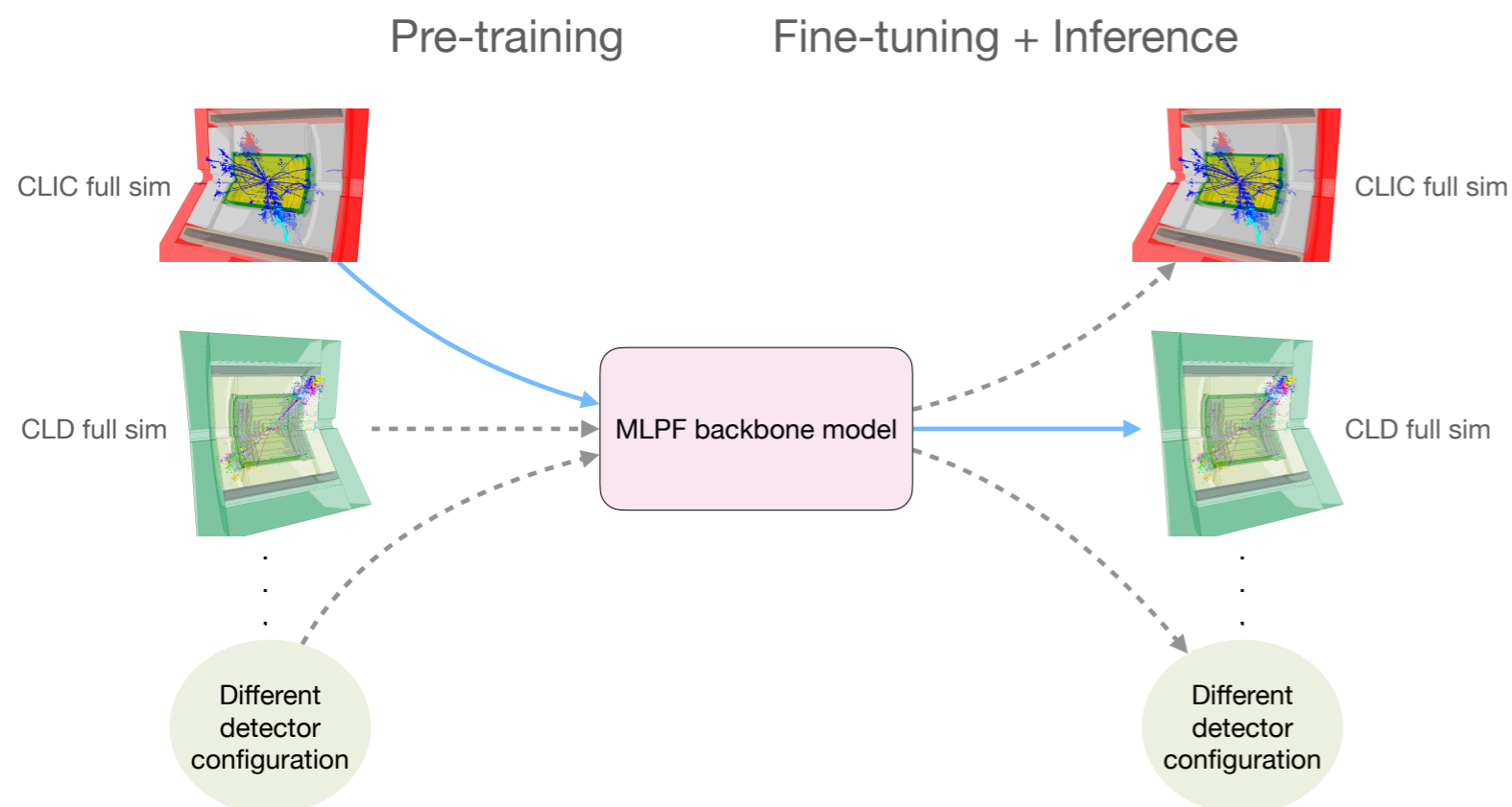
Particle-flow as a Machine-Learning task

- ▶ Can we formulate PF as an ML task (naturally “tunable” through re-training and portable to new hardware)?
- ▶ Learn a “set-to-set” function $f: X \rightarrow Y$, where $\{\text{tracks, energy clusters}\} \in X$ and $\{\text{particles}\} \in Y$



Cross-detector Transfer learning study

- ▶ Pre-train MLPF model on CLIC full sim. samples then fine-tune on full sim. samples from CLD at FCC-ee
- ▶ Fine-tuned model outperforms model trained from scratch for $<2\text{M}$ events, with 30% better jet resolution at 100k events, and outperforms PF baseline



(b) A pre-training and fine-tuning MLPF approach

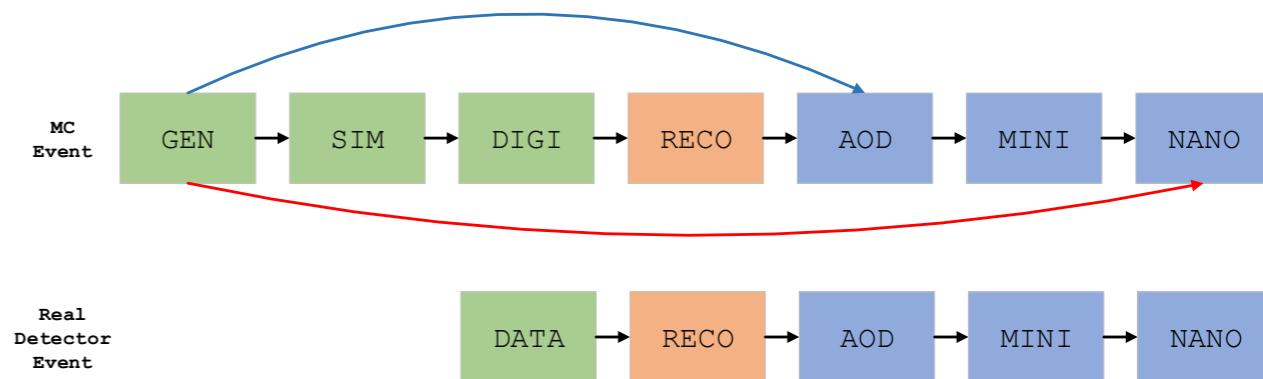
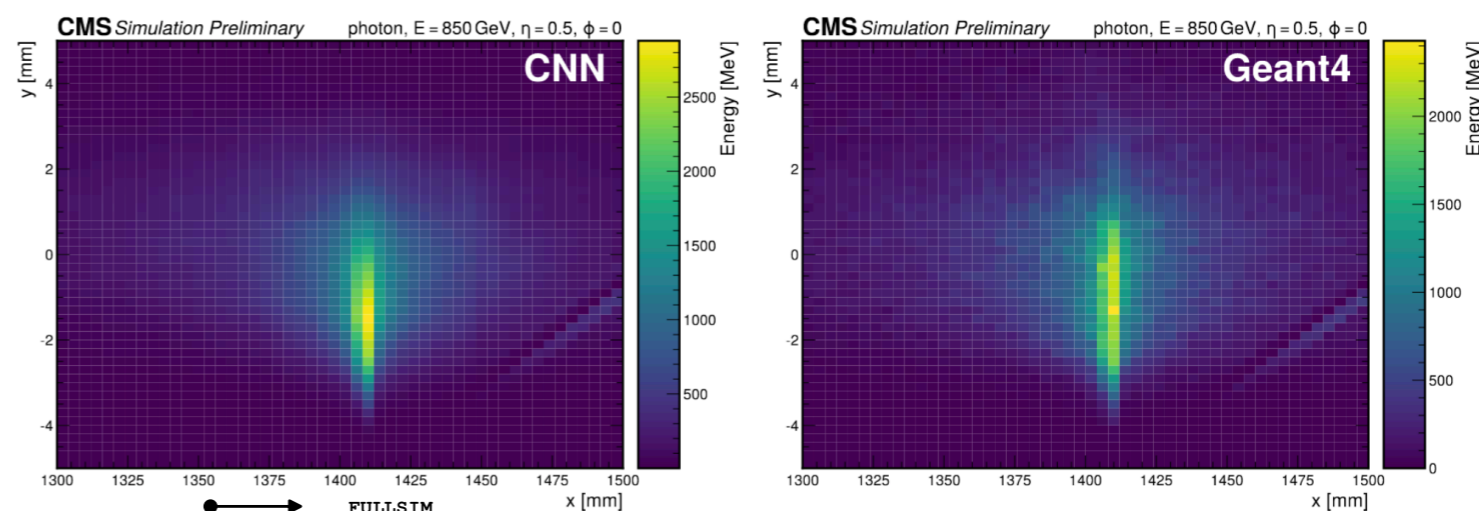
ML @ CMS: simulation

- Fully detailed simulation with generator+GEANT4 is computing intensive
- Generative models can though provide 1000x speed up
- Applicable at many levels: sampling, generation, detector model, analysis variables, etc.
- Careful study of statistical power of learned models over training samples

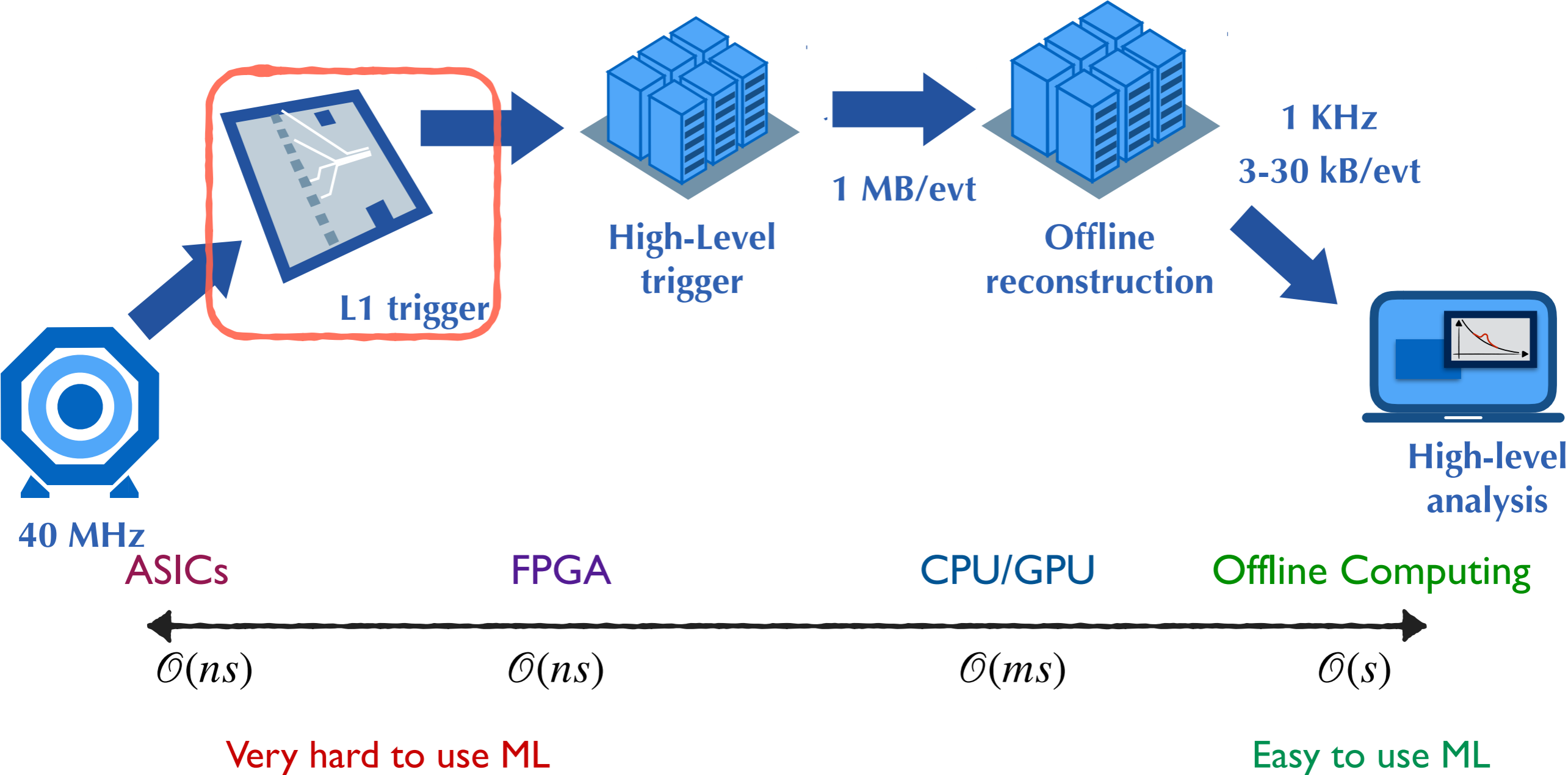
- A lot of R&D with examples in CMS: simulation of jets, photon and hadron showers in calorimeters

- A few ongoing use cases:

- Denoising with convolution neural networks
- FlashSim: end-to-end simulation using normalizing flows



Data reduction workflow @ LHC



The High-Luminosity LHC challenge

instantaneous luminosity

x 0.75

x 1-2

NOW

x 2.5

2029

x 5-7

LHC Run 1
 $\sqrt{s} = 7-8 \text{ TeV}$
30/fb

Long
Shutdown 1

LHC Run 2
 $\sqrt{s} = 13 \text{ TeV}$
150/fb

Long
Shutdown 2

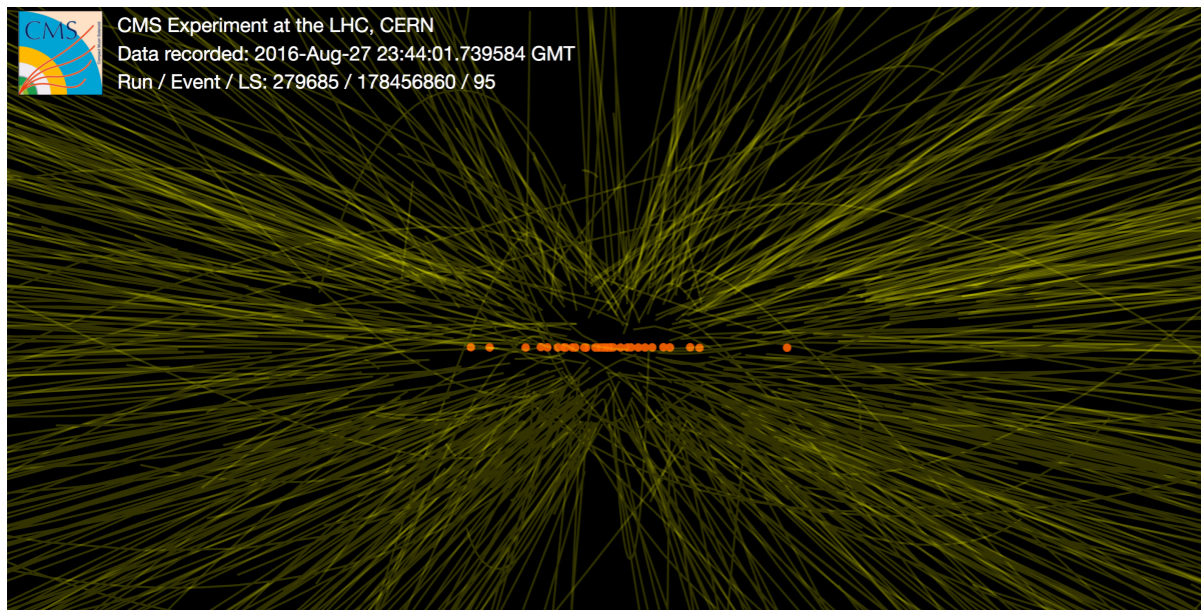
LHC Run 3
 $\sqrt{s} = 14 \text{ TeV}$
300/fb

Long
Shutdown 3

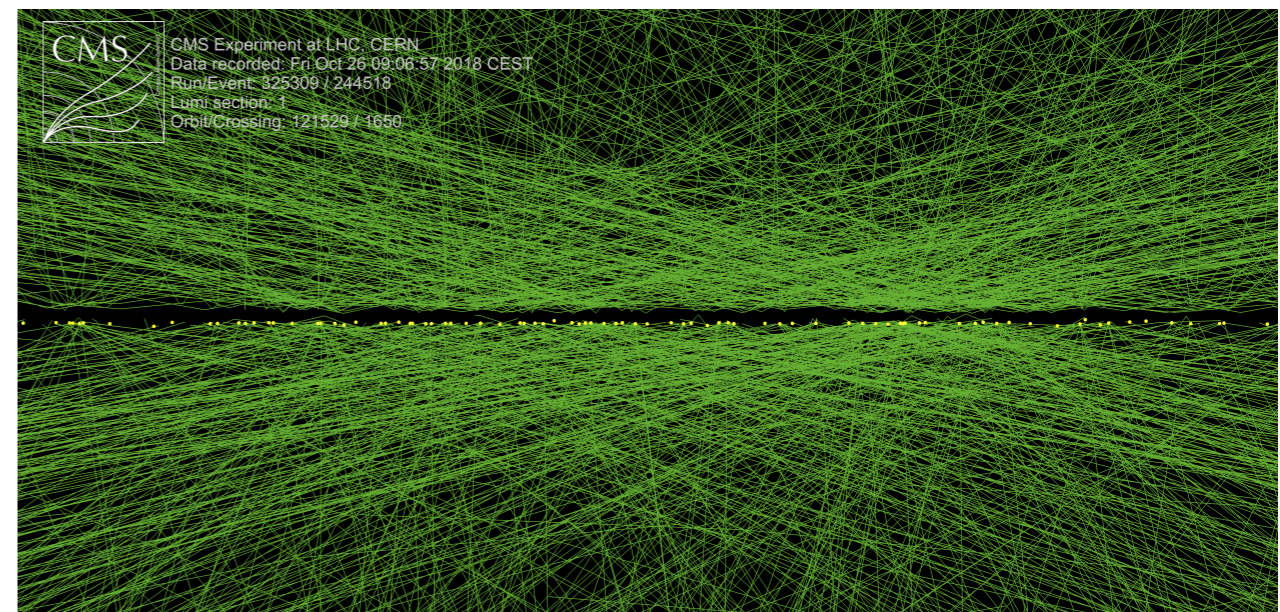
Run 4: HL-LHC
 $\sqrt{s} = 14 \text{ TeV}$
3000/fb

LHC TODAY

HL-LHC



40 simultaneous collisions
per bunch crossing



200 simultaneous collisions
per bunch crossing

+

more granular detector!

Intelligent Triggering

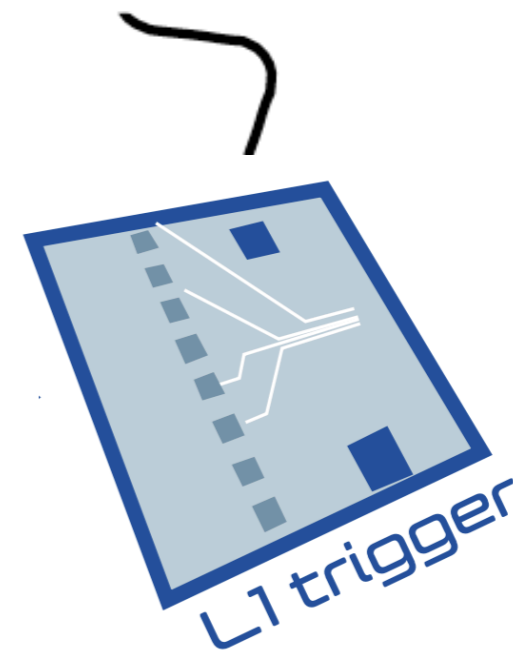
- Intelligent selection is most impact full at L1 trigger

Keep only BSM physics

make an smart decision!



Umm... What do I do then ?

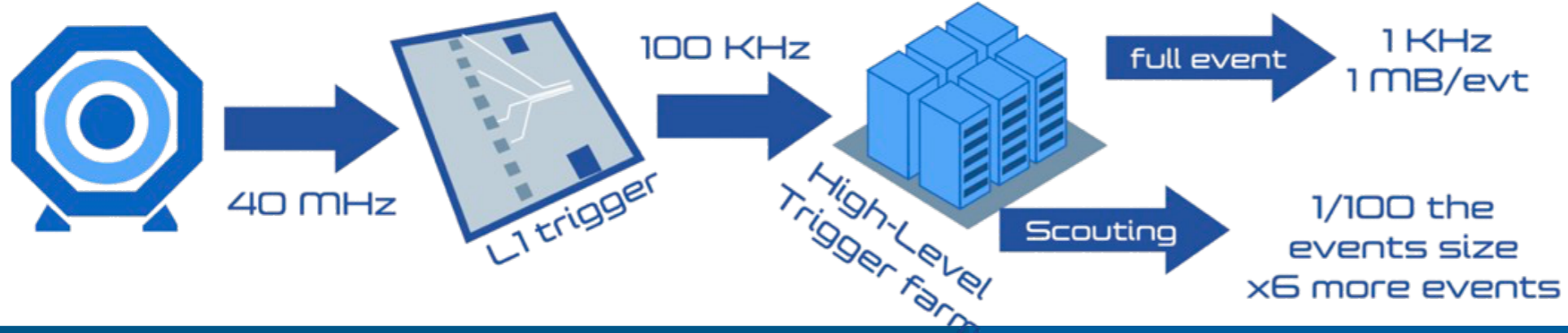
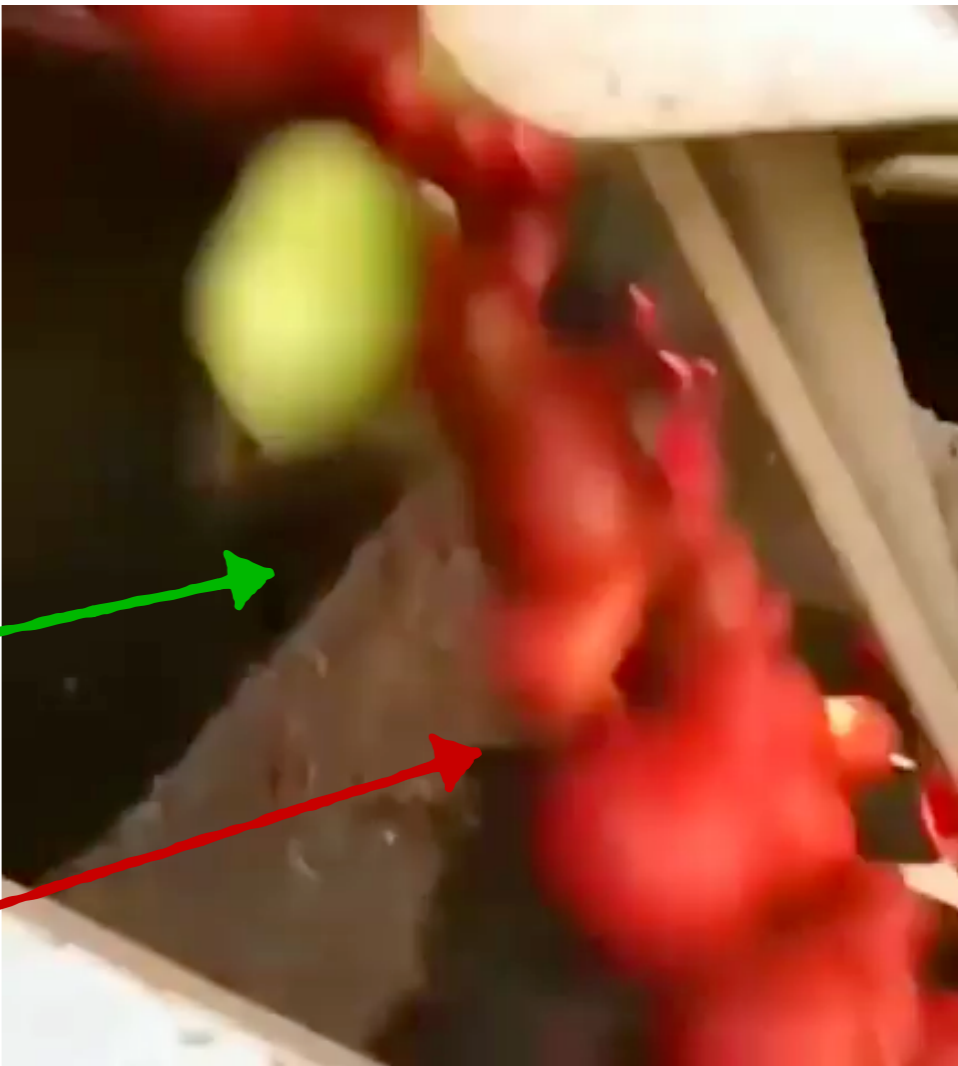


L1 Trigger in CMS

- CMS L1 trigger processes every collision
 - Need quick decisions w/ large data rate
- L1 runs on FPGAs
(w/ $3.8\mu s$ latency & selects 1 in 400 pp collisions)

Retain much of this
(Needs to be unbiased)

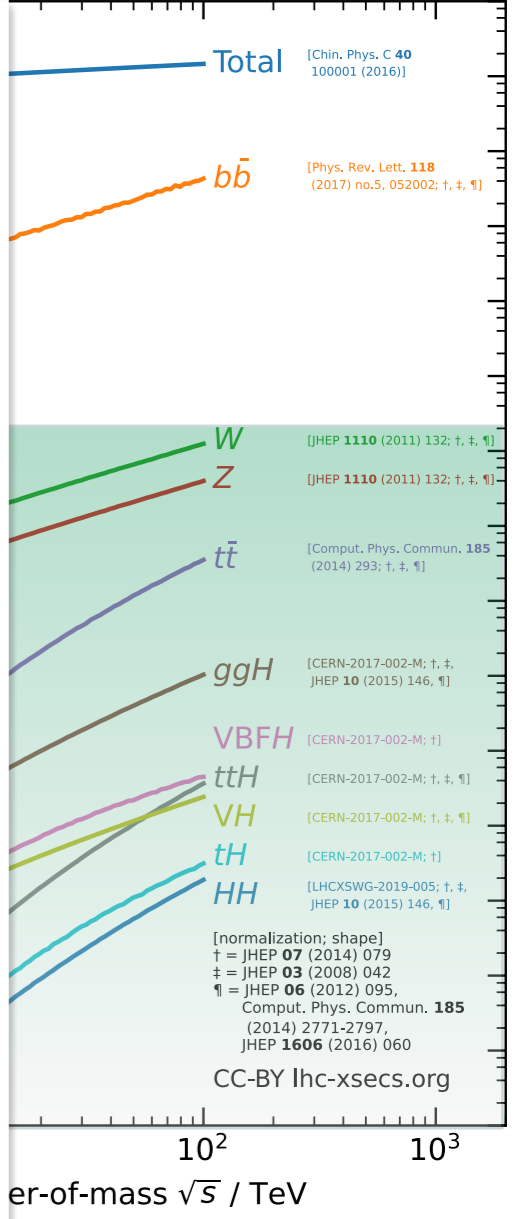
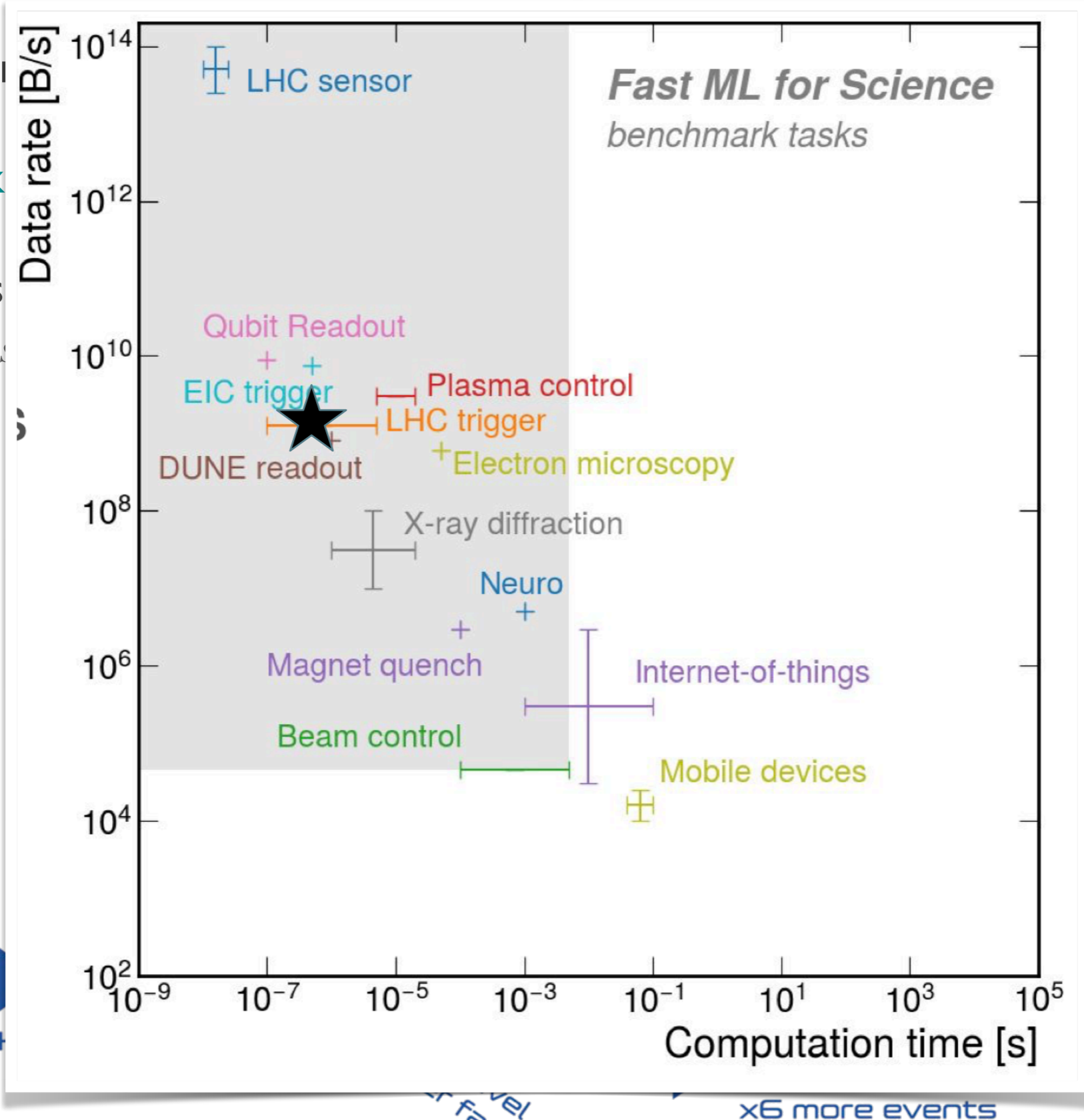
Discard a lot of this



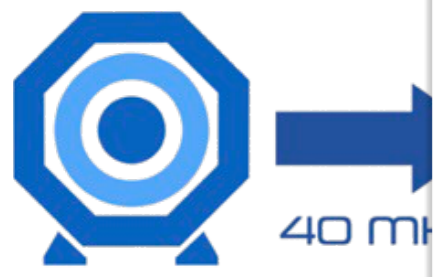
99.75% gets thrown away in first stage!

L1 Trigger in CMS

- CMS L1 trigger
 - Need quick
 - L1 runs (w/ 3.8μs)
 - selects 1

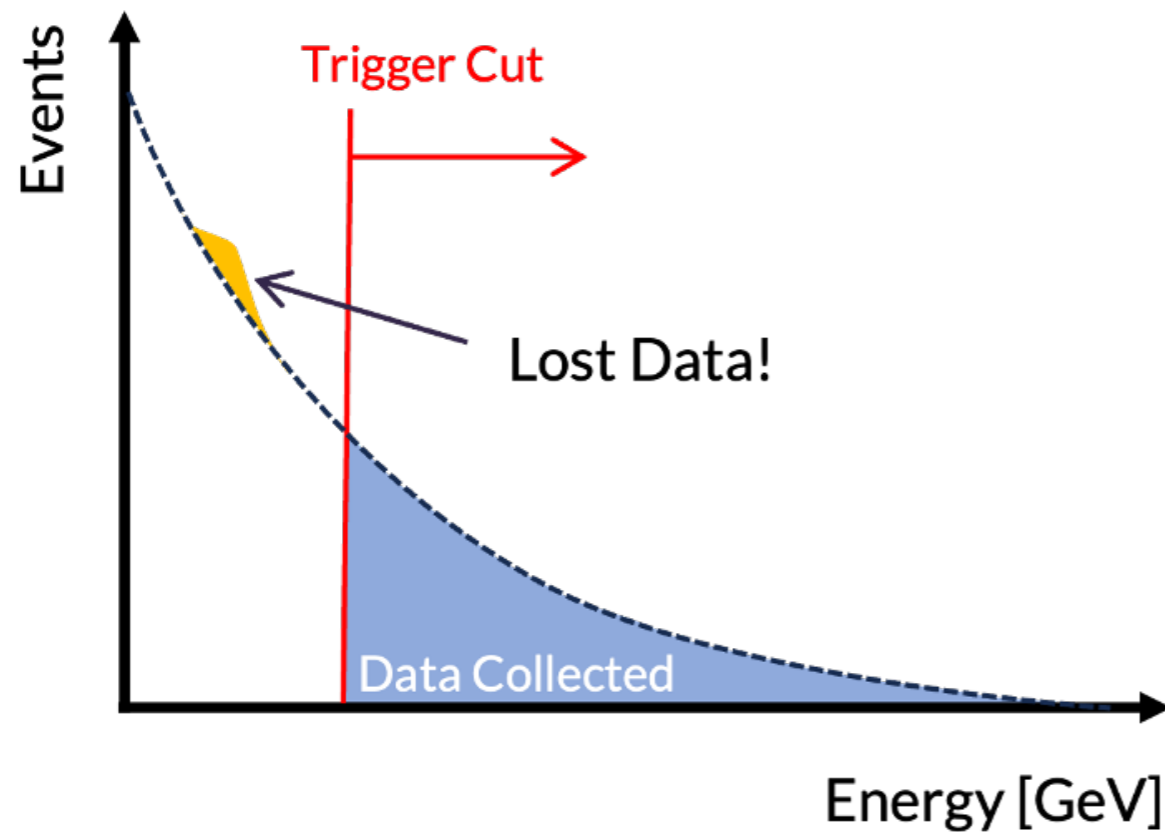


9.75% gets thrown away in first stage!



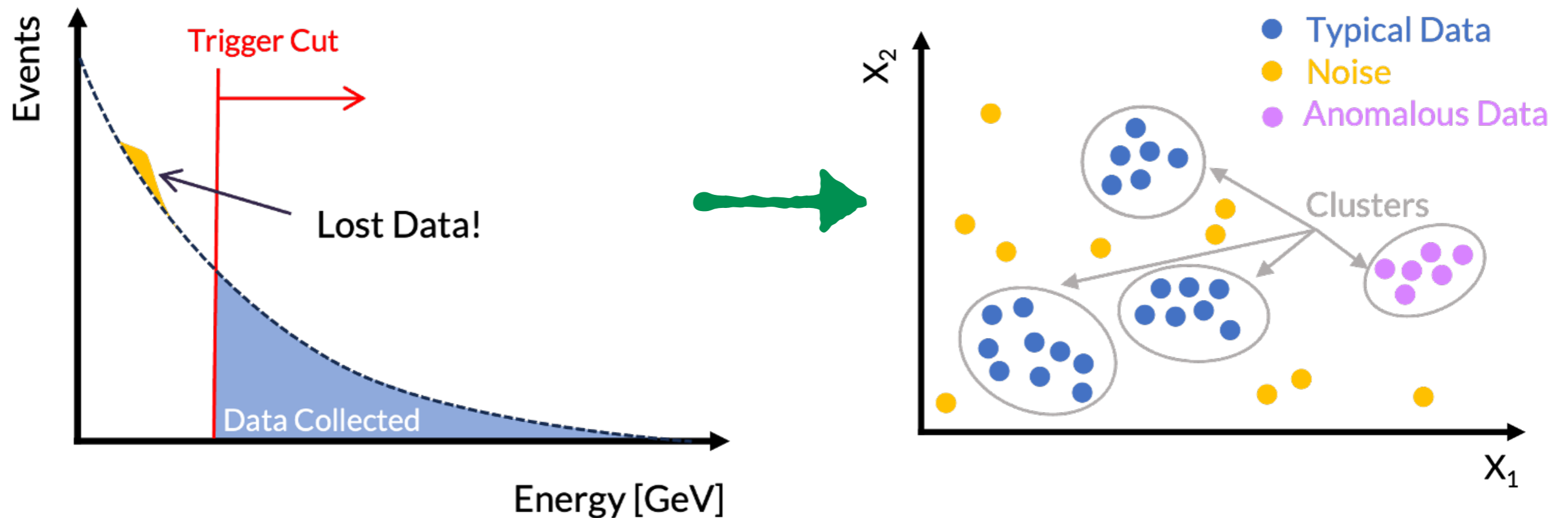
A new kind of data stream

- Simple kinematic selections to define trigger algorithms
 - **What if we miss it?** We might not know what we are looking for



L1 Anomaly Detection

- New Approach: **Anomaly Detection**
 - What if we select “**anomalous**” events based on decay topology
 - With **Machine Learning @ L1 trigger**

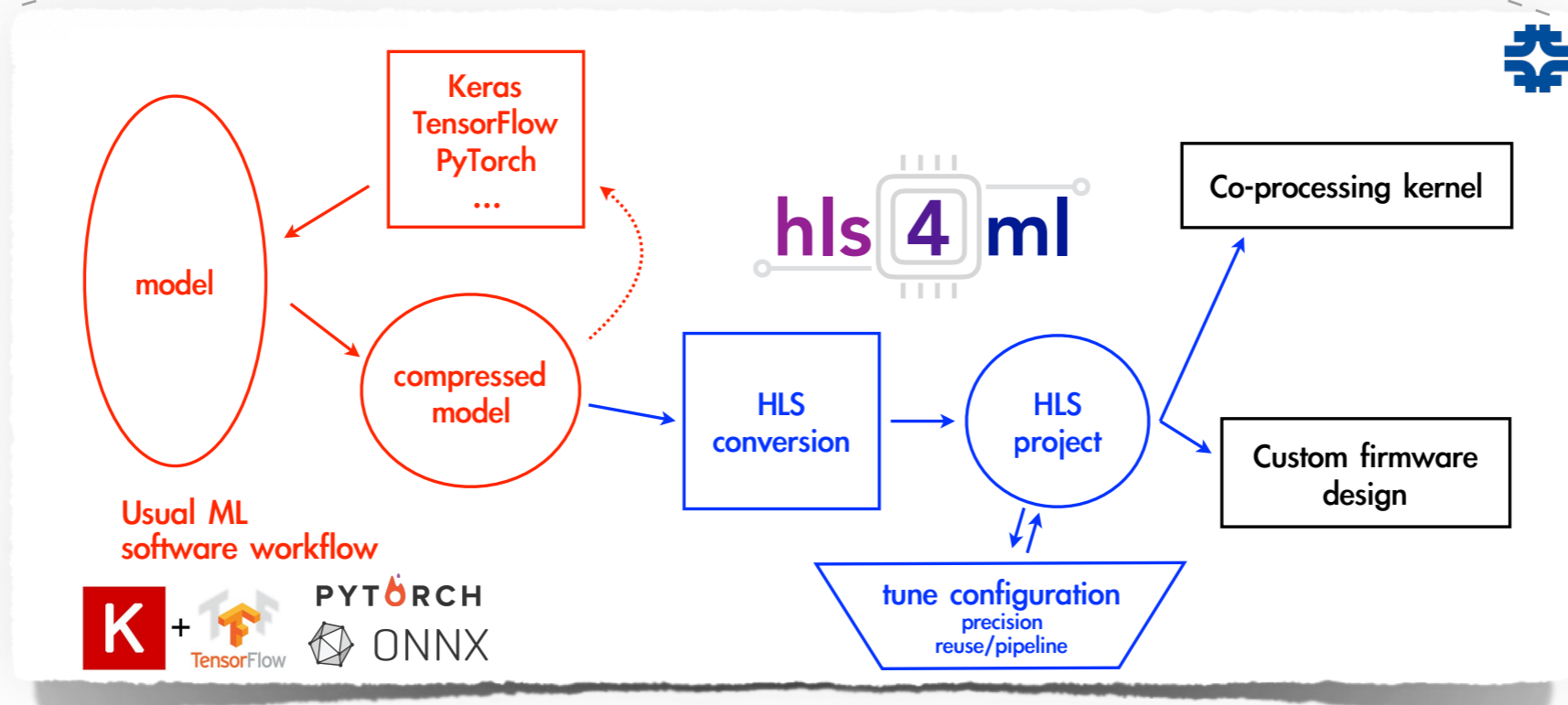


Enter  **AXOLITL**

Anomaly eXtraction Online L1 Trigger aLgorithm

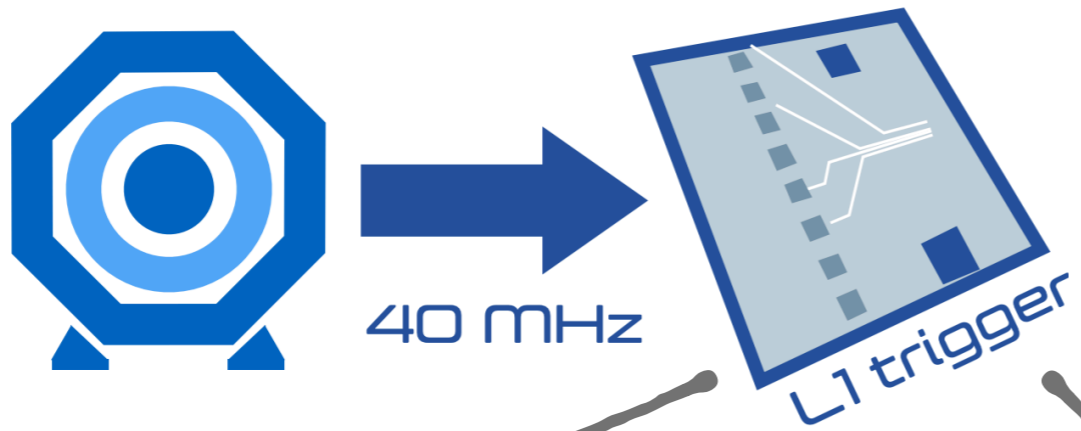
Fitting AXOL1TL on FPGA

- Instead of vanilla autoencoder, we use Variational Autoencoder
 - Use just the encoder part of the network for **AXOL1TL**
 - With **hls4ml** to deploy it on a FPGA

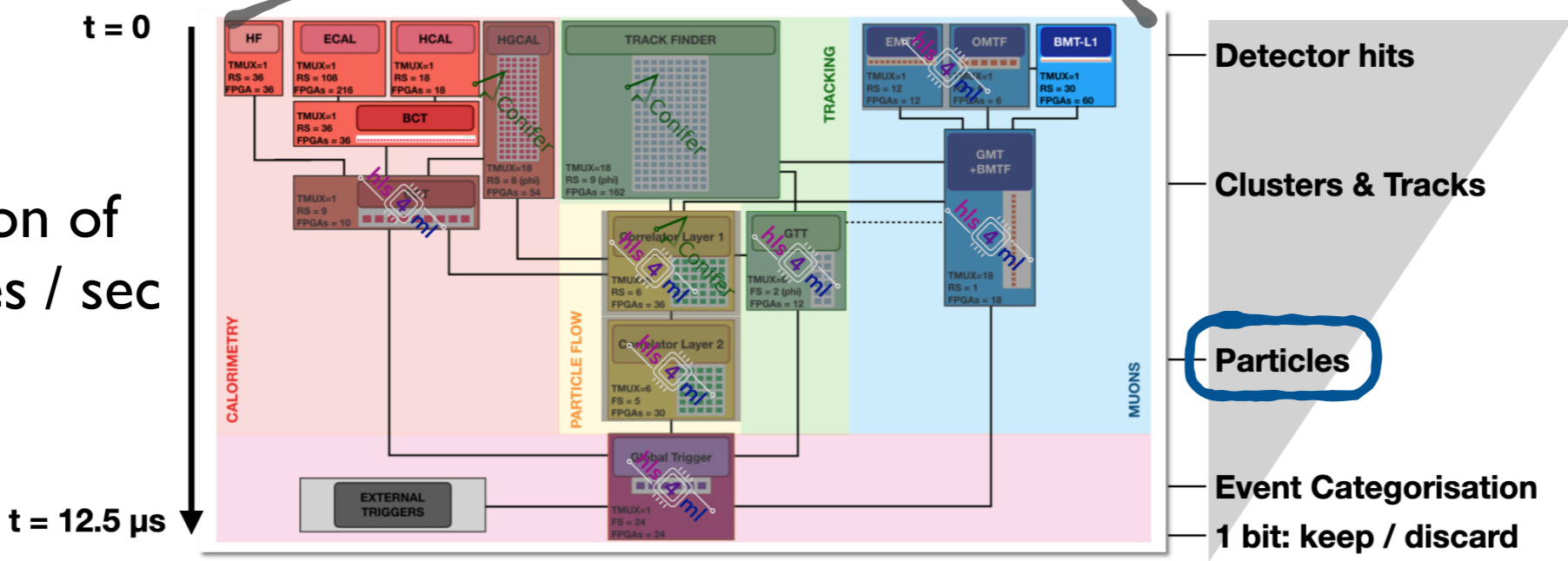


CMS @ HL-LHC

- CMS L1 Trigger @ HL-LHC is super charged with latest FPGAs + Algorithms
 - Can reconstruct as well as current High-Level Trigger @ 40 MHz
 - Even reconstruct individual particles for every single collision!

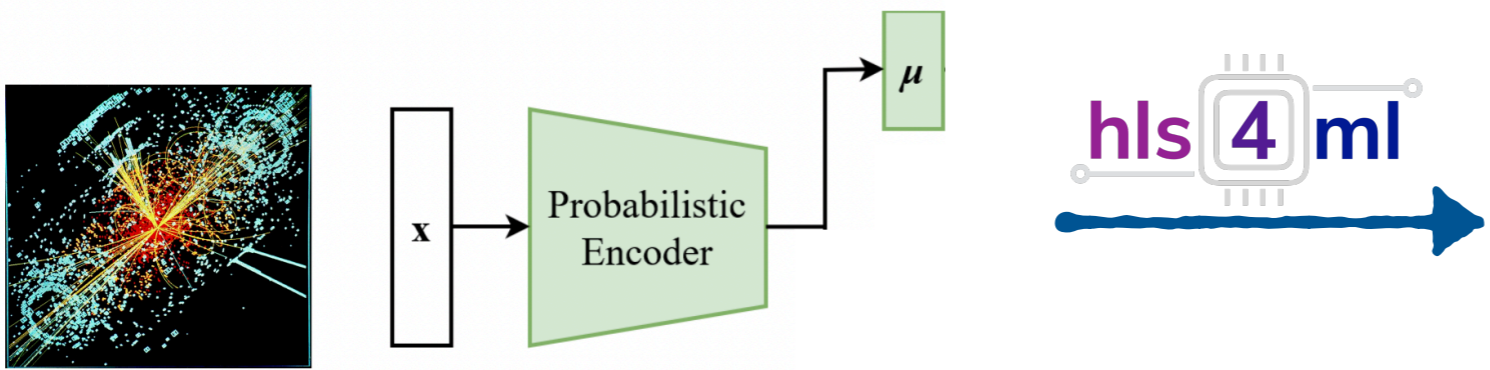


25 Billion of inferences / sec

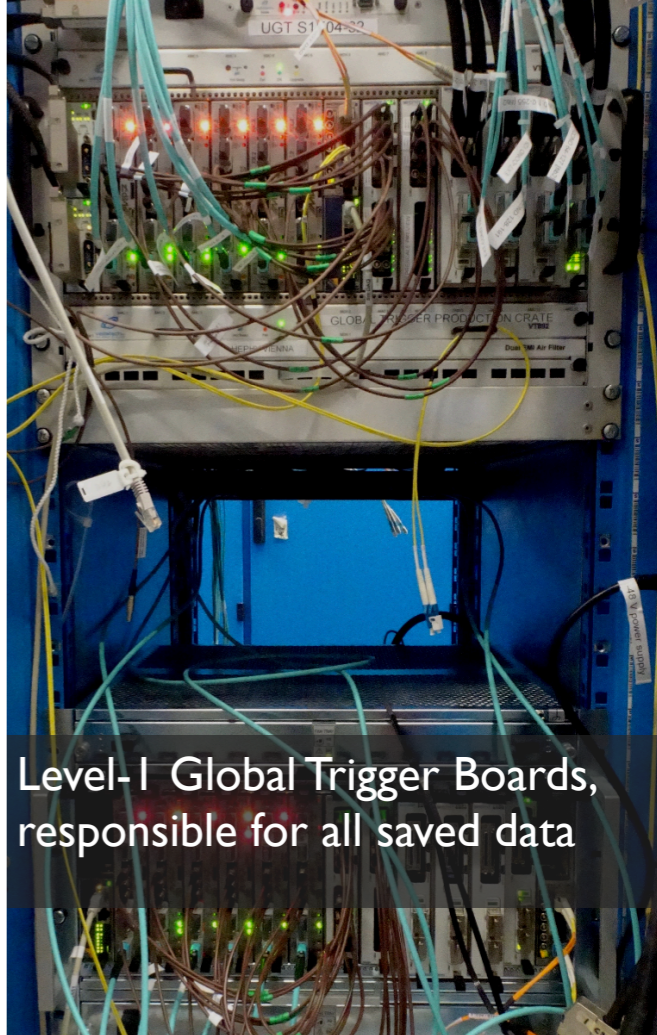
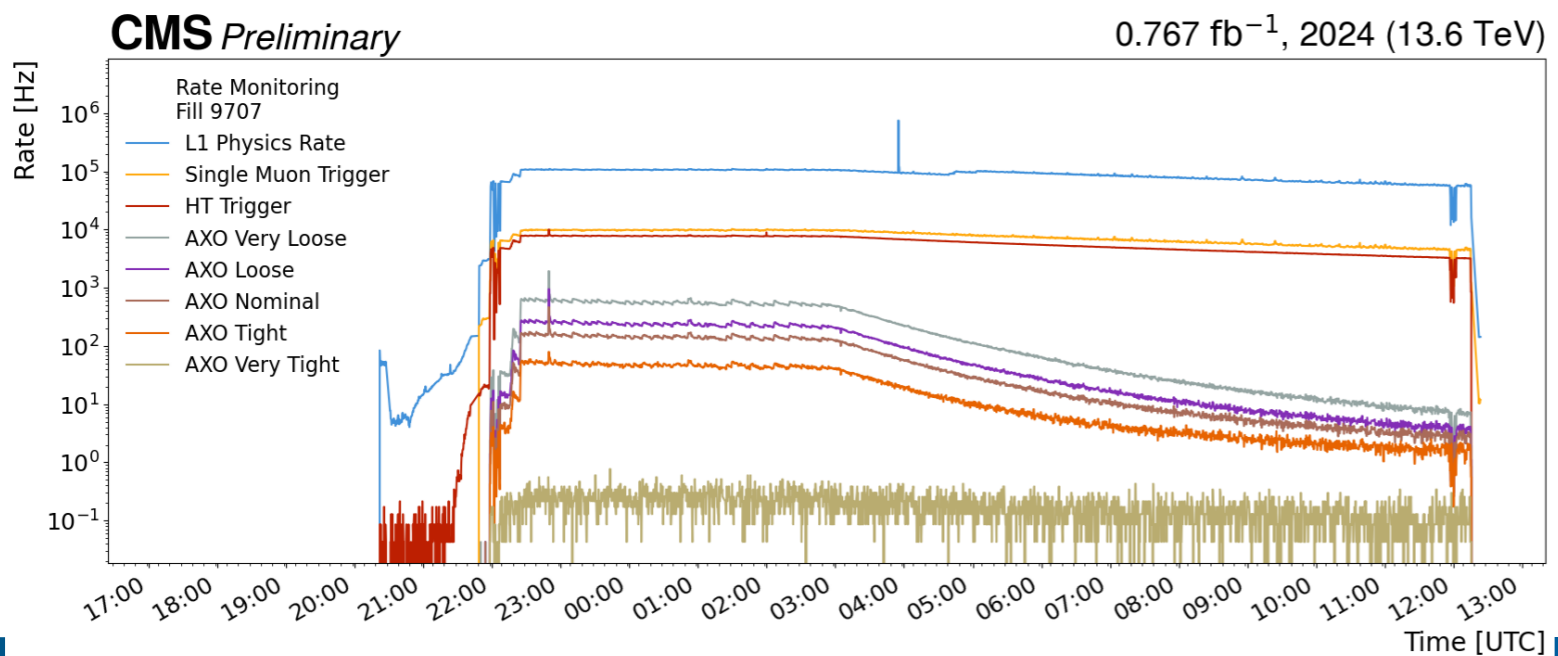


Fitting AXOL1TL on FPGA

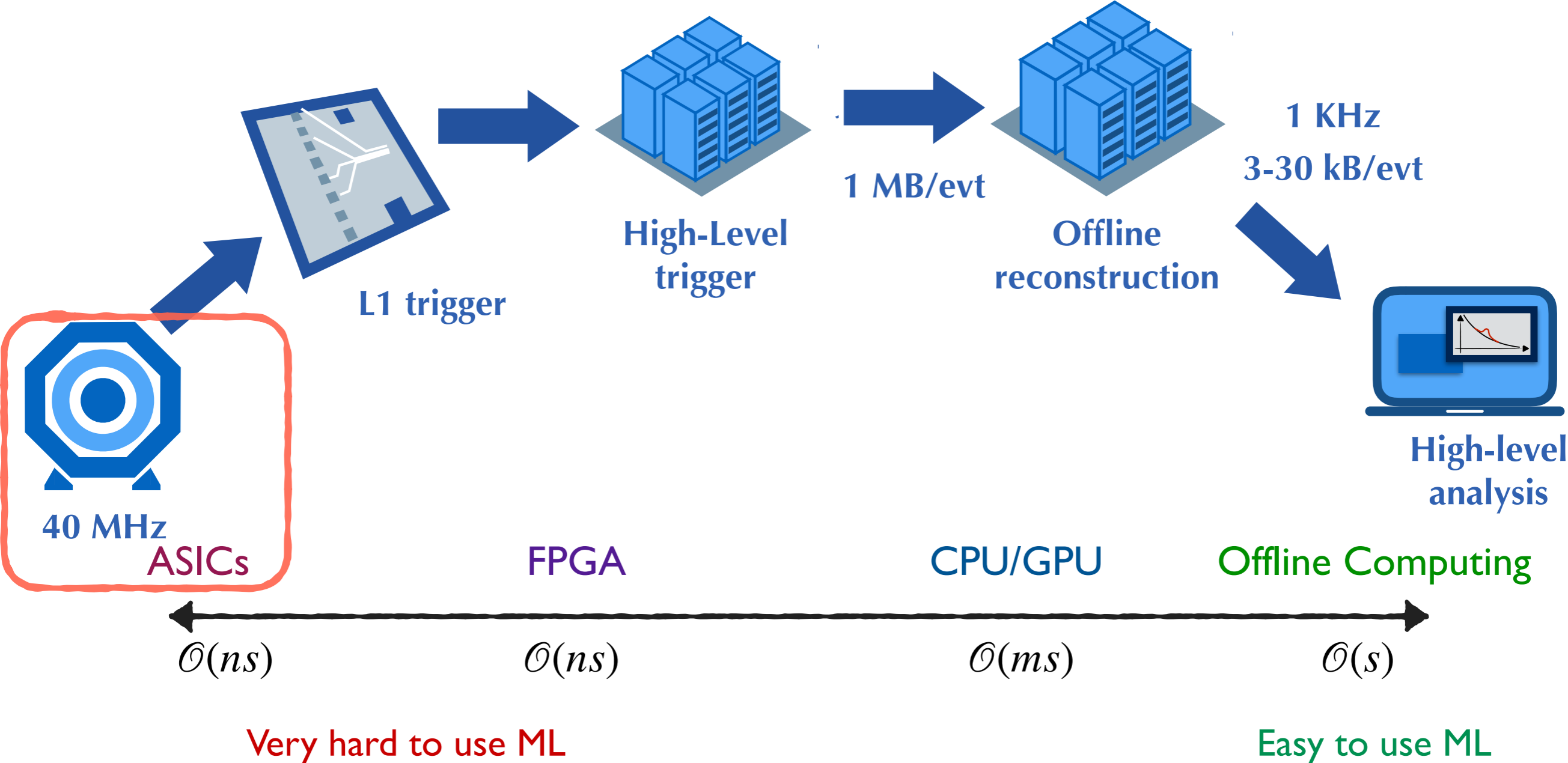
- Instead of vanilla autoencoder, we use Variational Autoencoder
 - With hls4ml to deploy it on a FPGA
 - Inference in just 50 nano seconds!



- Collecting data since 2024 @ LHC !



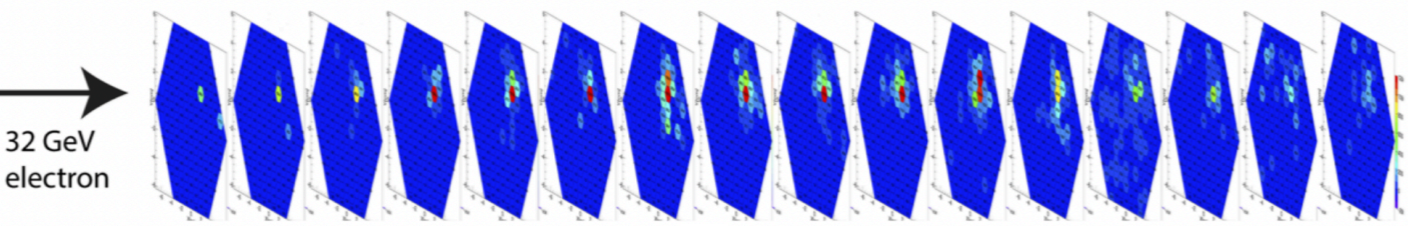
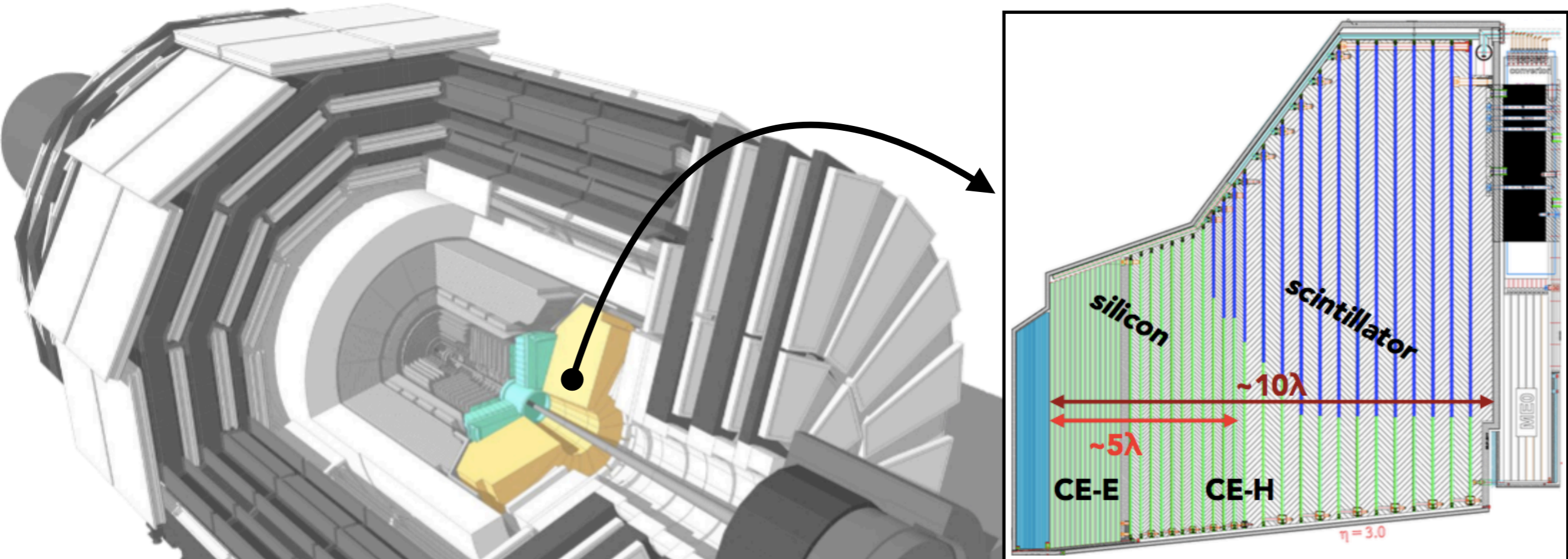
Data reduction workflow @ LHC



ML @ CMS: front-end electronics

The CMS High-Granularity Calorimeter

Novel technology for CMS endcap calorimeter:
50 layers with unprecedented number of readout channels (6M)!

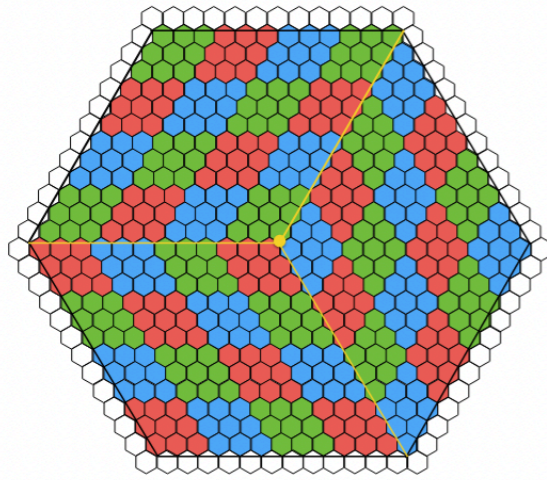


[CMS HGCal TDR](#)

ML @ CMS: front-end electronics

Input

HGCAL 8" hex module



432 silicon sensors → 48 trigger cells (TC) @ 7b per TC
(336b in total)

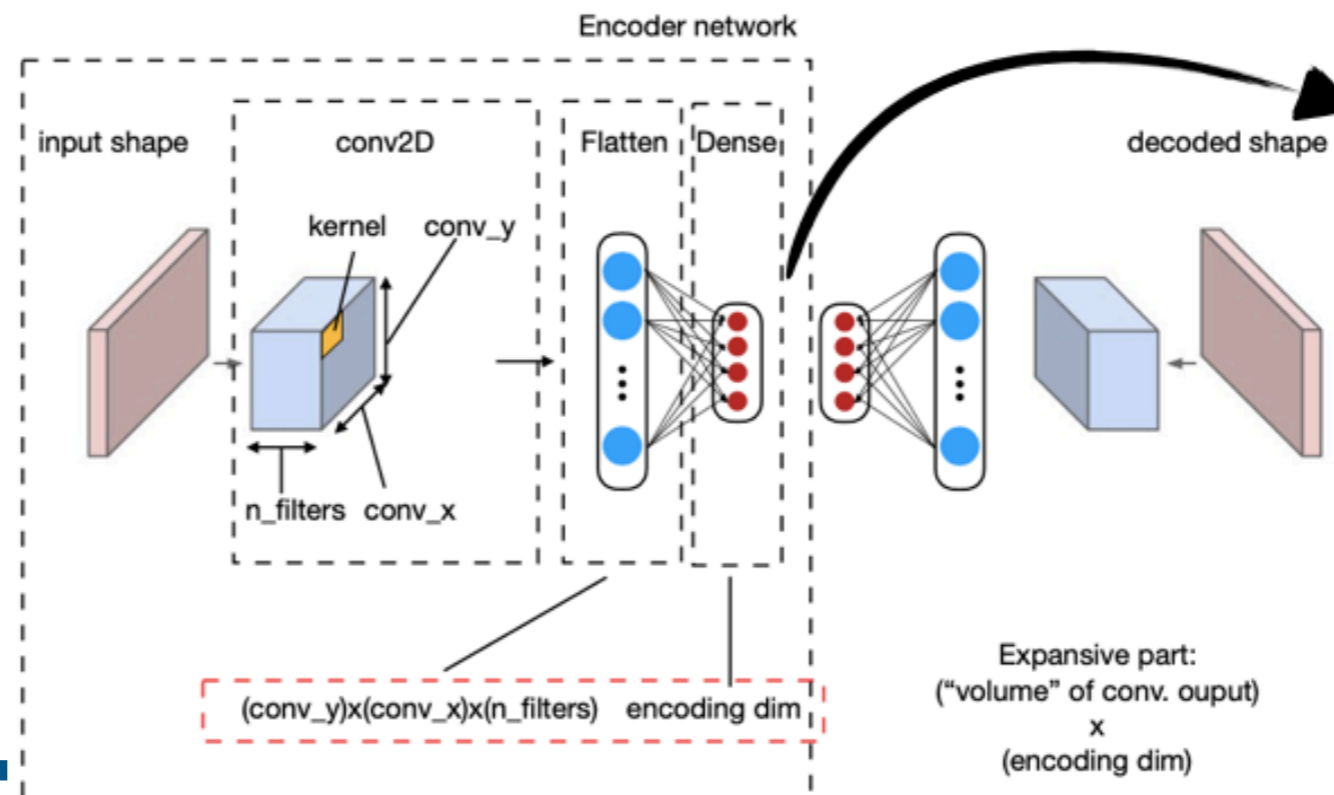
hls4ml used to obtain ASIC ML design!



Output:

“Super trigger cell” algo
3[16 TC sum] x 16–48 bits
= 48–144 bits
(depending on the position)

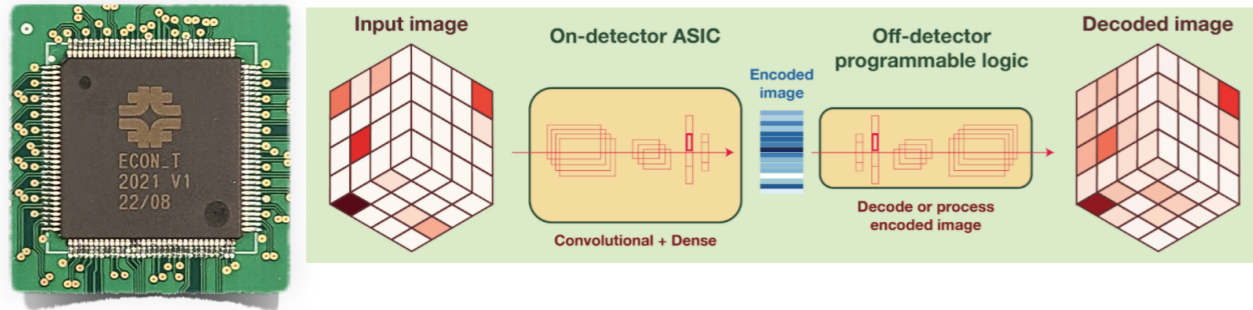
Can we do a better job of encoding the info in those bits w/o so much loss in granularity?



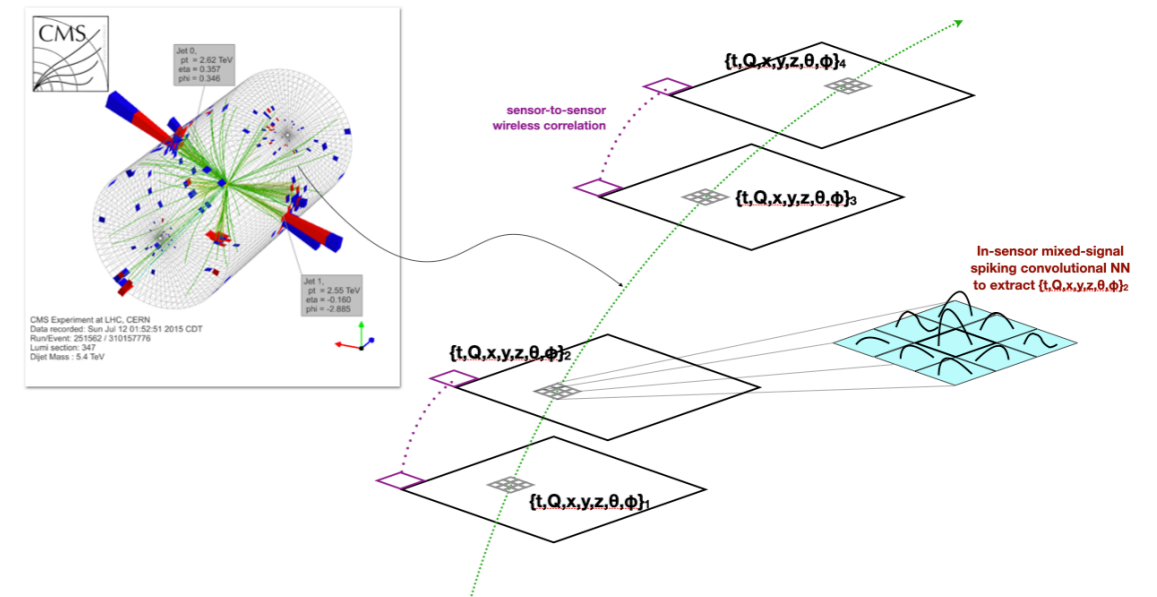
Really need quantized training here to optimize information encoding
Use QKeras!



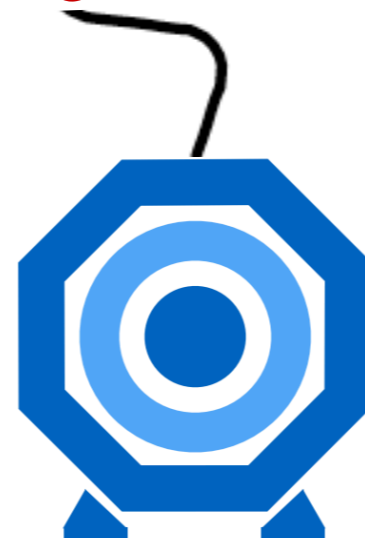
Data compression w/ Radiation hard ASICs



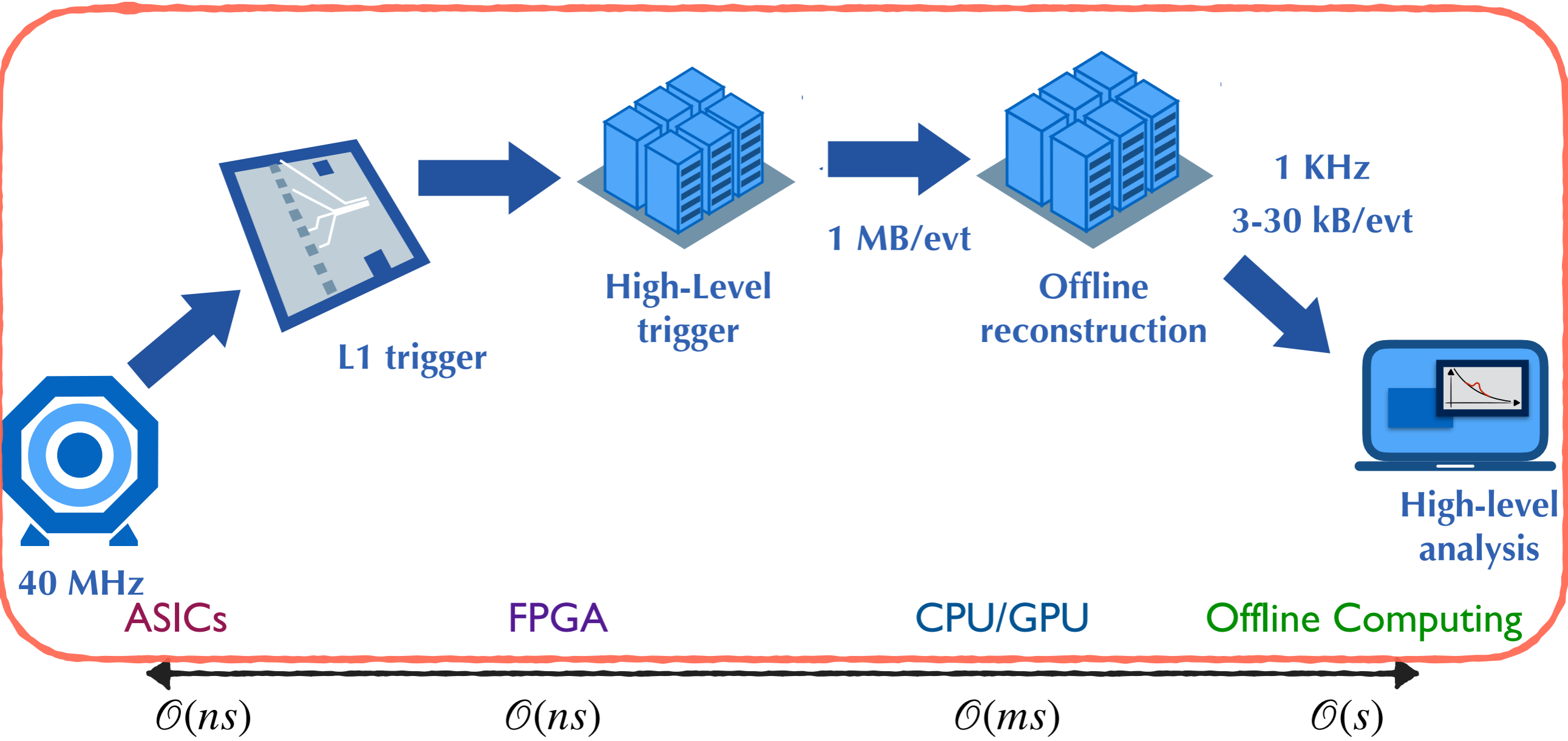
Smart pixels: Pixel sensors w/ AI on chip



Smart sensing w/ On chip AI



Data reduction workflow @ LHC

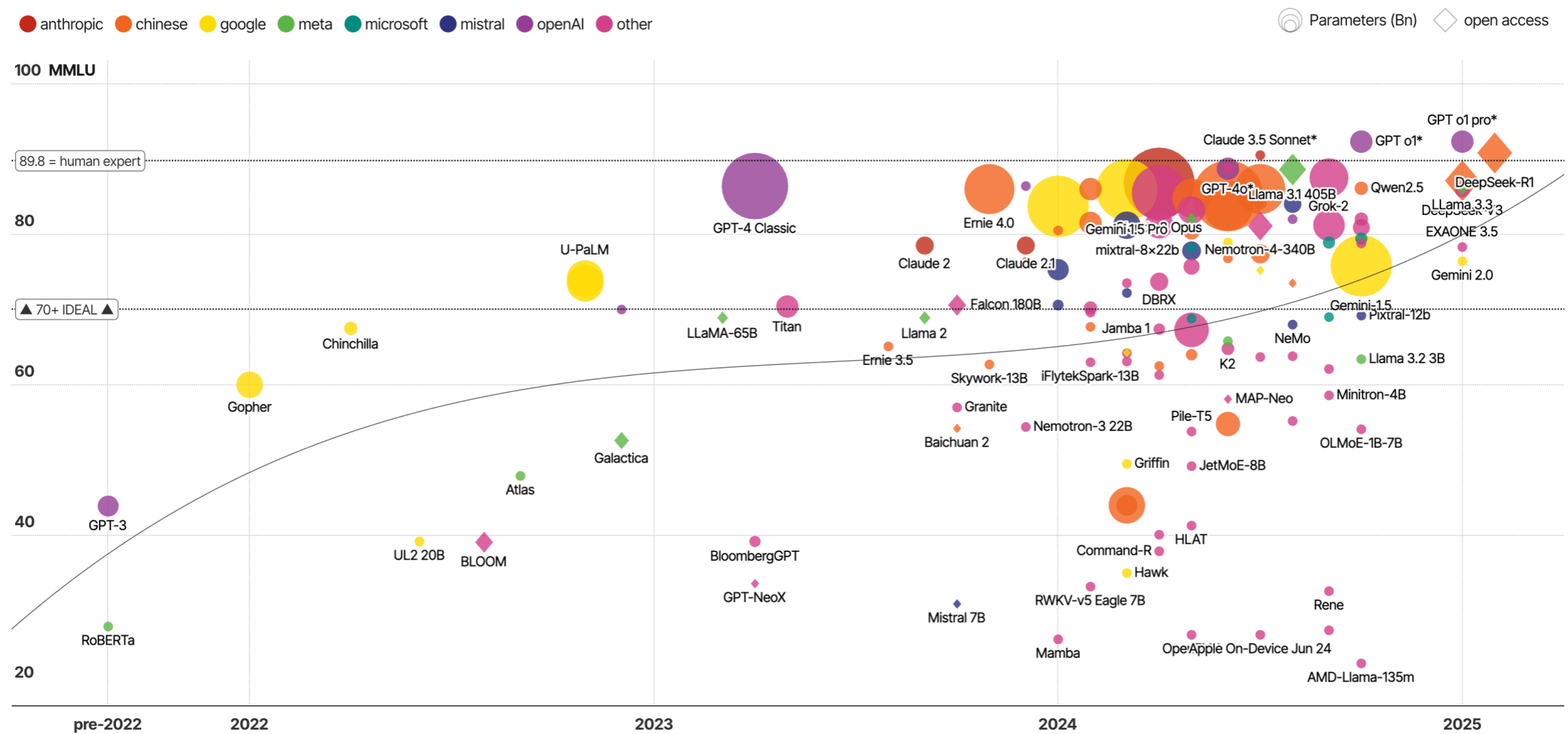


Very hard to use ML

Easy to use ML

The Rise of Large Language Models

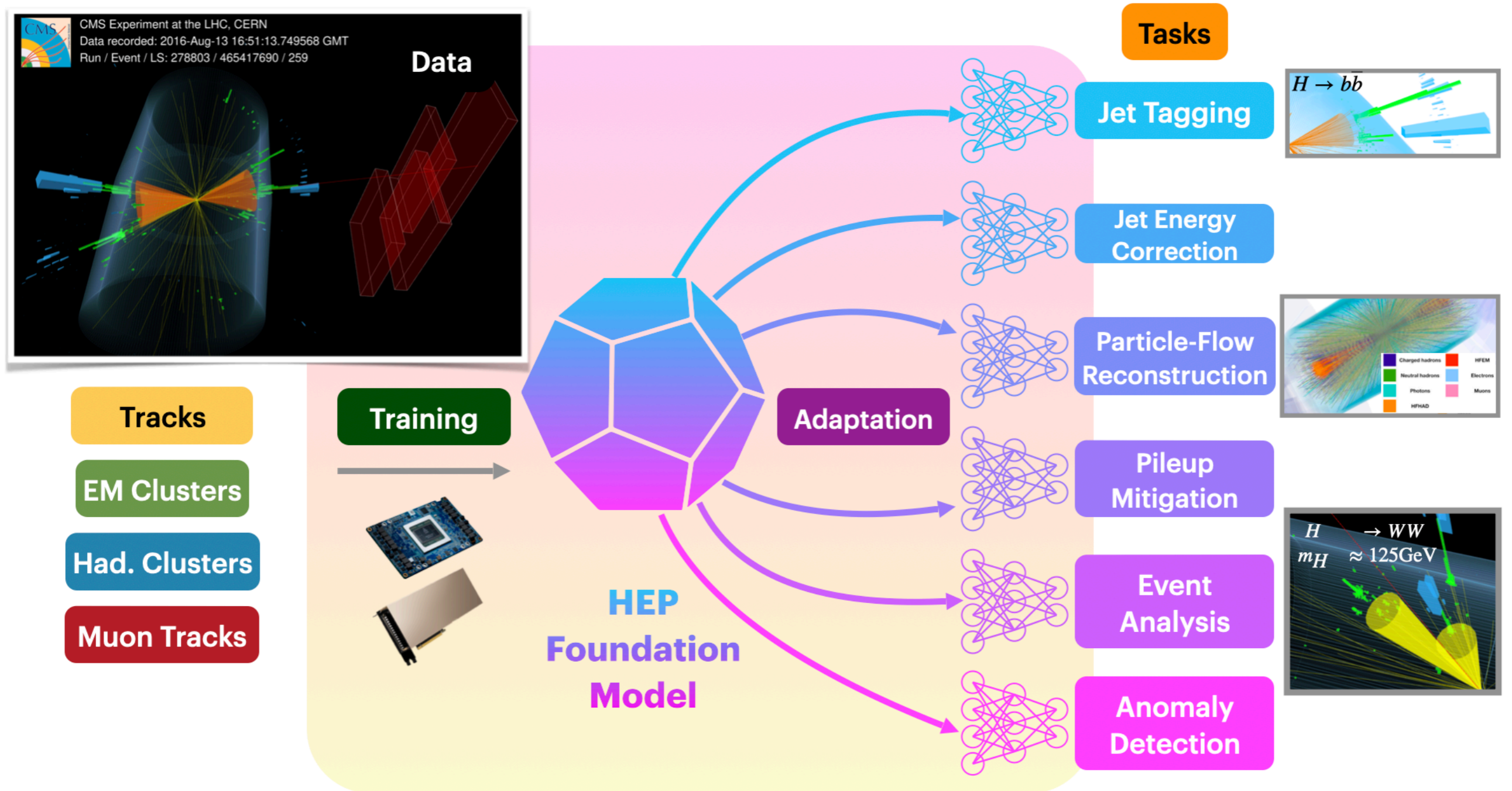
- ▶ Billion-parameter models pre-trained on massive unlabeled datasets can result in emergent capabilities surpassing human expertise in some areas



David McCandless, Tom Evans, Paul Barton
Informationisbeautiful // Jan 2024

<https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-llms-like-chatgpt/>
MMLU = benchmark for measuring LLM capabilities
* = parameters undisclosed // source: Life Architect // data

A Hep Foundation Model?



ML Group

- State-of-the art machine learning methods are increasingly being deployed in CMS
- Many opportunities to integrate/develop proof of concepts methods into the experiment with many challenges ahead to bring the most promising into production
- CMS ML group was created to supervise this effort: cms-conveners-ml@cern.ch

Our goal: enable, support, guide and foster ML developments in computing, POGs, PAGs



L2 Group: Machine Learning



Mia



Pietro

Forum (bi-weekly):

<https://indico.cern.ch/category/12412/>

L3 subgroup meeting (weekly, rotating):

<https://indico.cern.ch/category/12413/>

L3 topical groups

Knowledge



Daniel



Marius

Production



Hyejin

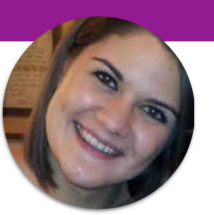


Yao

Innovation



Abhijith



Gaia

+ contacts to POG/PAGs and external initiatives

ML Subgroup: Knowledge

- Knowledge subgroup **collects, maintains, and disseminates knowledge** of ML algorithms in the CMS collaboration
 - Development and maintenance of CMS ML **benchmarks**, comparing and tracking the performance of algorithms, platforms, and ML frameworks on a set of benchmark CMS ML applications in reconstruction, simulation, trigger and computing
 - Consolidate **tutorials** and **lectures** on ML for CMS with the School Committee
 - Maintain list of **in-house experts** in various ML topics.
 - Prepares the ML part of the analysis questionnaire in collaboration with the Statistics Committee and documents good ML practices
- **Documentation** webpage: <https://cms-ml.github.io/documentation/>
 - Contributions welcome: <https://github.com/cms-ml/documentation/pulls>

Knowledge



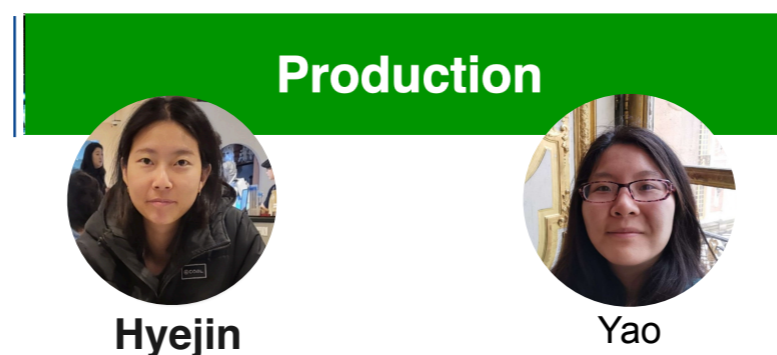
Daniel



Marius

ML Subgroup: Production

- Production subgroup delivers production-level training and inference for CMS ML algorithms
 - Development and maintenance of **ML training and inference workflows** in the CMS software stack, including model inference and support for external frameworks such as TensorFlow, PyTorch, ONNX, and MXNet
 - Work closely with CMS framework experts and all relevant O&C groups, for example, overseeing framework-related aspects and relevant software/computing groups
 - Development of training infrastructure to satisfy the needs of as many collaborators as possible



ML Frameworks

- **Keras/TensorFlow**

- Conventionally one of the most beginner-friendly formats
- Inference supported for production in CMSSW

- **PyTorch (Lightning)**

- More “pythonic” way of building models (especially models within models)
- Inference support in CMSSW forthcoming...

- **ONNX/ONNXRuntime**

- ONNX intended to be a universal exchange format (can convert from all libraries)
- For inference only; supported for CMSSW

- **Flax/JAX**

- Newest library based on “autograd”

- **XGBoost**

- For training boosted decision trees (fast)



ML Subgroup: Innovation

- Innovation subgroup identifies and applies new ML techniques to CMS challenges while working closely with CMS groups in areas where ML is expected to have a significant impact: reconstruction, trigger, simulation, and beyond
 - Discuss the relevance of new techniques and help with the adaptation and implementation of specific models
 - Develop specific methods for CMS that will lead to technical publications
 - Lead the organization of ML-oriented **hackathons** and **challenges** to help identify new applications and ML techniques in CMS
 - e.g. Upcoming hackathons in LLMs in CMS
 - Organize CMS internal **journal club** to discuss new and relevant ML results from inside and outside of particle physics

Innovation



Abhijith



Gaia

CMS ML Documentation



CMS Machine Learning Documentation

[Home](#)

Innovation

[ML Journal Club](#)

[ML Hackathons](#)

Resources

[Cloud Resources](#)

[Dataset Resources](#)

[FPGA Resource](#)

[GPU Resources](#)



[Ixplus-gpu](#)

[CERN HTCondor](#)

[SWAN](#)

[ml.cern.ch](#)

Guides

[Software environments](#)



[Optimization](#)



[General Advice](#)



Welcome to the documentation hub for the CMS Machine Learning Group! The goal of this page is to provide CMS analyzers a centralized place to gather machine learning information relevant to their work. However, we are not seeking to rewrite external documentation. Whenever applicable, we will link to external documentation, such as the iML groups [HEP Living Review](#) or their [ML Resources](#) repository. What you will find here are pages covering:

- ML best practices
- How to optimize a NN
- Common pitfalls for CMS analyzers
- Direct and indirect inferencing using a variety of ML packages
- How to get a model integrated into CMSSW

And much more!

If you think we are missing some important information, please contact the [ML Knowledge Subgroup](#)!

Last update: December 5, 2023

Summary

- Lots of R&D surrounding ML in CMS, which directly impact physics results
- Many activities coordinated by ML Group
 - We're here to help!
- Get involved!
 - Attend ML Forum (<https://indico.cern.ch/category/12412/>) and subgroup meetings (<https://indico.cern.ch/category/12413/>)
 - Contact conveners (cms-conveners-ml@cern.ch) or subgroup conveners for a slot to present your ML studies
 - Consider/discuss publishing your ML results in CMS
 - Release CMS datasets for ML studies through knowledge subgroup
 - Contribute to ML documentation (<https://cms-ml.github.io/documentation/>) by submitting pull requests or reporting issues
 - And more...

