# Accelerating Data Collection and Processing at the Large Hadron Collider

## Benjamin Nachman

*Lawrence Berkeley National Laboratory and the Simons Institute*
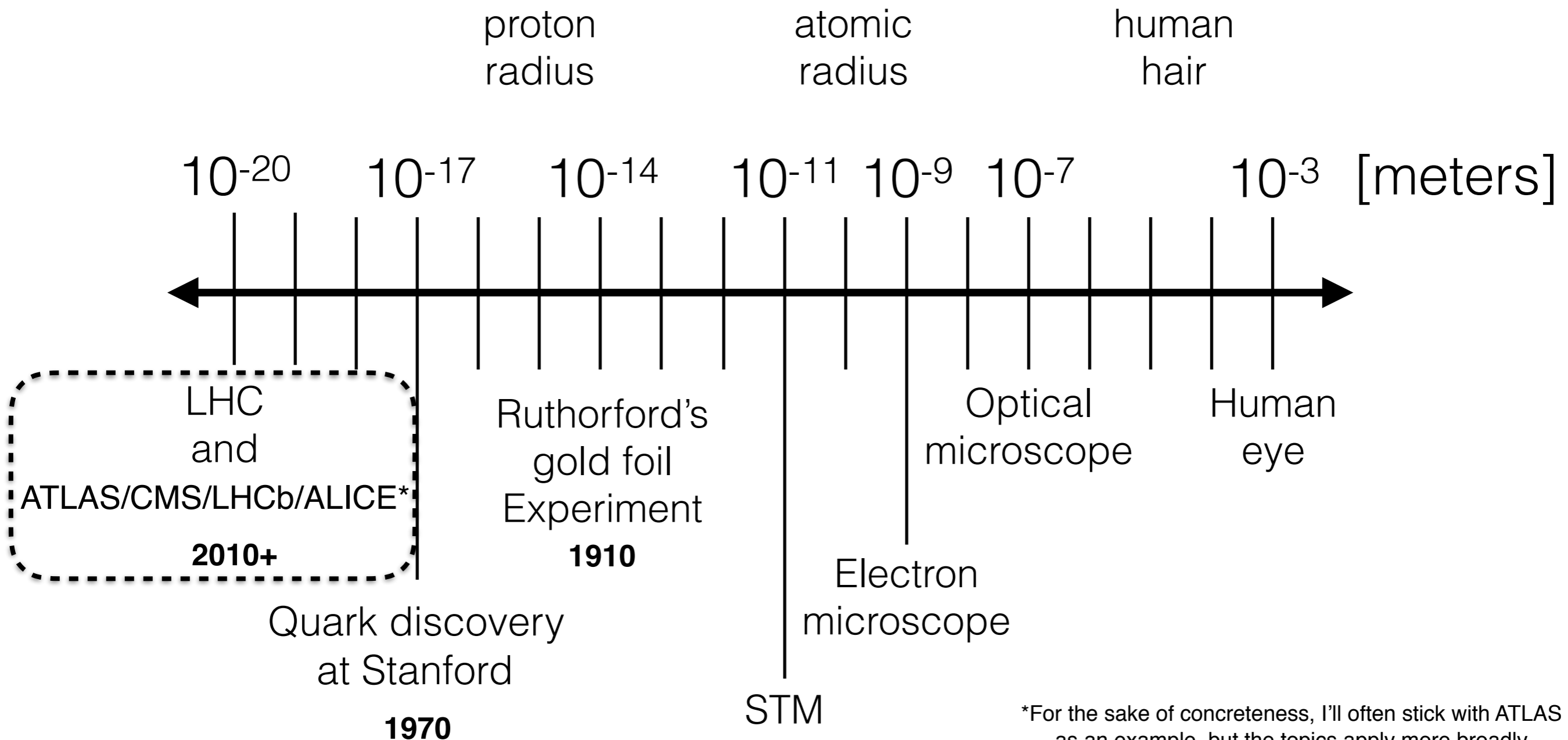
**Benjamin Nachman**

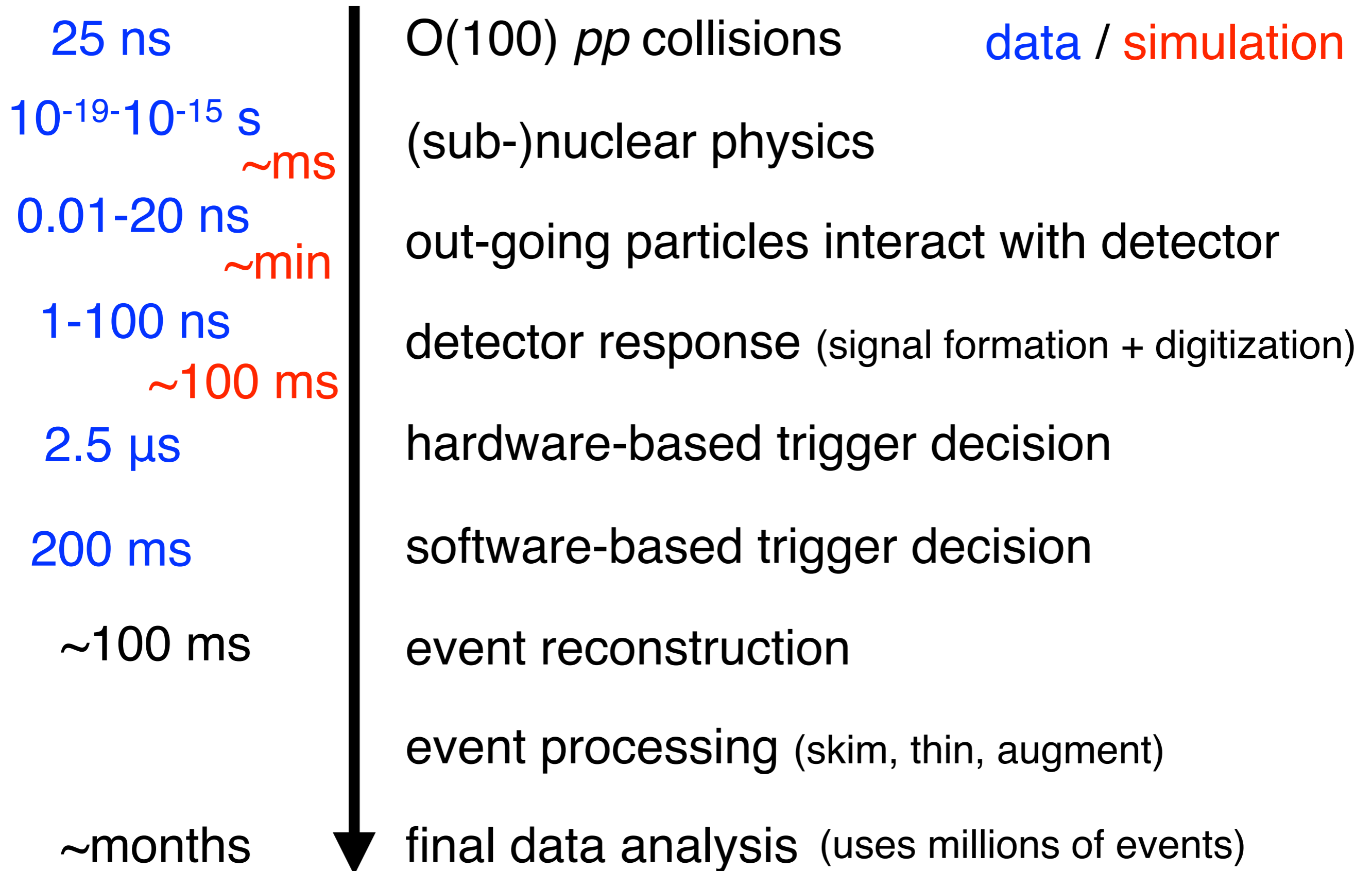*Lawrence Berkeley National Laboratory and the Simons Institute*

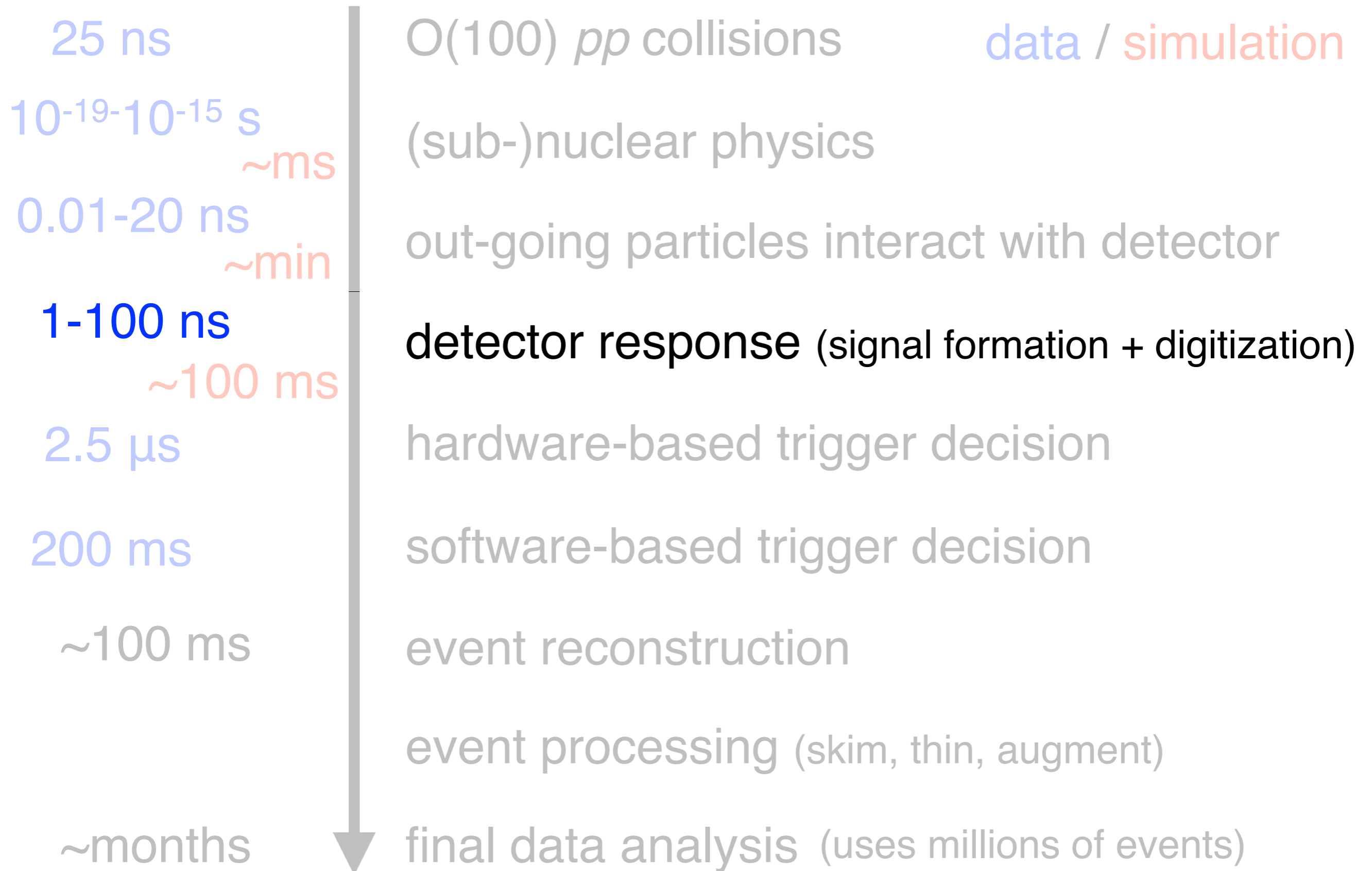*Real Time Decision Making for Applications in the Natural Sciences and Physical Systems, Feb. 28, 2018*

Goal: **We want to study the structure of the smallest building blocks of matter.** For this, we need the most powerful microscope ever built!
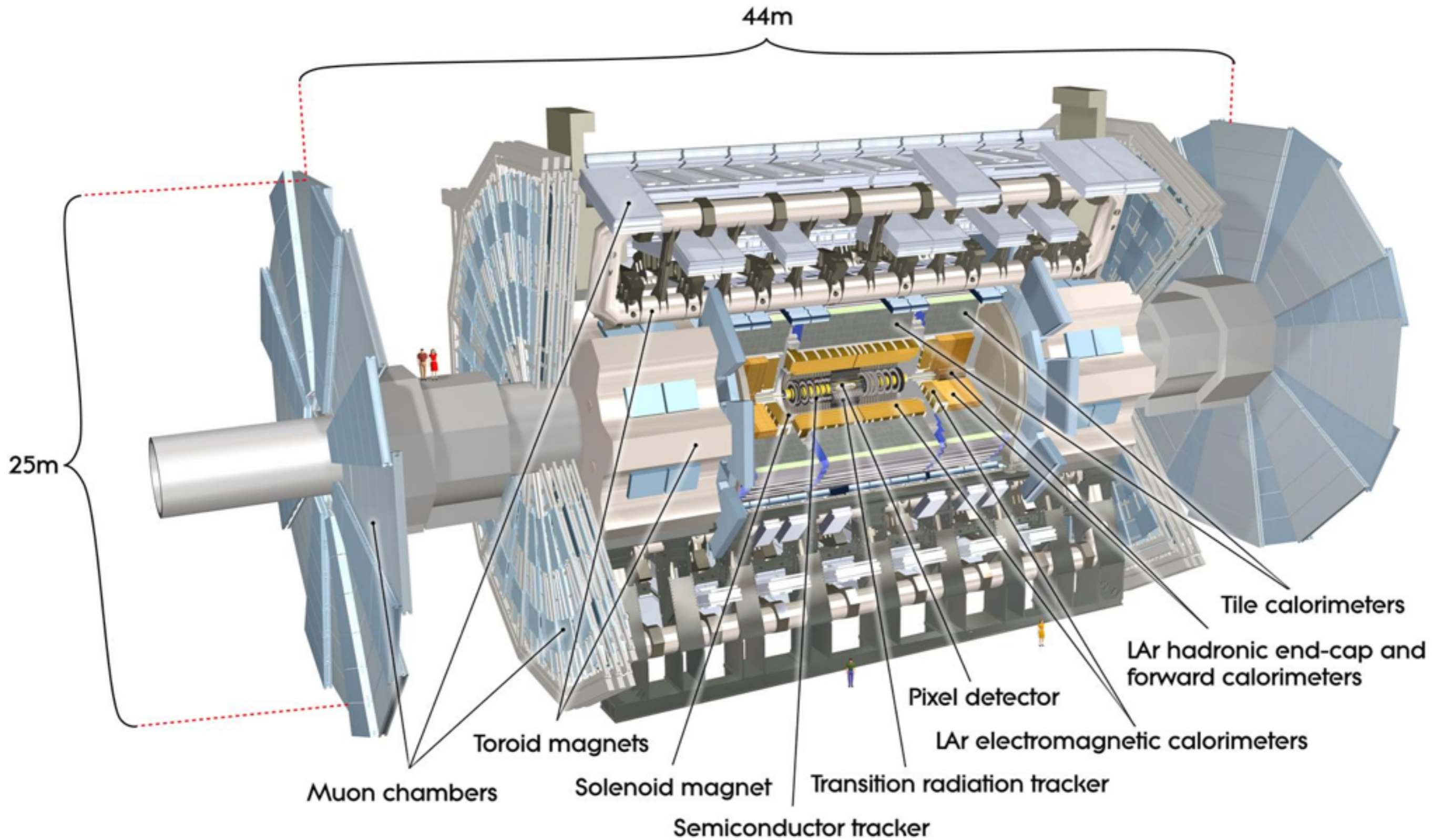
proton radius

atomic radius

human hair

$10^{-20}$    $10^{-17}$    $10^{-14}$    $10^{-11}$   $10^{-9}$   $10^{-7}$     $10^{-3}$   [meters]

LHC and ATLAS/CMS/LHCb/ALICE*

**2010+**

Ruthorford's gold foil Experiment

**1910**

Optical microscope

Human eye

Quark discovery at Stanford

**1970**

Electron microscope

STM

*For the sake of concreteness, I'll often stick with ATLAS as an example, but the topics apply more broadly
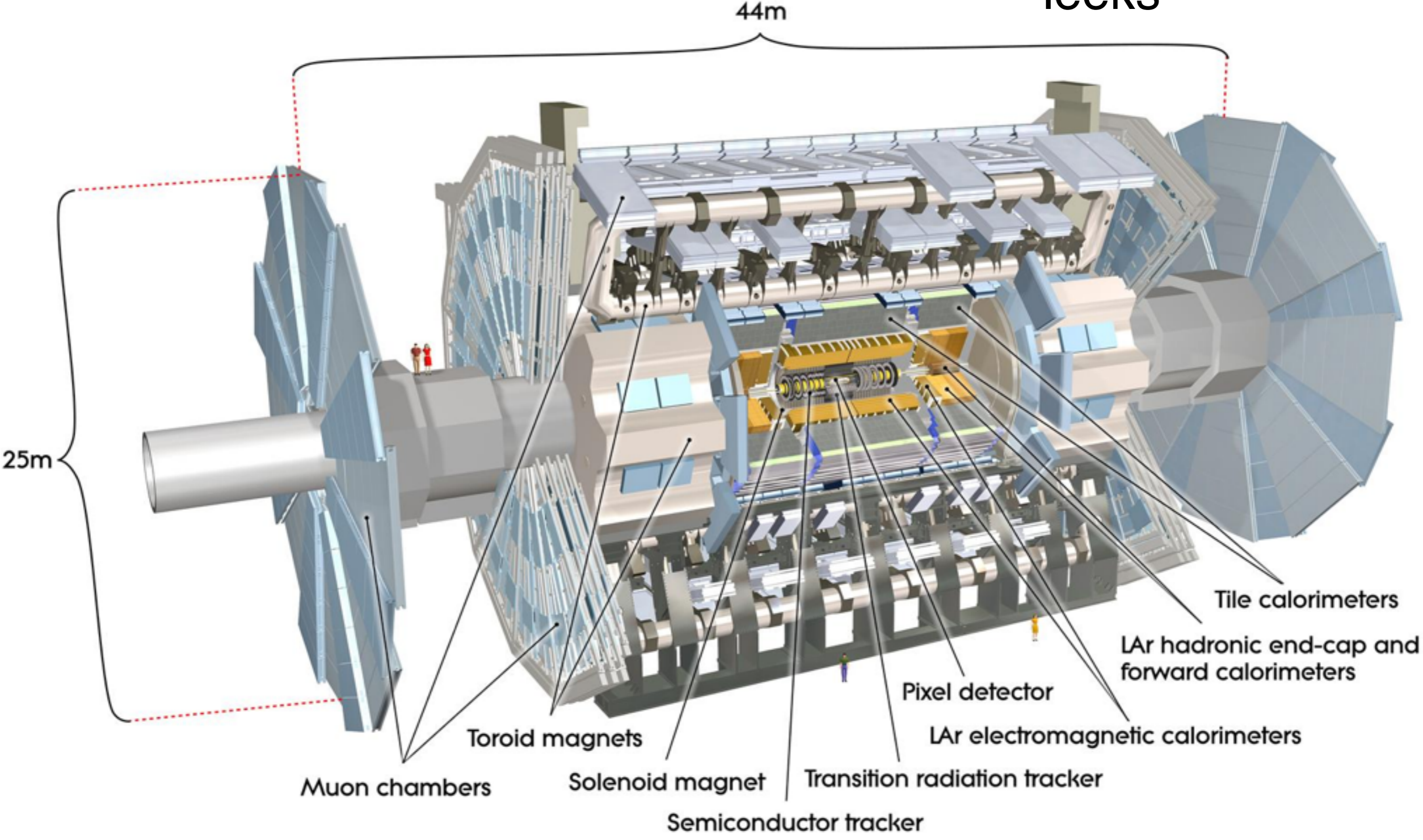
| | data / simulation |
|---|---|
| 25 ns | O(100) *pp* collisions |
| $10^{-19}$-$10^{-15}$ s<br>~ms | (sub-)nuclear physics |
| 0.01-20 ns<br>~min | out-going particles interact with detector |
| 1-100 ns<br>~100 ms | detector response (signal formation + digitization) |
| 2.5 μs | hardware-based trigger decision |
| 200 ms | software-based trigger decision |
| ~100 ms | event reconstruction |
| | event processing (skim, thin, augment) |
| ~months | final data analysis (uses millions of events) |

| | |
|---|---|
| 25 ns | O(100) *pp* collisions    data / simulation |
| $10^{-19}$-$10^{-15}$ s | (sub-)nuclear physics |
| ~ms | |
| 0.01-20 ns | out-going particles interact with detector |
| ~min | |
| **1-100 ns** | **detector response** (signal formation + digitization) |
| ~100 ms | |
| 2.5 µs | hardware-based trigger decision |
| 200 ms | software-based trigger decision |
| ~100 ms | event reconstruction |
| | event processing (skim, thin, augment) |
| ~months | final data analysis (uses millions of events) |

# Collider-based HEP detectors are like onions



44m

25m

Tile calorimeters

LAr hadronic end-cap and forward calorimeters

Pixel detector

LAr electromagnetic calorimeters

Transition radiation tracker

Semiconductor tracker

Solenoid magnet

Toroid magnets

Muon chambers

Biggest challenge for data volume is **innermost** layer

# Collider-based HEP detectors are like ~~onions~~ leeks



44m

25m

Tile calorimeters

LAr hadronic end-cap and forward calorimeters

Pixel detector

LAr electromagnetic calorimeters

Transition radiation tracker

Semiconductor tracker

Solenoid magnet

Toroid magnets

Muon chambers

# Biggest challenge for data volume is **innermost** layer

ATLAS EXPERIMENT

HL-LHC $t\bar{t}$ event in ATLAS ITK at $<\mu>=200$

**GHz/cm$^2$   ~0.1%/pixel/BC**

**Gbps/cm$^2$ ~streaming live audio from each pixel**

**1 Grad (TID) and $10^{16}$ $n_{eq}$/cm$^2$ (NIEL)**

| Generation | Run 1 (FEI3, PSI46) | Runs 2+3 (FEI4, PSI46DIG) | Runs 4+5 |
|---|---|---|---|
| Chip Size | 7.5 x 10.5 mm$^2$ <br> 8 x 10 mm$^2$ | 20 x 20 mm$^2$ <br> 8 x 10 mm$^2$ | > 20 x 20 mm$^2$ |
| Transistors | 3.5 M <br> 1.3 M | 87 M | ~1 G |
| Hit Rate | 100 MHz/cm$^2$ | 400 MHz/cm$^2$ | ~2 GHz/cm$^2$ |
| Hit Memory / Chip | 0.1 Mb | 1 Mb | ~16 Mb |
| Trigger Rate | 100 kHz | 100 kHz | 200 kHz - 1MHz |
| Trigger Latency | 2.5 µs <br> 3.2 µs | 2.5 µs <br> 3.2 µs | 6 - 20 µs |
| Readout rate | 40 Mb/s | 320 Mb/s | 1-4 Gb/s |
| Radiation | 100 Mrad | 200 Mrad | 1 Grad |
| Technology | 250 nm | 130 nm <br> 250 nm | 65 nm |
| Power | ~1/4 W/cm$^2$ | ~1/4 W/cm$^2$ | 1/2 - 1 W/cm$^2$ |

| Generation | Run 1 (FEI3, PSI46) | Runs 2+3 (FEI4, PSI46DIG) | Runs 4+5 |
|---|---|---|---|
| Chip Size | 7.5 x 10.5 mm² / 8 x 10 mm² | 20 x 20 mm² / 8 x 10 mm² | > 20 x 20 mm² |
| Transistors | 3.5 M / 1.3 M | 87 M | ≥1 G |
| Hit Rate | 100 MHz/cm² | 400 MHz/cm² | ~2 GHz/cm² |
| Hit Memory / Chip | 0.1 Mb | 1 Mb | ~16 Mb |
| Trigger Rate | 100 kHz | 100 kHz | 200 kHz - 1MHz |
| Trigger Latency | 2.5 µs / 3.2 µs | 2.5 µs / 3.2 µs | 6 - 20 µs |
| Readout rate | 40 Mb/s | 320 Mb/s | 1-4 Gb/s |
| Radiation | 100 Mrad | 200 Mrad | 1 Grad |
| Technology | 250 nm | 130 nm / 250 nm | 65 nm |
| Power | ~1/4 W/cm² | ~1/4 W/cm² | 1/2 - 1 W/cm² |

e.g. the camera in your phone on steroids, next to a nuclear reactor (unfortunately, Apple doesn't make one of these)

A significant component of the design and testing is happening right here at Berkeley!

11.8 mm ; 192 pixels

50 x 50 μm² pixels

20 mm ; 400 pixels

RD53 co-spokesperson Maurice Garcia-Sciveres + many others at Berkeley

Readout regions N x M pixel regions; helps to recover small charge hits. What is optimal?

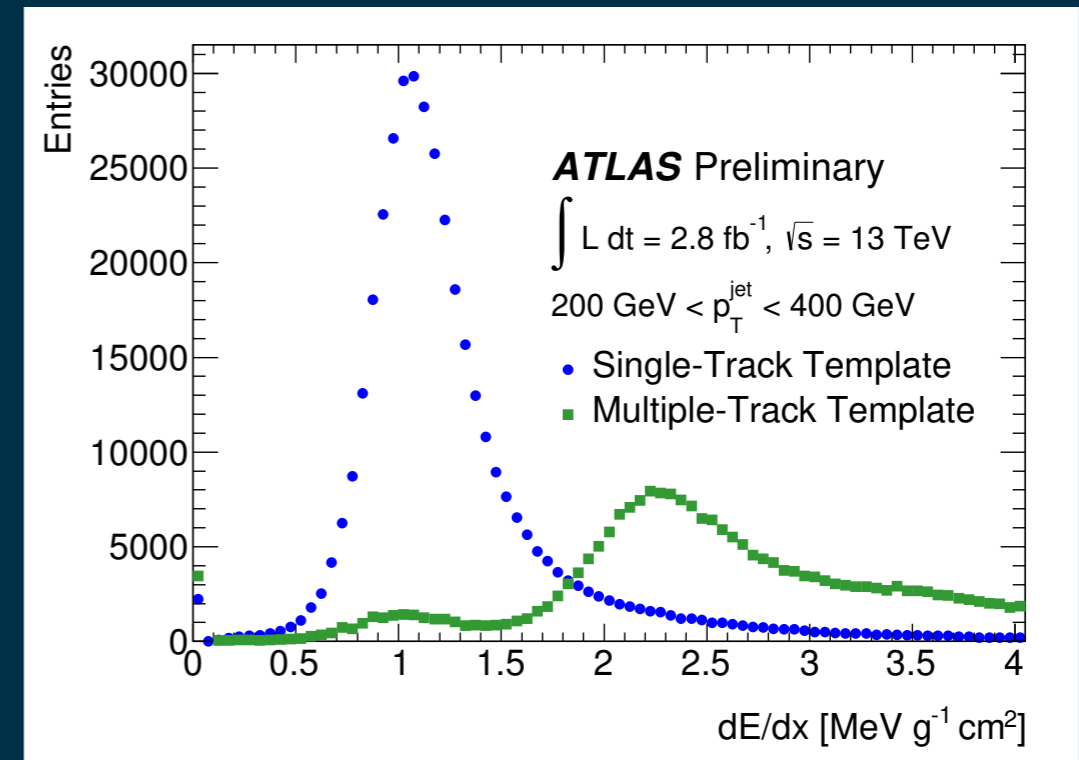This is just one of many choices we are currently studying!

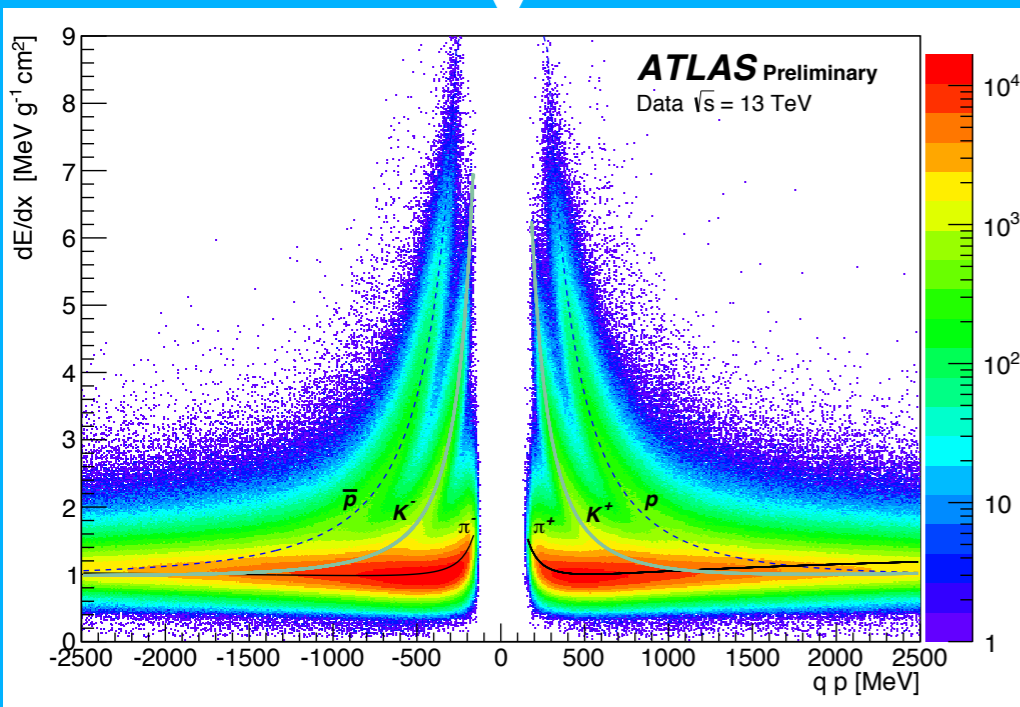Geant4 (Allpix) + Digitization (Truth) + Clustering (CCA)

50 x 50 x 150 µm³

- 2×2
- 4×1

**4 x 1** beats **2 x 2**

Mean Occupancy

η

ATLAS Simulation
$\sqrt{s}=7$ TeV

CCA Clustering
NN Clustering
3-pixel wide clusters

non-linear ToT interpolation

linear ToT interpolation

Arbitrary Units — Local x resolution [μm]

Position estimation

Particle identification

Particle multiplicity

Offline, the ToT is used for many purposes; however if not saved properly, performance may suffer!

ATL-PHYS-PUB-2016-007

ATLAS Preliminary
$\int L \, dt = 2.8 \text{ fb}^{-1}$, $\sqrt{s} = 13$ TeV
200 GeV $< p_T^{jet} <$ 400 GeV
• Single-Track Template
■ Multiple-Track Template

dE/dx [MeV g$^{-1}$ cm$^2$]

ATLAS Preliminary
Data $\sqrt{s} = 13$ TeV

$\bar{p}$  $K^-$  $\pi^-$  $\pi^+$  $K^+$  $p$

dE/dx [MeV g$^{-1}$ cm$^2$] — q p [MeV]

ATLAS-PIX-2015-002

Currently, ATLAS uses 4/8 bit ToT with a linear charge to ToT conversion.

ATLAS Simulation
$\sqrt{s}=7$ TeV

CCA Clustering
NN Clustering
3-pixel wide clusters

non-linear ToT interpolation

linear ToT interpolation

Offline, the ToT is used for many purposes; however if not saved properly, performance may suffer!

Position estimation

Particle identification

ATLAS Preliminary

$L_{int} = 2.8$ fb$^{-1}$, $\sqrt{s} = 13$ TeV

200 GeV < $p_T^{jet}$ < 400 GeV

Single-Track Template
Multiple-Track Template

dE/dx [MeV g$^{-1}$ cm$^2$]

ATL-PHYS-PUB-2016-007

ATLAS-PIX-2015-002

ATLAS Preliminary
Data $\sqrt{s}=13$ TeV

dE/dx [MeV g$^{-1}$ cm$^2$]

$\bar{p}$   $K^-$   $\pi^-$   $\pi^+$   $K^+$   $p$

q p [MeV]

**Question: what is the "best" way to store/utilize charge?**

Currently, ATLAS uses 4/8 bit ToT with a linear charge to ToT conversion.

Geant4 (Allpix) + Digitization (δ-ray veto), threshold = 600e
50 X 50 x 150 μm³, radius = 39 mm, |η| = 1

position resolution

worse

— N bit → N bit

···· 5 bit → N bit

**maximum achievable reduction**

Number of ToT bits

Can add digital logic so that N digitized bits are stored as M ≤ N bits.

There are $\binom{2^N - 2}{2^M - 2}$

possible functions mapping N to M bits.

We could gain space without much loss in performance
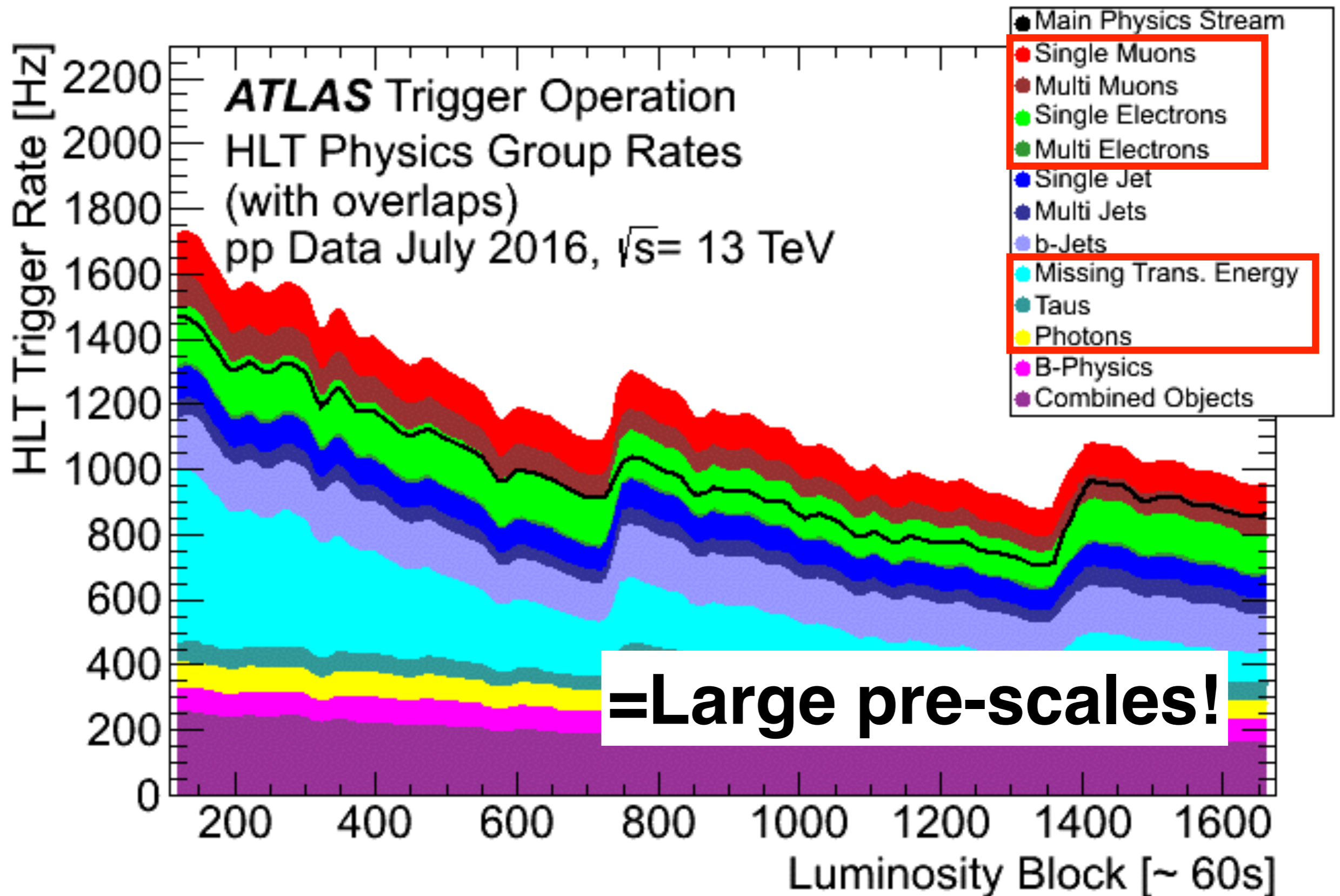
+ many more considerations for readout design!

| | |
|---|---|
| 25 ns | O(100) *pp* collisions          data / simulation |
| $10^{-19}$-$10^{-15}$ s | |
| ~ms | (sub-)nuclear physics |
| 0.01-20 ns | |
| ~min | out-going particles interact with detector |
| 1-100 ns | |
| ~100 ms | detector response (signal formation + digitization) |
| **2.5 μs** | **hardware-based trigger decision** |
| 200 ms | software-based trigger decision |
| ~100 ms | event reconstruction |
| | event processing (skim, thin, augment) |
| ~months | final data analysis (uses millions of events) |

# Cross-section is dominated by final states with no electrons, no muons, and no neutrinos



**Standard Model Production Cross Section Measurements** — Status: August 2016

ATLAS Preliminary
Run 1,2  $\sqrt{s}$ = 7, 8, 13 TeV

# Recorded events are dominated by final states with electrons, muons, neutrinos



**ATLAS** Trigger Operation
HLT Physics Group Rates
(with overlaps)
pp Data July 2016, $\sqrt{s}$= 13 TeV

Legend:
- Main Physics Stream
- Single Muons
- Multi Muons
- Single Electrons
- Multi Electrons
- Single Jet
- Multi Jets
- b-Jets
- Missing Trans. Energy
- Taus
- Photons
- B-Physics
- Combined Objects

**=Large pre-scales!**

HLT Trigger Rate [Hz]

Luminosity Block [~ 60s]

HLT = software trigger (L1 = hardware trigger)

**usual paradigm**

*If your favorite process cannot be triggered on, then it will **not be recorded** and there is **no way to analyze it**.*

**new paradigms**

*"If your favorite process cannot be triggered on inclusively, look for **associated production** with an object you can trigger."*

-I.S. Radiation, 2015

*"If your favorite process cannot be triggered on at HLT, make your analysis **faster and simpler** and do it after L1."*

-T.L. Analysis, 2012

(see Caterina's talk)

# usual paradigm

*If your favorite process cannot be triggered on, then it will **not be recorded** and there is **no way to analyze it**.*

# new paradigms

*"If your favorite process cannot ~~be triggered~~ ely, look for associated ~~production~~ you can trigger."*

Caveat: **large effective prescale** from the low cross-section from the associated production

-I.S. Radiation, 2015

*"If your favorite process cannot be triggered, make your analysis faster and do it after L1."*

Caveat: limited by the L1 rate (usually **harsher than HLT**)

-T.L. Analysis, 2012

ATLAS
EXPERIMENT

Run Number: 266904, Event Number: 25884352

Date: 2015-06-03 13:41:54 CEST

**trigger on pink**

**your favorite process could be in any of these ones**

L1 (TLA) or HLT (offline) rate

$$\left( \begin{array}{c} \text{"Zero bias} \\ \text{trigger" data} \end{array} \right) = \frac{h}{H} \times \left( \begin{array}{c} \text{Triggered} \\ \text{data amount} \end{array} \right)$$

40 MHz

| # *pp* interactions | L1 Rate | HLT Rate | ZBT | ZBT @HLT |
|---|---|---|---|---|
| 20 | 100 kHz | 100 Hz | $4 \times 10^5$ | 400 |
| 80 (~now) | 100 kHz | 1 kHz | $4 \times 10^4$ | 400 |
| 200 | 400 kHz | 10 kHz | $4 \times 10^3$ | 100 |

If you run a zero-background search and can't beat a trigger efficiency of ~0.02%, then you should be using the ZBT!

…and if you can do TLA, that number goes up to ~1%!

| # *pp* interactions | L1 Rate | HLT Rate | ZBT | ZBT @HLT |
|:---:|:---:|:---:|:---:|:---:|
| 20 | 100 kHz | 100 Hz | $4 \times 10^5$ | 400 |
| 80 (~now) | 100 kHz | 1 kHz | $4 \times 10^4$ | 400 |
| 200 | 400 kHz | 10 kHz | $4 \times 10^3$ | 100 |

Cross-sections with Leptophobic Z' in MG5

Effective Prescale

ZBT offline

ZBT @ HLT

~unprescaled photon trigger

~unprescaled quark/gluon trigger

Object $p_T$ [GeV]

Associated
photon
production

Associated
quark/gluon
production

There is a huge dataset that we are currently ignoring.

New physics may be hiding in these data
and we are collecting them anyway

Most powerful when combined with trigger-
level analysis (**so need to design ASAP**!)

**Takeaway message: the baseline is the ZBT >> 0!**

Think creatively about new possibilities…
the sky, and not the trigger, is the limit!

(also, remember ZBT offline has ~infinite time for processing)

| | |
|---|---|
| 25 ns | O(100) *pp* collisions      data / simulation |
| $10^{-19}\text{-}10^{-15}$ s | |
| ~ms | (sub-)nuclear physics |
| 0.01-20 ns | |
| ~min | out-going particles interact with detector |
| 1-100 ns | |
| ~100 ms | detector response (signal formation + digitization) |
| 2.5 µs | hardware-based trigger decision |
| **200 ms** | **software-based trigger decision** |
| **~100 ms** | **event reconstruction** |
| | event processing (skim, thin, augment) |
| ~months | final data analysis (uses millions of events) |

The extra pileup collisions add unwanted soft radiation on top of the event

ATLAS
EXPERIMENT

HL-LHC t$\bar{t}$ event in ATLAS ITK
at <μ>=200

This degrades trigger and offline performance

akin to image de-noising
… we can use ML for that!

# Pileup mitigation with machine learning



Leading vertex charged

Pileup charged

Total neutral

$\eta$

$\phi$

beam

Inputs to NN

10 filters $\times 2$

Leading vertex neutral

*Strange noise because we can measure ~2/3 of it (charged pileup)*

*…also a natural application of convolutional NNs!*

# Pileup mitigation with machine learning

"Pileup Mitigation with Machine Learning"

P. Komiske, E. Metodiev, **BPN**, and M. Schwartz, JHEP 12 (2017) 051

| | |
|---|---|
| 25 ns | O(100) *pp* collisions          data / simulation |
| $10^{-19}$-$10^{-15}$ s | |
| ~ms | (sub-)nuclear physics |
| 0.01-20 ns | |
| ~min | **out-going particles interact with detector** |
| 1-100 ns | |
| ~100 ms | detector response (signal formation + digitization) |
| 2.5 µs | hardware-based trigger decision |
| 200 ms | software-based trigger decision |
| ~100 ms | event reconstruction |
| | event processing (skim, thin, augment) |
| ~months | final data analysis (uses millions of events) |

25 ns

$10^{-19}$-$10^{-15}$ s

~ms

0.01-20 ns

~min

1-100 ns

~100 ms

2.5 μs

200 ms

~100 ms

~months

O(100) *pp* collisions          data / simulation

(sub-)nuclear physics

out-going particles interact with detector

detector response (signal formation + digitization)

hardware-based trigger decision

software-based tr...

Not needed for "real time" decisions, but we need **real time speed** to make this work

e... processing (skim, thin, augment)

final data analysis (uses millions of events)

*Not to scale!*

*Image inspired by
JHEP 02 (2009) 007*

**Spanning $10^{-20}$ m up to 1 m
can take O(min/event)**

This is only possible because of **factorization** (*Markov Property*): given the physics at one energy (~1/length) scale, the physics at the next one is independent of what came before.

**Spanning $10^{-20}$ m up to 1 m can take O(min/event)**

# Part I: "Hard-scatter"

We begin with equations of motion

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

*See this paper for adapting a ME to HPC*

*See this paper for ME integration with GNNs*



```
**************************************************************
*                                                            *
*                    W E L C O M E  to                       *
*            M A D G R A P H 5 _ a M C @ N L O               *
*                                                            *
*                                                            *
*                 *                        *                 *
*           *        * *          * *        *               *
*             * * * * * 5 * * * *                            *
*           *        * *          * *        *               *
*                 *                        *                 *
*                                                            *
**************************************************************
```

## Many tools exist for automating this highest energy step

*For many cases, this is slow but not limiting (yet)*

# Part II: Quarks → protons (+ friends)

Makes heavy use of MCMC



Not a limiting factor in terms of computing time.

State-of-the-art for material interactions is Geant4.

Includes electromagnetic and hadronic physics with a variety of lists for increasing/decreasing accuracy (at the cost of time)

*This accounts for O(1) fraction of all HEP competing resources!*

Geant 4

It is important to mention that **after** Geant4, each experiment has custom code for *digitization*

this can also be slow; but is usually faster than Geant4 and reconstruction



deposited charge

MIP

Preamplifier output

40 MHz clock

analog threshold

Time

ToT = 2

It is important to mention that **after** Geant4, each experiment has custom code for *digitization*

N.B. **calorimeter energy deposits** **factorize** (sum of the deposits is the deposit of the sum) but **digitization** **(w/ noise) does not!**

deposited charge

Preamplifier output

analog threshold

Time

ToT = 2

Goal: replace (or augment) simulation steps with a faster, powerful generator based on state-of-the-art machine learning techniques

To start: attack the most important part:
**Calorimeter Simulation**

# Why should **you** care?
## Many analyses are forced to use a Geant4-based simulation as current fast sim. is not good enough.



**Standard Model Production Cross Section Measurements**   *Status: July 2017*

# Why should **you** care?
## Many analyses are forced to use a Geant4-based simulation as current fast sim. is not good enough.

**Standard Model Production Cross Section Measurements**   *Status: July 2017*

$\sigma$ [pb]

$10^{11}$  △ ○ total (x2)
□ △ ○ inelastic

$10^6$   *incl.*

**ATLAS** Preliminary
Run 1,2 $\sqrt{s}$ = 7, 8, 13 TeV

Theory

LHC pp $\sqrt{s}$ = 7 TeV

# If we don't do something, the HL-LHC won't be possible.  If we do something now, we can save O($10 million/year).

$n_j \geq 3$  $n_j \geq 4$  $n_j \geq 3$  $n_j \geq 4$  $ZZ$ $ZZ$  $H \rightarrow WW$  $W\gamma$
$n_j \geq 4$  $n_j \geq 5$  *s-chan*  $Z\gamma$
$n_j \geq 5$  $n_j \geq 4$  $n_j \geq 6$  $H \rightarrow \tau\tau$
$n_j \geq 6$  $n_j \geq 5$  $n_j \geq 7$  *Zt*

1

$n_j \geq 6$ $n_j \geq 5$  VBF
$H \rightarrow WW$
$n_j \geq 7$ $n_j \geq 6$

$10^{-1}$

$n_j \geq 7$  $n_j \geq 8$
$H \rightarrow \gamma\gamma$

$10^{-2}$  $n_j \geq 7$  $H \rightarrow ZZ \rightarrow 4\ell$

$10^{-3}$  $W^{\pm}W^{\pm}$

$WZ$

| pp | Jets $R=0.4$ | $\gamma$ | W | Z | $t\bar{t}$ | t | VV | $\gamma\gamma$ | H | WW | $W\gamma$ | $t\bar{t}W$ | $t\bar{t}Z$ | $t\bar{t}\gamma$ | Wjj EWK | Zjj EWK | WW | $Z\gamma\gamma$ | $W\gamma\gamma$ | $WW\gamma$ | $Z\gamma jj$ EWK | VVjj EWK |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | fid. | fid. | fid. | fid. | tot. | tot. | fid. | fid. | fid. | fid. | tot. | tot. | fid. | fid. | Excl. tot. | fid. | fid. | fid. | fid. | fid. | fid. |

Training NN's is slow,
but evaluation is **fast**

Physics-based
simulations are **slow**



What if we can learn to
simulate calorimeter
showers with a NN?

φ

z

η

Generative Adversarial Networks (GAN):
*A two-network game where one **maps noise to images** and one **classifies images as fake or real**.*
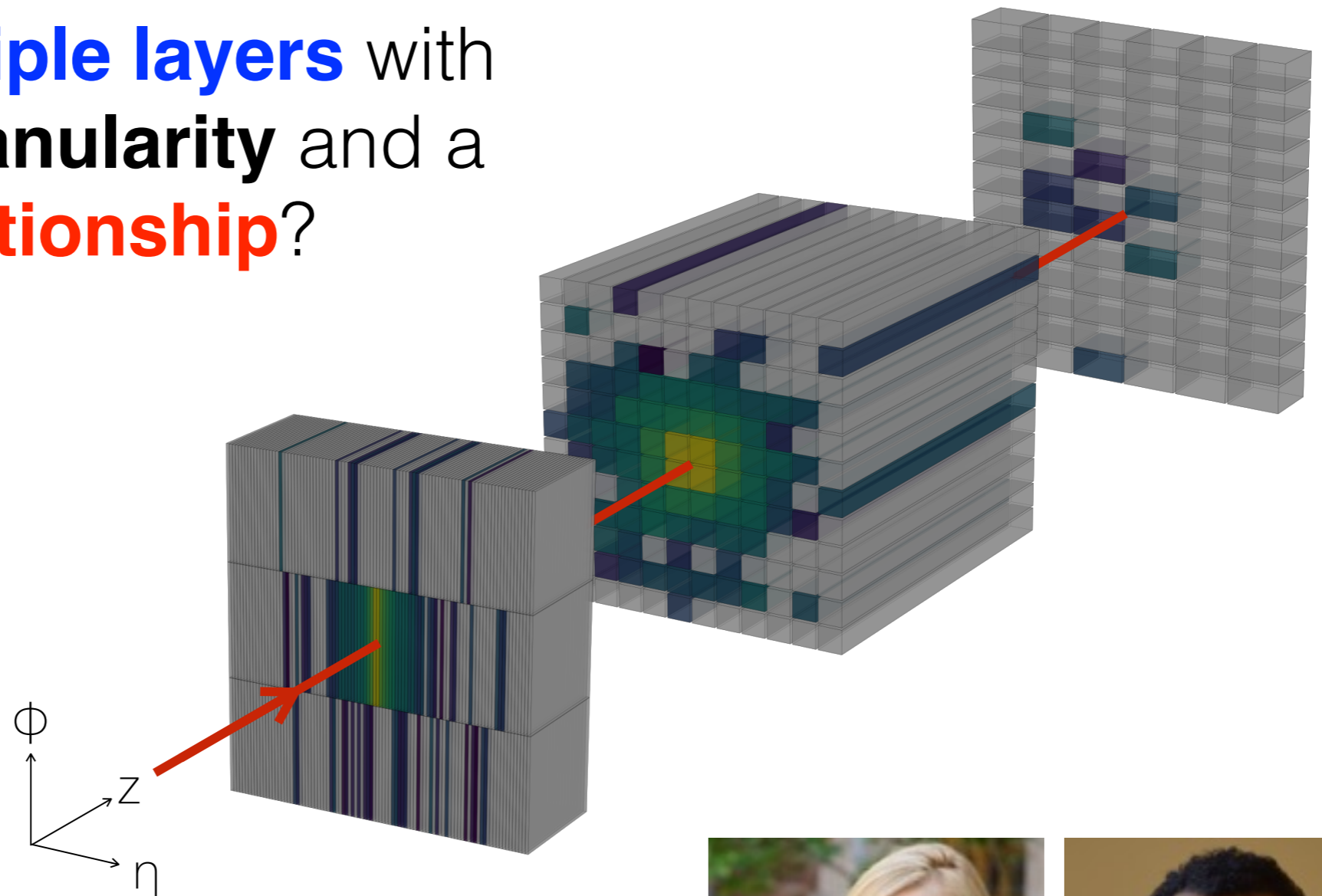
noise

GAN

{real,fake}

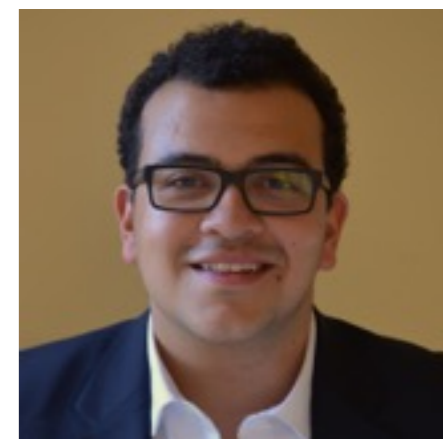When **D** is maximally confused, **G** will be a good generator

Physics-based simulator


Average generated signal image


Average Pythia signal image

What about **multiple layers** with **non-uniform granularity** and a **causal relationship**?



$\phi$

z

$\eta$

# Generator Network

One image
per calo layer

One network per particle type;
input particle energy

INPUTS

OUTPUTS

particle
energy
$E$

1

rescale

Scalar
multiplication

latent
space
$z$

1024

**LAGAN**
*G*

**LAGAN**
*G*

**LAGAN**
*G*

Resize

LCN

$W_{01}$
$W_{11}$

Linear
Combination

Resize

LCN

$W_{12}$
$W_{22}$

Linear
Combination

use layer i as
input to layer i+1

ReLU to
encourage sparsity

# Average Images

Geant4



CaloGAN

# Shower shapes



Qualitative agreement; clearly also room for improvement.

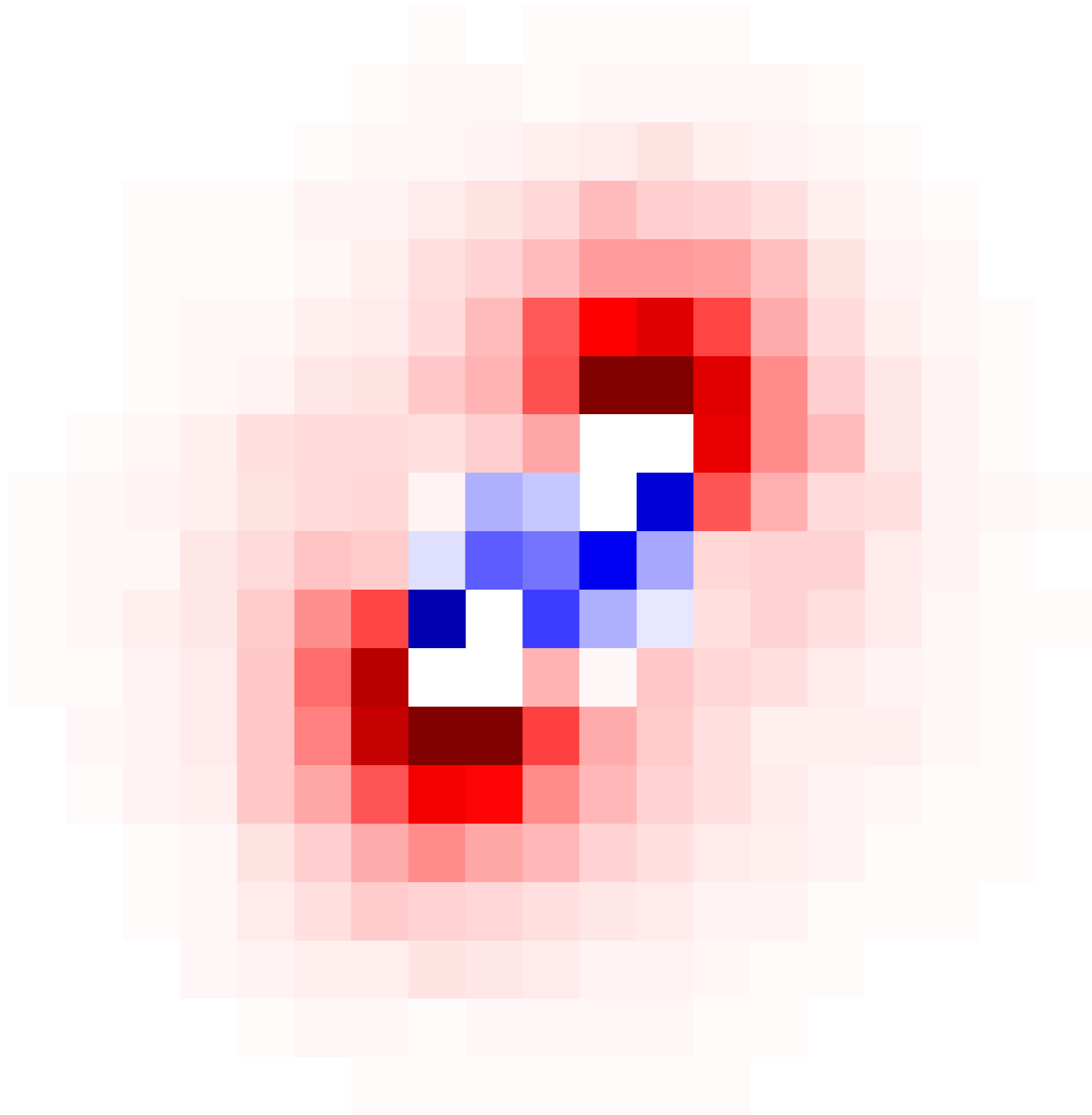In fact, one could add this into the training; for now held out for validation.

# Timing

| Generation Method | Hardware | Batch Size | milliseconds/shower |
|---|---|---|---|
| GEANT4 | CPU | N/A | 1772 ← |
| CALOGAN | CPU | 1 | 13.1 |
| | | 10 | 5.11 |
| | | 128 | 2.19 |
| | | 1024 | 2.03 |
| | GPU | 1 | 14.5 |
| | | 4 | 3.68 |
| | | 128 | 0.021 |
| | | 512 | 0.014 |
| | | 1024 | 0.012 ← |

*M. Paganini, L. de Oliveira,*
***BPN**, PRL 120 (2018) 042003*

The LHC is a unique science tool with extreme challenges related to the data rate: real time / ultra fast algorithms are required.





**There are many exciting opportunities and ideas for fully exploiting our data** we must make sure no stone is left unturned !
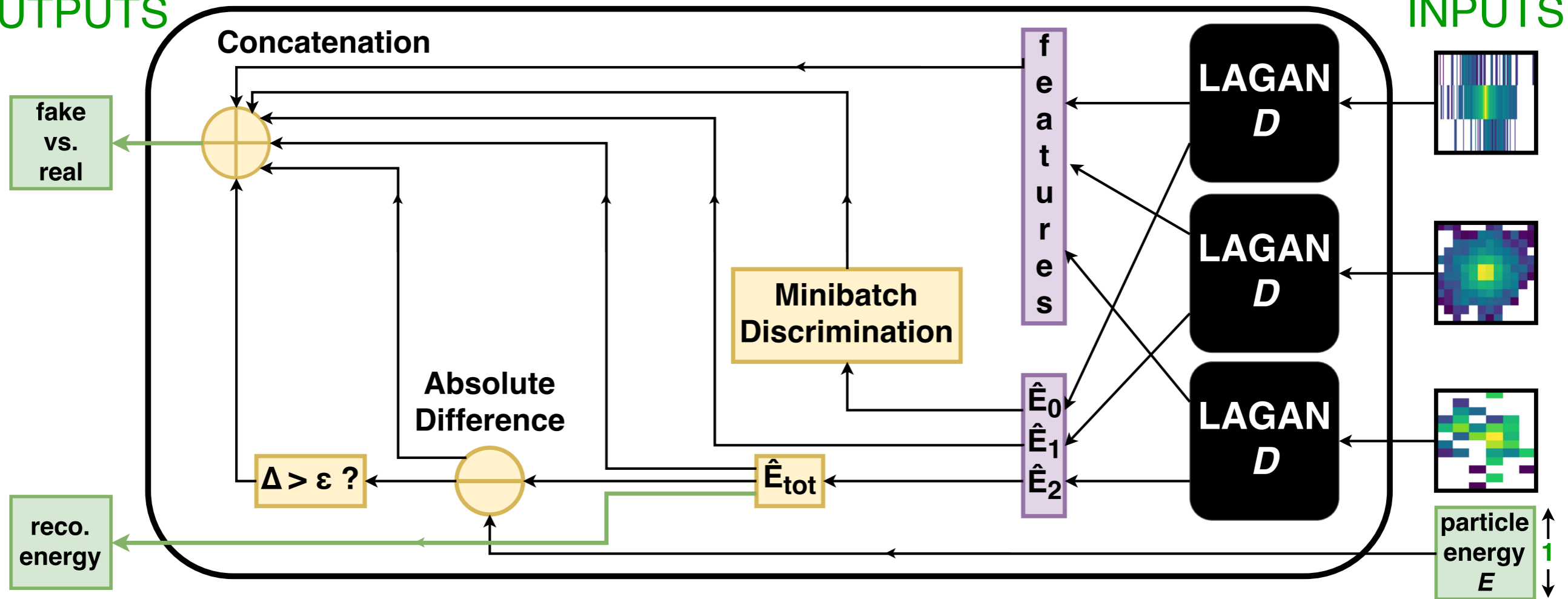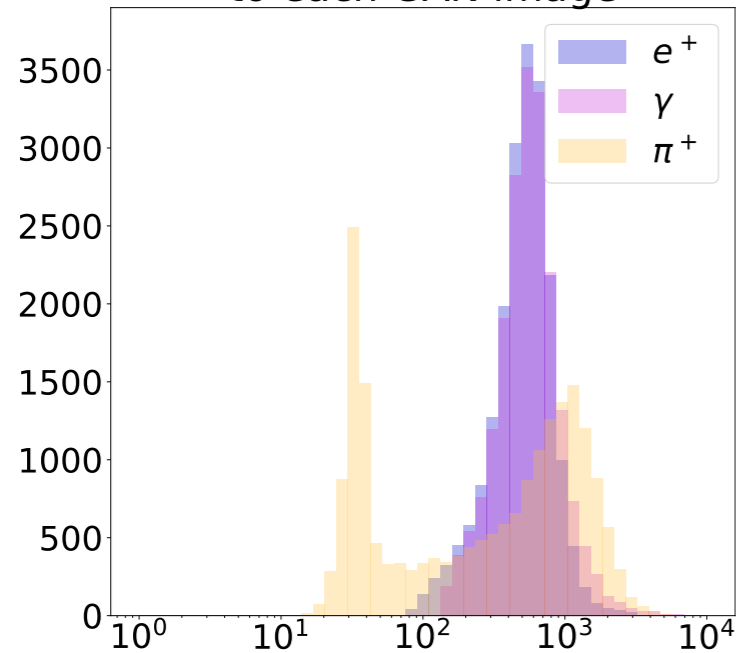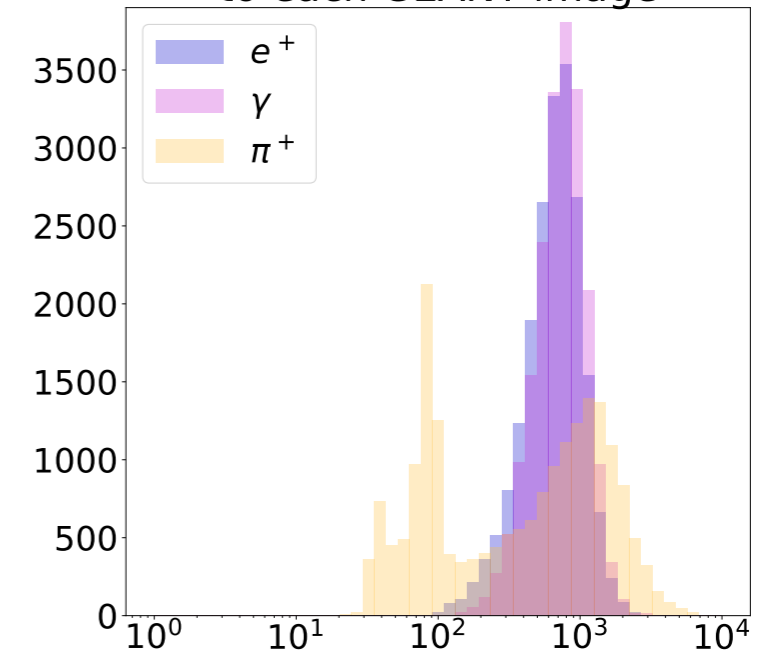
Fin.

# "Overtraining"

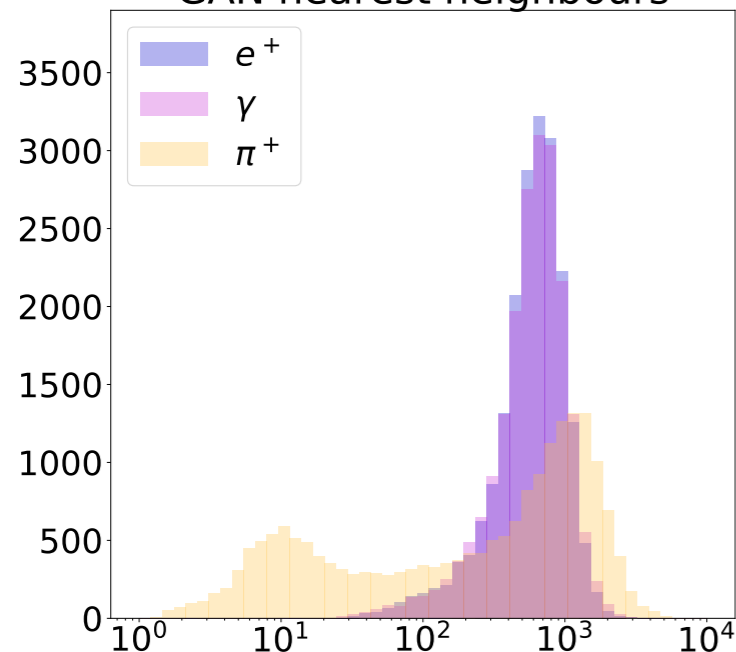Nearest GEANT neighbour to each GAN image


Nearest GAN neighbour to each GEANT image

not memorizing

A key challenge for GANs is the diversity of generated images. This does not seem to be a problem for CaloGAN.


GAN nearest neighbours


GEANT nearest neighbours

no mode collapse