# Detector-embedded tracking with the 'RETINA algorithm'

**Giovanni Punzi**
*University & INFN-Pisa*
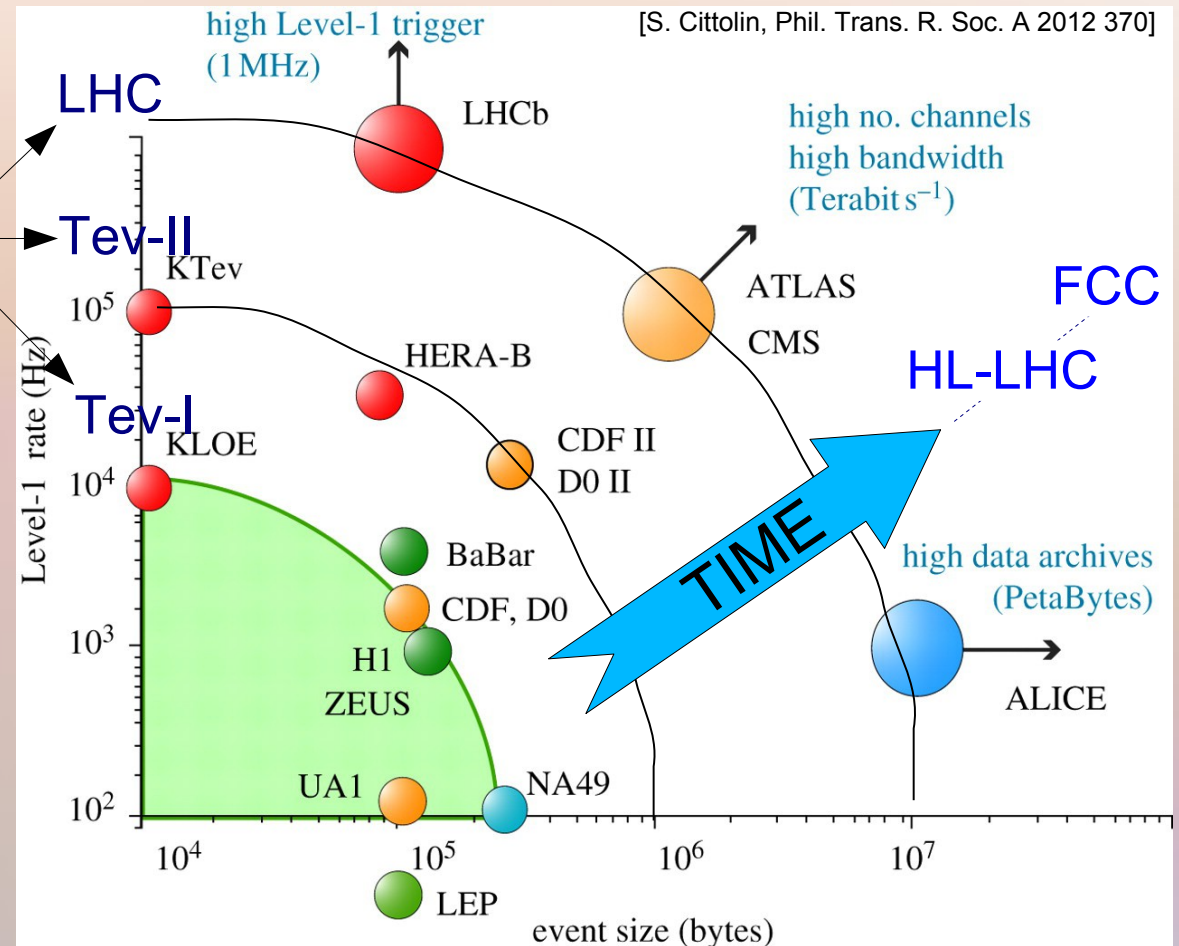
IFAE 2017
*Trieste, April 19-21, 2017*

# History of HEP Data-Processing in a plot

**Experiment generations**

LHC

Tev-II

Tev-I

Complexity and computational load kept increasing, in pace with electronics advancements
→ Data processing still a major cost item of experiments
→ Often a major technical constraint



[S. Cittolin, Phil. Trans. R. Soc. A 2012 370]

high Level-1 trigger (1 MHz)

high no. channels high bandwidth (Terabit s$^{-1}$)

FCC

HL-LHC

LHCb

KTev

HERA-B

ATLAS CMS

KLOE

CDF II D0 II

TIME

BaBar

CDF, D0

high data archives (PetaBytes)

H1 ZEUS

ALICE

UA1

NA49

LEP

Level-1 rate (Hz)

$10^5$

$10^4$

$10^3$

$10^2$

$10^4$      $10^5$      $10^6$      $10^7$
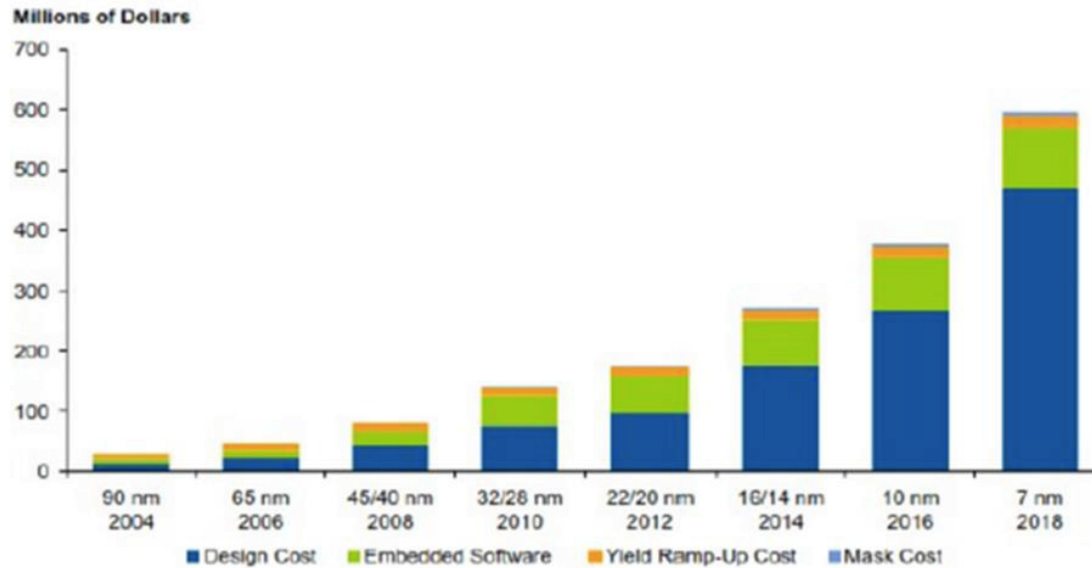
event size (bytes)

But:  - Physics landscape now asking for more ***precision***
        - Moore's law ***slowing down***
        …  symptoms that HEP will face a computing roadblock

# NOT simply 'business as usual'
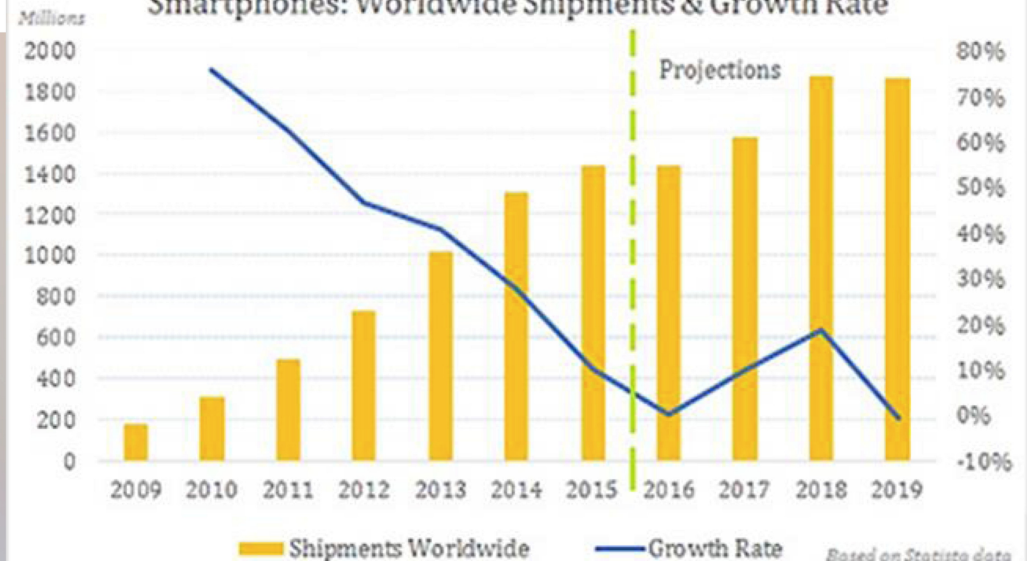


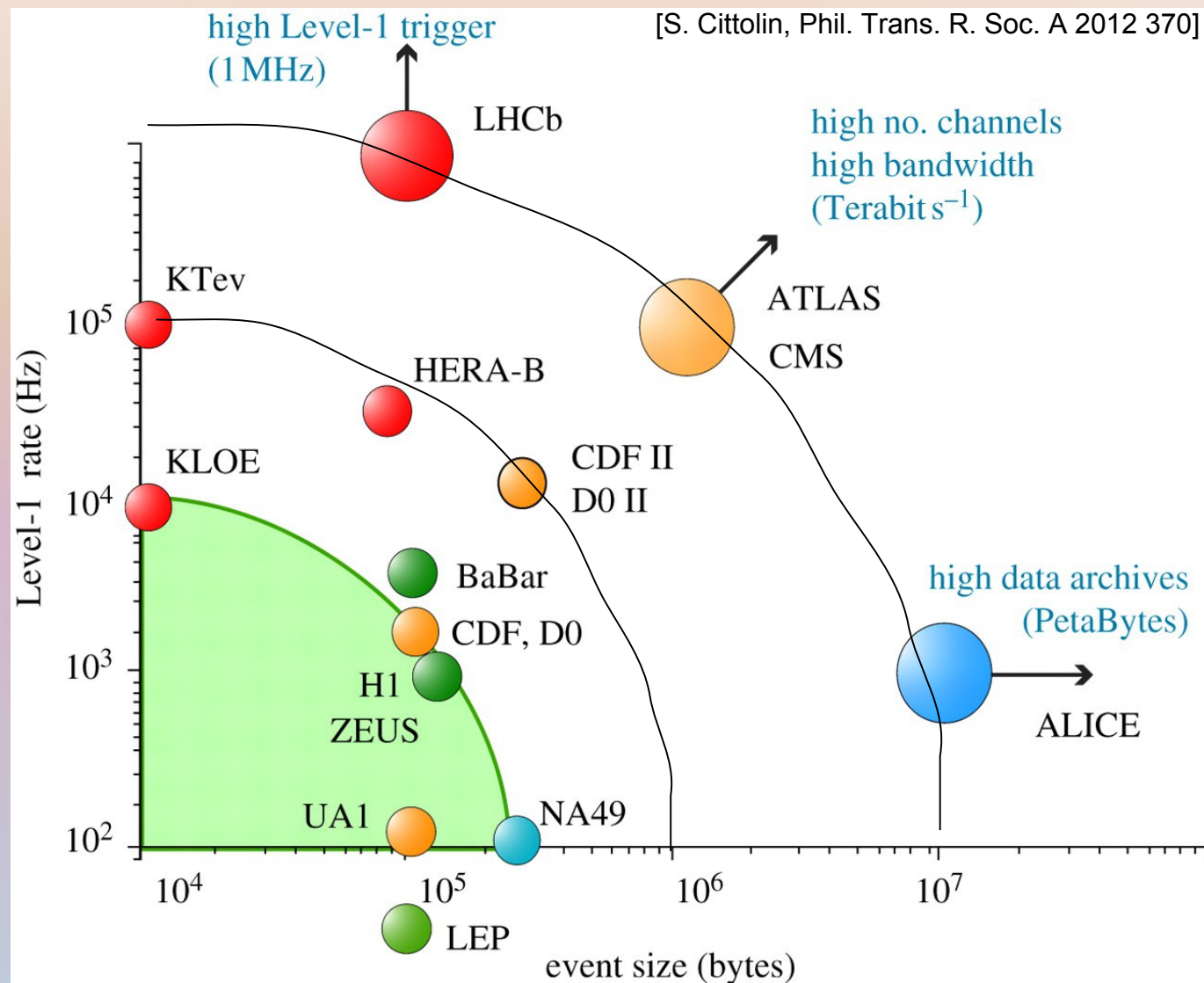**Estimated Cost of Developing Lower Node Chips**

Market Realist

- Increasing development costs...

... and decreasing returns



Smartphones: Worldwide Shipments & Growth Rate
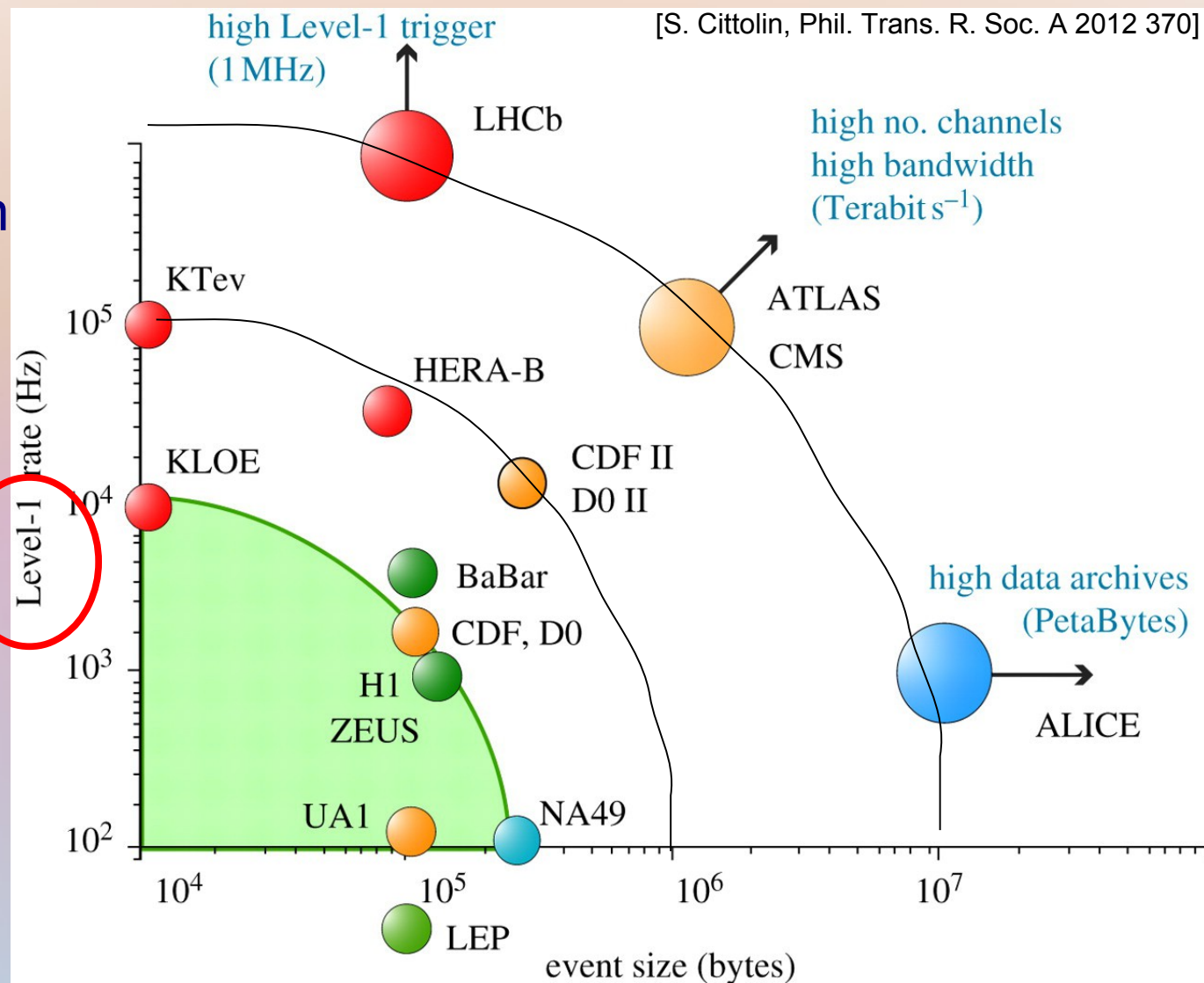
# More facts hidden in this plot...



[S. Cittolin, Phil. Trans. R. Soc. A 2012 370]

# More facts hidden in this plot

This is NOT
The full rate !
$>10^2$ reduction
by Level-1
pre-selection

**Level-1**
is a <u>major</u>
<u>challenge</u>
for the future



[S. Cittolin, Phil. Trans. R. Soc. A 2012 370]

high Level-1 trigger
(1 MHz)

high no. channels
high bandwidth
(Terabit s$^{-1}$)

LHCb

KTev

HERA-B

ATLAS
CMS

KLOE

CDF II
D0 II

$10^5$

$10^4$

BaBar

CDF, D0

high data archives
(PetaBytes)

$10^3$

H1
ZEUS

ALICE

$10^2$

UA1

NA49

LEP

Level-1 rate (Hz)

event size (bytes)

$10^4$    $10^5$    $10^6$    $10^7$

# The Level-1 challenge

- Traditionally based on simple quantities, cheap and fast to calculate

- Future: more complex, "precision physics", large pileup...
  => No easily-extracted, small portion of the event allowing to reduce data for later processing.

  $\rightarrow$ need *detailed processing (tracking)* of data in each crossing

  - ATLAS: tracking trigger (FTK) at almost the full crossing rate

  - CMS: need tracker readout at 40 MHz in order to do L1 decision

  - LHCb: "signal" in every collision, full event analysis at 40 MHz

  - FCC: *all SM physics will be "low-Pt" physics*

    - rate of top events ~3 kHz

    - Pile-up ~$10^3$

  **All are hard problems – drain resources from HLT**

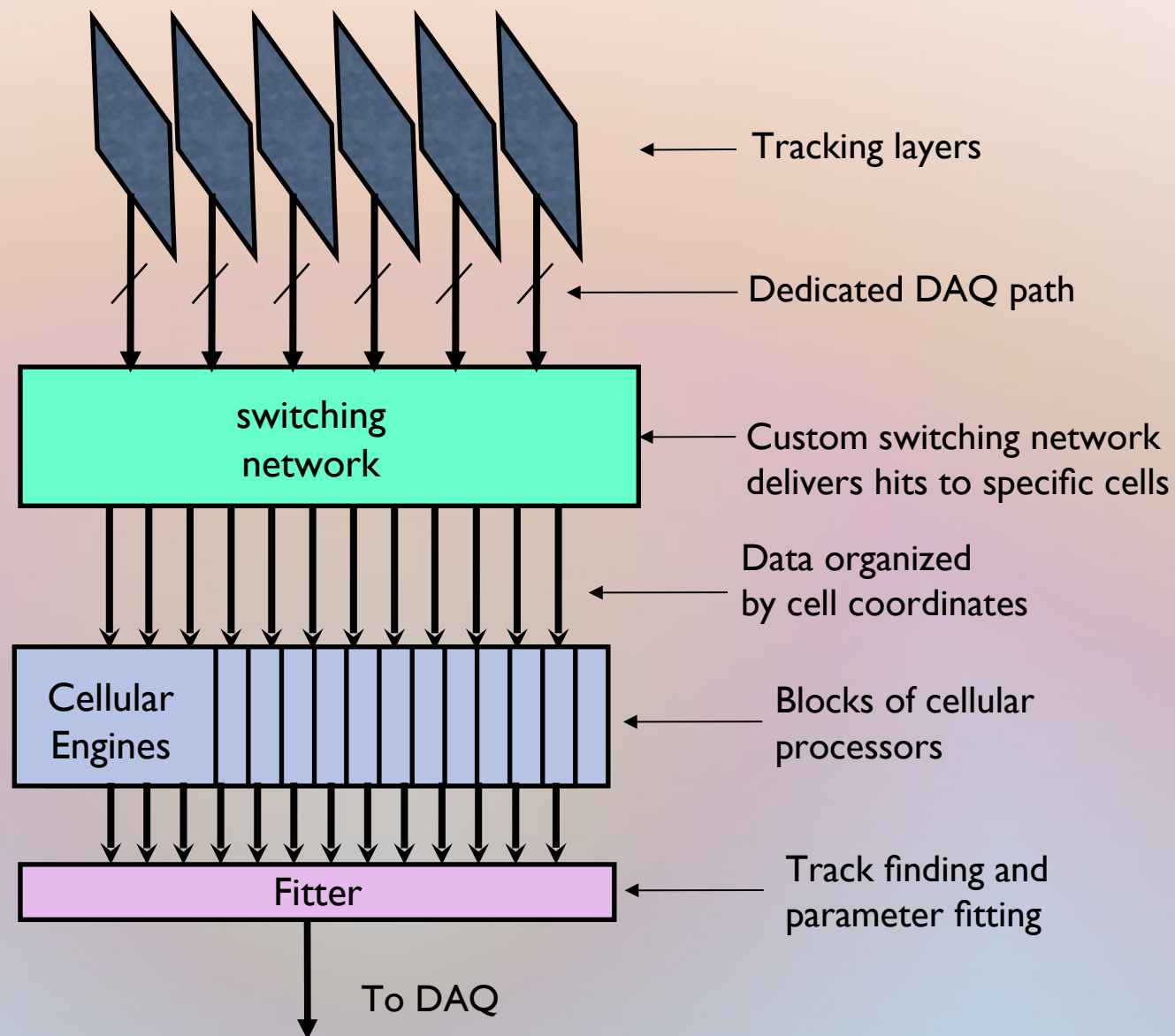# Technologies for Level-1 tracking

- Fastest approach to tracking up to now: direct matching to a bank of **stored templates**

- First large system to use this method has been CDF, at the Tevatron, where a real-time processor named SVT was capable of reconstructing quality tracks @30kHz in ~10µs.

- Based on custom ASICs implementing content-addressable memory (Associative Memory [NIM A278, (1989), 436-440])

- It actually worked (allowed CDF to discover Bs oscillations)

- Same approach continuing in FTK for ATLAS and in the planned Phase 2 upgrade for CMS  ( O(MHz) event rate)

- 'RETINA' approach an updated version of this idea, aiming at even better performance → eventually **embedded tracking** *( = tracking detectors producing tracks, rather than raw hits)*
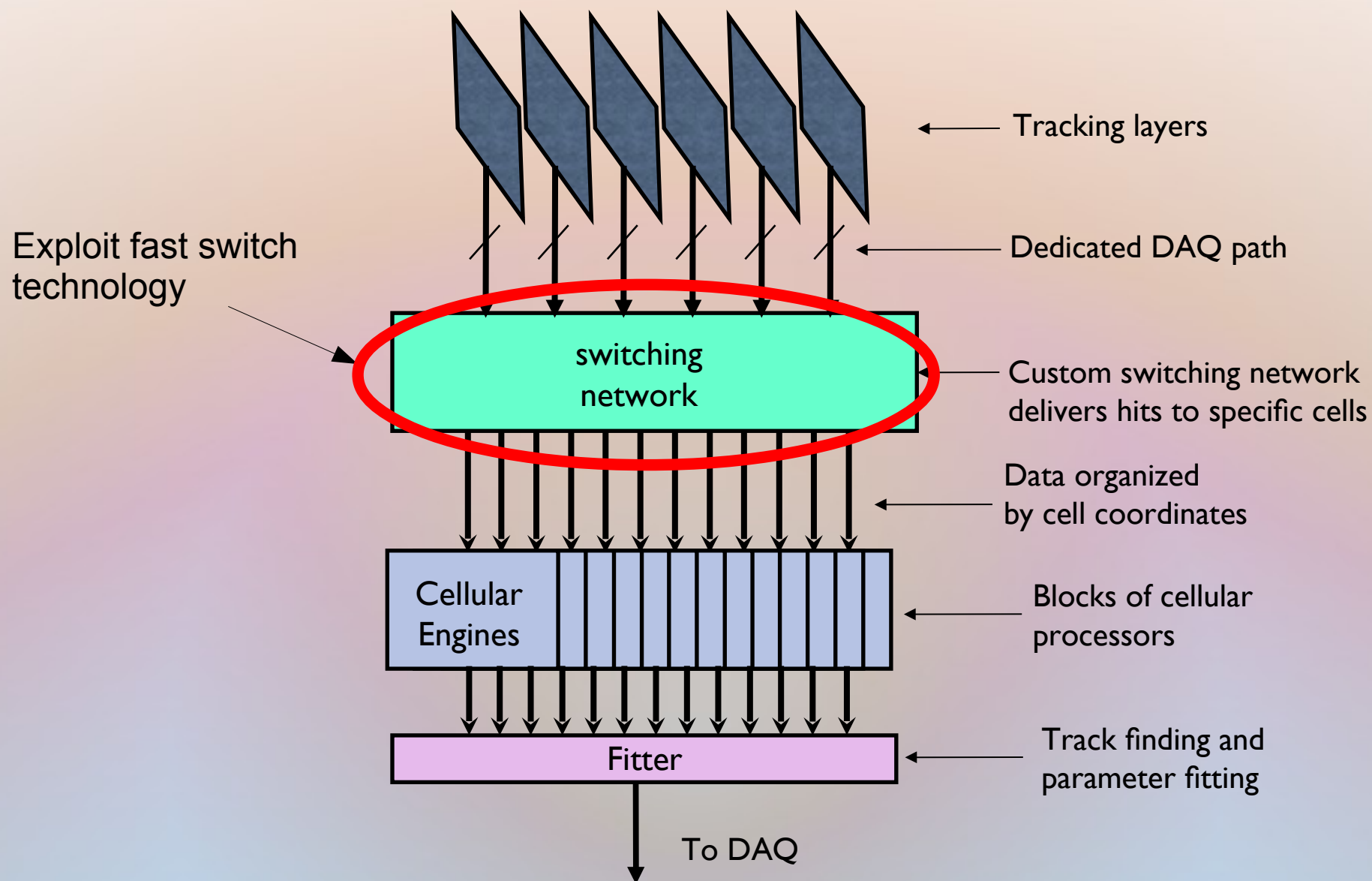
# Basic ideas behind RETINA

1) <u>Specialization:</u> build a dedicated system. Remember that the success of GPUs stems from specialization for a narrow purpose. Aim to build something that does for Tracking what the GPU did for Graphics (just with a smaller market...) (a "TPU").

2) <u>Template matching:</u> Inherit from the Associative Memory idea of parallel template matching and push it even further

3) <u>Inspiration from natural vision:</u> Natural neural systems are capable of pattern recognition in ~30 'time units' vs ~2000 for the best artificial systems up to now → analog weighting, extreme parallelism, and *"bandwidth over calculation"*.

4) <u>Opportunistic technology:</u> exploit telecom industry: COTS, optical fibers, FPGA – industry's choice for complex projects with small productions (CT scanners, high-end radars...)
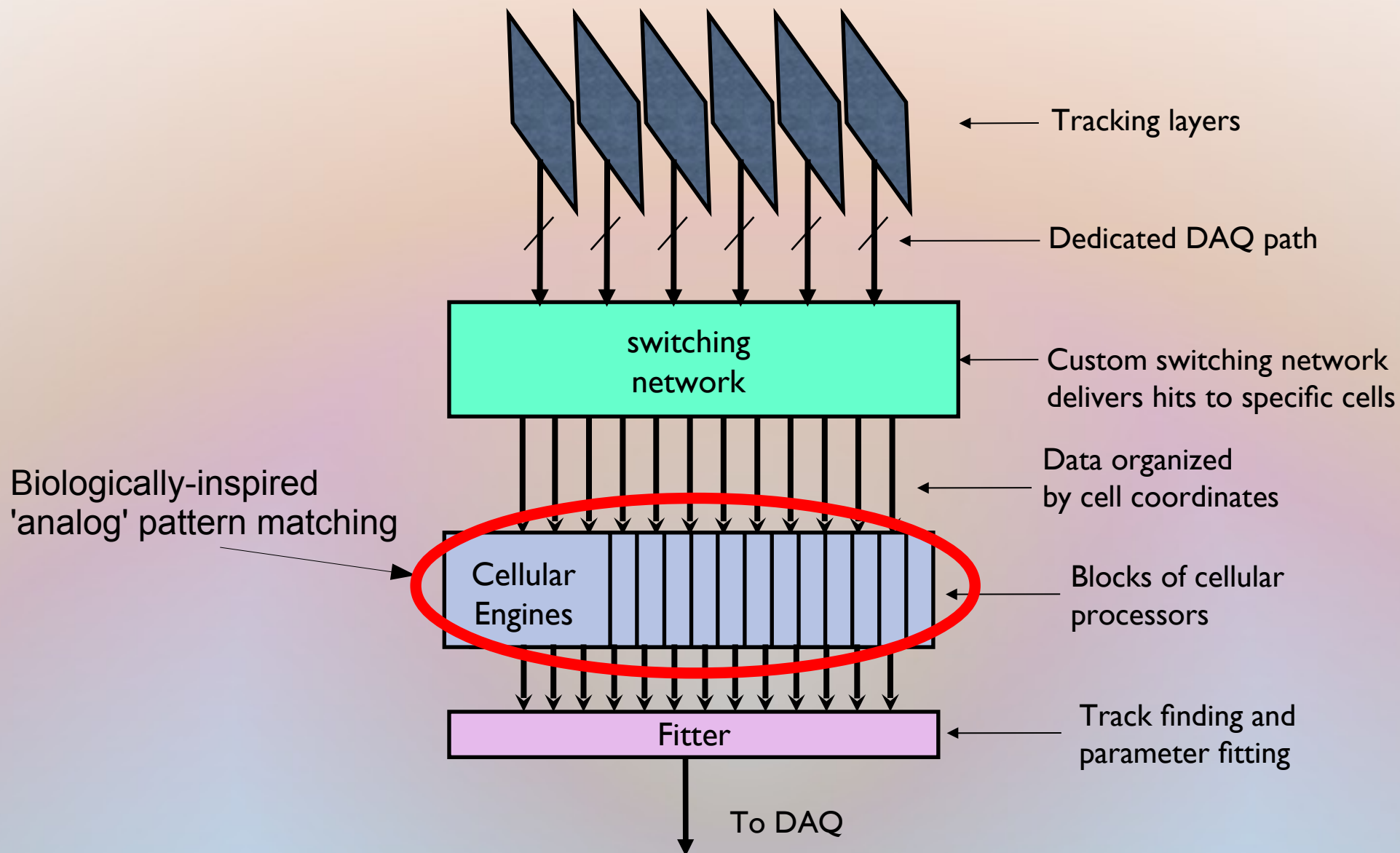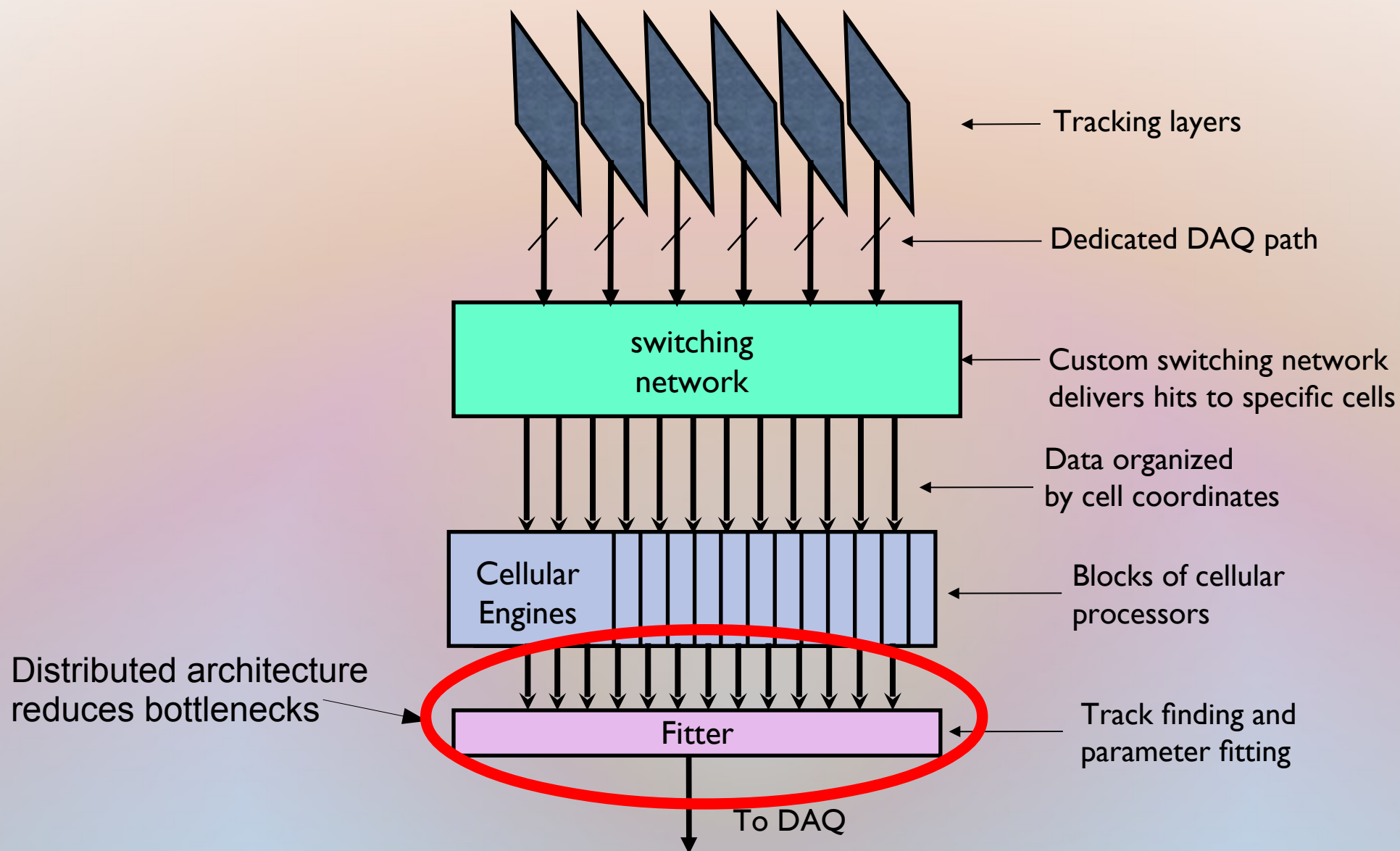
# RETINA system architecture



Tracking layers

Dedicated DAQ path

switching network

Custom switching network delivers hits to specific cells

Data organized by cell coordinates

Cellular Engines

Blocks of cellular processors

Fitter

Track finding and parameter fitting

To DAQ

# RETINA system architecture



Tracking layers

Exploit fast switch technology

Dedicated DAQ path

switching network

Custom switching network delivers hits to specific cells

Data organized by cell coordinates

Cellular Engines

Blocks of cellular processors

Fitter

Track finding and parameter fitting

To DAQ

# RETINA system architecture



Tracking layers

Dedicated DAQ path

switching network

Custom switching network delivers hits to specific cells

Data organized by cell coordinates

Biologically-inspired 'analog' pattern matching

Cellular Engines

Blocks of cellular processors

Fitter

Track finding and parameter fitting

To DAQ

# RETINA system architecture



Tracking layers

Dedicated DAQ path

switching network

Custom switching network delivers hits to specific cells

Data organized by cell coordinates

Cellular Engines

Blocks of cellular processors

Distributed architecture reduces bottlenecks

Fitter

Track finding and parameter fitting

To DAQ

# Hit delivery via smart switching

- Hits must be delivered only to the cells that need them (can be more than one)
- Switch network "knows" where to deliver
- Information is *embedded in the network* via distributed LUTs

**Data processing occurs *while* data is being moved - not *afterwards***

one hit

**Pay the price of a (temporary) bandwidth increase**

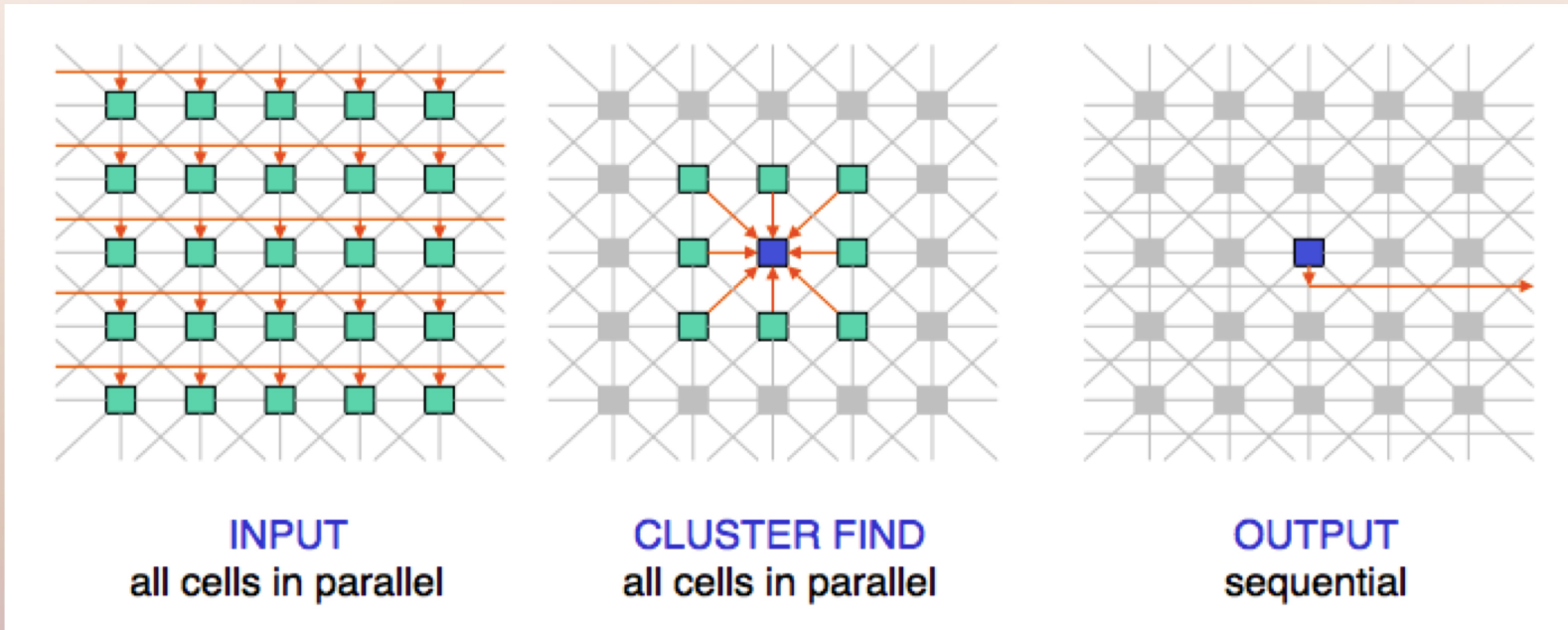# Custom switching network built from uniform elementary blocks

# "neural-like" tracking algorithm



- Map track parameter space into an array of cells, implemented in hardware
- Each cell performs a weighted sum of hits near to the track trajectory (inspired by biological receptive fields of visual cortex)
- A valid track appears as a cluster of cell responses – parameters can then be determined by interpolation of nearby cells → save on hardware size

- First work in this direction in year 2000 [L. Ristori, "An Artificial retina for Fast Track Finding" NIM A453(2000),425]  (historical reason for the name)

- Mathematically related to "Hough transform"[P.V.C.Hough,Conf.Proc. C590914(1959),554]
  – but the essence of the retina approach is *architectural implementation*

# Cellular engine working principle



**INPUT**
all cells in parallel

**CLUSTER FIND**
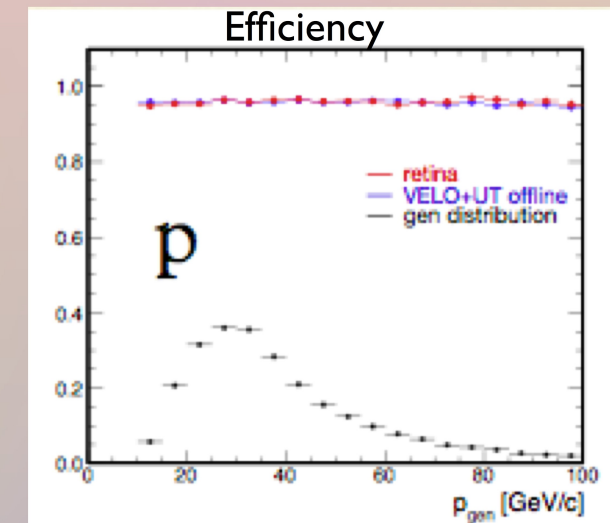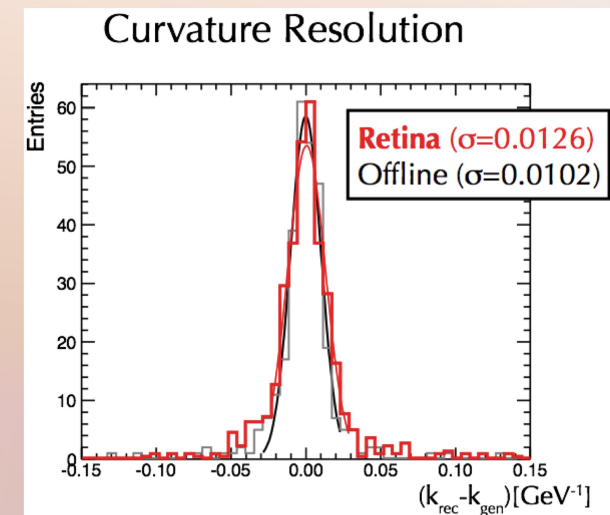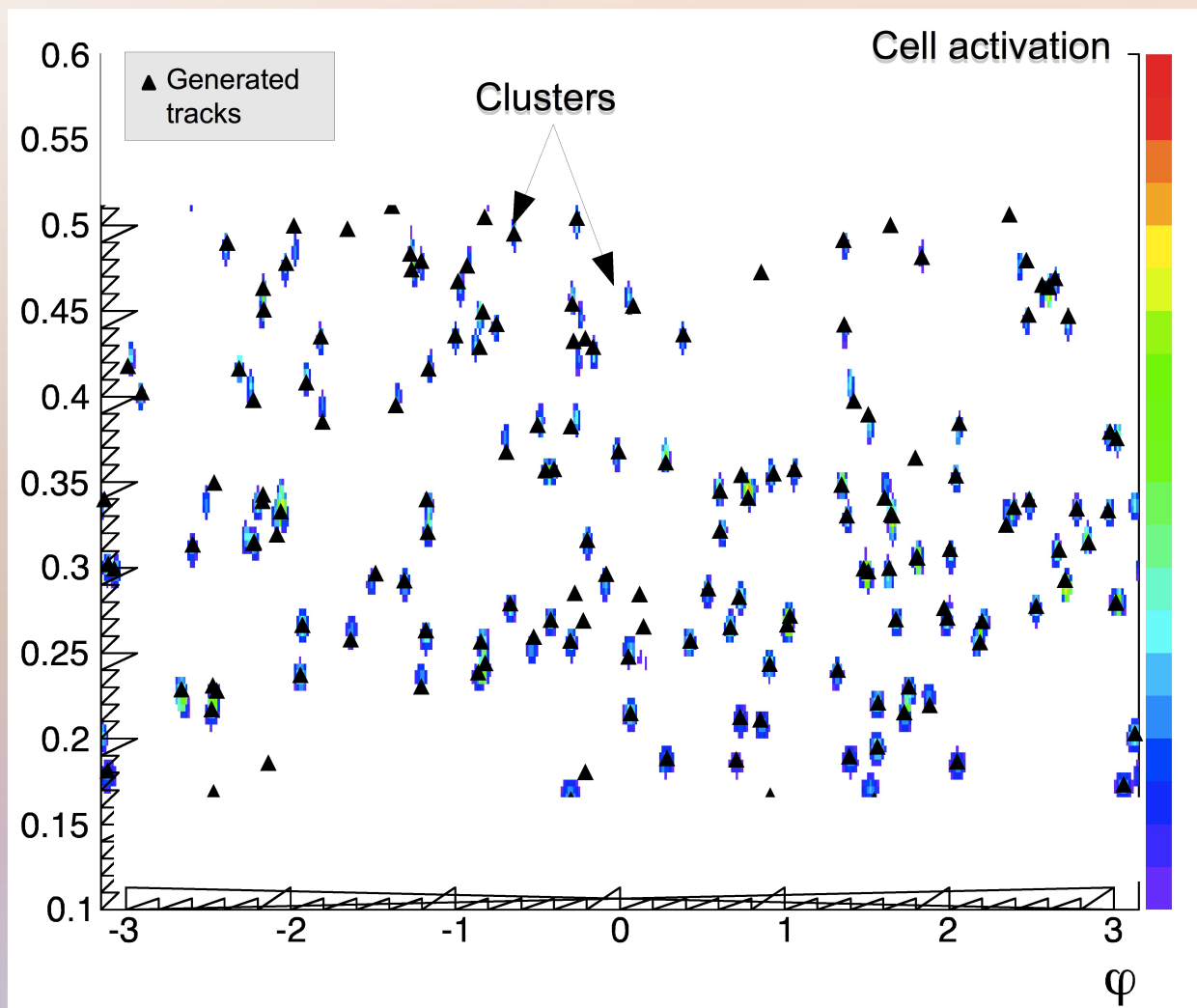all cells in parallel

**OUTPUT**
sequential

1. Computing engine in each cell computes weighted sum in parallel
2. Each node deals with nearby cells → local clustering
3. If cluster center, output result to next stage (not shown)
4. Output of several nodes input to local fitting logic to finalize track

Everything happens in pipeline without wait states (data-flow)
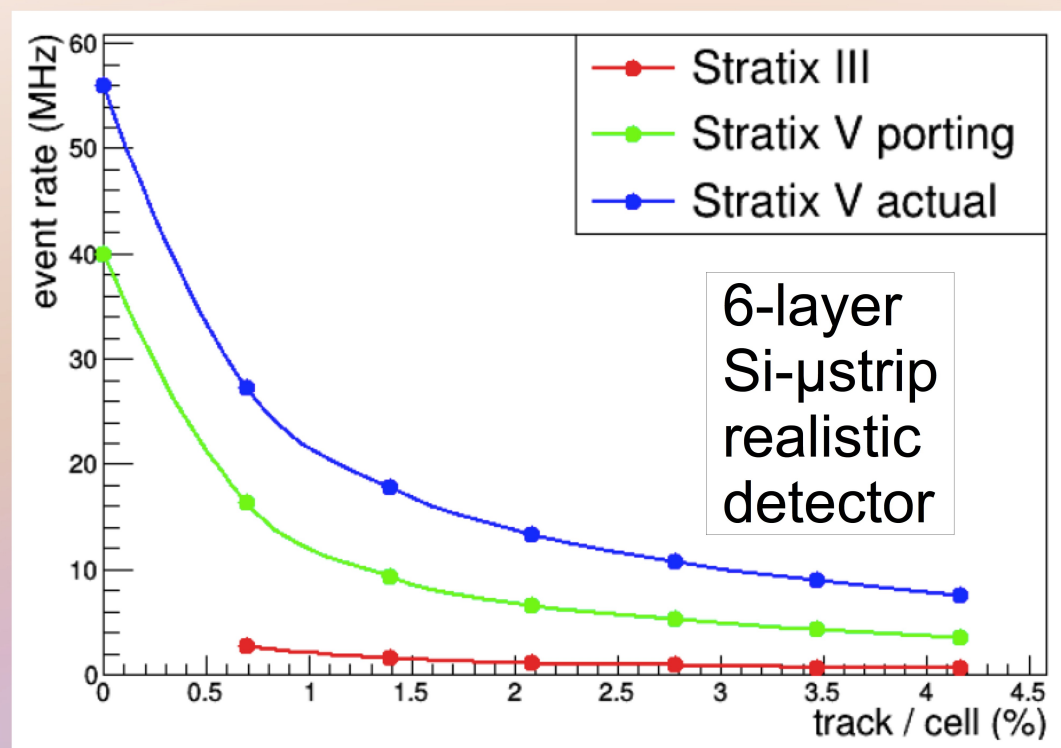
# Retina Tracking algorithm works



- Example simulation of 6-layer pixel detector [CERN-LHCb-PUB-2014-026]
- Shows offline-like efficiency/resolution possible with Retina algorithm

# Hardware Prototype



- Thanks to INFN-CSN5 RETINA project

- 2 Stratix-V (1MLE, high speed grade) **1.2 Tb/s bidirectional bandwidth** up to 700 MHz clock

- On-board CPU, ample DDR memory, 96 inter-FPGA LVDS connections

- 96 high-speed SerDes I/O (12 Gb/s)

- With optical links, buffer memories, disks, CPU rack etc, for high-rate tests

- Can be used as "building block" for an entire high-performance tracker → results on this prototype readily extrapolate to real systems

- Ample choice of interfaces allow accommodating a variety of applications and connect to different systems

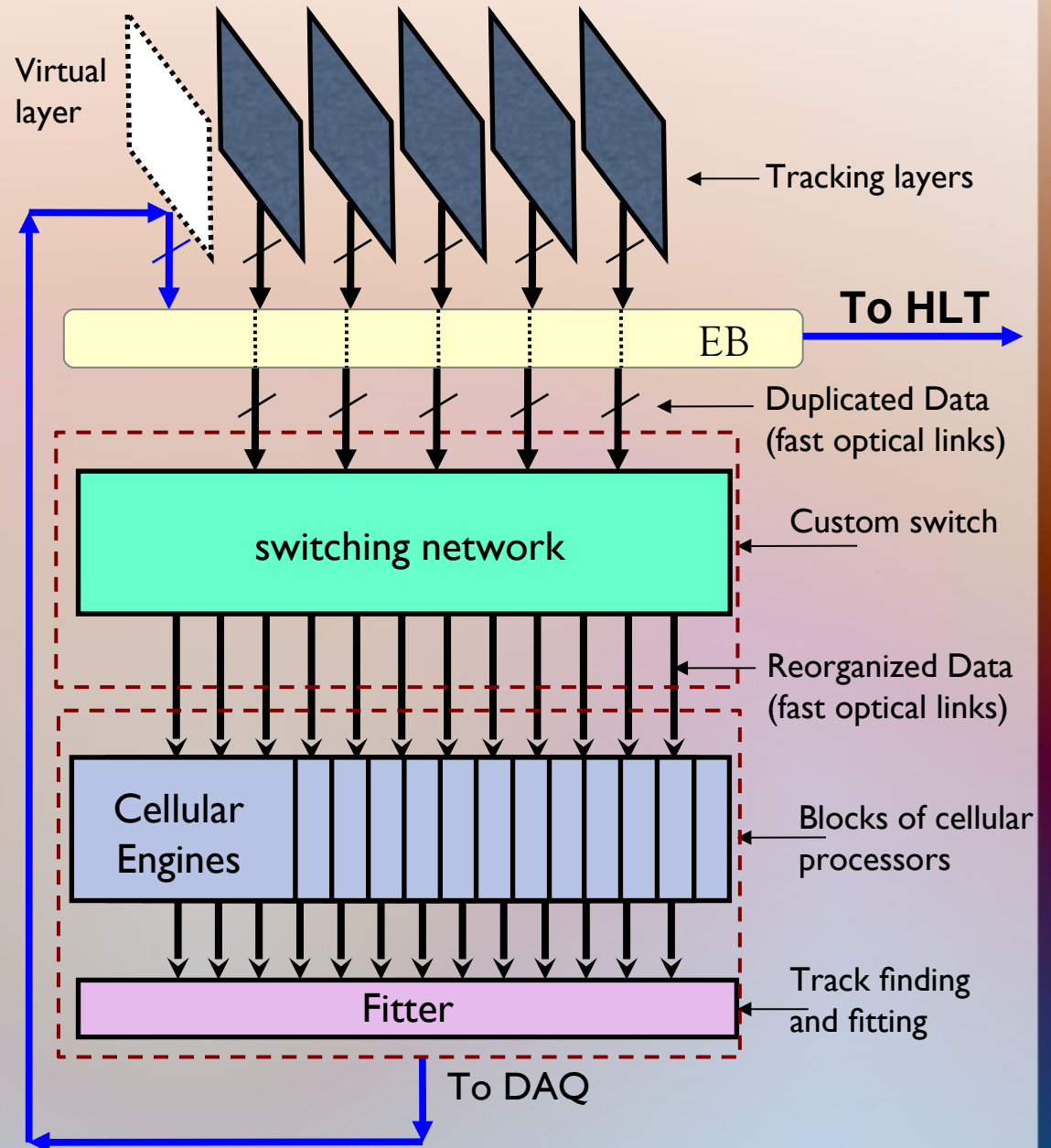# Prototype Measured performances



- Prototype achieved LHC crossing rate (40 MHz)
- Hardware cost   <0.1 € /kHz of tracks (prototype)
   - power cost: 0.2 mW /kHz of tracks
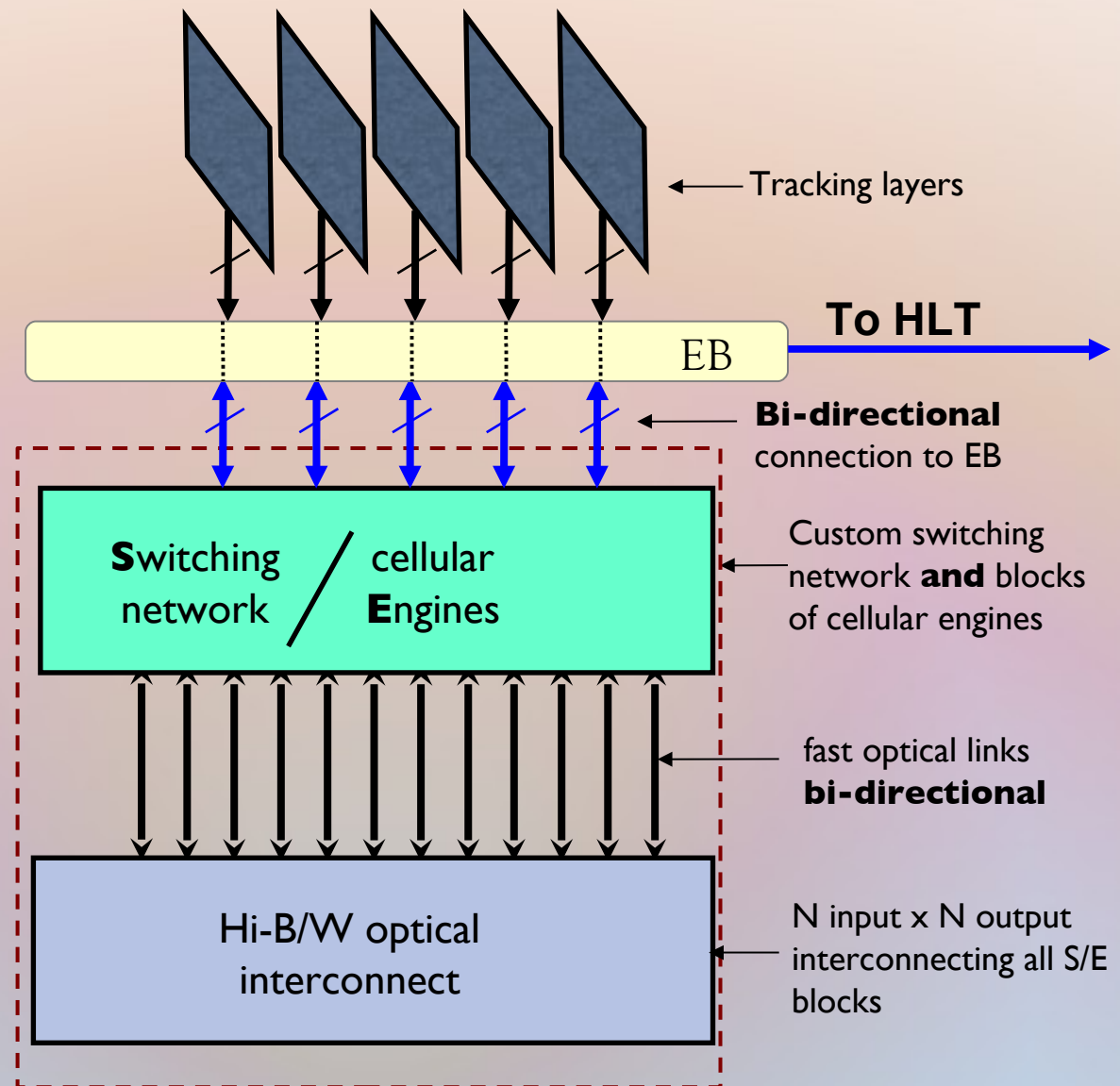- Very short latency <0.5 µs <u>facilitates embedding</u>

# Embedded RETINA

- Consider a configuration where detector data enters an Event Builder (EB) before going to HLT

- Data can be duplicated inside EB and sent to RETINA system

- RETINA output can be fed back to the EB, appearing as an extra "Virtual layer" of the detector, producing reconstructed tracks, ready for up-front use by HLT
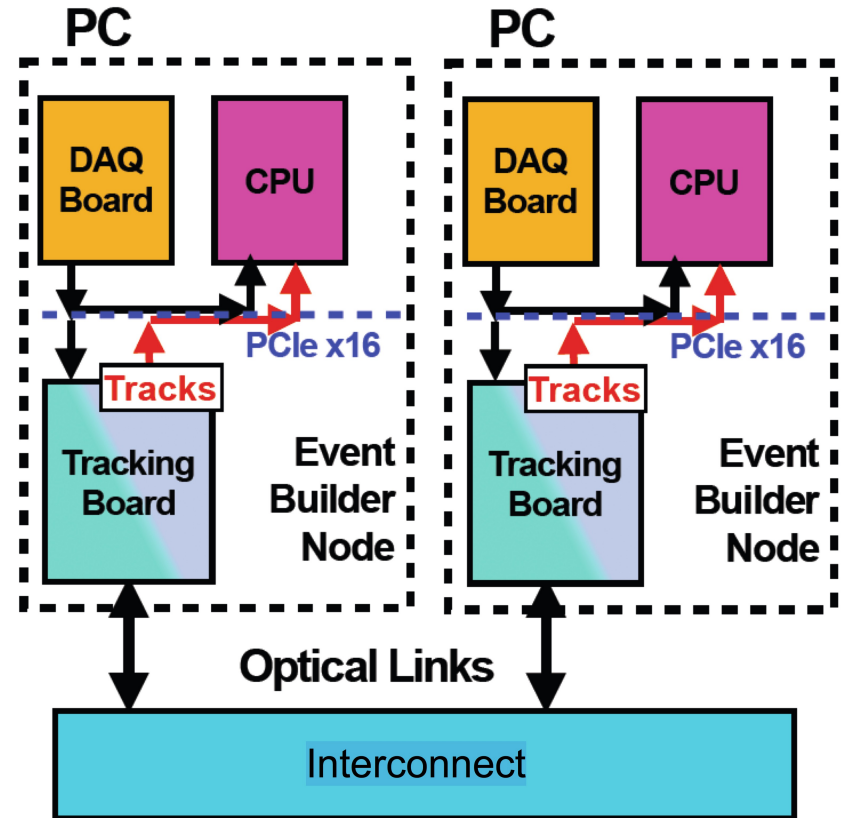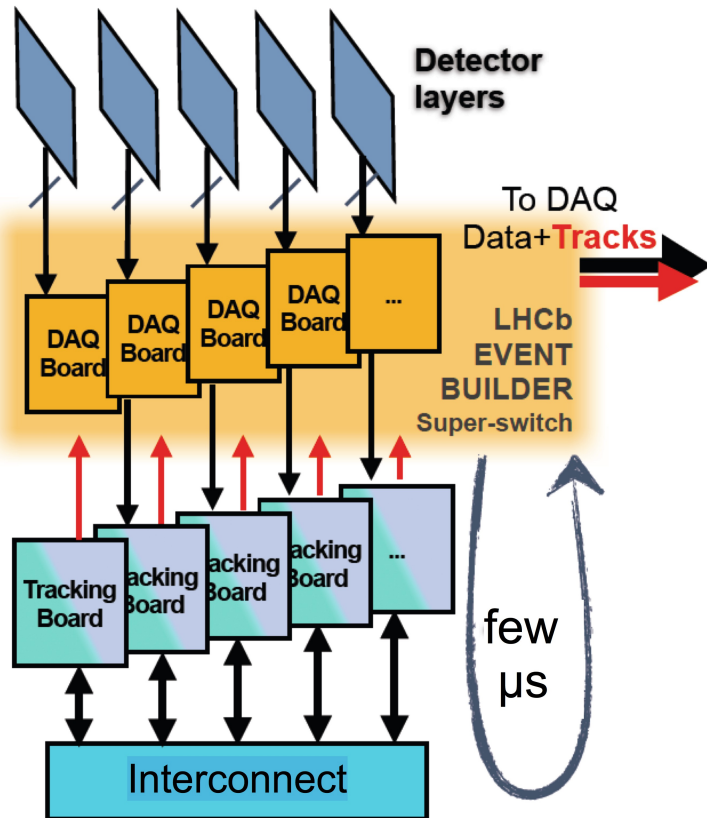
- Requires short enough latency

Virtual layer

Tracking layers

To HLT

EB

Duplicated Data (fast optical links)

switching network

Custom switch

Reorganized Data (fast optical links)

Cellular Engines

Blocks of cellular processors

Fitter

Track finding and fitting

To DAQ

# Improved bi-directional configuration



Tracking layers

EB

**To HLT**

**Bi-directional** connection to EB

**S**witching network / cellular **E**ngines

Custom switching network **and** blocks of cellular engines

fast optical links **bi-directional**

Hi-B/W optical interconnect

N input x N output interconnecting all S/E blocks

# Distributed-embedded RETINA



- A single tracking board performs **both** switching and clustering
  - Reads small detector portion, outputs small parameter space
- Easier to implement large global bandwidths
- Allows use of standard commercial PCIe FPGA boards
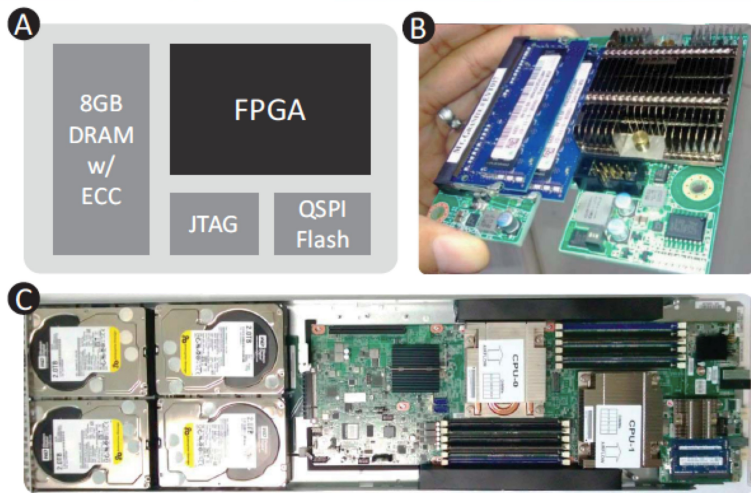
# Microsoft's CATAPULT architecture



Figure 1: (a) A block diagram of the FPGA board. (b) A picture of the manufactured board. (c) A diagram of the 1 U, half-width server that hosts the FPGA board. The air flows from the left to the right, leaving the FPGA in the exhaust of both CPUs.
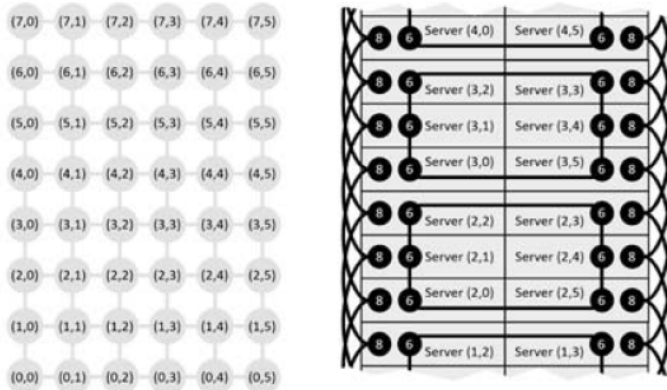


Figure 2: The logical mapping of the torus network, and the physical wiring on a pod of 2 x 24 servers.
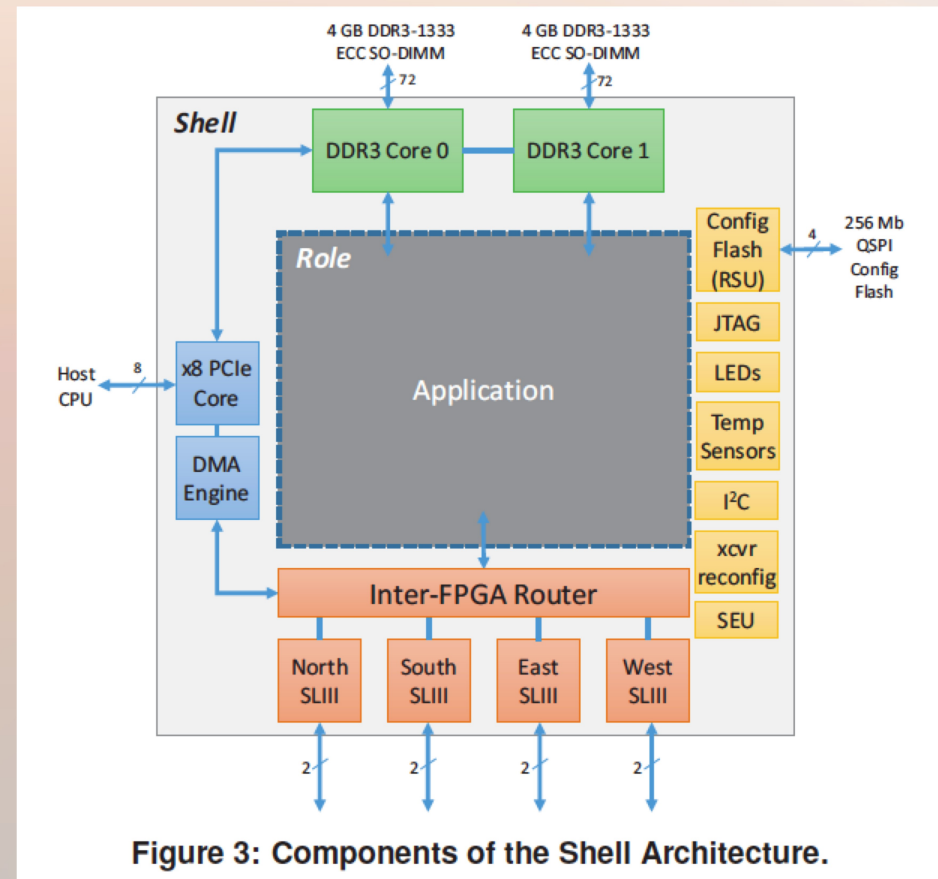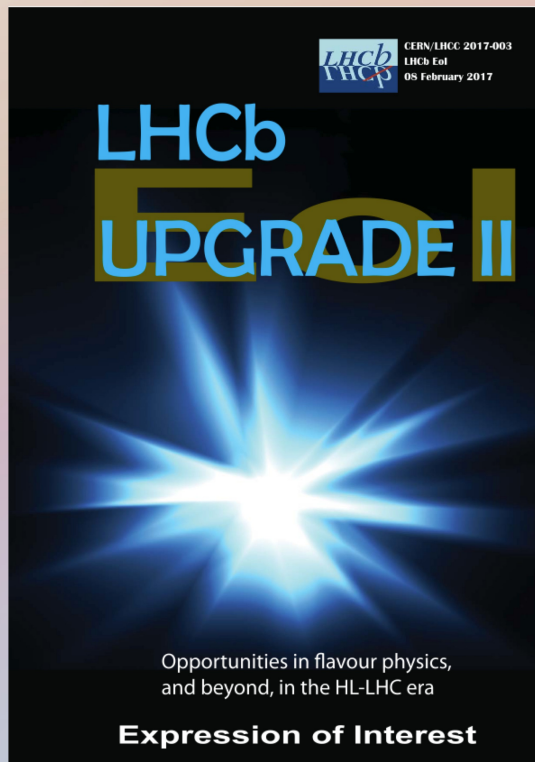


Figure 3: Components of the Shell Architecture.

- Interesting structural analogy with independently-developed 'Catapult' system
- Distributed, inter-connected FPGA boards (powering Bing in clouds)
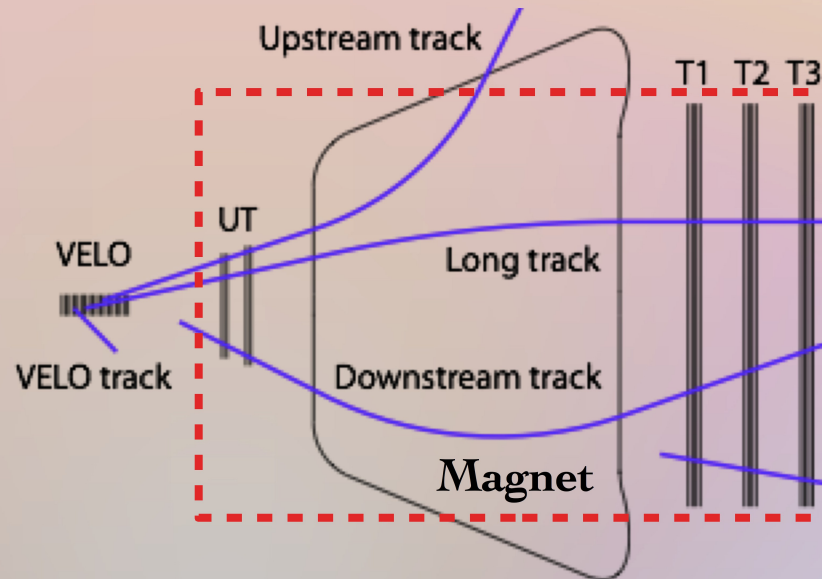- Larger latencies, but similar issues

# Application: LHCb Phase-2 upgrade

- LHCb recently published an EOI for upgrading to $L=2*10^{34}$ to better exploit the HL-LHC

- Includes proposal of a real-time RETINA tracker for long-lived tracks, embedded in the Event Builder



CERN/LHCC 2017-003 (8/2/2017)

- 16 layers w/ stereo and B field
- Phase IB: ~10 Gtracks/s @$2*10^{33}$
  = *all CMS tracks >2GeV @$5*10^{34}$*
  *(challenging!)*

# Summary

- A key to future HEP experiments will be the capability of real-time reconstruction by special-purpose processors.

- RETINA project aimed at designing better real-time tracking processors, using architectures inspired by natural vision, and technologies from telecom industry

- Encouraging first steps towards a future with *detector-embedded data reconstruction*

# Bibliography

G.Punzi and L.Ristori, Triggering on heavy flavors at hadron colliders, Ann.Rev.Nucl.Part.Sci. 60 (2010) 595-614

L. Ristori, An artificial retina for fast track finding, NIM A 453 (425-429), http:// inspirehep.net/record/539203

M. M. Del Viva, G. Punzi, The brain as a trigger system, http://arxiv.org/abs/ 1410.5123

M. M. Del Viva, G. Punzi, D. Benedetti, Information and Perception of Meaningful Patterns [PDF]

N. Neri et al., First results of the silicon telescope using an 'artificial retina' for fast track finding, Talk at ANIMMA15 conference, http://www.ipfn.ist.utl.pt/ANIMMA2015/

S. Stracka, A specialized processor for track reconstruction at the LHC crossing rate, Talk at Connecting The Dots Workshop 2015, https://indico.physics.lbl.gov/ indico/getFile.py/access? contribId=14&sessionId=2&resId=0&materialId=slides&confId=149

R. Cenci, Artificial retina processor for track reconstruction, Talk at Connecting The Dots Workshop 2015, https://indico.physics.lbl.gov/indico/getFile.py/access? contribId=2&sessionId=9&resId=0&materialId=slides&confId=149

A. Abba et al., Progress Towards the First Prototype of a Silicon Tracker Using an 'Artificial Retina' for Fast Track Finding, Poster at TWEP-P14, https://indico.cern.ch/ event/299180/session/7/contribution/64

A. Piucci, Reconstruction of tracks in real time at high luminosity environment at LHC, Master thesis, https://etd.adm.unipi.it/theses/available/etd-06242014-055001/.

D. Ninci, Real-time track reconstruction with FPGA at LHC, https://etd.adm.unipi.it/ theses/available/etd-11302014-212637/.

F. Spinella et al., The TEL62: A real-time board for the NA62 Trigger and Data Acquisition. Data flow and firmware design, IEEE Nucl. Sci. Symp. Conf. Rec., 1 (2014).

A. Abba et al., The artificial retina for track reconstruction at the LHC crossing rate, arXiv:1411.1281 [ICHEP 2014], https://inspirehep.net/record/1326137.

N. Neri, First prototype of a silicon tracker using an 'artificial retina' for fast track finding , PoS TIPP2014 (2014) 199 [TIPP2014], https://inspirehep.net/record/ 1315951.

A. Abba, The artificial retina processor for track reconstruction at the LHC crossing rate, JINST 10 (2015) 03, C03018 [WIT2014], https://inspirehep.net/record/1315154.

A. Abba et al, Simulation and performance of an artificial retina for 40 MHz track reconstruction, JINST 03-10 (C03008) [WIT2014], https://inspirehep.net/record/ 1314984.

A. Abba et al., A Specialized Processor for Track Reconstruction at the LHC Crossing Rate, JINST 9 (C09001) 2014 [INSTR14], https://inspirehep.net/record/ 1303542.

A. Abba et al., *The Readout Architecture for the Retina-Based Cosmic Ray Telescope*, Real Time Conference (RT), 2014 19th IEEE-NPSS, [IEEE-RT 2014] http://dx.doi.org/10.1109/RTC.2014.7097516.

A. Abba, et al., A retina-based cosmic rays telescope, Real Time Conference (RT), 2014 19th IEEE-NPSS, [IEEE-RT2014], http://dx.doi.org/10.1109/RTC. 2014.7097515, https://inspirehep.net/record/1367442.

A. Abba et al., A specialized track processor for the LHCb upgrade, CERNa-LHCb- PUB-2014-026 https://cds.cern.ch/record/1667587.