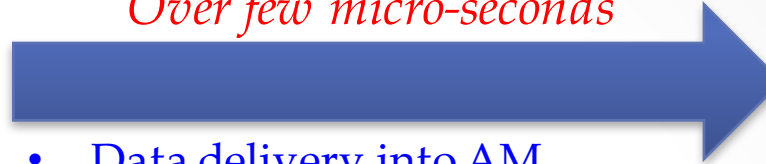
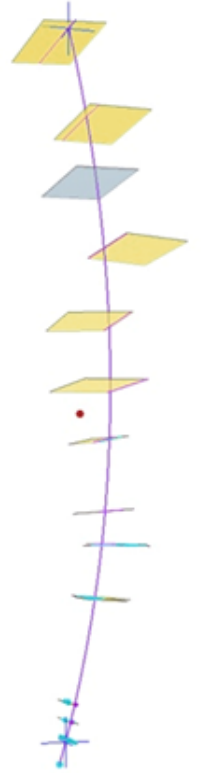


Associative Memory for HEP

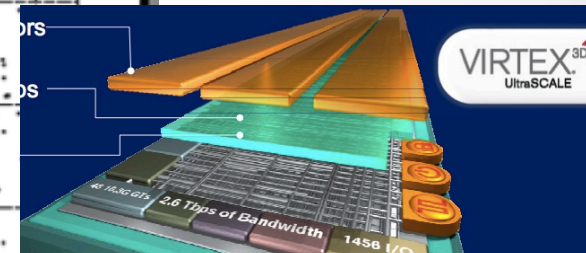
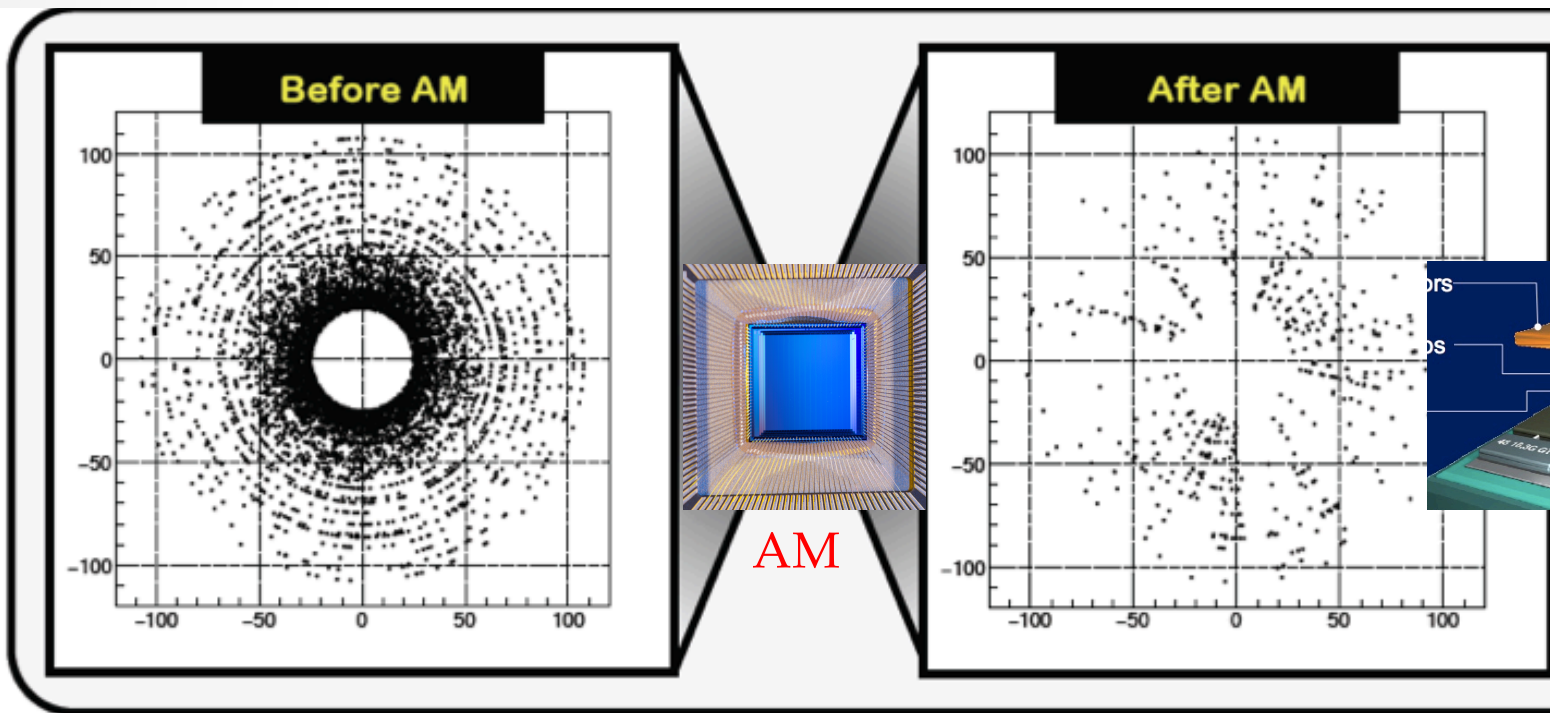
Over few micro-seconds



- Data delivery into AM
- Pattern Recognition (AM+FPGA)
- Track Fitting (FPGA)



Ted Liu, Fermilab (Oct 5th, 2017)



FPGA

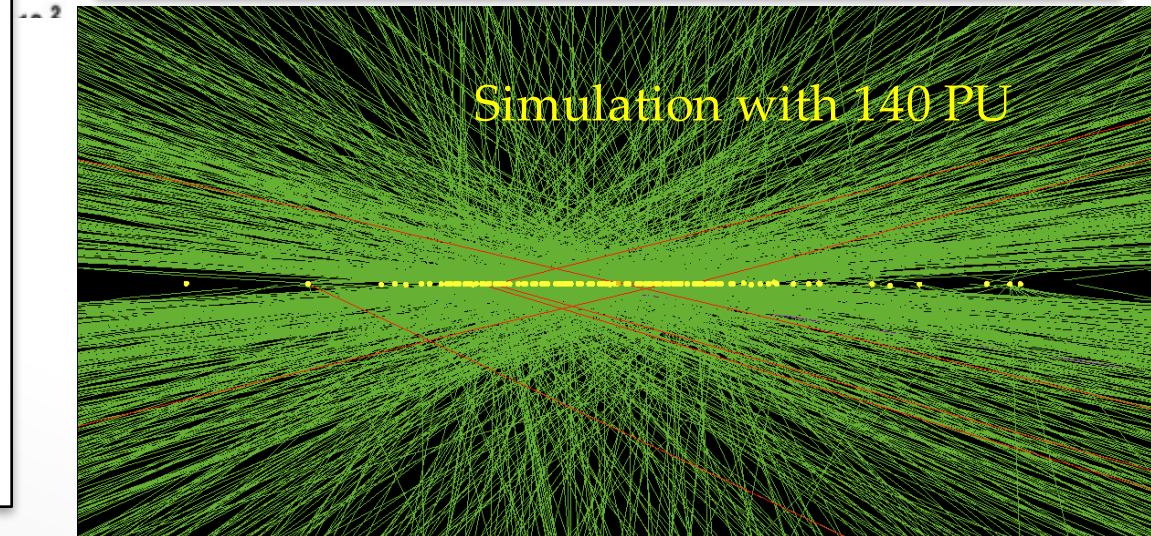
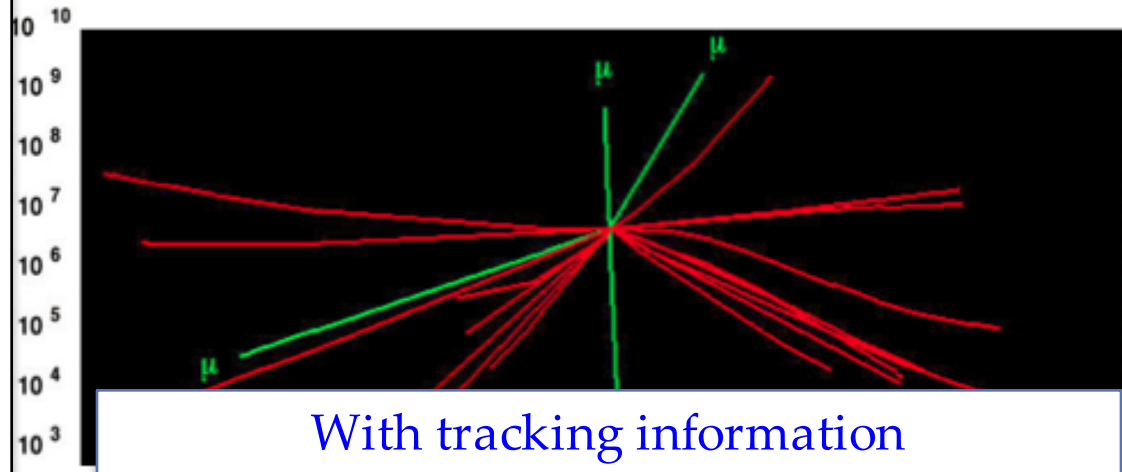
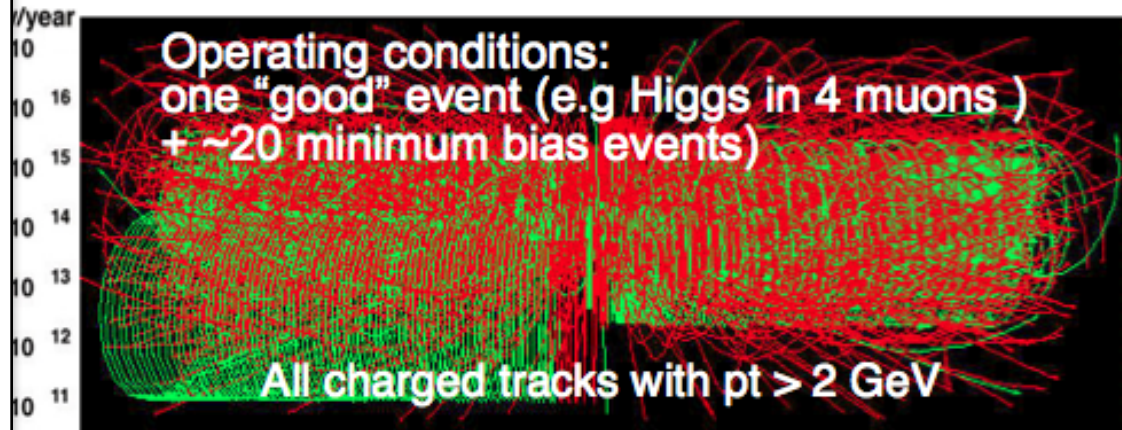
CMS L1 Tracking Trigger:

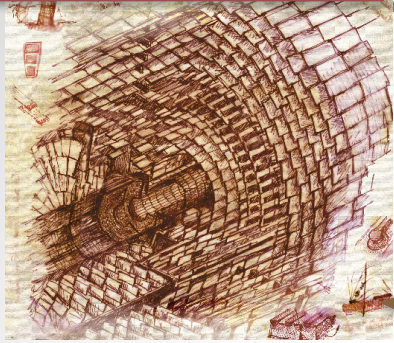
Will need to reconstruct charged particle trajectories “on-the-fly” for every beam crossing (25 ns, or 40 Million beam crossings per second), from an ocean of input data (bandwidth required to transfer data ~ 100Tb/s)

This requires extremely fast high bandwidth data communication as well as massive pattern recognition and track fitting power, with near zero latency (~ few μ s total)

This is challenging!

Pileup at HL-HC: ~ 200 (only 20 shown here)





Data transfer

Data
formatting

Partition detector into
trigger towers/sectors

Pick your favorite method:

- Associative Memory Approach*
- Hough Transformation*
- Tracklet-based*
- your choice here... (any new ideas?)*

• *Challenging issues*

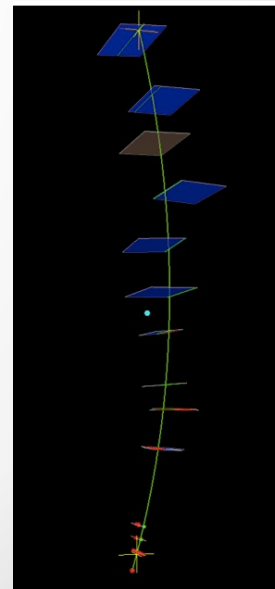
- (1) Data Reduction at detector/sensor stage*
- (2) Data Transfer (rad hard, high bandwidth, low power link)*
- (3) Data Formatting/Delivery*
- (4) Pattern Recognition (using AM as a filter)*
- (5) Track Fitting ... (track fitting is done in FPGA)*

Pattern
Recognition

Track
Fitting

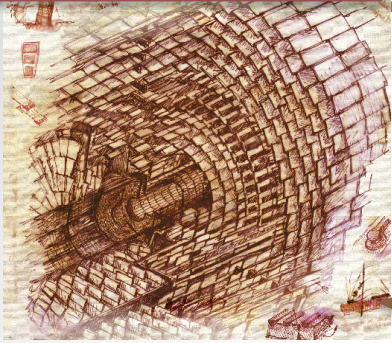
Finer pattern recognition

Associative Memory approach is a known HEP technique for silicon based tracking trigger for hadron collider, but so far only at Level 2 trigger. The challenge here is in the very short latency at Level 1. This has never been done in HEP before...



Detector design for triggering

AM + FPGA based Tracking Trigger



Data transfer

Data formatting

Partition detector into trigger towers/sectors

Associative Memory Approach

Pattern Recognition

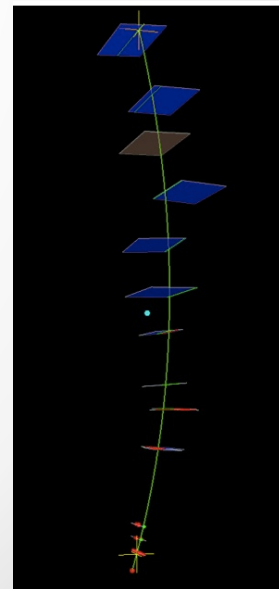
Finer pattern recognition

Track Fitting

Challenging issues

- (1) Data Reduction at detector/sensor stage
 - (2) Data Transfer (rad hard, high bandwidth, low power link)
 - (3) Data Formatting/Delivery
 - (4) Pattern Recognition
 - (5) Track Fitting ...
- } Tracking Trigger

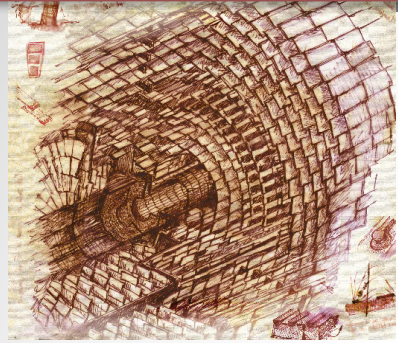
Core R&D for AM + FPGA approach:
 High performance Data Formatting/Delivery
 High performance Associative Memory (ASIC)
 system architecture development to board and chip specifications



Detector design for triggering

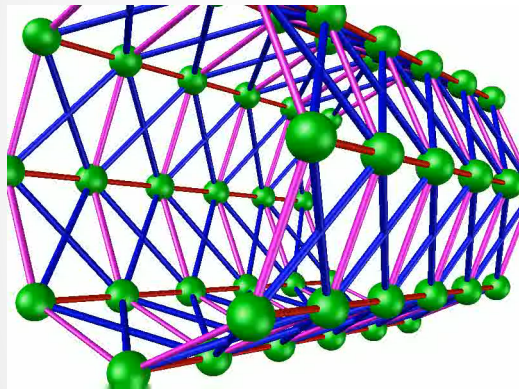
AM-based Tracking Trigger R&D for HL-LHC

VIPRAM in 3D

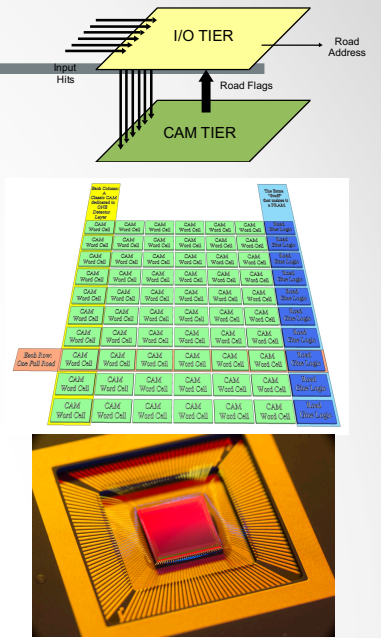
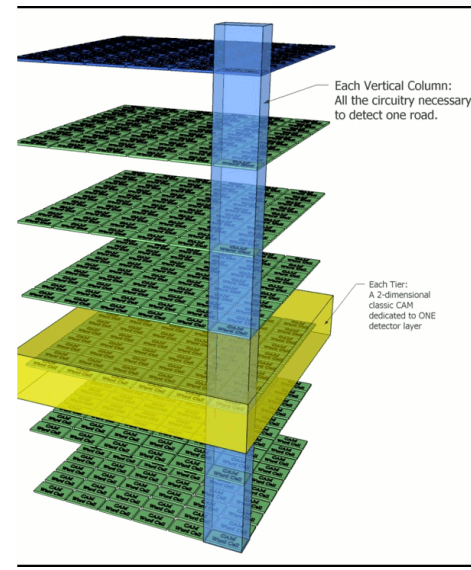


Data transfer

L1 Track Trigger architecture



Custom ASIC



Data formatting

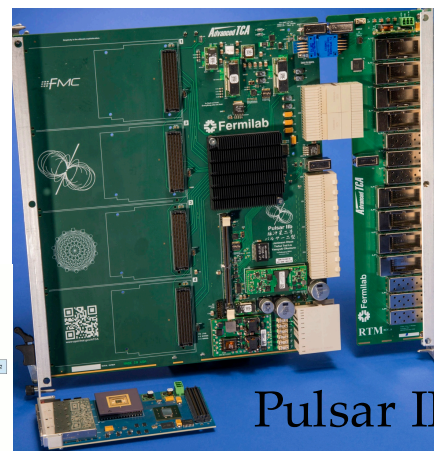
Pattern Recognition

VIPRAM: Vertically Integrated Pattern Recognition Associative Memory

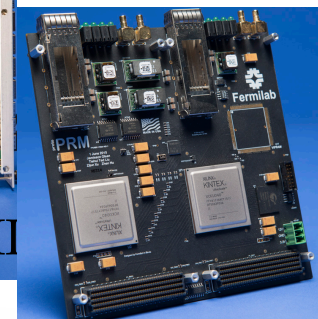
Track Fitting

Extensive simulation work

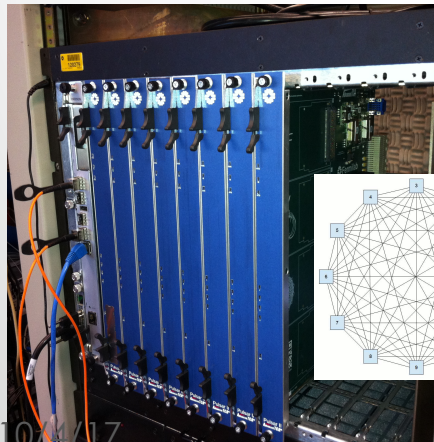
ATCA



Pulsar II

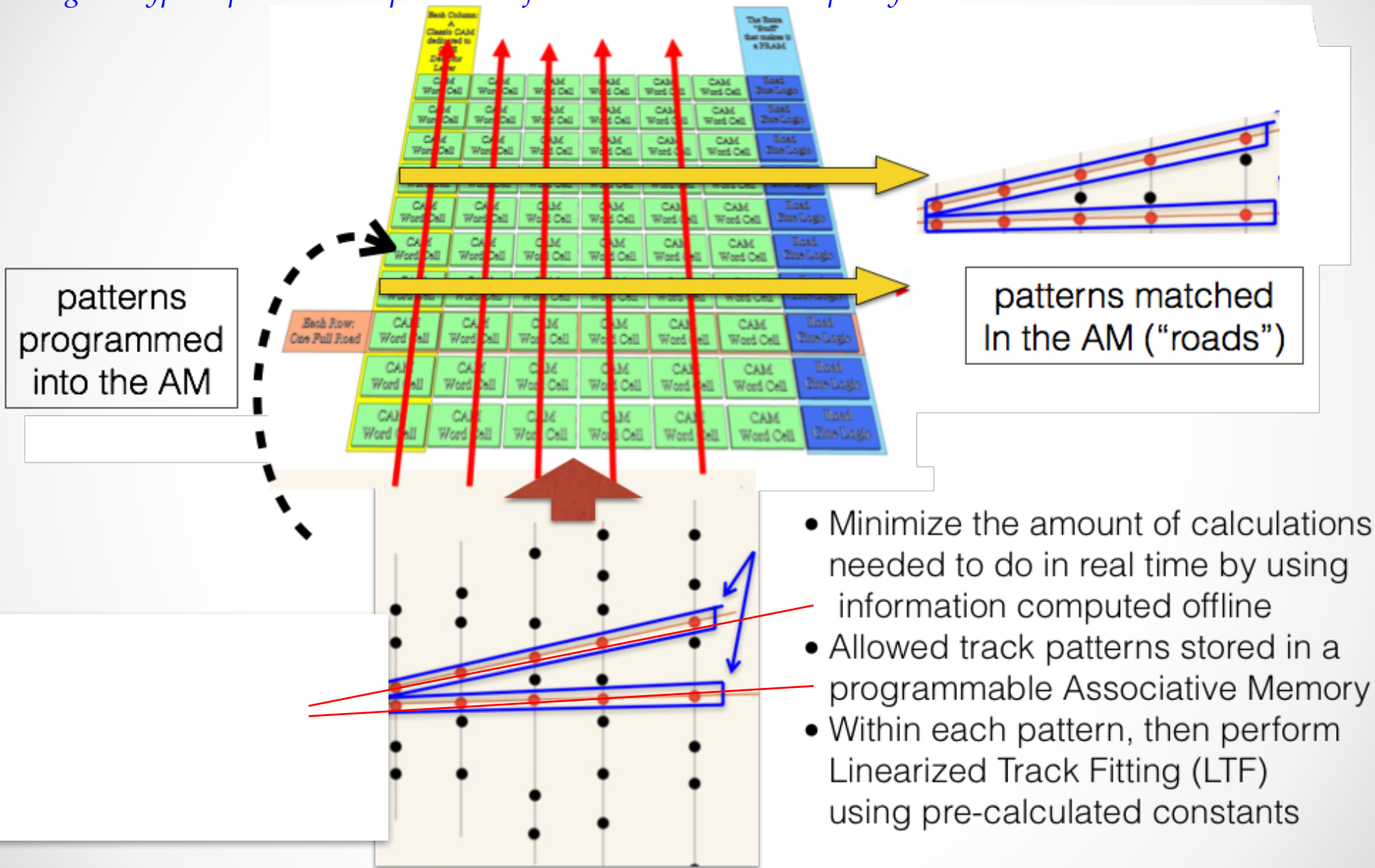


AM+FPGA approach
Pattern Recognition Mezzanine (PRM)

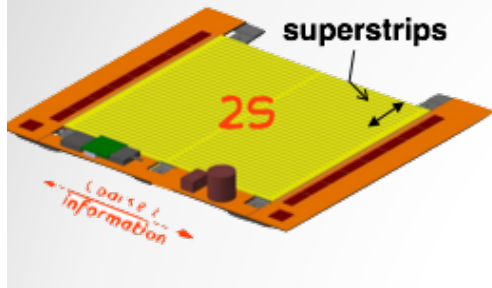


AM concept

Massive parallel processing to tackle the intrinsically complex combinatorics of track finding algorithms, *avoiding the typical power law dependence of execution time on occupancy*

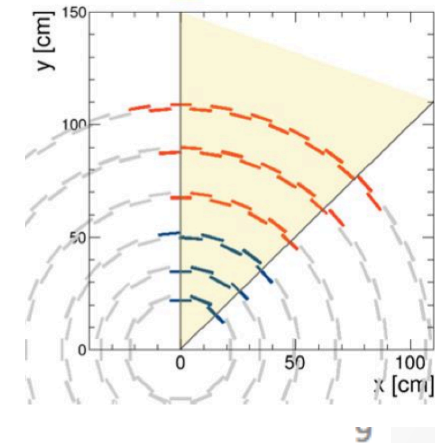
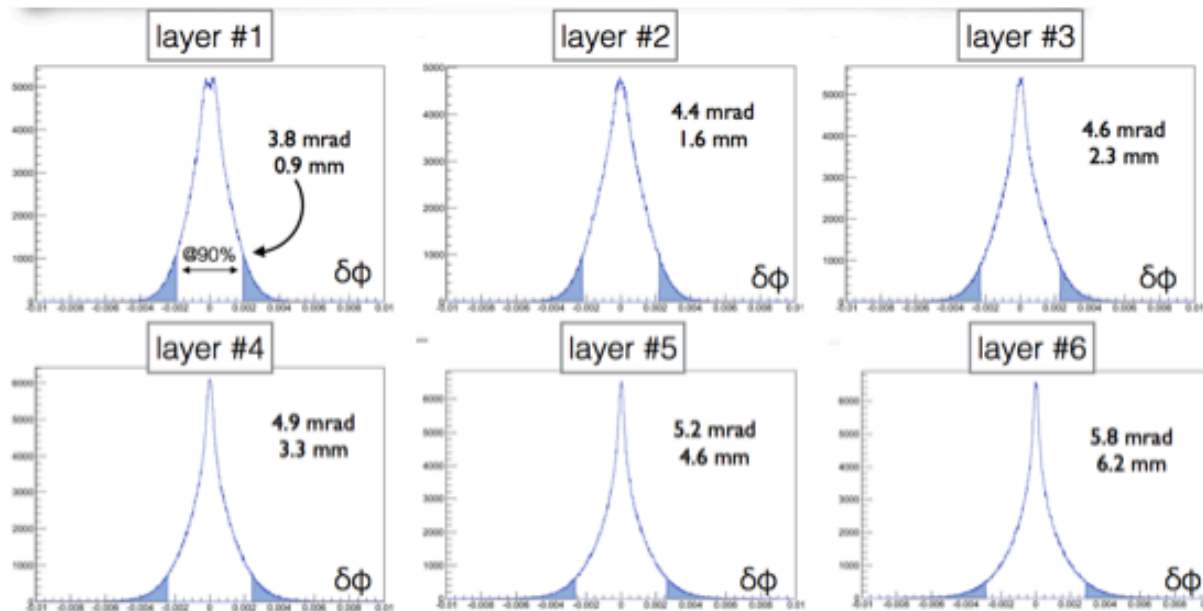
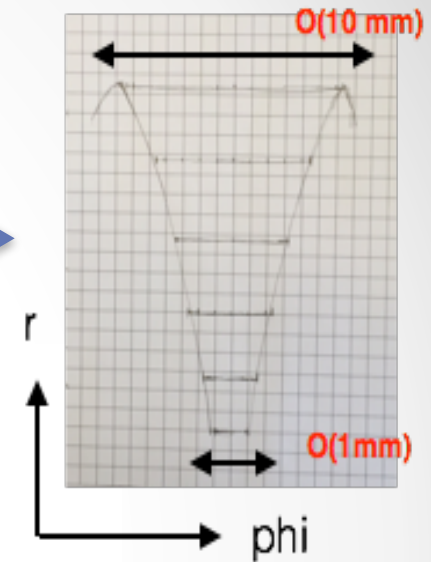


Tailoring the AM approach to CMS phase 2 tracker



What is the natural bin size (superstrip) for pattern recognition?

- use single muon gun in $\{\eta, \phi, P_T\}$ phase space
- find natural spread in ϕ due to multiple scattering etc
- "fountain-like" pattern configuration
- then study slicing in z



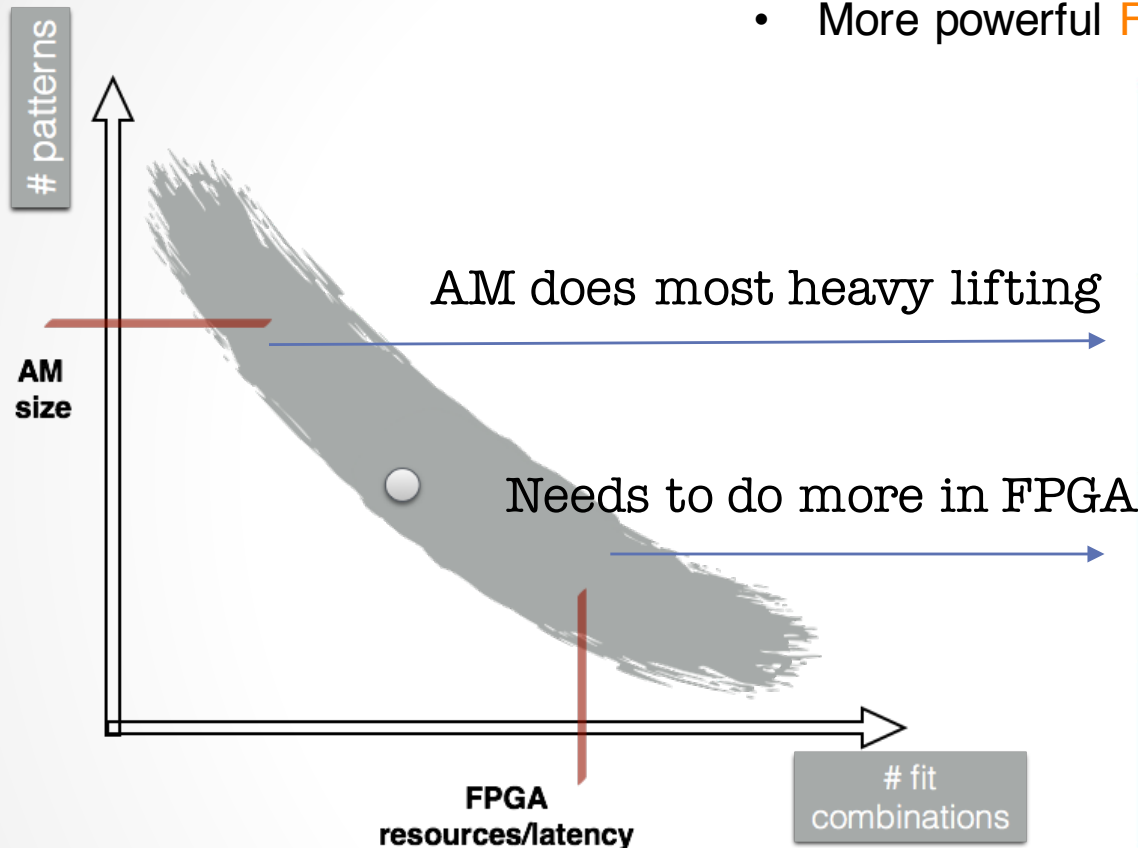
→ This sets the scale of the size of AM patterns, as well as the size of total pattern bank

AM+FPGA system optimization

Pattern Recognition

Interplay between the AM and the FPGA:

- More powerful AM => less demand on the FPGA
- More powerful FPGA => less demand on the AM



| AM size | <roads/evt> | <fits/evt> |
|-------------------|-------------|------------|
| ~1M sf1_nz8 | 19.9 | 57 |
| 256K sf0.5_nz1 | 182 | 790 |
| 64K sf0.8_nz1 | 237 | 1431 |

Track Fitting

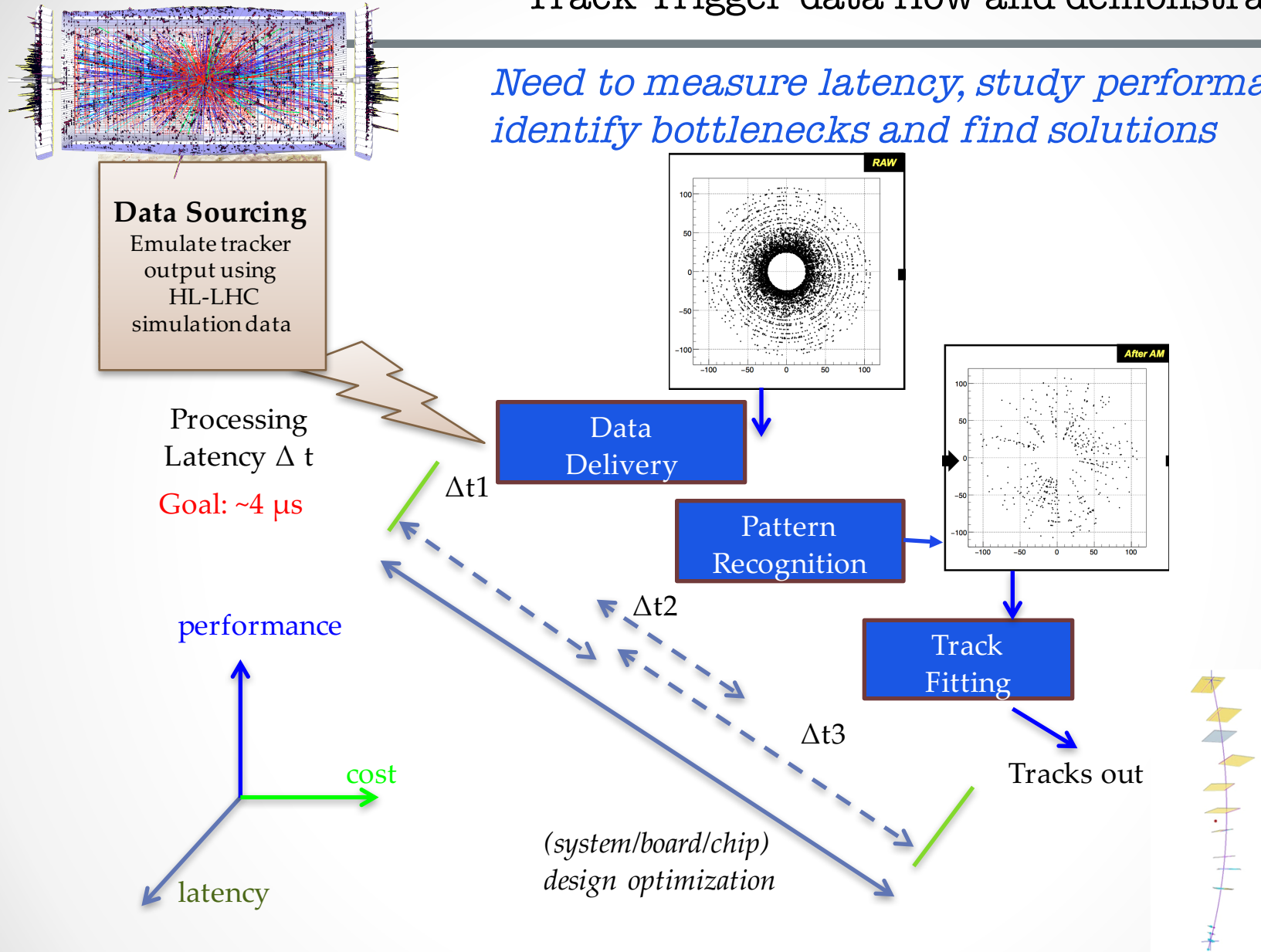
A good working point: as a balance between the size of the AM and the number of roads and fits to perform per beam crossing

The chip specification needs to be studied extensively ... from system point of view

ttbar + PU200
Trigger tower:
1/8 in phi
1/6 in eta

Track Trigger data flow and demonstration

Need to measure latency, study performance and identify bottlenecks and find solutions



The AM chip specification has to be determined from system point of view

New design required: from system to chip

- System architectural design (2012-2013) for CMS L1
 - *System Architecture fully studied/established in 2013*
 - *Concept of demonstrator developed in 2013*
- System level performance requirements simulation (2013 – present)
 - *Preliminary results on system requirements:
AM patterns <~1M patterns/Tower, or ~128K/chip
Roads fired per trigger tower ~100, stub combinations/tower ~300*
- High performance ATCA board R&D (2011 - present)
 - The initial R&D goal: ~1Tbps ATCA board with large mezzanines for AM+FPGA
 - ***This goal has been achieved in 2014 with Pulsar 2b design***
 - ***High performance Pattern Recognition Mezzanine cards working in 2015
(UltraScale FPGA based with VIPRAM emulated in FPGA)***
 - ***System level demonstration to fully corner the specifications of the chip (2016)***
- High performance Associative Memory R&D at FNAL
 - Vertically Integrated Pattern Recognition Associative Memory) 2011- present
 - The first 2D VIPRAM prototype fully tested in 2014
 - *The first 2-tier VIPRAM_L1CMS dedicated for L1CMS (submitted 2016)*
 - *The FPGA implementation/emulation of VIPRAM_L1CMS ASIC (2015)*
 - *The first 3D VIPRAM-3D design submitted 2016*

Top down



Bottom up

VIPRAM R&D Status Summary

Initial Idea: ~2010

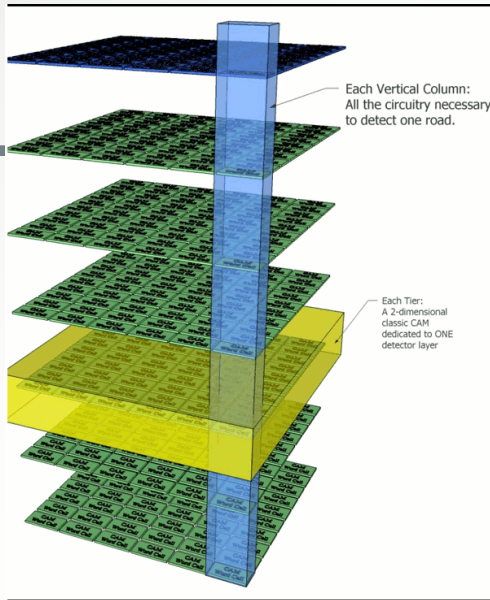
VIPRAM concept paper: 2011

CDRD award: 2012

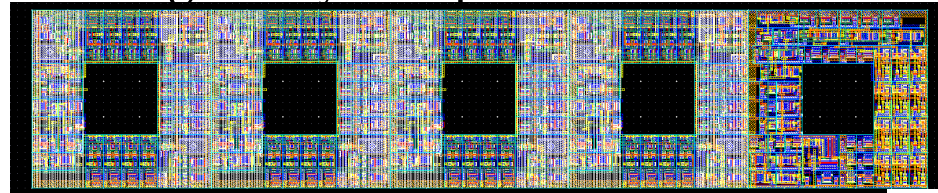
First pure 2D Design submission: 2013

First ProtoVIPRAM00 successfully tested: 2014

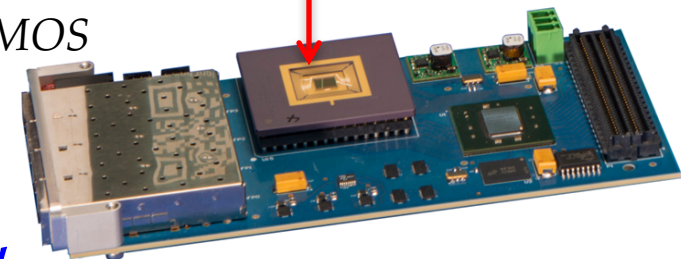
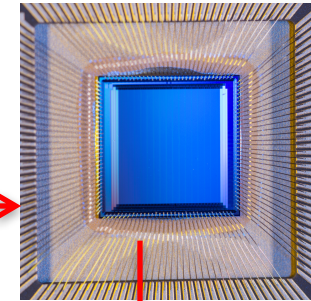
→ **Design building blocks are ready for 3D stacking**



2D design fully compatible with 3D stacking



CAM cell size: 25 μm x 25 μm , @ 130nm GF CMOS



Main work in 2015:

The generic 3D multi-tier design (protoVIPRAM01) ready

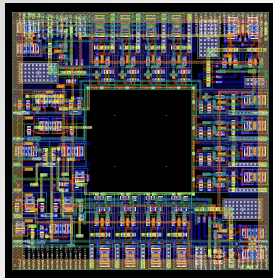
The 2-tier design for CMS L1 tracking trigger (protoVIPRAM02) ready

Both designs submitted in 2016, wafers received later 2016, 3D processing on going.

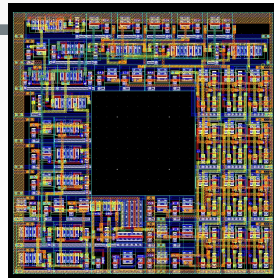
Main ASIC designer: Jim Hoff

Extensive simulation and optimization work done

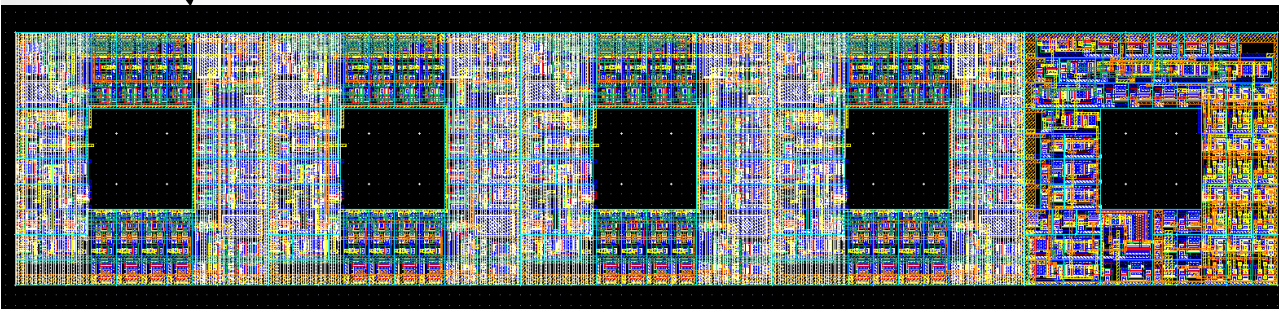
much of the work done by engineering students



CAM cell

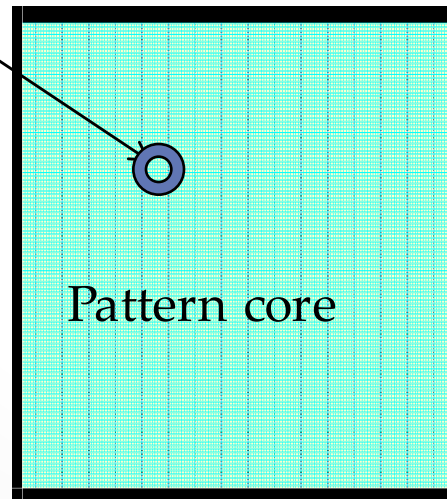


Majority Logic cell

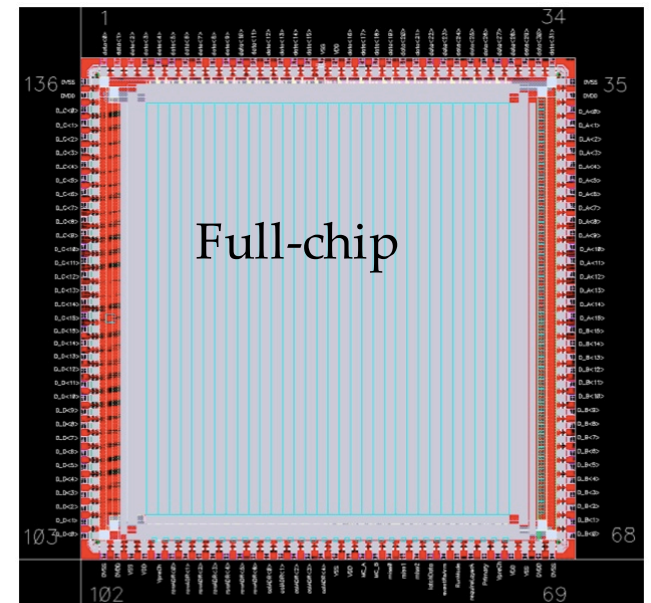


Single Pattern

Power delivery
Signal timing ...



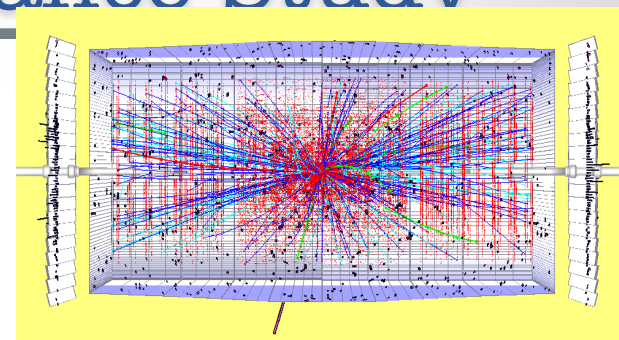
Pattern core



Full-chip

Chip Testing and Performance Study

- Use data and patterns from 140 Pile Up CMS AM simulation for HL-LHC to study the chip performance in close to “battle field” conditions

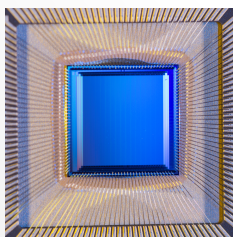
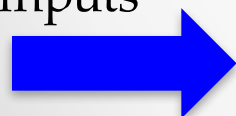


| Freq(MHz) | Found/Expected Avg. |
|-----------|---------------------|
| 50 | 100 |
| 66 | 100 |
| 71 | 100 |
| 76 | 100 |
| 83 | 100 |
| 90 | 100 |
| 100 | 100 |
| 111 | 100 |
| 125 | 98.85057471 |

| layer | <Nstubs> |
|-------|----------|
| 1 | ~90 |
| 2 | ~60 |
| 3 | ~50 |
| 4 | ~40 |
| 5 | ~40 |
| 6 | ~40 |

protoVIPRAM00 shows 100% performance around ~111 MHz, design target was 100 MHz
Power consumption ~ 300mW, as expected

stub inputs



Fired patterns



Extensive stress test performed to study chip performance in extreme conditions.

AM+FPGA approach

STUDY OF CAM CELL TOPOLOGIES AND DESIGN AND SIMULATION OF PROTOVIPRAM

Submitted in fulfillment of the requirements of First Degree BITS C422T
August 2012 - November 2012

By

Siddhartha Joshi

URA VISITING SCHOLARS AT FERMILAB PROJECT REPORT

GRANTEE INFORMATION

Name (last, first, MI): Xia, Wenbo

Date of the Award: 09/21/2012 Date of (expected) Completion:

Home Institution: Southern Methodist University

Fermilab Sponsor: Particle Physics Department

Title of Project: Power & Thermal Modeling & Management in 3D IC Technologies

Project Accomplishments:

- Power measurement for the proto-VIPRAM chip by circuit simulation
- Power model for the proto-VIPRAM chip; AM+FPGA approach
- Thermal model for the proto-VIPRAM chip

DESIGN, SIMULATION AND TESTING OF PROTOVIPRAM

Submitted in fulfillment of the requirements of First Degree BI
January 2013 - April 2013

By

Siddhartha Joshi

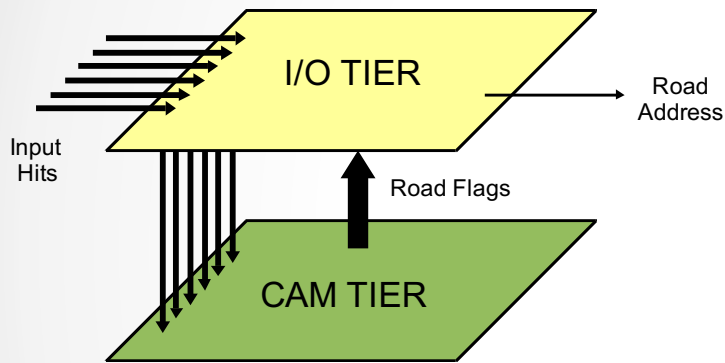
The project has attracted engineering school students. Few engineering school thesis, work done at Fermilab as intern or URA fellowship.

Second URA fellowship for VIPRAM project: Northwestern EE school Ph.D student: Dawei Li (2014-2015)

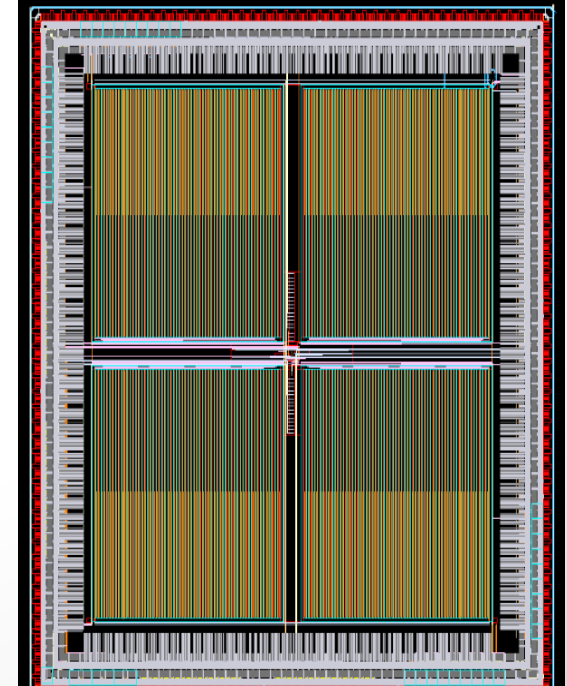
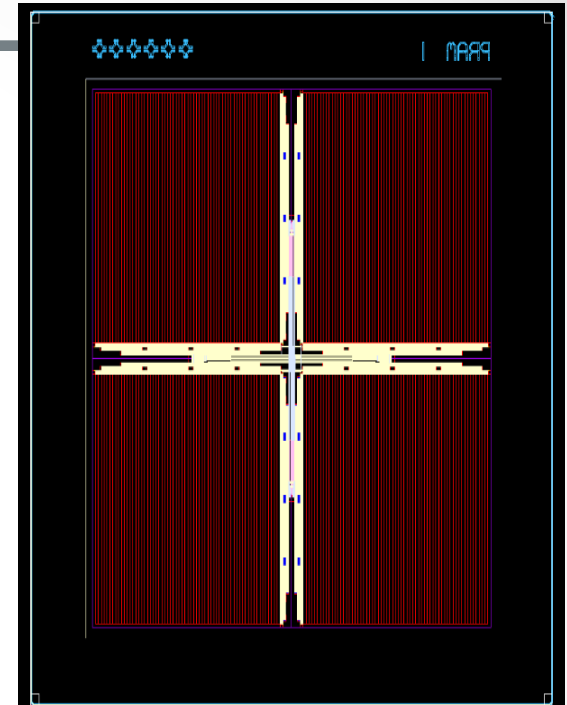
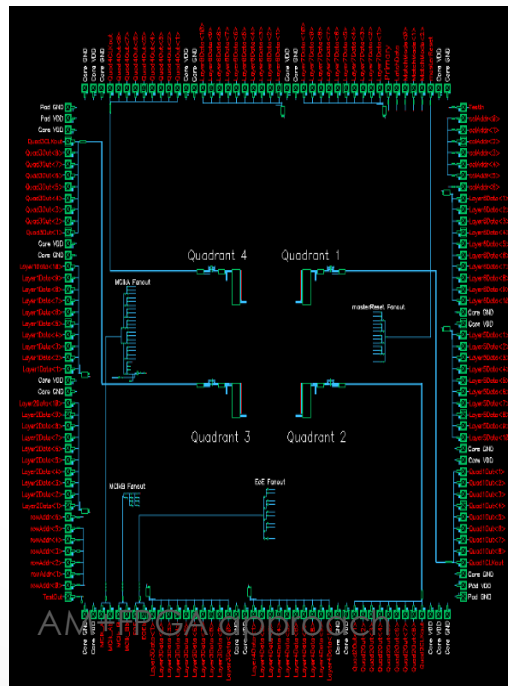
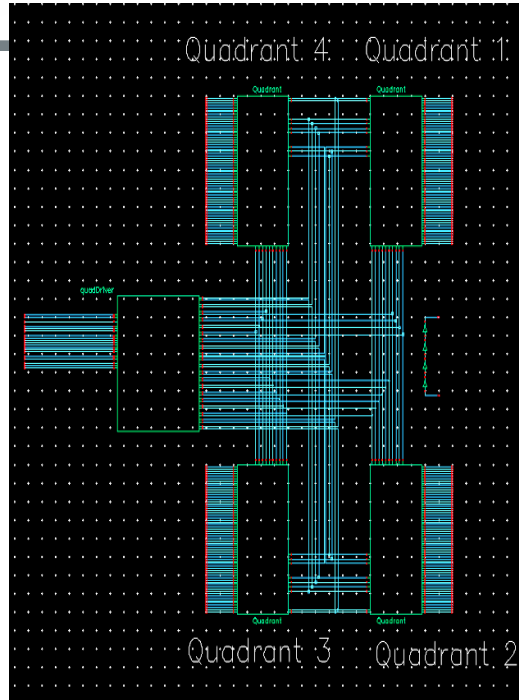
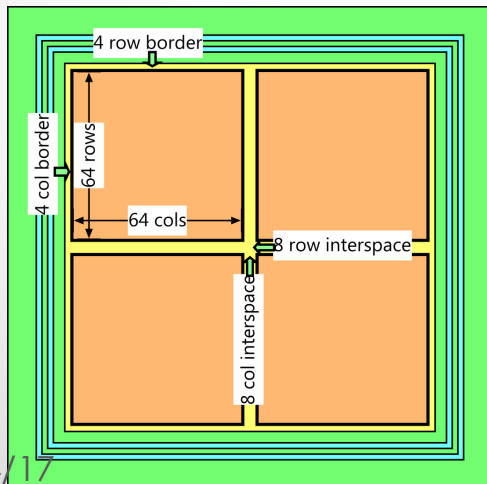
Li, Dawei et al, "A methodology for power characterization of associative memories," in *Computer Design (ICCD), 2015 33rd IEEE International Conference on*, vol., no., pp.491-498, 18-21 Oct. 2015

The L1CMS “CAM Tier and I/O Tier” Schematic and Layout

A new 2-tier design for CMS L1 tracking trigger



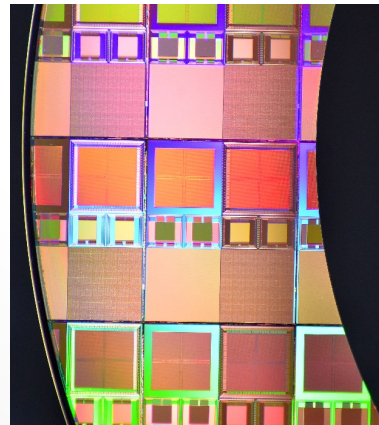
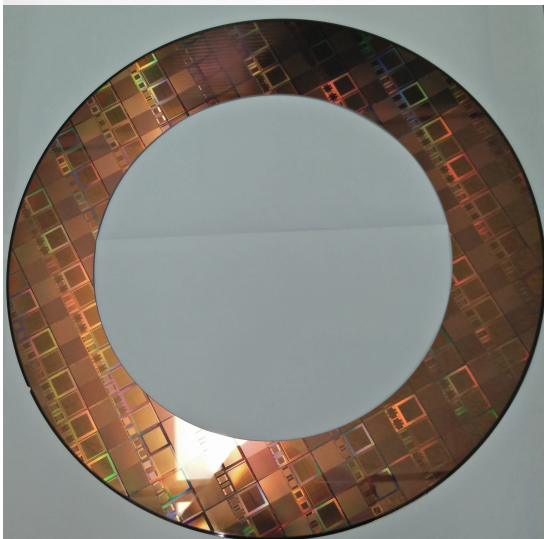
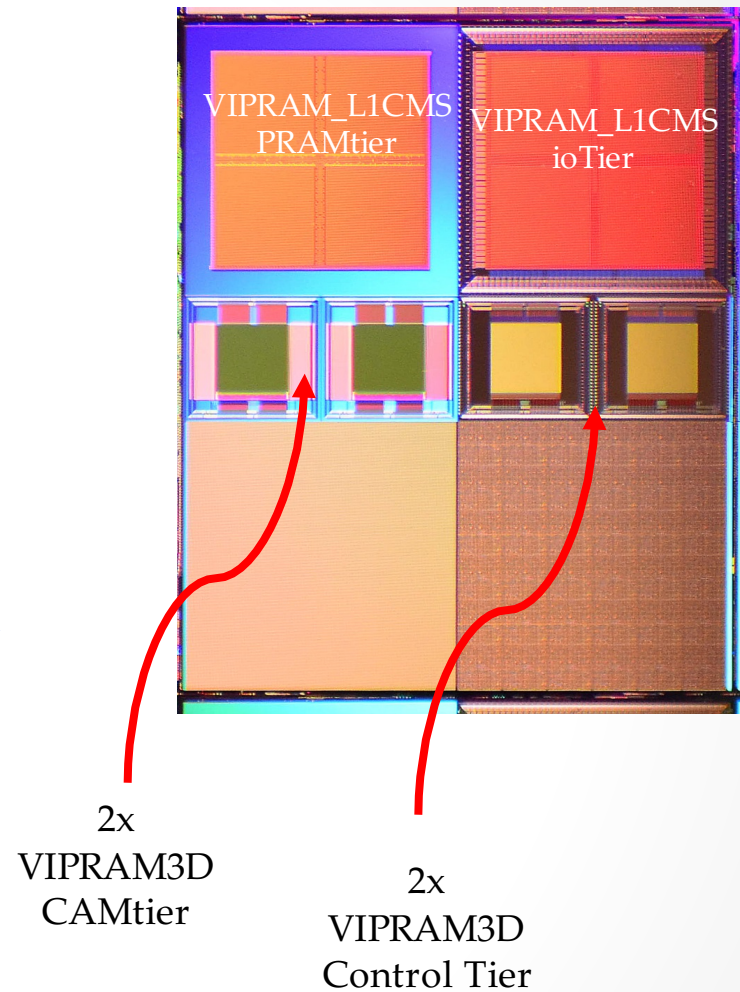
submitted in 2016



10/4/17

Status

- March 2016: the designs submitted (MPW)
- Aug. 2016: wafers were received from Global Foundries
- Sept. 2016: 3D fabrication began at Novati.
- final chip delivery is delayed, due to new 3D assembly line



Annulus remaining after 8 inch center was cored for 3D fabrication

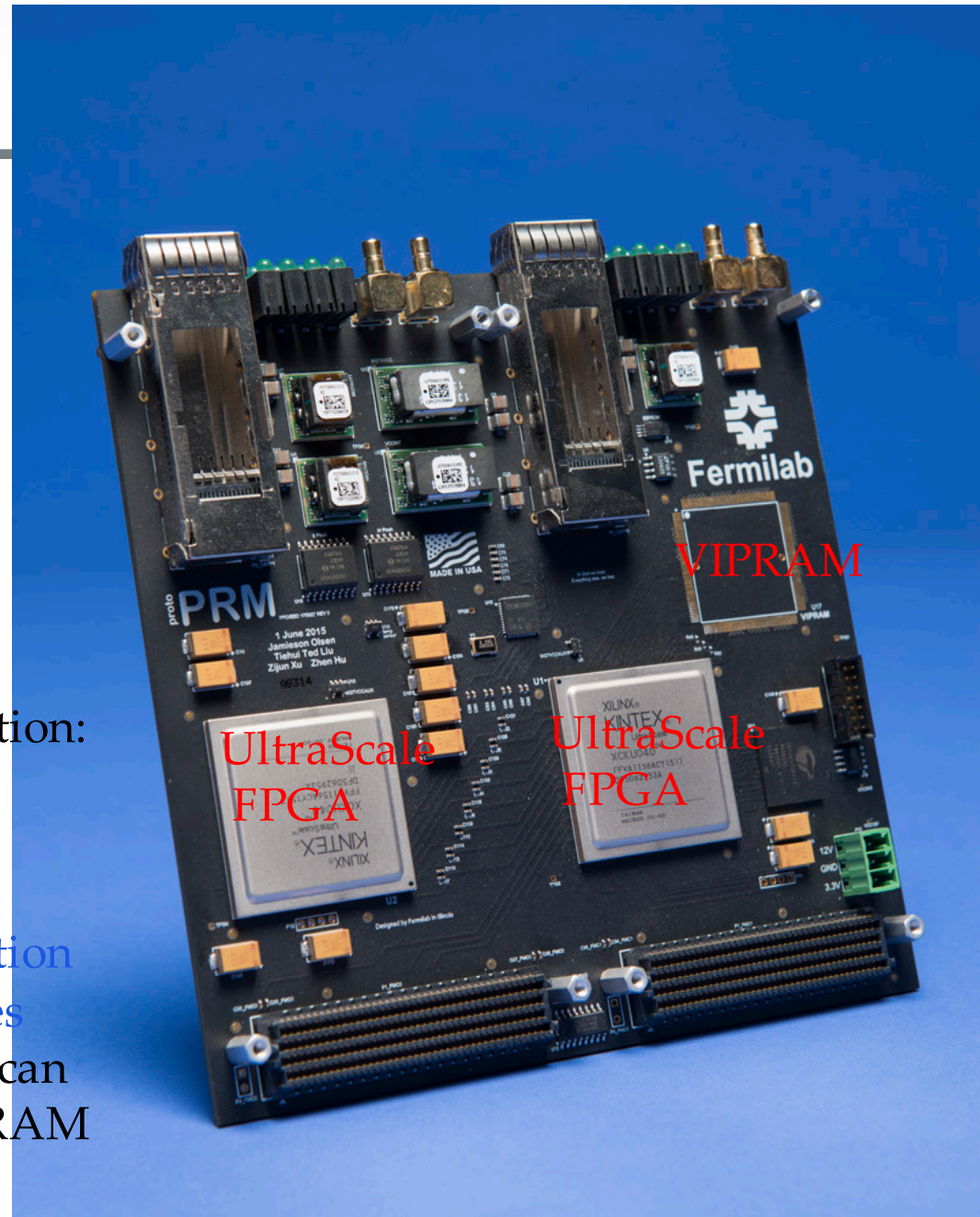
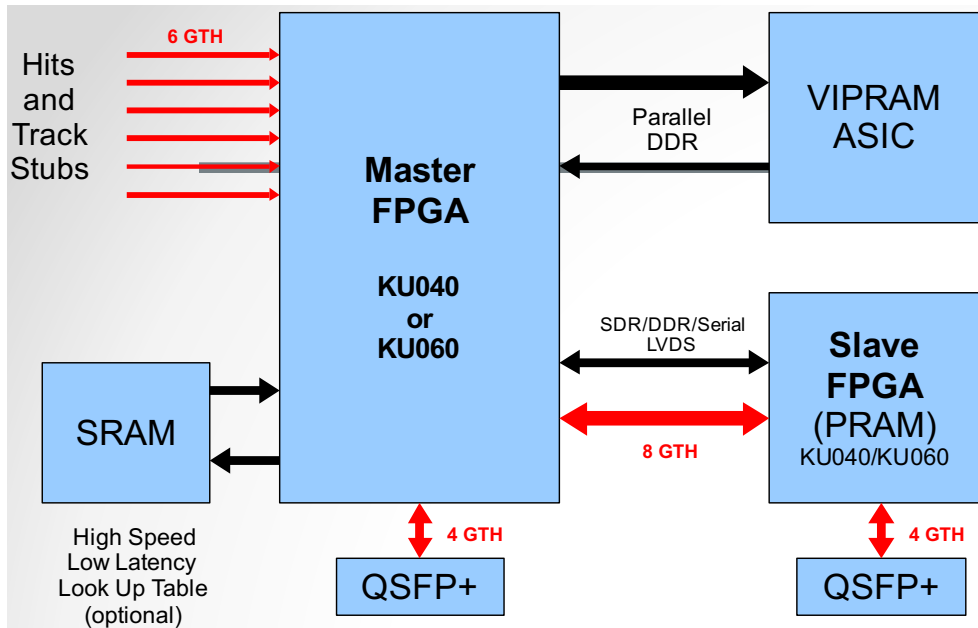
Collaboration with SMU engineering school

- W. Xia, "Thermal Analysis of the proto-VIPRAM2D chip", SMU MS Thesis (URA, 2014)
- W. Xia, T. Zhang, P. Gui, T. Liu and J. Hoff, "Thermal Analysis for the proto-VIPRAM00 chip", **Journal of Instruments, JINST_10_C03015, 2015.**

- *SMU also works with FNAL on other ASIC projects:*
- T. Liu, X. Wang, R. Wang, G. Wu, T. Zhang and P. Gui, "Temperature Compensated Triple-Path PLL with KVCO Non-Linearity Desensitization Capable of Operating at 77 K," **IEEE Transactions on Circuits and Systems I: Regular Papers (TCAS-I), in press.**
- J. R. Hoff, G. W. Deptuch, G. Wu, and P. Gui, "A Cryogenic Lifetime Studies of 130 nm and 65 nm nMOS Transistors for High-Energy Physics Experiments", **IEEE Transactions on Nuclear Science, Vol. 62, No. 3, June 2015.**
- Wu, G. W. Deptuch, J. Hoff, and P. Gui, "Degradations of threshold voltage, mobility and drain current and the dependence on transistor geometry for stressing at 77 K and 300 K", **IEEE Transactions on Device and Materials Reliability, Vol. 14, Iss. 1, pp. 477-483, March 2014.**

Collaboration with Northwestern engineering school

- "A Content Addressable Memory with Multi-Vdd Scheme for Low Power Tunable Operation", to be published in IEEE MWSCAS 2017, by Siddhartha Joshi, Dawei Li, Seda Ogrenci-Memik, Grzegorz Deptuch, James Hoff, Sergo Jindariani, Tiehui Liu, Jamieson Olsen, Nhan Tran.
- "A methodology for power characterization of associative memories," 33rd IEEE Int. Conf. Comput. Des, 2015. by Dawei Li, Siddhartha Joshi, Seda Ogrenci-Memik, James Hoff, Sergo Jindariani, Tiehui Liu, Jamieson Olsen, Nhan Tran.
- "Design and Testing of the first 2D prototype VIPRAM", T. Liu, G. Deptuch, J. Hoff, S. Jindariani, S. Joshi, J. Olsen, N. Tran and M. Trimpl, JINST 10, no. 02, C02029 (2015), doi:10.1088/1748-0221/10/02/C02029.
- J. Hoff et al, "VIPRAM_L1CMS: a 2-Tier 3D Architecture for Pattern Recognition for Track Finding", Fermilab Technical Publication CONF-16-690-PPD, submitted to IEEE NSS proceedings, 2016.
- G. Deptuch et al, "Performance Study of the First 2D Prototype of VIPRAM", Fermilab-PUB-17-385-CMS-E-PPD, submitted to TNS, <https://arxiv.org/abs/1709.08303>
- Both students, Sid Joshi and Dawei Li, got Intel job offers after working on VIPRAM.



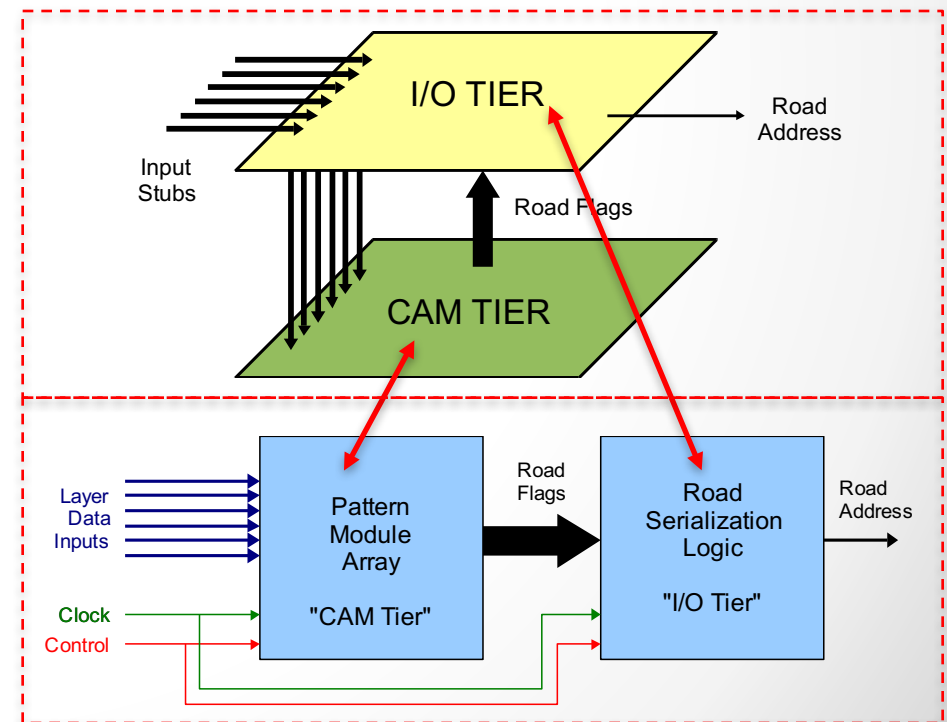
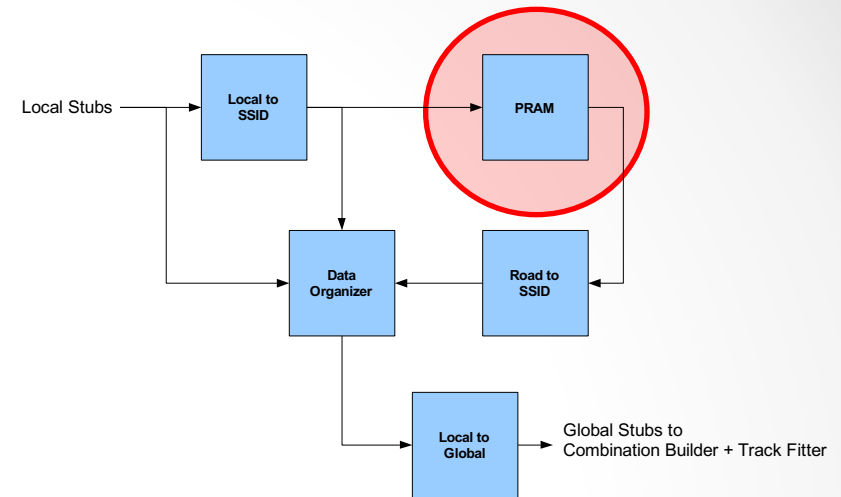
Pattern Recognition Mezzanine (protoPRM)

Designed as a development platform for demonstration and VIPRAM optimization:

- The slave FPGA will act as AM first
 - For initial demonstration
 - To fully optimize the AM specification to minimize the ASIC design cycles
- With VIPRAM loaded, the slave FPGA can
 - Act as AM and compare with VIPRAM performance side by side
 - Or be dedicated for track fitting

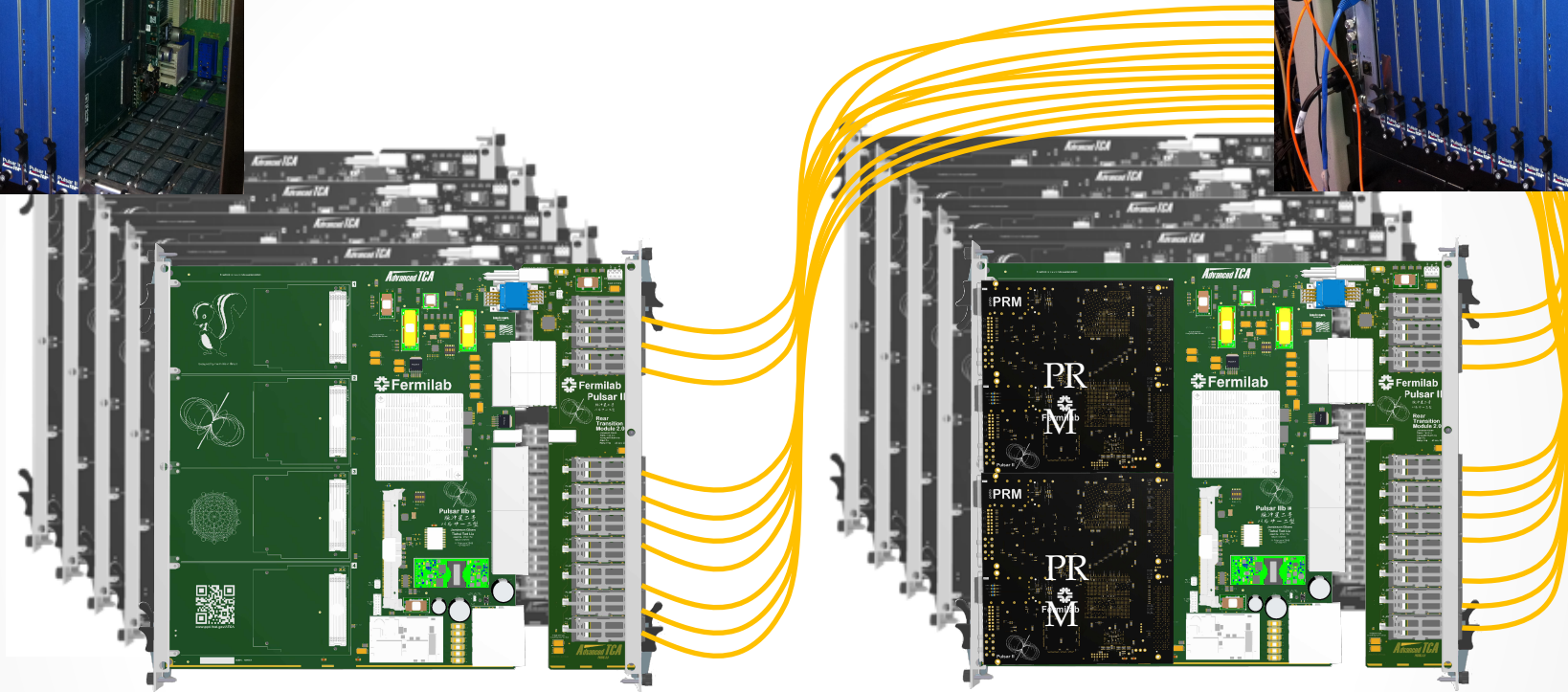
AM in FPGA: Overview

- AM in FPGA: very **closely follows** the AM ASIC (chip) design
 - **Match two silicon tiers** in ASIC with two modules in FPGA firmware
 - CAM Tier -> a 2D array of Patterns
 - I/O Tier -> input and output of fired roads
 - **Pipelined operation**
 - CAM tier: processes pattern matching with stubs for current event N
 - I/O tier: outputs road addresses for event N-1 at the same time
 - With KU60 FPGA, can implement 4k patterns closely follow the ASIC functionalities and operation
 - Can implement much more if memory W/R operation is not included



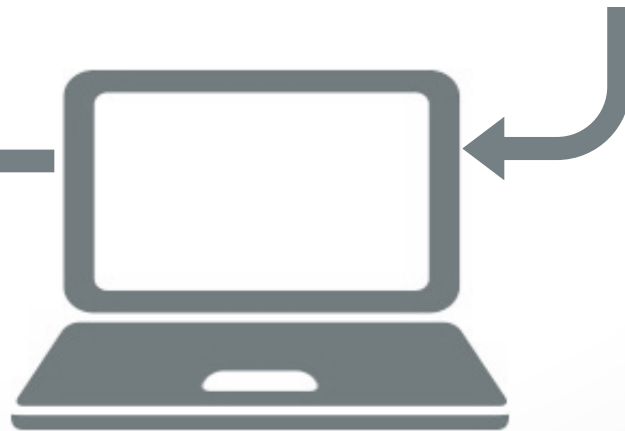
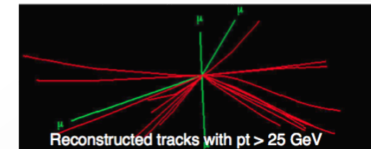
Trigger Tower/crate system level

vertical slice demo (2016-)



Data Source Boards
CMS silicon tracker
data emulation

Pattern Recognition Boards



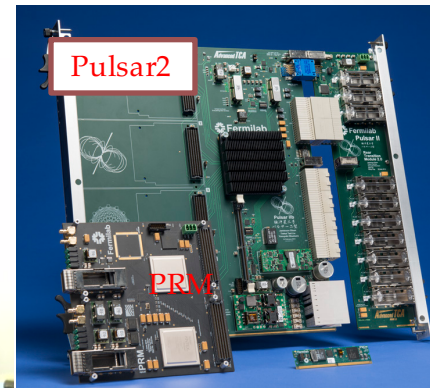
AM+FPGA approach

CMS L1 Tracking Trigger System Demonstration at FNAL/LPC

4.8 Tbps
fiber connections



TTCCi



One shelf contains Pulsar 2b
Data Source Boards (DSB)
Emulating ~400 detector modules
(one trigger tower)

One shelf contains Pulsar 2b as
Pattern Recognition Boards (PRB)

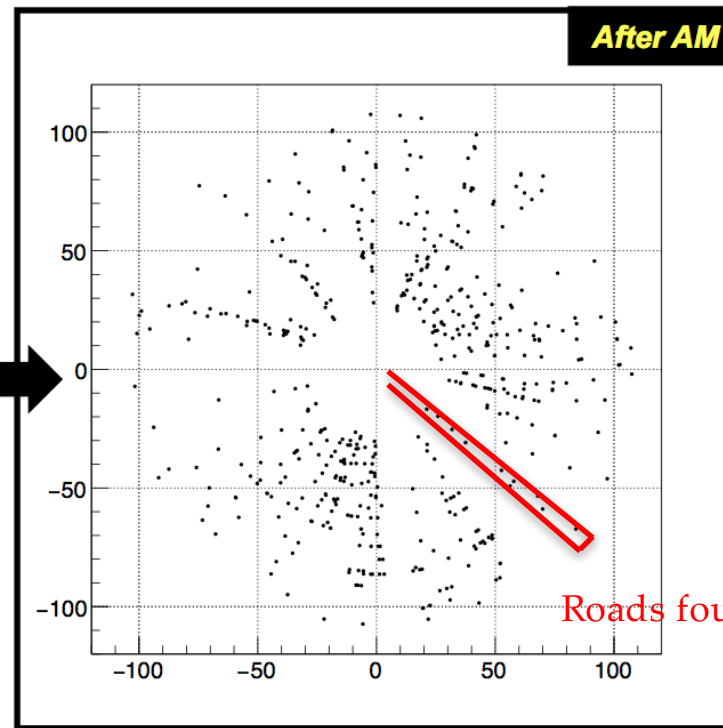
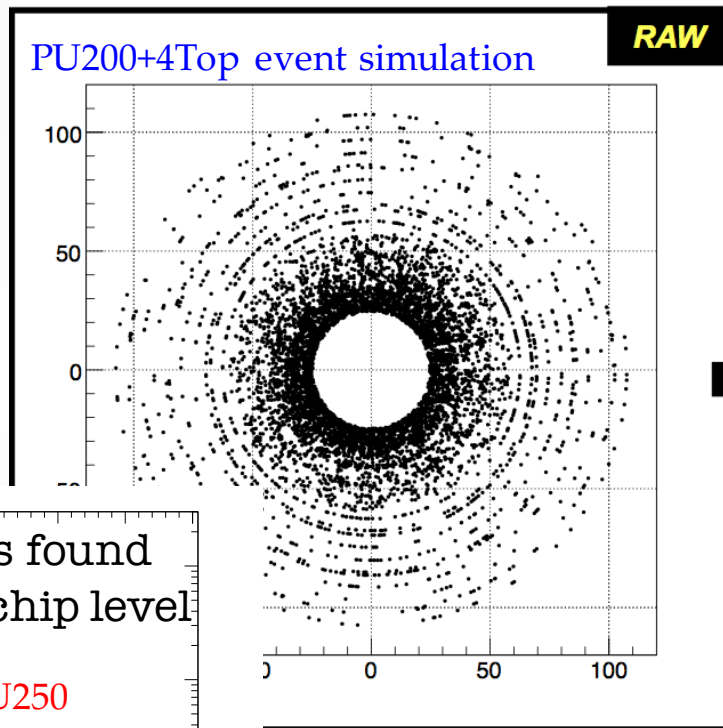
Each PRB can host two
Pattern Recognition
Mezzanines (PRM)

VIPRAM_L1CMS emulated in FPGA

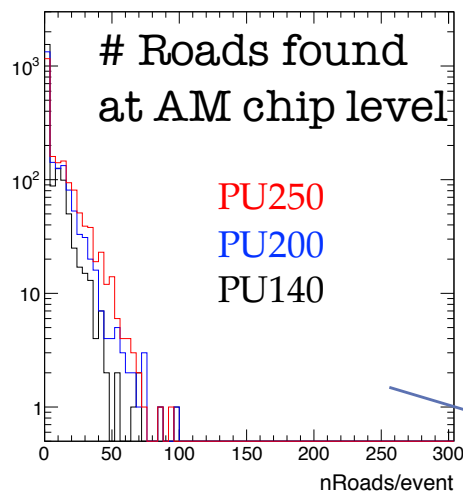
Close collaboration:
FNAL, Northwestern, U. Florida, Texas A&M, Brazil(SPRACE/UERJ) and China (Peking)
With FULL support of LPC

Summary of Data Delivery and PR latency

From DTC input ($t=0$), at $t=1850$ ns the first road with full resolution stubs are ready for downstream track fitting.
in other words, the picture on the right is emerging ...



Roads found ready for TF

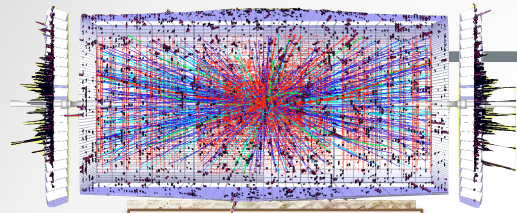


A factor of ~ 10 reduction of occupancy, and *more importantly*, stubs/hits are organized in found roads ("hits of interest"), making the track fitting task inside FPGA much easier...

At AM chip level, at PU250, on average $\sim 10/250k$ patterns fires, $< 10^{-4}$

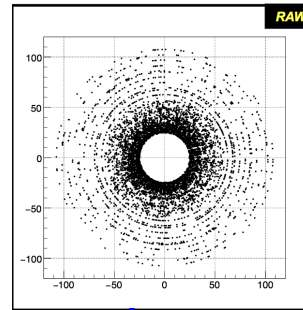
Track Trigger data flow and demonstration

Measured latency at different stages:



Data Sourcing
Emulate tracker output using HL-LHC simulation data

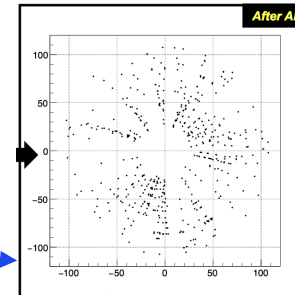
Processing Latency Δt
Goal: $\sim 4 \mu s$



Data delivery latency measured: $1.2 - 1.7 \mu s$

Data Delivery

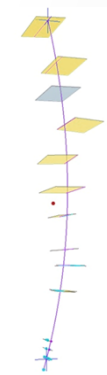
Pattern Recognition



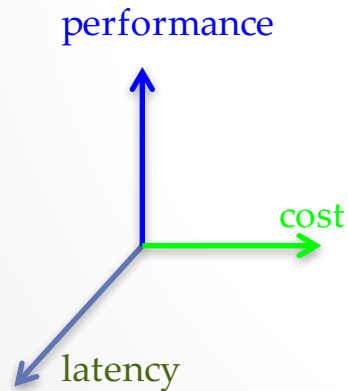
Pattern Recognition: 1st road out at $1.85 \mu s$

Track Fitting

Tracks out



Track Fitting finishes at $2.53 \mu s$



(system/board/chip) design optimization

Group photo taken at weekly meeting at FNAL/LPC



***Close collaboration:
FNAL, Northwestern, U. Florida, Texas A&M,
Brazil (SPRACE/UERJ) and China (Peking).
With FULL support of LPC***

One Comment

*For associative memory approach, it should be possible to naturally include **timing information** in the future to potentially and significantly improve the pattern recognition power (**pattern recognition in both space and time**). In addition, for certain applications where the pattern density requirement is not high, AM can be implemented in FPGAs.*

Summary

- Developing a high performance pattern recognition chip for L1 tracking trigger at HL-LHC is much more involved than just ASIC design itself
 - Much of the work is in the specifications/demonstration of the chip
 - This has to be done from (new) system point of view
- What has been achieved for chip specifications
 - System architecture fully developed
 - Extensive simulation for system optimization to define the pattern density requirement for the AM chip
 - System level demonstration to fully develop and corner the system interface specifications of the AM chip (applies to both 2D and 3D designs)
- 3D based VIPRAM chip designs have been extensively simulated and optimized:
 - 2D prototype chips have been successfully tested and work as expected
 - The new 3D designs have been submitted and are in 3D processing (delayed)
 - Extensive studies on 3D related power /thermal analysis done and still on going, knowledge & experience to be gained will benefit other future 3D projects

***The R&D effort has benefitted significantly
from working closely with many university groups...***

Backup slides

High Performance AM chip development road map

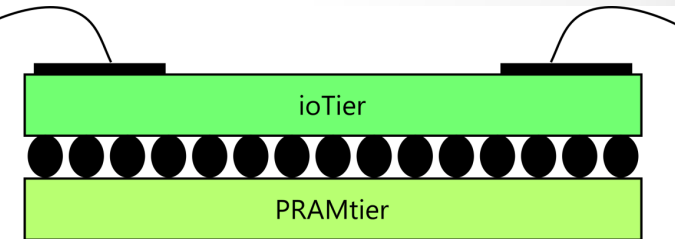
AM chip in 2D:
INFN AMchip06 to AMchip08 (65nm scale down to 28nm),
FNAL first 2D mini-ASIC 130nm @100MHz (scale down if needed)

Three-prong
developments

FNAL VIPRAM-L1CMS:

Simple 2-tier design, scale from current 130nm to 40 nm

~200K patterns/chip @ 200MHz



FNAL generic R&D exploring multi-tier VIPRAM-3D

