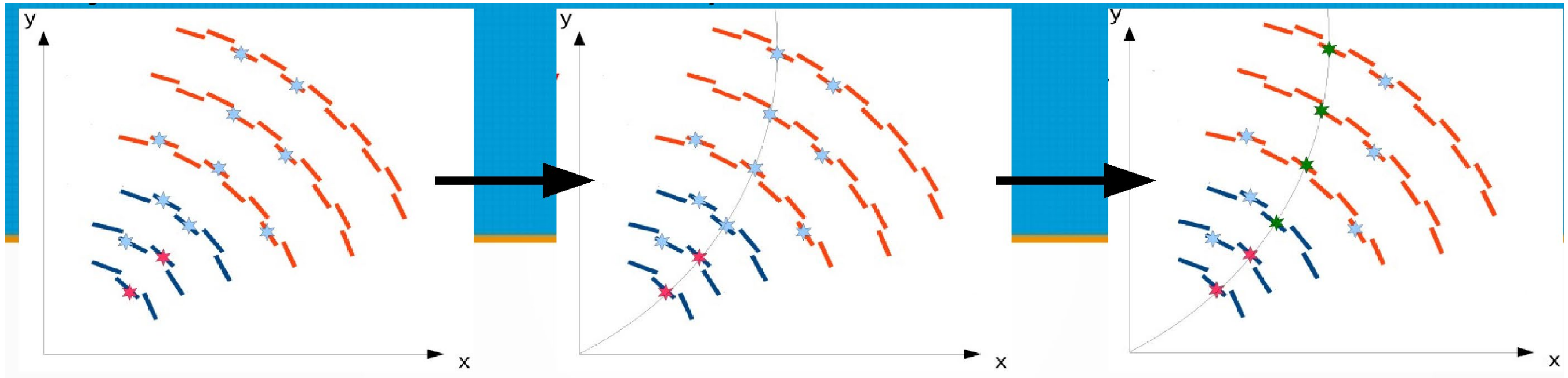# A 3D FPGA Track Segment Seeding Engine Based on the Tiny Triplet Finder

Jinyuan Wu

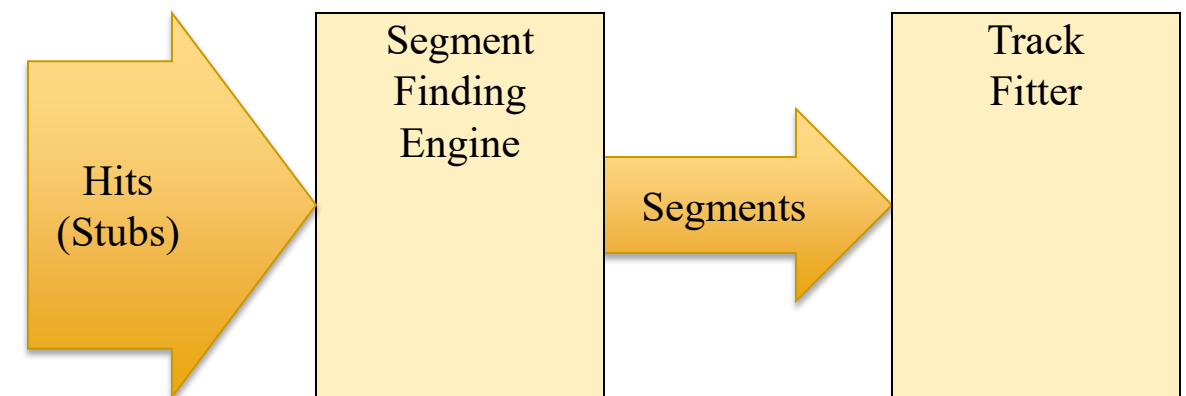Fermilab

Sept. 2021

# Introduction

- An exercise of implementing a 3-hit 3-D track segment seeding engine based on the Tiny Triplet Finder is presented.

- It is intended to group hits into track segments and to feed into later Kalman filtering stage.

- The silicon resource consumption of the seeding engine is small, and it can be put into a low-cost FPGA.

- *In our work, we use CMS L1 tracking Hybrid Algorithm as the example because it is challenging. But this work is not in official CMS project*
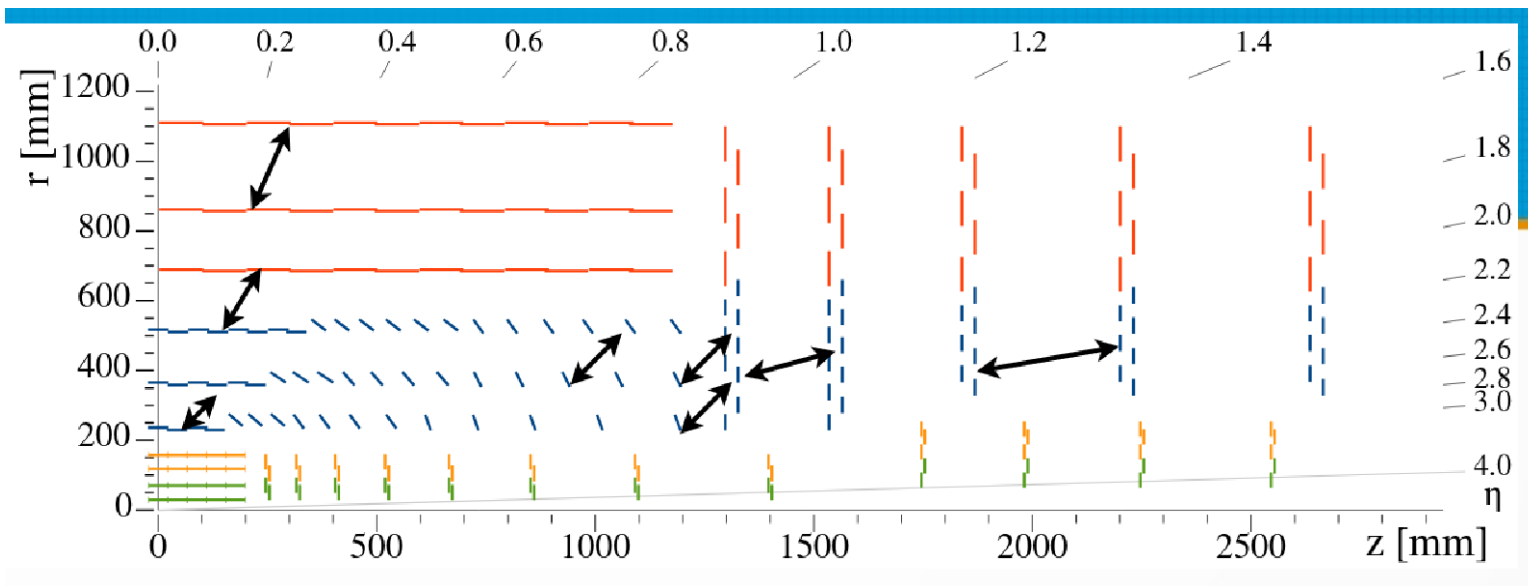
# Tracking in Multilayer Detectors



Andrew Hart et. al. 2020

- Well accepted tracking approaches take the following stages:
  - Hits or stubs are first grouped into track segments.
  - Track segments are fitted into tracks.



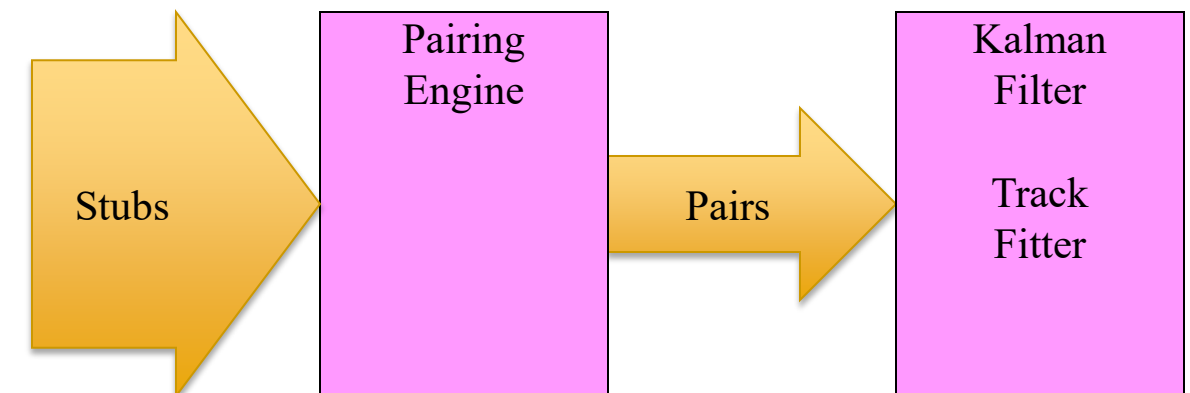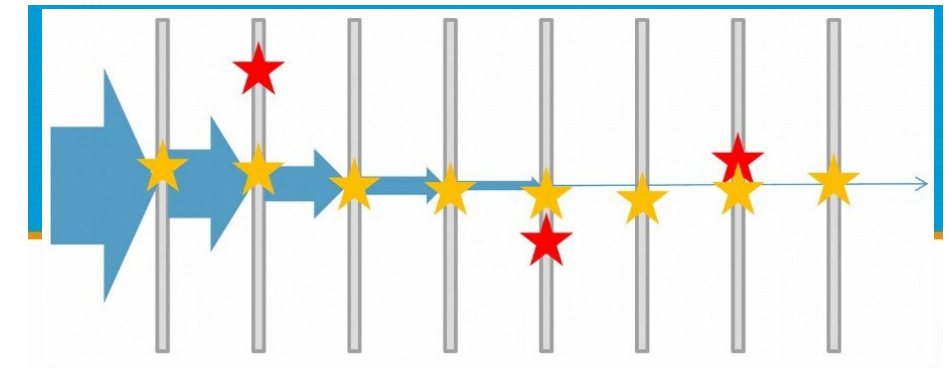Hits (Stubs) → Segment Finding Engine → Segments → Track Fitter

# The CMS L1 Tracking Hybrid Algorithm



Andrew Hart et. al. 2020
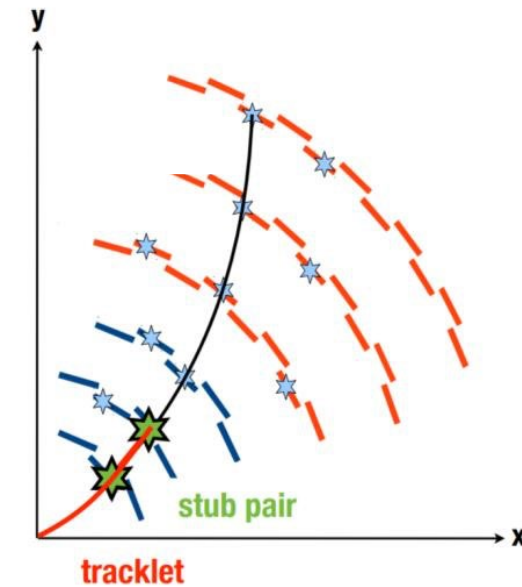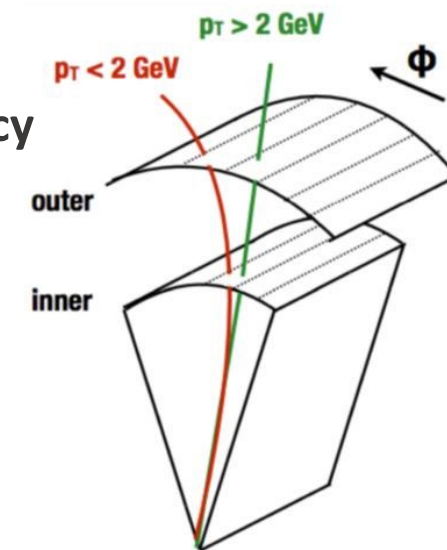


- Hits from different layers are first paired
- Track segments with paired hits are extended.
- Hits are fitted into tracks.



Stubs → Pairing Engine → Pairs → Kalman Filter / Track Fitter
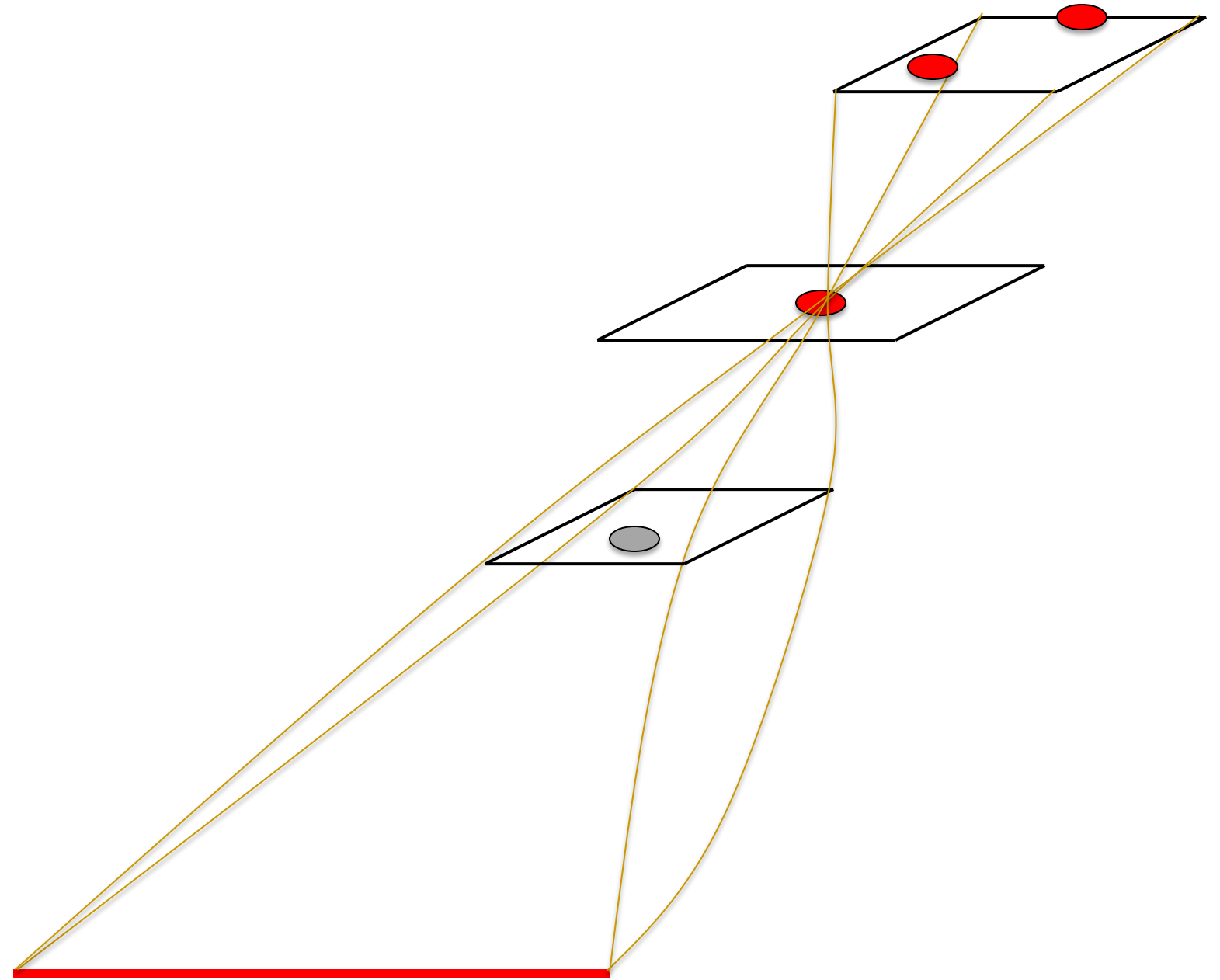
# Pairing

**Stub Organisation & Seeding**
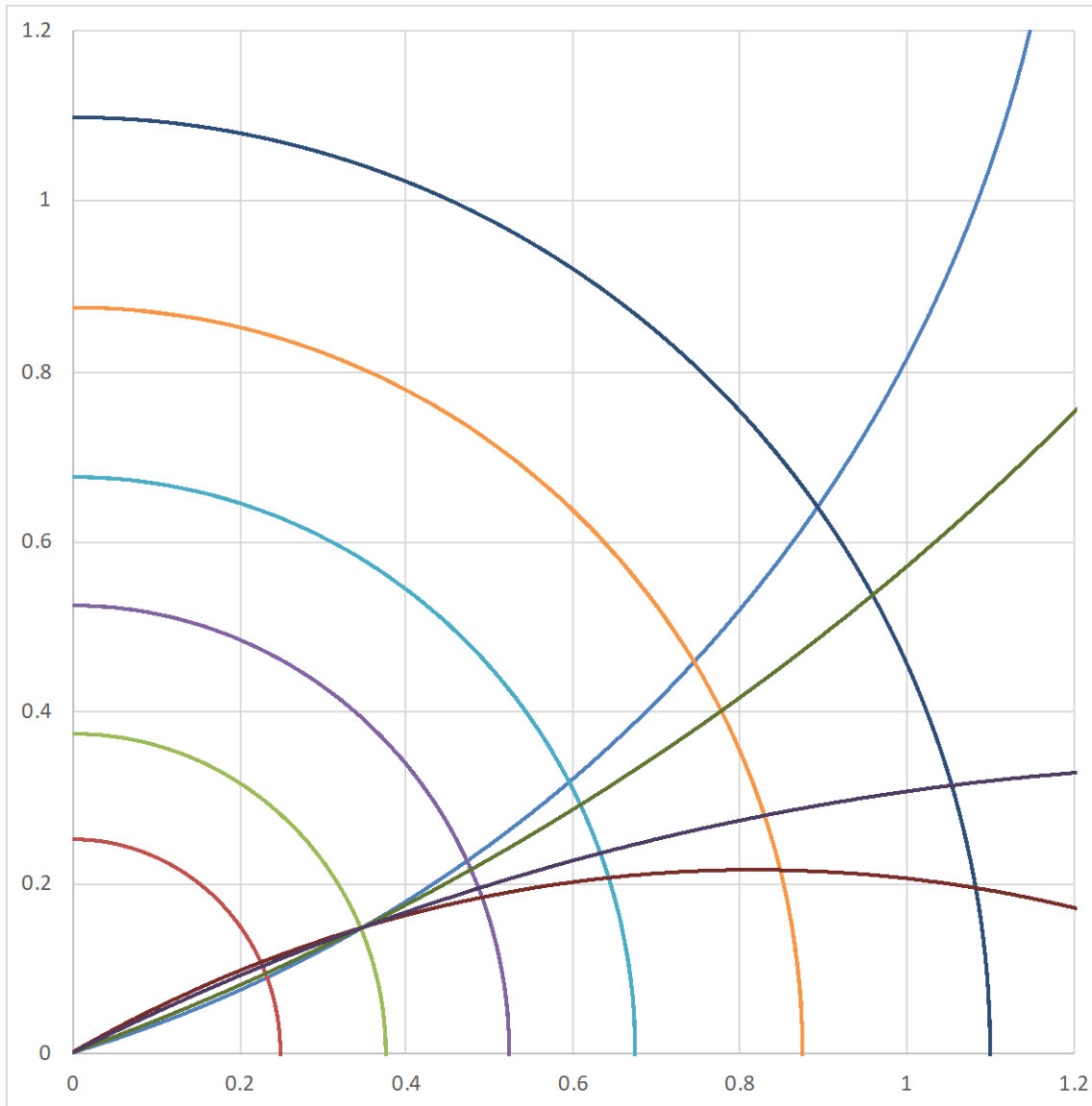
- Seeds (tracklets) are generated from **pairs of layers**
- Layer pairs selected to give full coverage in η, including **redundancy**
- Seeding step **massively parallelised** by internal geometrical partitioning & stub organisation (Virtual Modules - VMs)

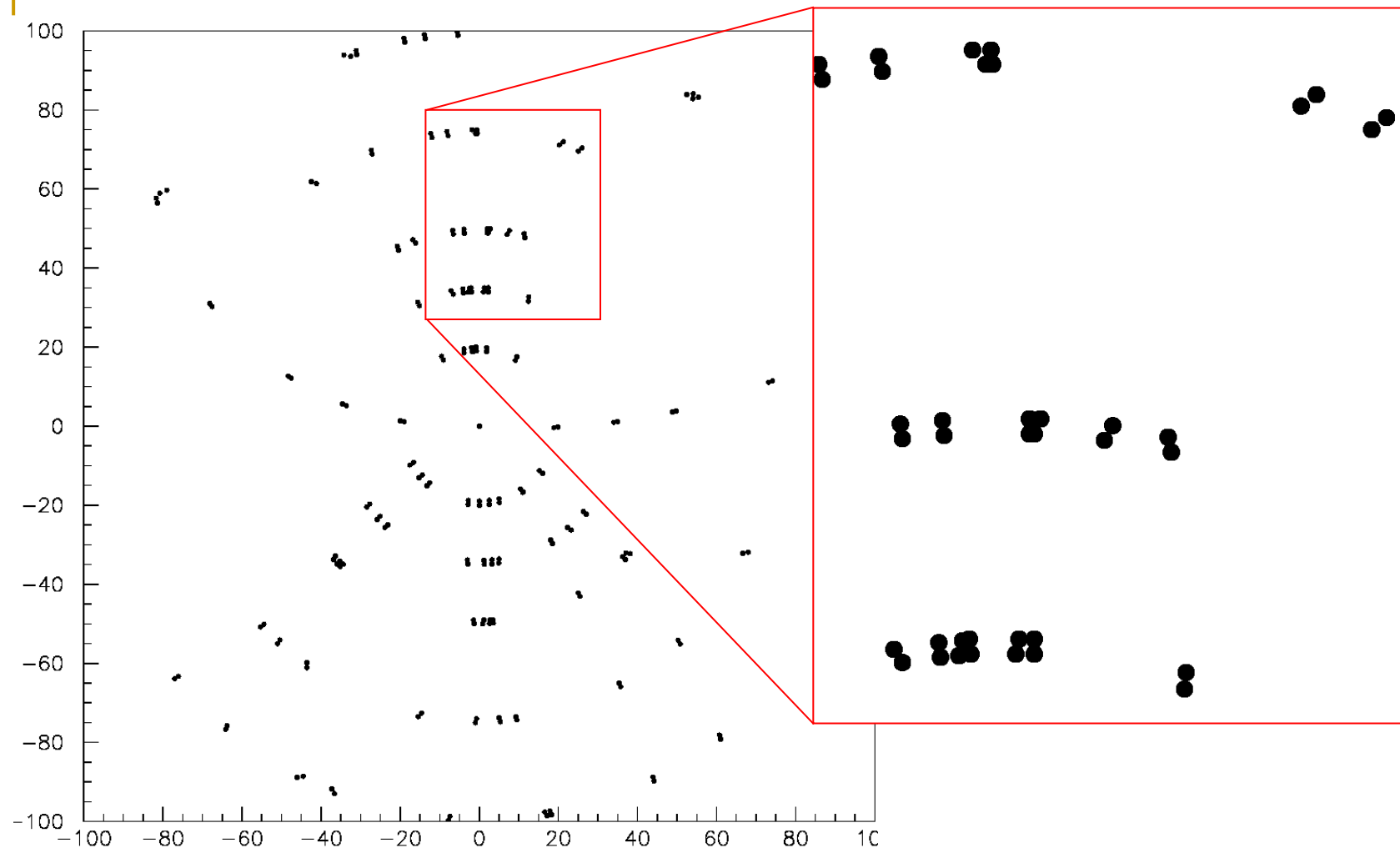Mark Pesaresi et. al. ACES 2020

- The CMS L1 Tracking Hybrid Algorithm: find a pair => find third => fit and find others

# Pairing Stubs



- For any stub in Layer 2, a range in phi and z must be searched in Layer 1 or/and Layer 3.
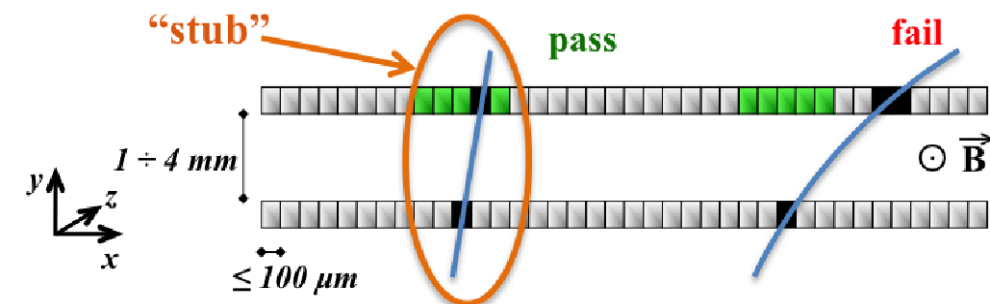- Fake pairs may occur.

# Pairing or Triplet Finding



- If the separation of silicon sensor layers were large (~ 10 mm), then stub **pairing** would be sufficient.

- Given the separation of ~ 4 mm in PS and 2S module, correlating **triplets** of stubs becomes necessary.

Wu et. al. 2004
https://indico.cern.ch/event/418639/contributions/1018451/attachments/868833/1216631/lowpt_lecc2004p.pdf

# From Pairing to Triplet Seeding

- Simple pairing will creates more fake hit pairs, which must be eliminated in later stages.

- When three or more hits (stubs) are picked up together in triplet seeding stage, less fake segments are expected.

**fake**

**fake**

**Seeding with Pairs**

**Seeding with Triplets**

# Triplet Finding

- If three hits can be correlated together in a single action, would that be useful in the Hybrid Algorithm?

- It seems that the triplet finding is more natural for:
    - Displaced track finding
    - Beam-spot constraint tracking

# Triplet Seeding



- Stub pairing won't eliminate stubs not associated with a high Pt track.
- Using three layers or more to most eliminate fake stubs in the seeding engine.

# Triplet seeding is good.
# The algorithm is straightforward.
# The challenge is how to fit it into FPGA.

# Tiny Triplet Finder:
# -- A low resource usage scheme

# The Tiny Triplet Finder for Cylindrical Detector Geometry



- The Tiny Triplet Finer is a scheme to reuse a single set of track road coincidence logic for whole detector.
- It uses shift (or rotate) to bring bit arrays to feed the track road coincidence logic.

# General Structure of the Tiny Triplet Finder

Shifter

Bit-wise Coincident Logic

Shifter

Bit Array

Bit Array

*R1/R3

*R2/R3

Triplet Map Output To Decoder

- A typical Tiny Triplet Finder includes:
  - Bit arrays for different layers of the detector.
  - Shifters
  - Bitwise coincidence logic.

# An Implementation Exercise

# An Implementation Exercise Using CMS L1 Tracking as Example



- The detector in the exercise have the same radii and lengths as CMS Outer Tracker

# System Division



Louise Skinnari et. al. 2020

- System Division:
  - 18 time slices
  - 9 regions
- Each TFP serves:
  - 40 degrees.
  - 25 ns x 18 = 450 ns or 225 clock cycles @500 MHz

# Serving Range of Seeding Engine

- In each TFP, 4 seeding engines are implemented, each serving 2.4 m (full length) x 10 degrees plus overlap.

- Each engine sees <100 hits/layer per event.

- The layer 2 is divided with exact 10-degree boundaries.

- Overlapping in:
  - Layer 1: +- 2.16 degrees, 14.3 degrees total
  - Layer 3: +- 2.62 degrees, 15.2 degrees total

# Generated Events: r-z View



- In 240 cm (z) x 10 degrees (phi) range, up to 112 hits per layer are generated, out of which, 10 are from good tracks.
- The z0 of the vertices spread in +- 10 cm along beam axis.

- Minimum PT: 2GeV/c.
- The hits on Layer 2 is limited to 10 degrees, but hits on Layer 1 and 3 spread up to 16 degrees.
- The impact parameters of the tracks: +- 2 mm.

# The 3D Seeding Engine for SP Layers



- The parameters chosen in r-z plane are z0 and z375 (z375 close to z of Layer 2).
- For any stub in Layer 2, the phi pattern in Layer 1 and Layer 2 with matching (z0, z375) are checked.
- The triplets matching not only in r-z, but also in r-phi are reported for further track fitting.

# Bin Numbers

| | Width | Number |
|---|---|---|
| z0 | 2 cm | 10 |
| z375 | 1 cm | 240 |
| phi | 0.125 degrees (0.55 mm at 25 cm) | 128 (/16 degrees) |

**phi**

**z375**

**z0**

- Neither phi bins nor z bins can be too fine.
- Just checking r-phi or r-z would have too high fake rate.
- The 3-D track seeding (i.e., checking both r-phi and r-z) should provide sufficient fake stub rejection.

# The Track Seeding Engine

# Structure of the Seeding Engine: Two RAM Banks

# Structure of the Seeding Engine: Single Physical RAM View

# General Timing

# Operation: (1) Filling Hits

# Pattern Shifting and checking

- The bit patterns with the same z0 from Layer 1 and Layer 3 are shifted to align with the coincidence road array.

- The amount of shift is determined by phi2.

- Matching bits are output for further process.

# The Seeding Engine & Clustering Buffer



- Full data of hits are stored in the clustering buffer.
- After triplet finding, the full data are sent to the Kalman filter stage for further track fitting.

# Implementation in Cyclone 5 FPGA

# Home Test Stand and the Evaluation Module





- The seeding engine core firmware is implemented in Altera C5G evaluation module (179 USD ea.).
- The FPGA is Cyclone V GX 5CGXFC5C6F27C7N (171 USD ea.).

# The Seeding Engine Block I/O



- The hit data from Layer 1, 2 and 3 are fed into the engine, one hit per clock cycle.
- The primary coincidence outputs are in port DDYY[63..0] which represents up to 64 possible 1/PT values of the track segment.

# Structure of the seeding engine



- The entire engine is organized as a pipeline so that it can run at high speed (250 MHz in Cyclone V).
- Blocks: Address Generation; Hugh Trans. Cells; Shifter Bank; Road Coincidence

# The Seeding Engine

# Address Conversion for Hugh Transform Cells



AddrGenL1

EvData1x[127..0]     HitCmd[1..0]
EvData2x[127..0]     HitPos[7..0]
SEL2                 HitAddr[79..0]
Refresh
CK250

SEL2
Refresh
CK250

inst5

Pipeline Steps:3

Pipeline Steps:9

HughTransCells21c4

HitCmd[1..0]     HitPtn9a[1
HitPos[7..0]     HitPtn8a[1
HitAddr[79..0]   HitPtn7a[1
TRST             HitPtn6a[1
CK250            HitPtn5a[1
                 HitPtn4a[1
                 HitPtn3a[1
                 HitPtn2a[1
                 HitPtn1a[1
                 HitPtn0a[1

TRST
CK250

inst2

- For Layer 1 or Layer 3 Hugh Transform Cells have 10 RAM blocks corresponding 10 bins in $z0$.

- The z-coordinate $z1$ or $z3$ is converted to $z375$ for each $z0$ by checking the lookup table ROM.

- The ROM contents depend on detector geometry and are generated using simulation data.

MemInit1.txt

```
--Format: Quartus .mif

WIDTH=80;|
DEPTH=256;

ADDRESS_RADIX=HEX;
DATA_RADIX=HEX;

CONTENT BEGIN
00 : 00000000000000000000;
01 : 00000000000000000000;
02 : 00000000000000000000;
03 : 00000000000000000000;
04 : 00000000000000000000;
05 : 00000000000000000000;
06 : 00000000000000000000;
07 : 00000000000000000000;
08 : 00000000000000000000;
09 : 00000000000000000000;
0a : 00000000000000000000;
0b : 00000000000000000000;
0c : 00000000000000000000;
0d : 00000000000000000000;
0e : 00000000000000000000;
0f : 00000000000000000000;
10 : 00000000000000000000;
11 : 00000000000000000000;
12 : 00000000000000000000;
13 : 00000000000000000000;
14 : 00000000000000000000;
15 : 00000000000000000000;
16 : 00000000000000000000;
17 : 00000000000000000000;
18 : 00000000000000000000;
19 : 00000000000000000000;
1a : 00000000000000000000;
1b : 00000000000000000000;
1c : 00000000000000000000;
1d : 00000000000000000000;
1e : 00000000000000000000;
1f : 00000000000000000000;
20 : 00000000000000000000;
21 : 00000000000000000000;
22 : 00000000000000000000;
23 : 00000000000000000000;
24 : 00000000000000000000;
25 : 00000000000000000000;
26 : 00000000000000000000;
27 : 00000000000000000000;
28 : 01000000000000000000;
29 : 02010000000000000000;
2a : 04030201000000000000;
2b : 05040302010000000000;
2c : 07060504030201000000;
2d : 08070605040302010000;
2e : 0a090807060504030201;
2f : 0b0a0908070605040302;
30 : 0d0c0b0a090807060504;
31 : 0e0d0c0b0a0908070605;
32 : 100f0e0d0c0b0a090807;
33 : 11100f0e0d0c0b0a0908;
```
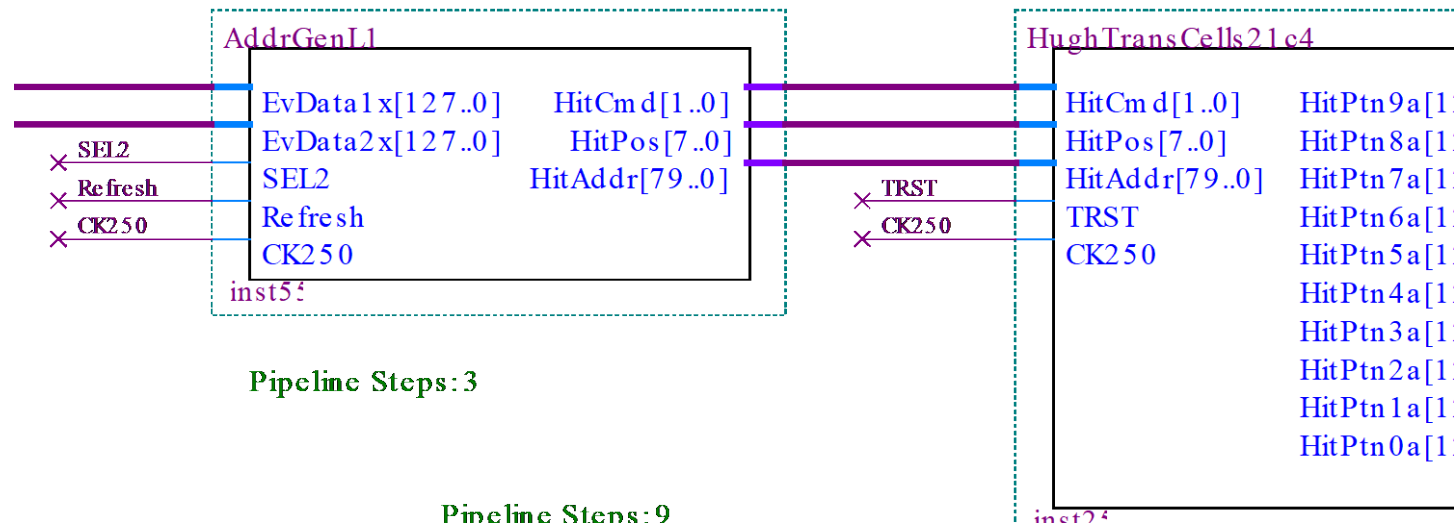
MemInit3.txt

```
--Format: Quartus .mif

WIDTH=80;
DEPTH=256;

ADDRESS_RADIX=HEX;
DATA_RADIX=HEX;

CONTENT BEGIN
00 : 22222324242525262627;
01 : 23232424252526272728;
02 : 23242425262627272828;
03 : 24252526262727282929;
04 : 2525262627282829292a;
05 : 25262727282829292a2b;
06 : 2627272828292a2a2b2b;
07 : 27272829292a2a2b2b2c;
08 : 282829292a2a2b2c2c2d;
09 : 2829292a2b2b2c2c2d2d;
0a : 292a2a2b2b2c2c2d2e2e;
0b : 2a2a2b2b2c2d2d2e2e2f;
0c : 2a2b2c2c2d2d2e2e2f30;
0d : 2b2c2c2d2d2e2f2f3030;
0e : 2c2c2d2e2e2f2f303031;
0f : 2d2d2e2e2f2f30313132;
10 : 2d2e2e2f303031313232;
11 : 2e2f2f30303131323333;
12 : 2f2f3030313232333334;
13 : 2f303131323233333435;
14 : 30313132323334343535;
15 : 31313233333434353536;
16 : 32323333343435363637;
17 : 32333334353536363737;
18 : 33343435353636373838;
19 : 34343535363737383839;
1a : 34353636373738383939;
1b : 3536363737383839393a;
1c : 36363738383939393a3a3b;
1d : 3737383839393a3b3b3c;
1e : 373838393a3a3b3b3c3c;
1f : 3839393a3a3b3b3c3c3d3d;
20 : 39393a3a3b3c3c3d3d3e;
21 : 393a3b3b3c3c3d3d3e3f;
22 : 3a3b3b3c3c3d3e3e3f3f;
23 : 3b3b3c3d3d3e3e3f3f40;
24 : 3c3c3d3d3e3e3f404041;
25 : 3c3d3d3e3f3f40404141;
26 : 3d3e3e3f3f4040414242;
27 : 3e3e3f3f404141424243;
28 : 3e3f4040414142424344;
29 : 3f404041414243434444;
2a : 40404142424343444445;
2b : 41414242434344444546;
2c : 41424243444445454646;
2d : 42434344444545464747;
2e : 43434444454646474748;
2f : 43444545464647474849;
30 : 44454546464748484949;
31 : 4545464747484849494a;
```

# Register-Like Block RAM for Hit Patterns



Event N
225 CLK

Event N+1
225 CLK

Event N+2
225 CLK

L1, L3
112 CLK

L2
112 CLK

R

Single Clock Writing

Single Clock Reading

Single Clock Refreshing

- The hit patterns are stored in the register-like BRAM.

- Each hit is processed in a single clock cycle.

- Single clock cycle refreshing allows fast preparation for next event.

## Register-Like Block RAM: Implementation, Testing in FPGA and Applications for High Energy Physics Trigger Systems

Jinyuan Wu

*Abstract*— In high energy physics experiment, trigger block memories are utilized for various algorithm. The indexed searching algorithm, often demanded to reset all memory location between different event, is a feature not supported in regular block memory. Another common demand is to be able to update the contents in any memory location in a single clock cycle. These demands can be fulfilled with registers but the implementation with registers is unaffordable. In this paper, a register-like memory with such schemes is designed which allows single clock cycle reading and single clock cycle writing for any memory location and also allows refreshing the entire memory within a single clock. The prototype and test results are presented.

*Index Terms*— Trigger System, FPGA, Applications

### I. INTRODUCTION

IN high energy physics experiments, trigger systems perform essential roles for reaching the physics goals of the experiments. As luminosity in collider experiments increases, sophisticated trigger systems utilizing nearly full detector resolutions are demanded. Algorithms with offline analysis complexity such as associated memories or "artificial retina" implementations in hardware/firmware environment are attempted [1-2]. In these algorithms, a common building element is a block memory capable of fast updating and fast global refreshing.
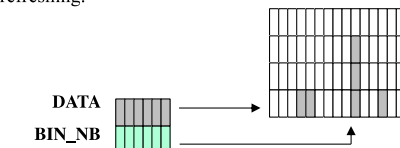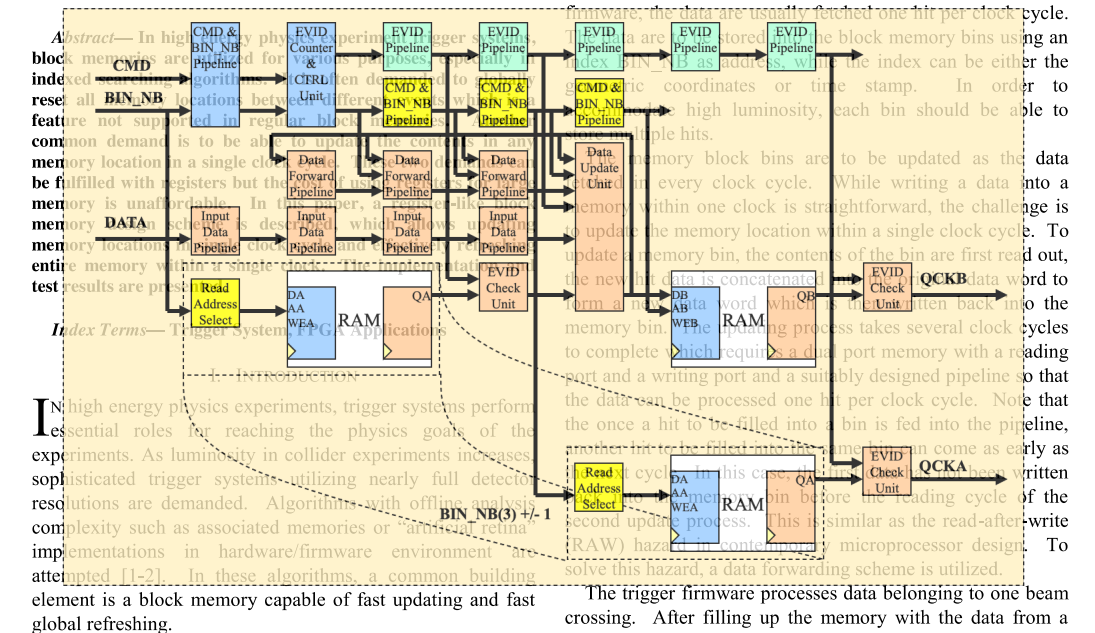


Fig. 1. An application of block memory bins used for in high energy physics experiment trigger system.

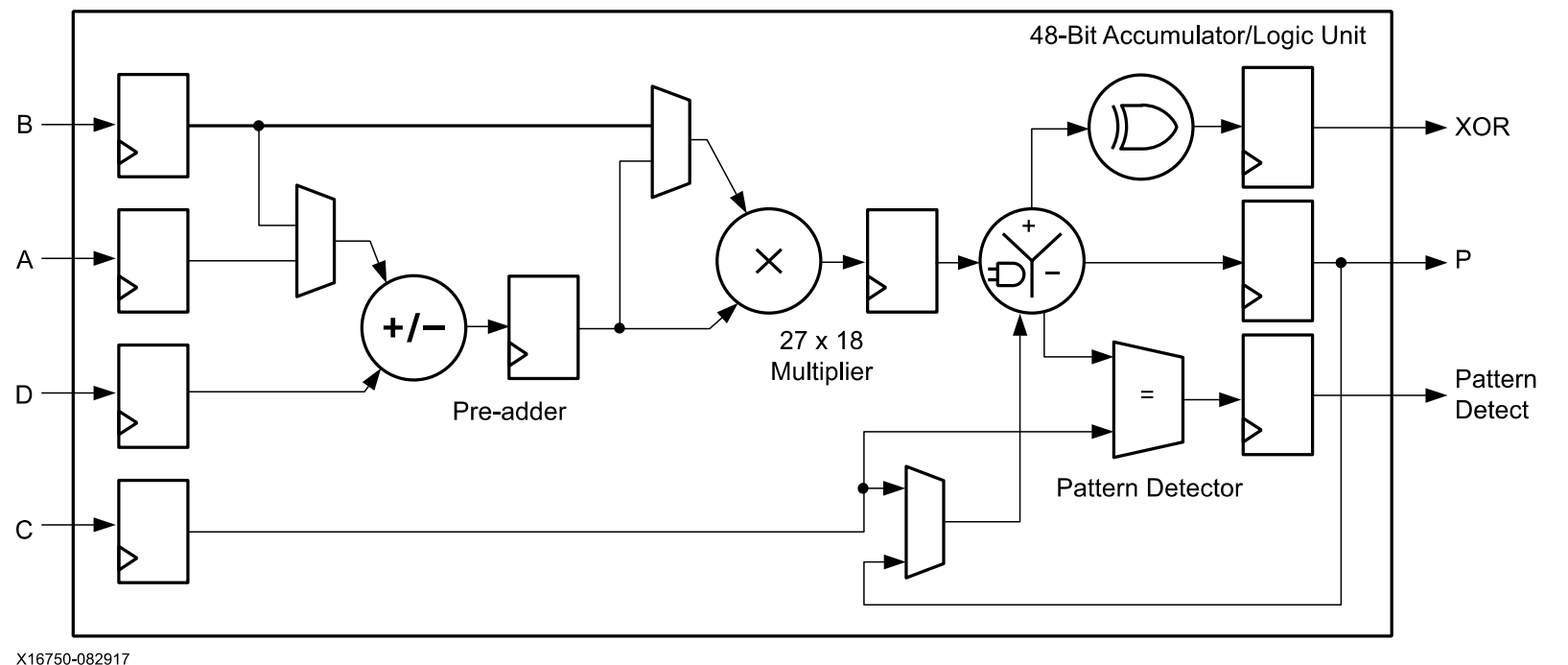A typical application of the block memories organized in memory bins is shown in Fig. 1. The detector data representing hit coordinates and timing information are fed into the trigger system and in the processing stages of the

firmware, the data are usually latched one hit per clock cycle. The hits are to be stored in the block memory bins using an index BIN_NB as address, while the index can be either the geometric coordinates or time stamp. In order to support high luminosity, each bin should be able to store multiple hits. The memory block bins are to be updated as the data come in every clock cycle. While writing a data into a memory within one clock is straightforward, the challenge is to update the memory location within a single clock cycle. To update a memory bin, the contents of the bin are first read out, the new data is concatenated to form a new data word which is then written back into the memory bin. The updating process takes several clock cycles to complete which requires a dual port memory with a reading port and a writing port and a suitably designed pipeline so that the data can be processed one hit per clock cycle. Note that the once a hit to be filled into a bin is fed into the pipeline, another hit to be filled into the same memory bin can come as early as the next clock cycle. In this case, the data to be written into the memory bin before the reading cycle of the second update process. This is similar as the read-after-write (RAW) hazard in contemporary microprocessor design. To solve this hazard, a data forwarding scheme is utilized.

The trigger firmware processes data belonging to one beam crossing. After filling up the memory with the data from a beam crossing, the algorithm will search the data based on the index of the bins. Note that the searching process may not address all bins containing the hit data and it may also address empty bins. After reading process, all memory bins will be effectively cleared to prepare for the next beam crossing. It is well known that regular block memories do not support global reset. To fulfill this requirement, an event ID tagging scheme is used.

The single clock updating and global refreshing schemes developed in our previous work [3] are combined into a unified scheme in this work. In this paper, global refreshing scheme is first discussed in Section II, followed by the pipeline structure with data forwarding support in Section III. The implementation and test results of the entire scheme are presented in Section IV.

### II. GLOBAL REFRESH SCHEME FOR BLOCK MEMORY

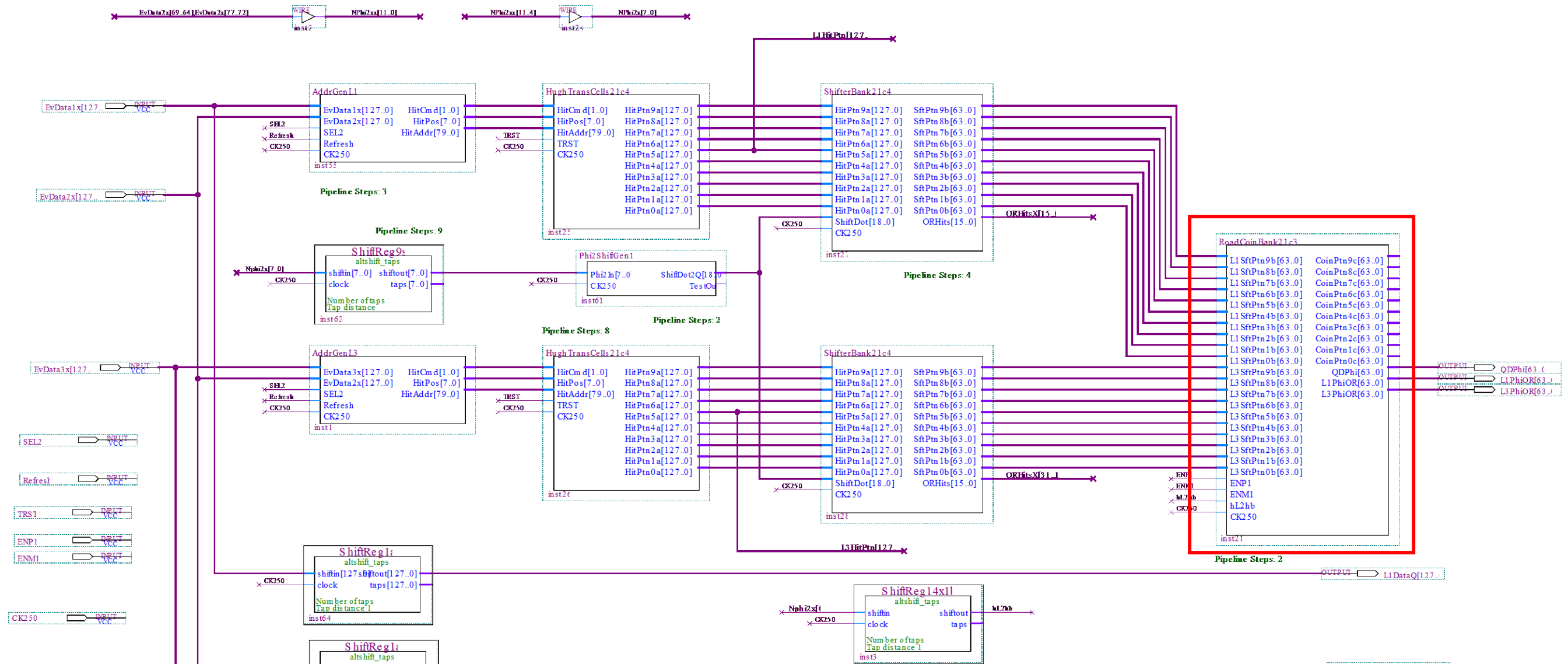Intrinsically, the block memories can only be accessed one

# The Seeding Engine

# Shifters



- The multipliers in FPGA DSPs are used to implement shifters.

# The Seeding Engine

# Road Coincidence
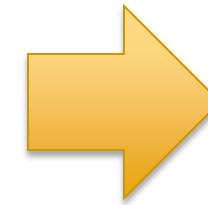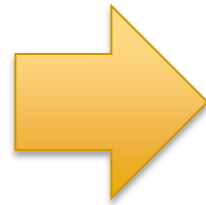


- Shifted bit patterns are checked by coincidence road array.
- The r-phi logic equations are generated with simulation.

# Test Results

# Raw Input Data and Output Data



- Hits (stubs) are sent into the seeding engine.
- Each line is a hit in a detector layer.
- Typical data: 100 hits/layer with 10 good tracks with 90 random hits.
- The seeding engine outputs a 64-bit pattern, and the non-zero bits represent a potential track segment.

# Good Track Acceptance and Fake Rejection



- Good track acceptance is better than 99%.
- Increasing number of bins in phi causes better fake rejection.

# Fake Tracks due to High Occupancy

- When number of total hits per event (BX) is low, fake rate is low.

- Event beyond occupancy of 3000 stubs/layer for the whole detector (the AvMulti3K line), total number of coincidence is still manageable.

- The dash line Limit4 indicate total coincidence allowed in fitting stage if each coincidence is processed with 4 clock cycles.
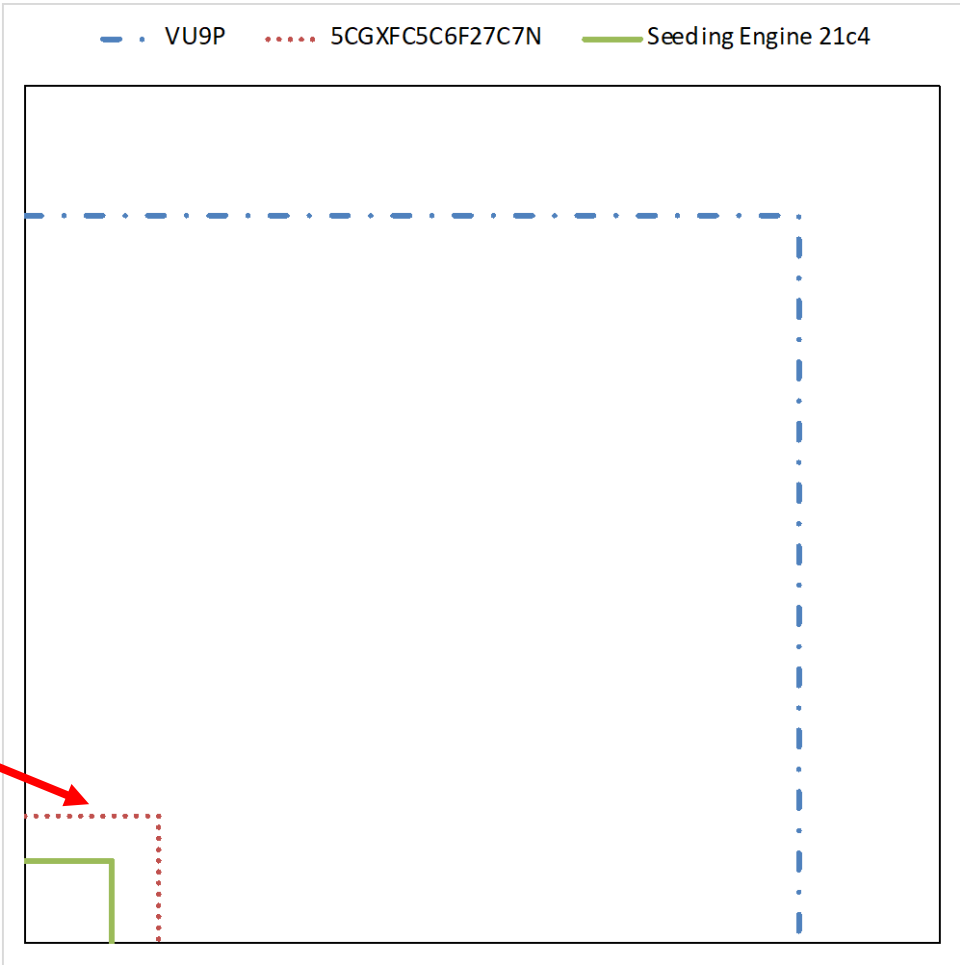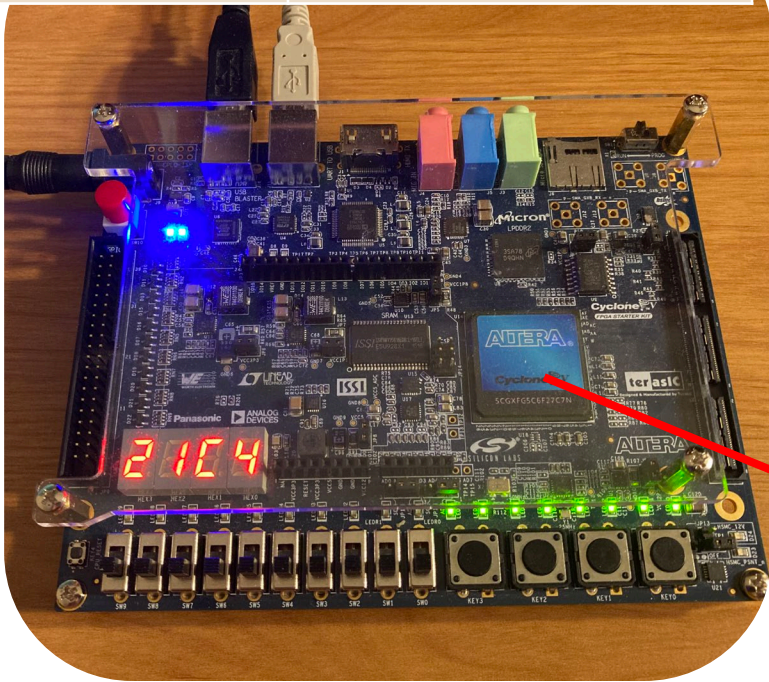
# Performance of the Seeding Engine



- Most of fake hits are rejected while good tracks are preserved.
- Data volume sent to the fitting stage is much smaller.

# Resource Usage:

| | Used in a Seeding Engine core | Available in Cyclone 5 5CGXFC5C6F27C7N | Available in Ultra-scale Plus VU9P |
|---|---|---|---|
| Logic Cells | 12446 (ALM) Eq. 32.9K (LC) | 29080 (ALM) | 2586K (LC) |
| Block RAM | 150 (M10K) 743K (bits) | 446 (M10K) 4567K (bits) | 75.9M (bits) |
| DSP | 100 (16x16 Multipliers) | 150 (DSP 27x27) Eq. 300 (16x16) | 6840 (DSP 18x27) |

- The seeding engine core fits the Cyclone 5 device with ~45% resource usage.
- Four cores are needed in VU9P and it looks possible.





Legend: VU9P ···· 5CGXFC5C6F27C7N ── Seeding Engine 21c4

# Summary

- **Algorithms**: well-known in the field.

  - ❑ The fake rate of track segment seeding with triplets is lower.

- **Implementations**: exist in the field.

  - ❑ The challenge is how to fit into a reasonable sized FPGA.

- **Methods** of solving the implementation challenge: this work.

  - ❑ Using Tiny Triplet Finder so that only a small number of coincidence roads (product terms) are implemented.

  - ❑ Hit patterns are implemented with "register-like RAM" which allows single clock update, read and refresh operations.

  - ❑ Lookup tables and coincidence roads are auto initiated using simulation data to accommodate detector geometry or alignment variations.
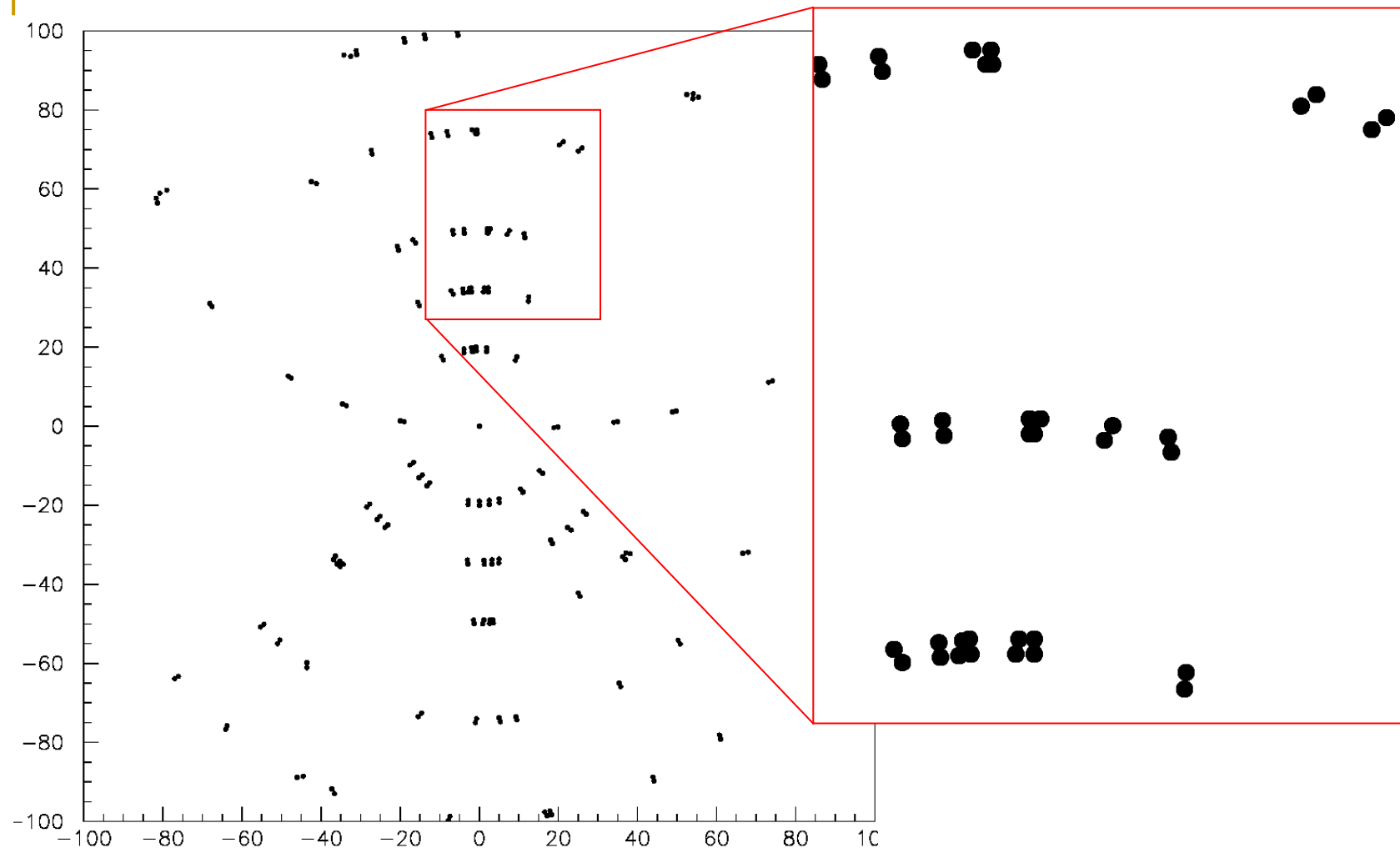
The End

Thanks

# Introduction

- For high-luminosity operation, 3D track segment seeding becomes necessary:
    - Usable bin sizes in z or phi direction can't be too fine due to:
        - Multiple scattering
        - Non-zero impact parameter
    - The 2D segment seeding would yield too many fake segments.
    - Constraints in both r-z and r-phi planes help to reduce fake segments.
- Using the Tiny Triplet Finder to implement 3D track segment seeding engines
    - The Tiny Triplet Finder is a low-resource usage scheme for track segment finding in FPGA.
    - It becomes affordable to implement 3D seeding engines since the resource usage is low.
- A true 3D track segment seeding engine has been implemented and tested.
    - Detector partition, data timing etc. are based on the seeding-fitting approach in recent CMS presentations: [Mark Pesaresi, ACES 2020] and [Luigi Calligaris RT2020].
    - Stubs are generated on idealized Outer Tracker 3 PS layers.
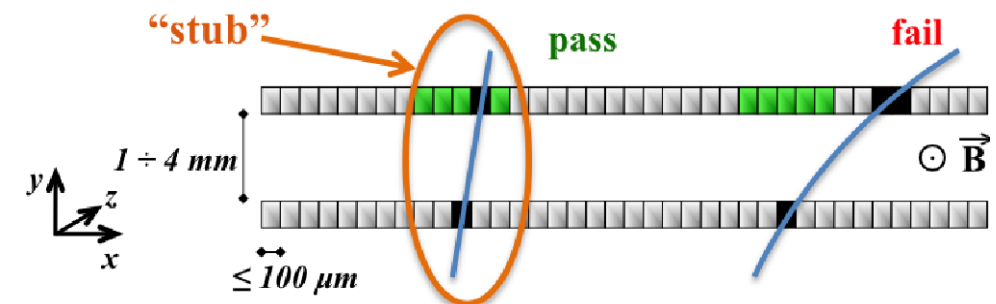    - Operates as intended.

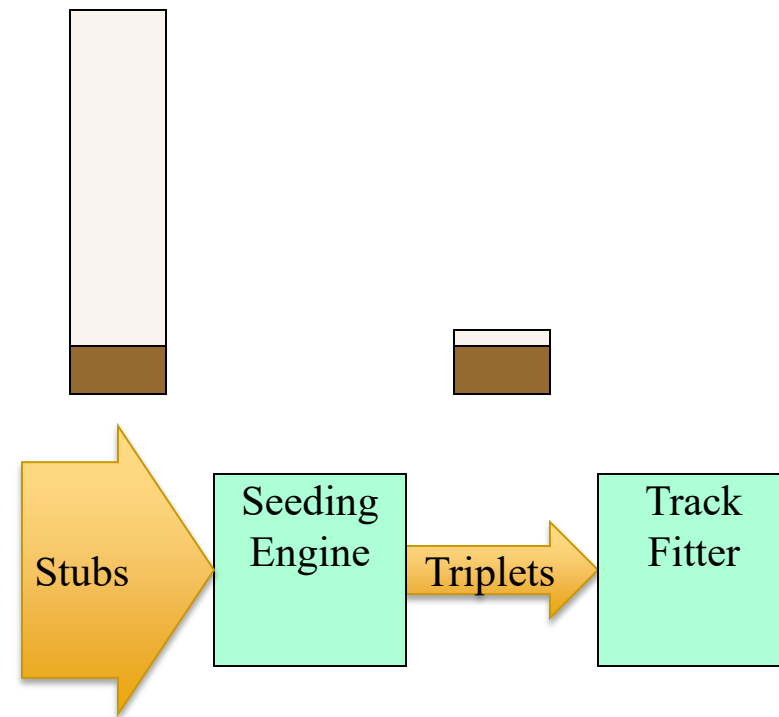# Pairing or Triplet Finding



- If the separation of silicon sensor layers were large (~ 10 mm), then stub **pairing** would be sufficient.

- Given the separation of ~ 4 mm in PS and 2S module, correlating **triplets** of stubs becomes necessary.

Wu et. al. 2004
https://indico.cern.ch/event/418639/contributions/1018451/attachments/868833/1216631/lowpt_lecc2004p.pdf
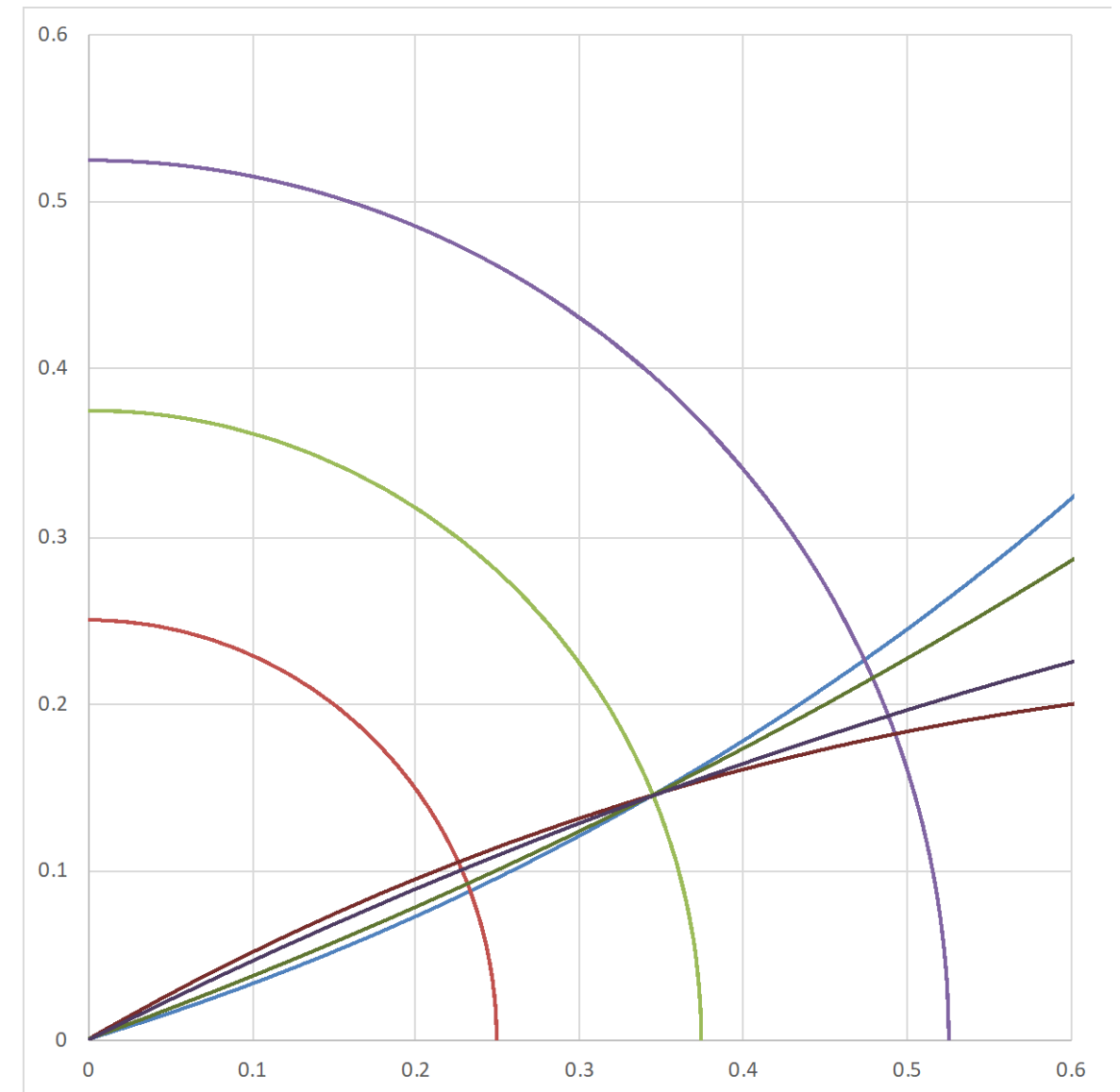
# Why Seeding

Mark Pesaresi et. al. ACES 2020

- At 200 pile-up, the tracker sees 15k-20k stubs => ~ 3000 stubs/layer.

- There are ~300 high Pt tracks (Pt > 2 GeV).

- About 90% stubs are not associated with any high Pt track. => more than 90% stub pairs are fake.

- Requesting additional constraints from 3 or more stubs (triplet) will pre-exclude fake pairs.

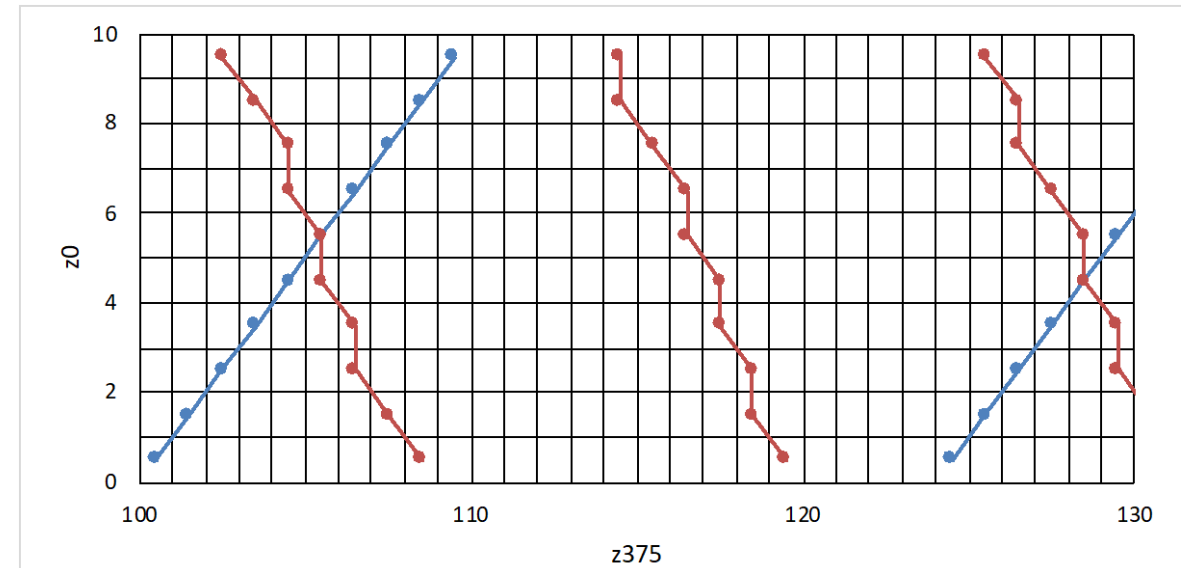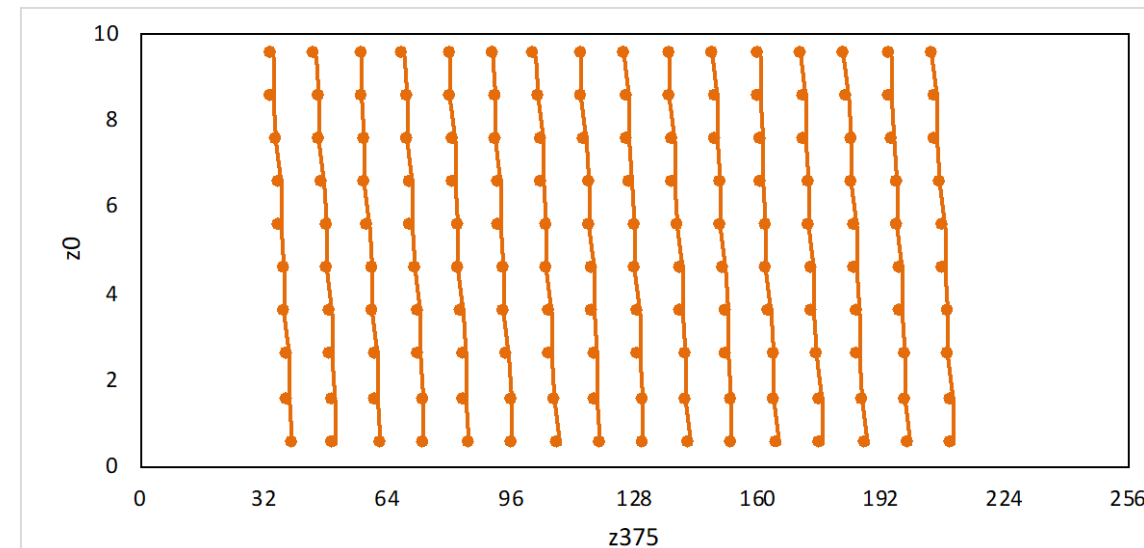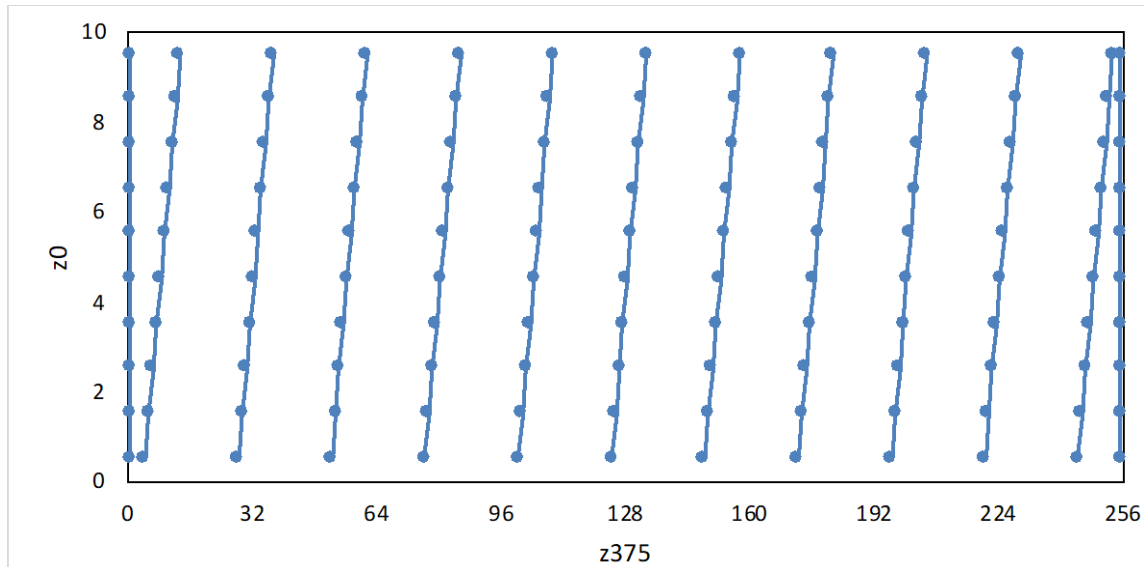# The Searching Range: Layer 2 to 1 and 3

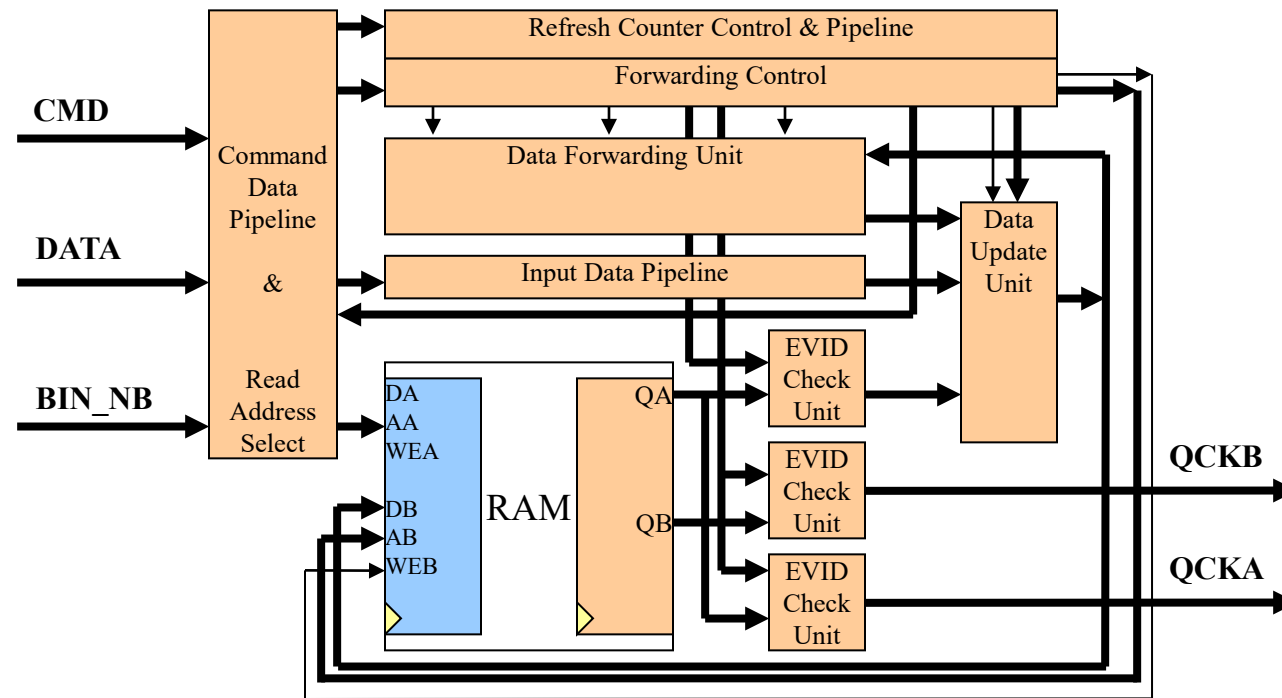| | | | |
|---|---|---|---|
| PT(GeV/c) | 2.000 | | |
| B(T) | 4.000 | | |
| R0(m) | 1.667 | | |
| Z0(m) | 0.200 | | |
| Detector Length(m) | 2.400 | | |
| Number of stubs per layer | 3000.000 | | |
| | Layer 1 | Layer 2 | Layer 3 |
| R(m) | 0.250 | 0.375 | 0.525 |
| Phi(rad) | 0.075 | 0.113 | 0.158 |
| Searching range in phi (rad) | 0.075 | | 0.091 |
| (deg.) | 4.316 | | 5.205 |
| (mm) | 18.834 | | 47.691 |
| Search range in Z (mm) | 66.667 | | 80.000 |
| Total Layer Area (m^2) | 3.770 | | 7.917 |
| | | | |
| Number of Hits in searching area | 0.999 | | 1.446 |



- For almost all 3000 stubs in Layer 2, on average a hit can be found in Layer 1 or Layer 3.

- About 90% of these pairs are not part of true tracks.

- When 3000 stubs/layer are sent to the stub pairing stage, there will be about 3000 pairs output to the tracking stage without significant reduction.
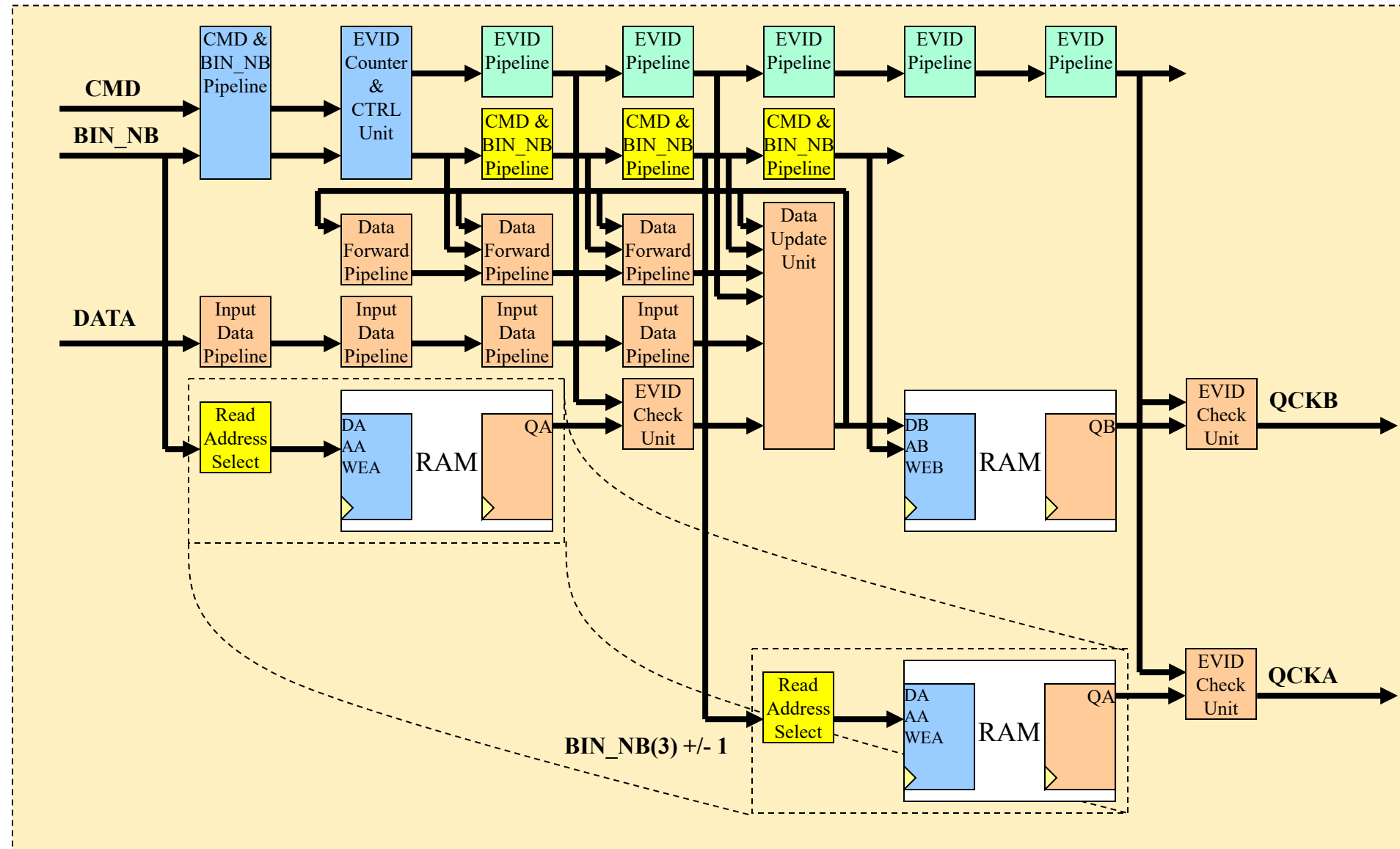
# Hits in the Hugh Transform Cells







- A hit is written into 10 z0 RAM blocks.
- The hits in Layer 1 and Layer 3 have different angles.
- The z coordinate of Layer 2, z375 is used to address all RAM block simultaneously.
- If an intersection of Layer 1 and Layer 3 hits are met, it might be a potential track segment.

# Register-Like Block RAM Block Diagram



- The

# Detailed Block Diagram



- Boundary coverage is also supported during reading operations.