

The new variable resolution Associative Memory for Fast Track finding aka AMchip04

WIT May 3-5, 2012, Pisa



Alberto Annovi

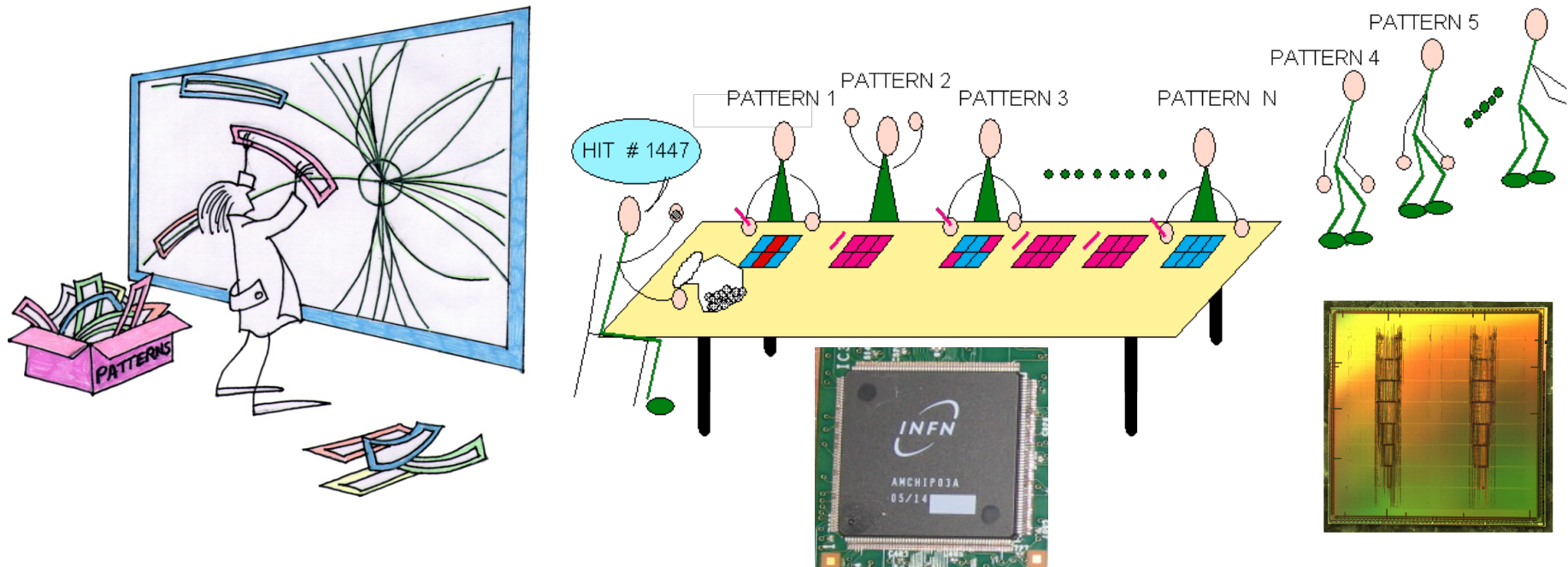
Istituto Nazionale di Fisica Nucleare
Laboratori Nazionali di Frascati



FTK algorithm: Pattern recognition & Track fitting

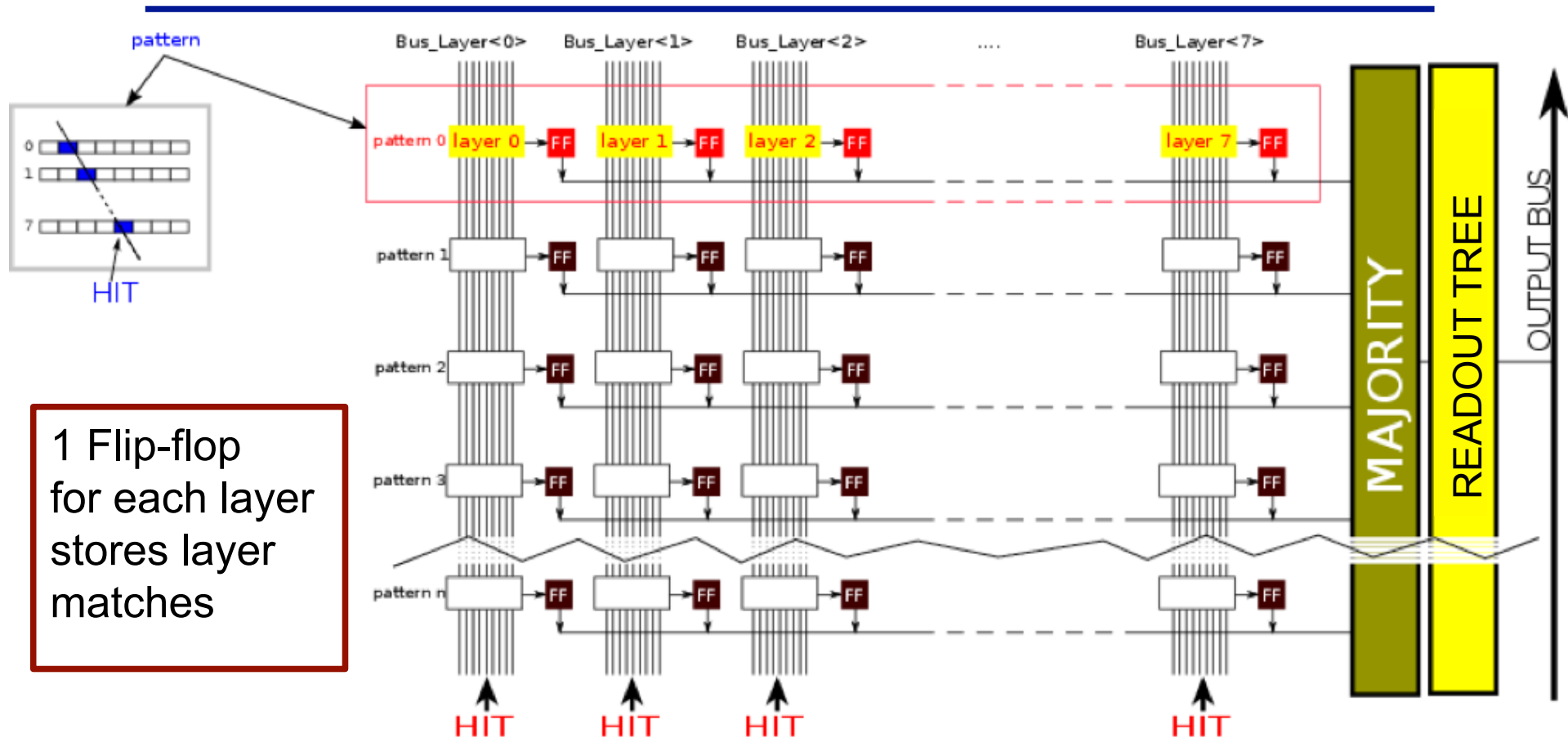
See Thu B. Penning's talk

- Pattern recognition – find track candidates with enough Si hits



- $O(10^9)$ prestored patterns simultaneously see the silicon hits leaving the detector at full speed.
- Based on the **Associative Memory** chip (content-addressable memory) initially developed for the CDF Silicon Vertex Trigger (**SVT**).

AM working principle



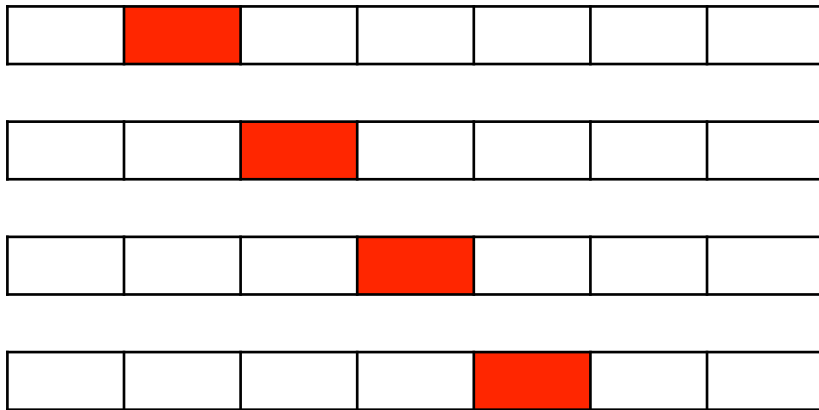
1 Flip-flop for each layer stores layer matches

All patterns compared in parallel with incoming data. Look for correlation of data received at different times. (Feature unique to AMchip)

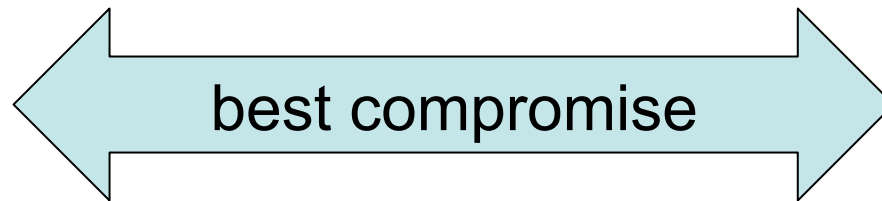
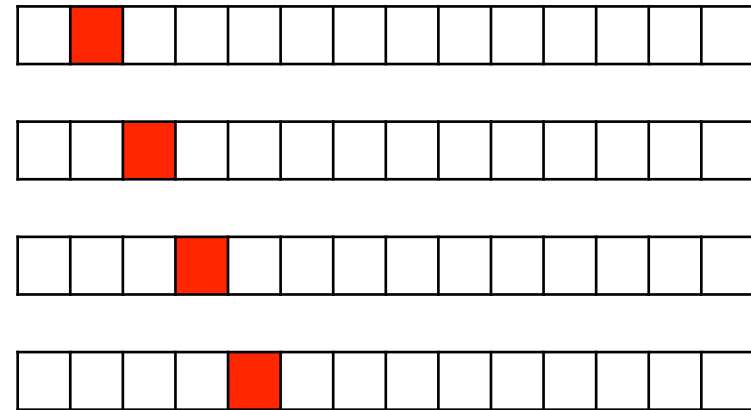
Fast pattern matching. Flexible input: position, time, objects...

Generatig the pattern bank

Wide patterns



Thin patterns



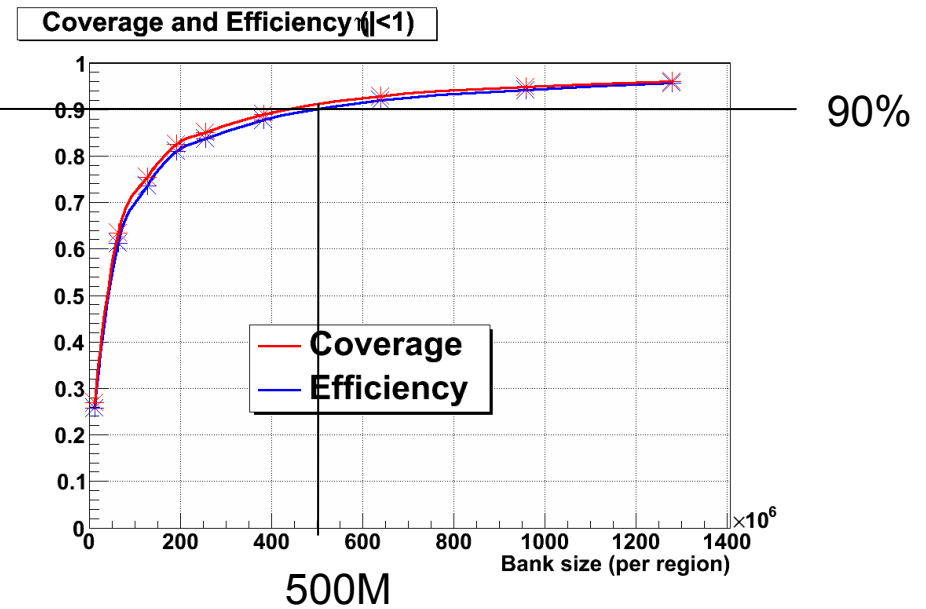
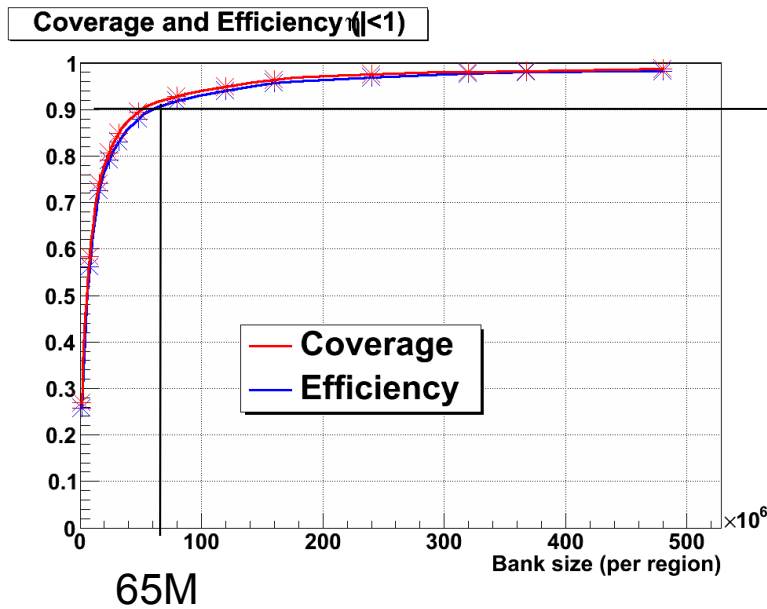
High efficiency
with less patterns (hardware)
BUT more fakes

More patterns (hardware)
for same efficiency
less fakes

Pattern efficiency

Pattern size
r- ϕ : 24 pixels, 20 SCT strips
z: 36 pixels

Pattern size (half size)
r- ϕ : 12 pixels, 10 SCT strips
z: 36 pixels



of patterns in Amchips (barrel only, 45 ϕ degrees)

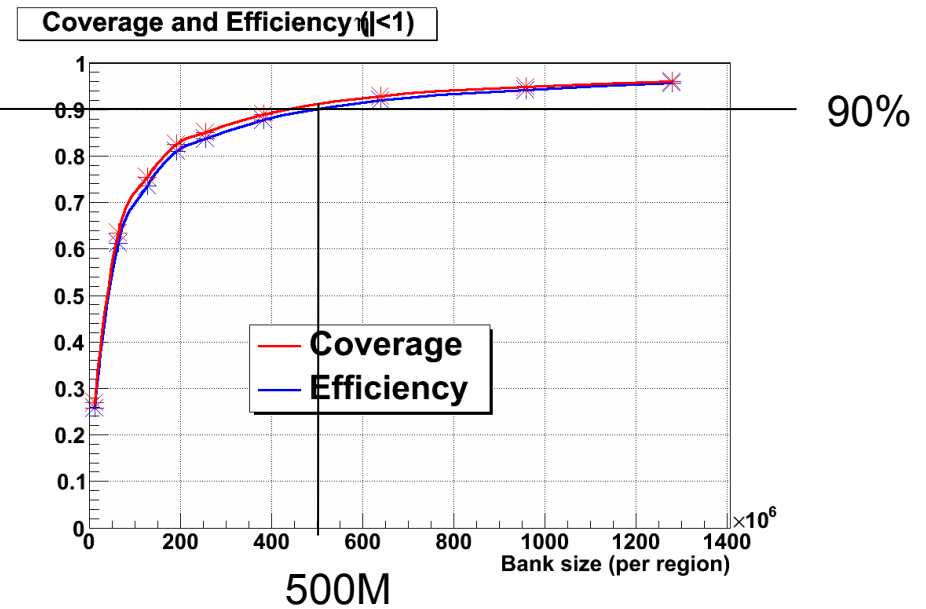
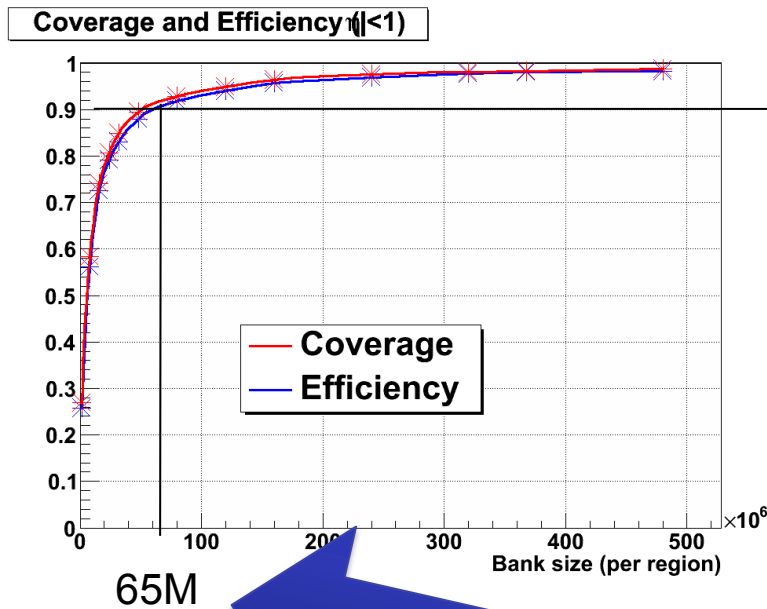
<# matched patters/event @ 3E34> = 342k

<# matched patters/event @ 3E34> = 40k

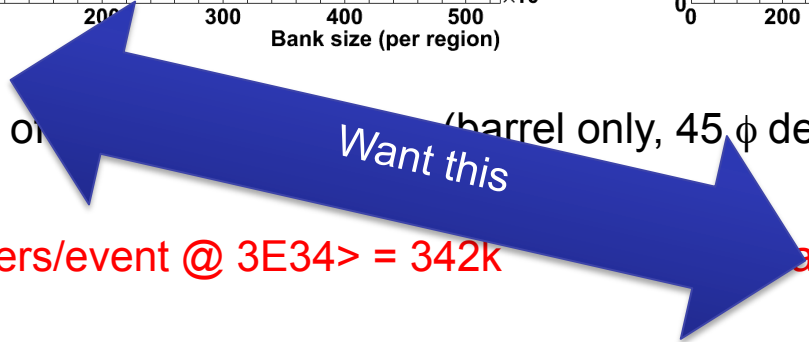
Pattern efficiency

Pattern size
 r- ϕ : 24 pixels, 20 SCT strips
 z: 36 pixels

Pattern size (half size)
 r- ϕ : 12 pixels, 10 SCT strips
 z: 36 pixels



of (barrel only, 45 ϕ degrees)

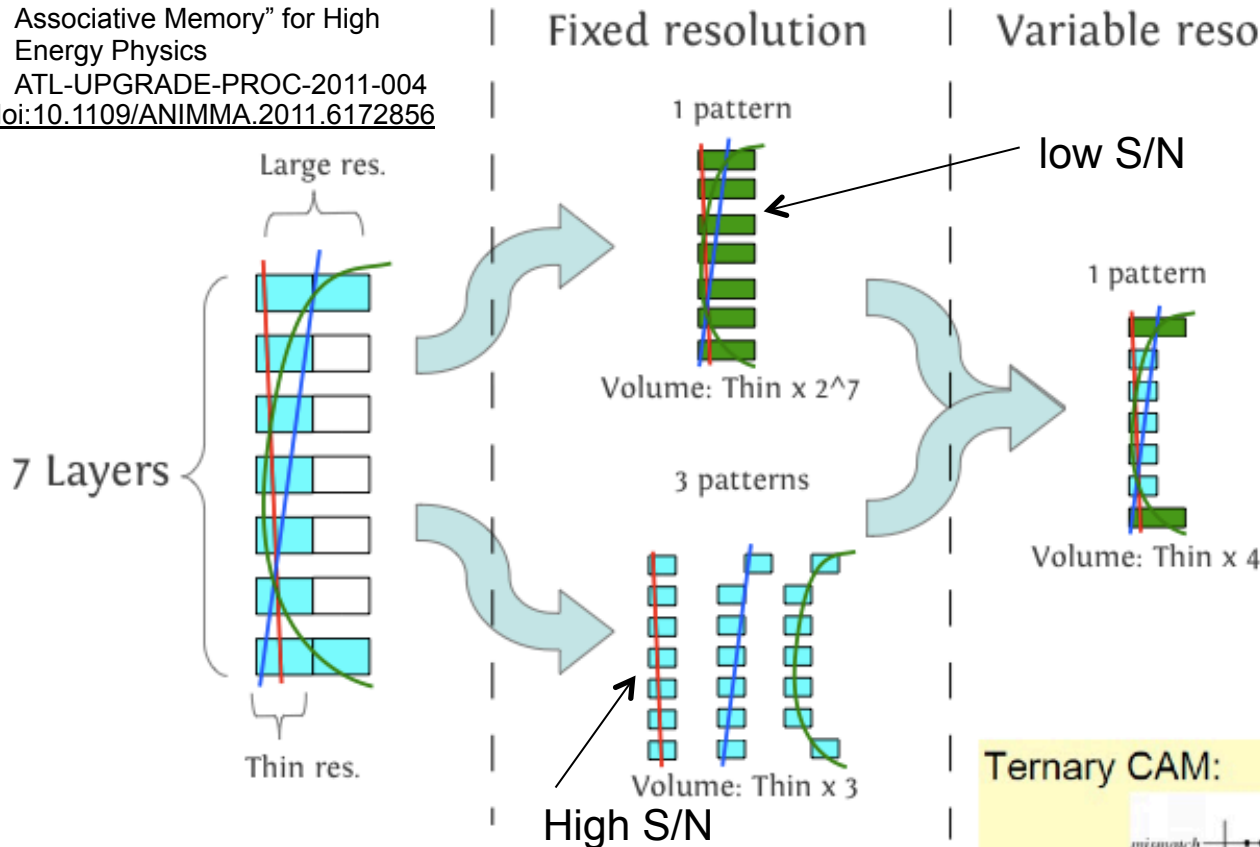


<# matched patters/event @ 3E34> = 342k

<# matched patters/event @ 3E34> = 40k

AMCHIP04: VARIABLE RESOLUTION

A new "Variable Resolution Associative Memory" for High Energy Physics
 ATL-UPGRADE-PROC-2011-004
 doi:10.1109/ANIMMA.2011.6172856



Good rejection and occupy only one pattern location.

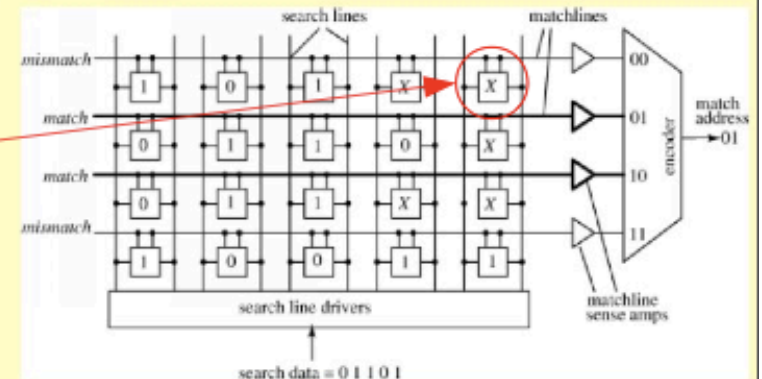
Per-pattern choice of optimal resolution.

We can use **don't care** on the least significant bit when we want to match the **pattern layer @ Large resolution** or use all the bits to match it **@ Thin resolution**

Coincidence window is programmable layer by layer and pattern by pattern

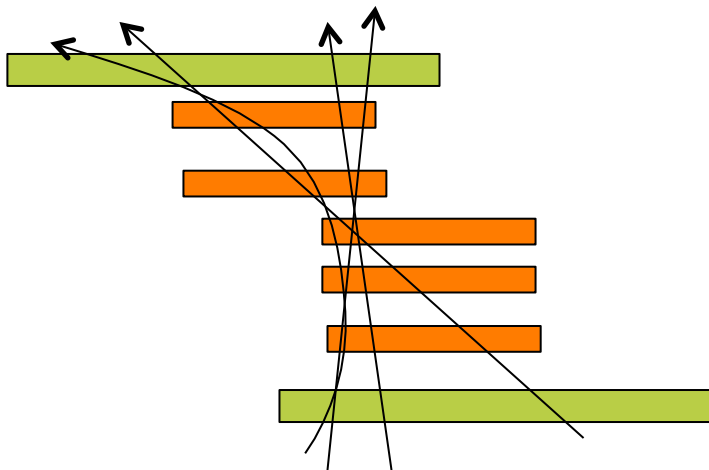
Ternary CAM:

feature: Don't Care Bits



Many bits variable resolution

1 bit variable resolution

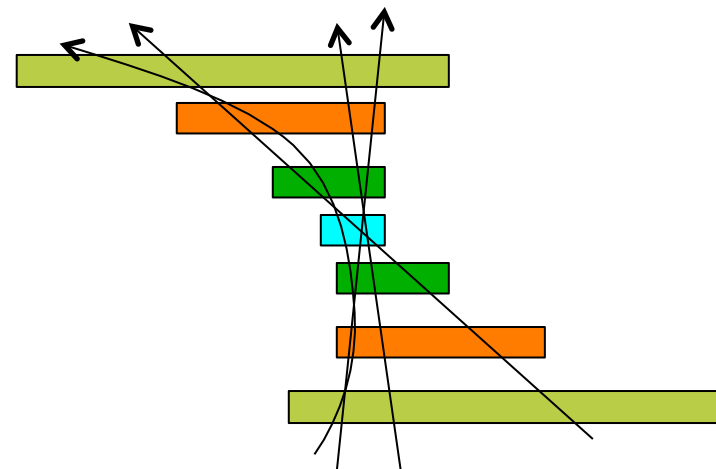


1 pattern

Volume 4^* 

Volume $2^{(7*2)*4^*}$  = 2^{16} 

3 bit variable resolution



1 pattern

Volume $1/4^*$ 

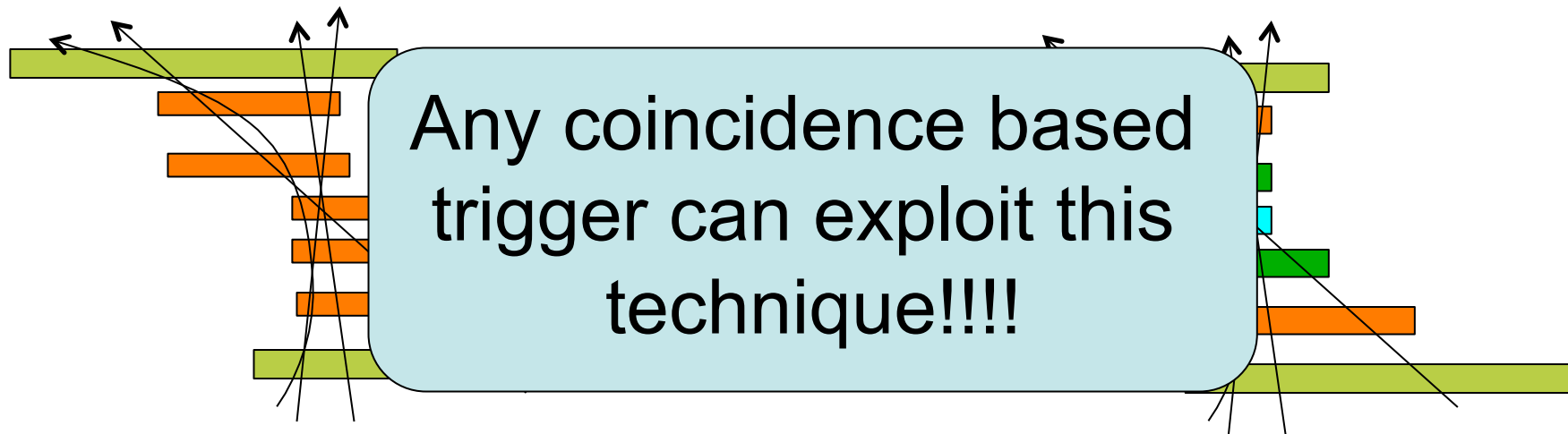
Volume 2^{12} 

**1/16 less volume
1/16 less fakes!!!**

Many bits variable resolution


1 bit variable resolution

3 bit variable resolution



1 pattern

Volume 4^* 

Volume $2^{(7*2)*4^*}$  = 2^{16} 

**1/16 less volume
1/16 less fakes!!!**



1 pattern

Volume $1/4^*$ 

Volume 2^{12} 

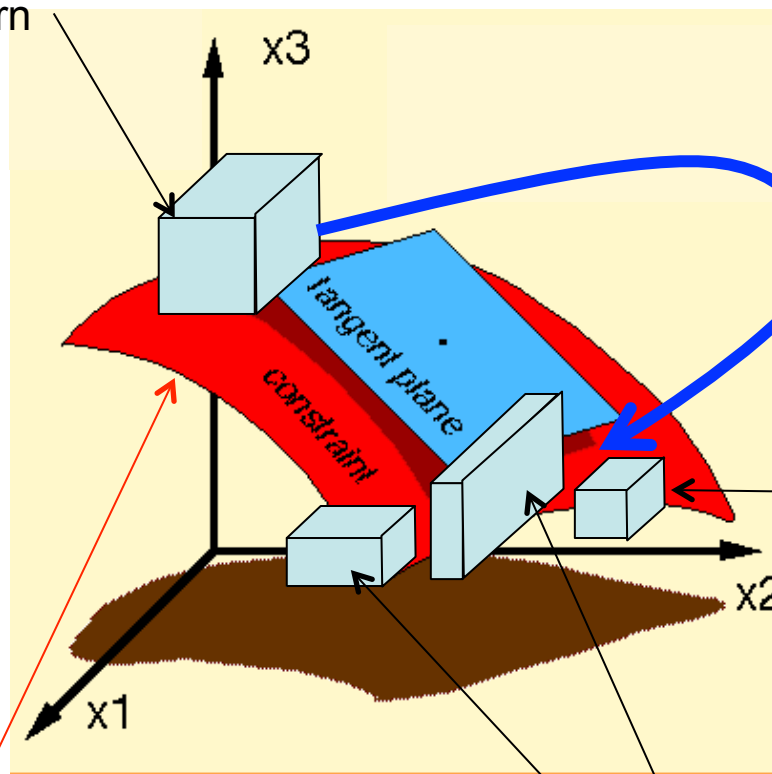
The patterns: a different point of view

5 strip + 3 pixel layers
→ 11 coordinates
→ 11D hit coord. space

A factor of 2 on each side
→ a factor 2^{11} less volume
→ $O(1/2048)$ less fakes!!

The pattern bank:
• cover the track manifold with patterns.
• covered space outside manifold → fakes.
• variable resolution → dramatically improves S/N

Large pattern

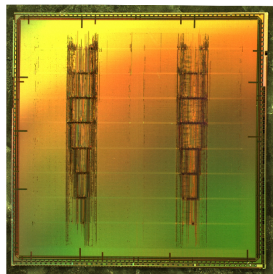
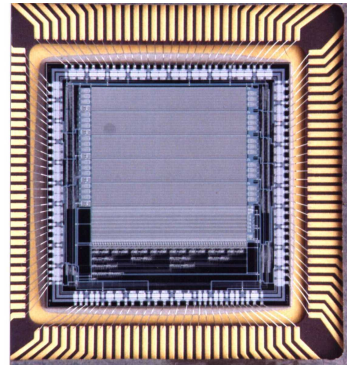


Thin pattern

5D track manifold

Variable resolution patterns

AM chips from 1992 to 2005



- (90's) **Full custom VLSI chip** - $0.7\mu\text{m}$ (INFN-Pisa)
- **128 patterns, 6x12bit words each**
- **384k patterns (SVT total)**

F. Morsani et al., “The AMchip: a **Full-custom** MOS VLSI Associative memory for Pattern Recognition”, IEEE Trans. on Nucl. Sci., vol. 39, pp. 795-797, (1992).

On the opposite side: **FPGA** for the same AMchip

P. Giannetti et al. “A Programmable Associative Memory for Track Finding”, Nucl. Instr. and Meth., vol. A413/2-3, pp. 367-373, (1998).

G Magazzu' I progetto standard cell presented @ LHCC (1999)

In the middle: **Standard Cell $0.18\mu\text{m}$** (INFN-Pisa-Ferrara) → **5000 pattern/chip** Amchip
SVT upgrade total: 6M patterns

L. Sartori, A. Annovi et al., “A VLSI Processor for Fast Track Finding Based on Content Addressable Memories”, **IEEE Transactions on Nuclear Science**, Volume 53, Issue 4, Part 2, Aug. **2006** Page(s):2428 - 2433

AMchip03 array (AM board)



See the FTK processing unit poster

AMchip Comparison

	AMchip03	AMchip04	Effect
Technology	180nm	65nm	x8 pattern density
Clock freq.	50MHz	100MHz	faster, higher power cons.
Die size	10x10mm ²	12x12mm ²	x1.5 patterns prototype 3.5x4 mm ²
Core voltage	1.8V	1.2V	lower power consumption
Core power	1.3W	2W	at 40MHz and 100MHz respectively
Selec. Prech.	No	Yes	~80% power saving
Full custom	No	Yes	x2 pattern density
Layers	6 (or 12)	8	¾ pattern density
Patters/chip	5k	80k	8k in AMchip04 prototype
Bits/layer	up to 18	up to 15	
Ternary/layer	N/A	3 to 6	better S/N with variable resolution

2 event "buffers": readout 1st, load 2nd event

The hard part: **FTK goal 1 billion pattern for LHC phase I**
 push pattern density to the limit, **keep power under control despite**
x16 patterns x8/6 layers and 40MHz --> 100MHz
 would mean x50 power consumption with same design & technology

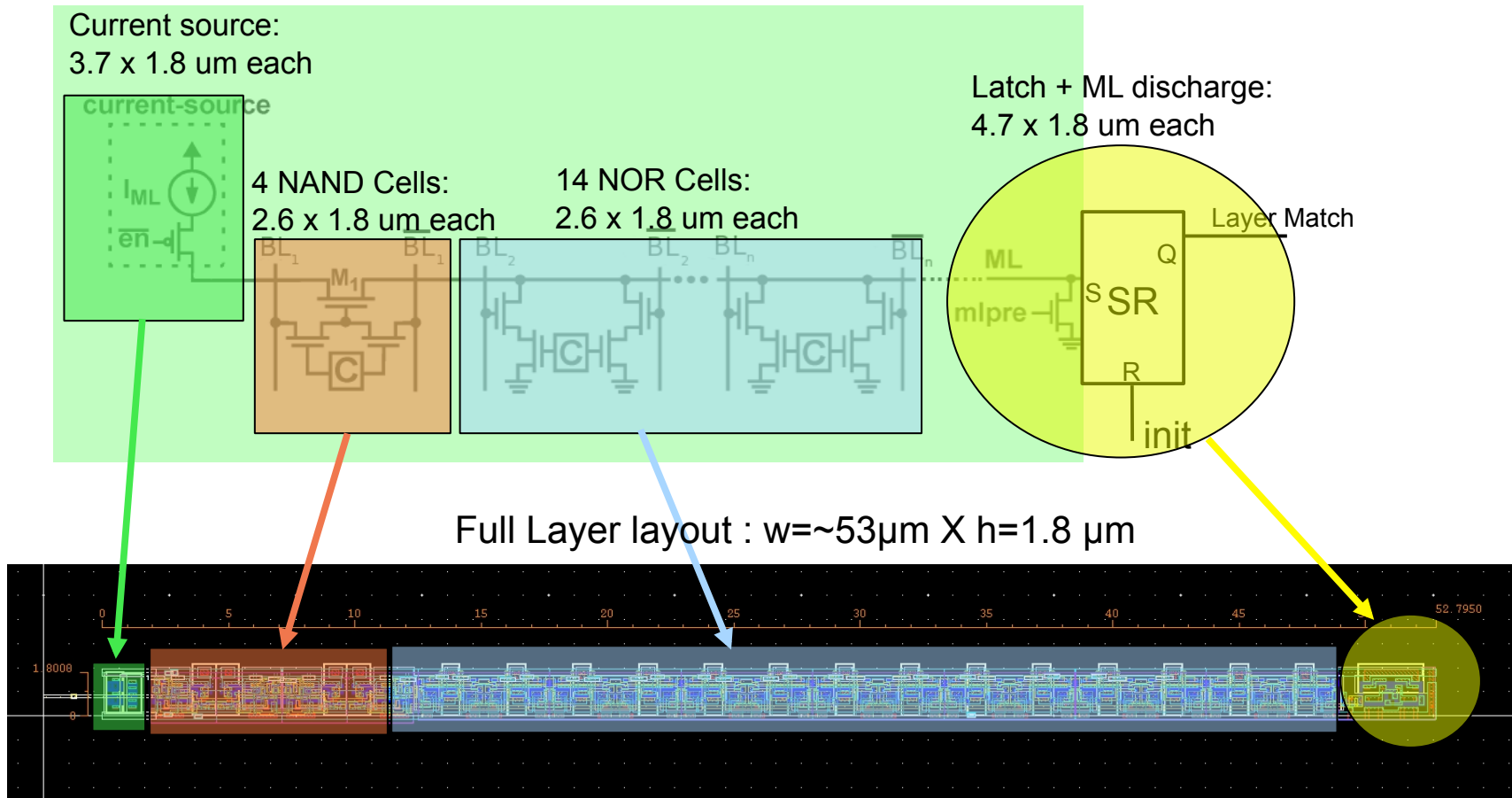
AM chip04 functions / specs

- Store pre-calculated trajectories (patterns)
 - Each pattern: 8 positions (numbers or words) one for each layer
- Compare patterns with incoming data
 - Detectors hits for one event
- For each event readout patterns
 - with enough hits 8/8, 7/8 or 6/8
- For each pattern readout:
 - Pattern address (ID) + bitmap of fired layers
- Configuration and pattern loading through JTAG interface

Associative Memory Layer

To save power we have used two different match line driving scheme:

- *Current race scheme (dummy layer timing)*
- *Selective precharge scheme*



Current race and selective - precharge schemes

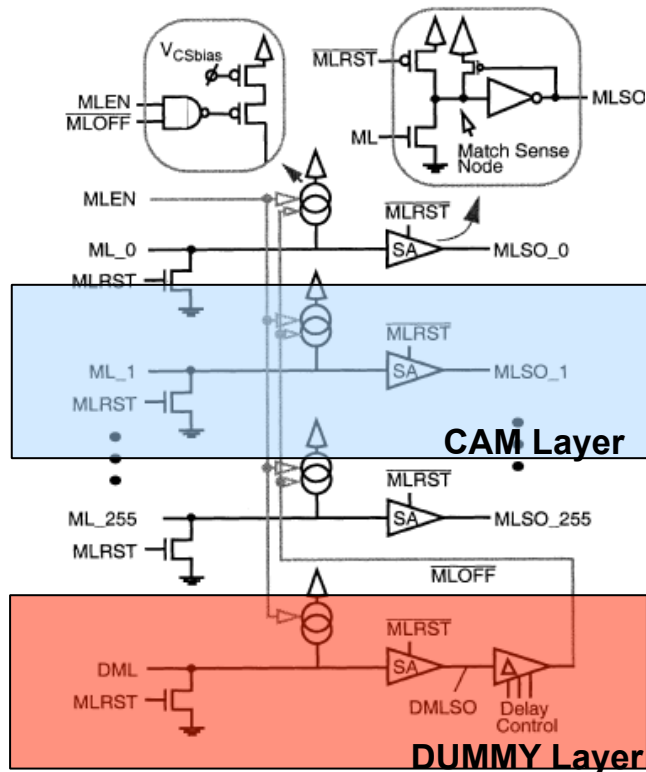


Fig. 5. Current-race ML sensing scheme.

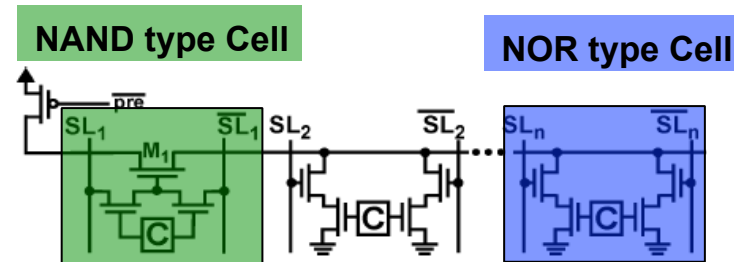
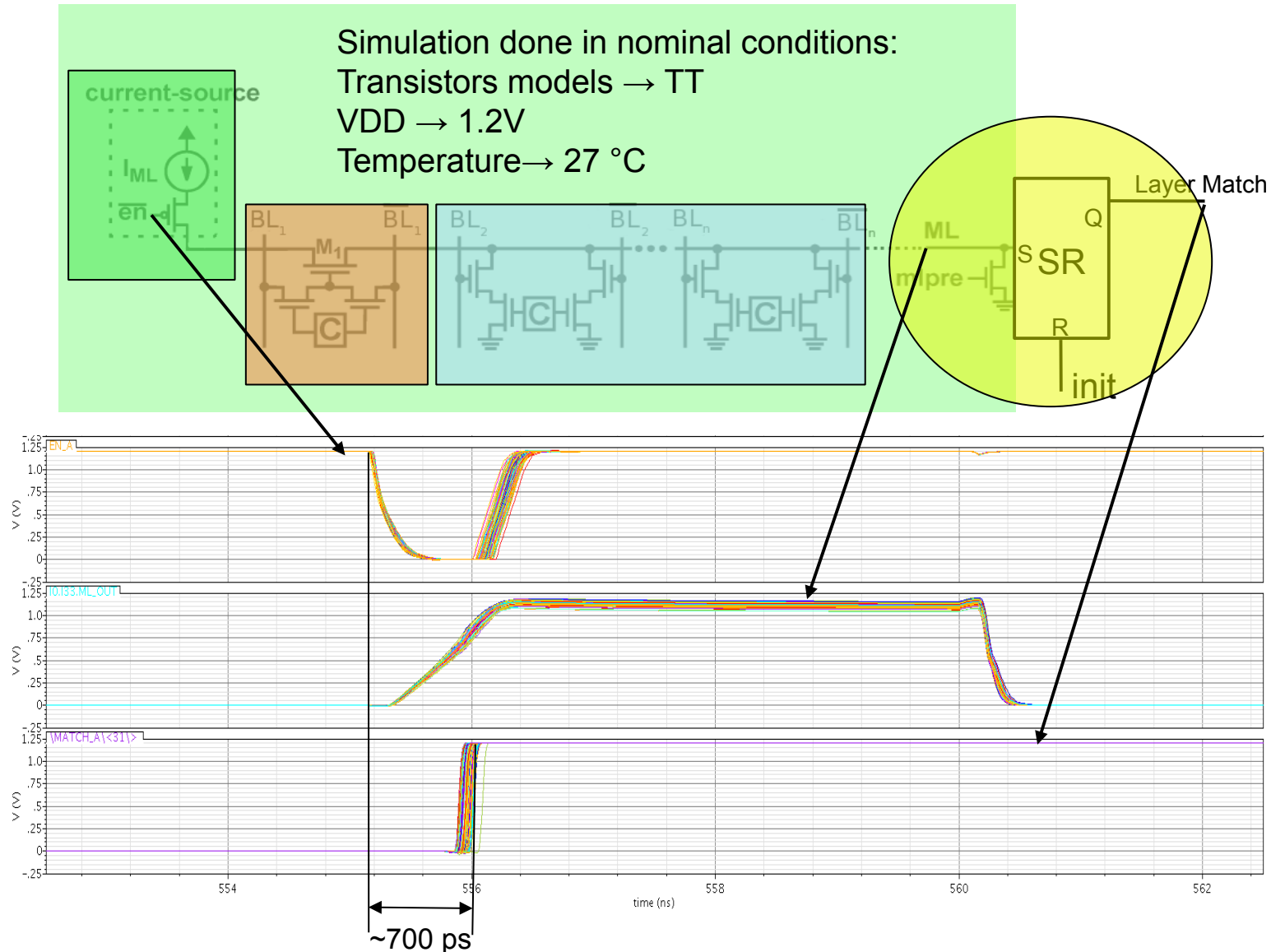


Fig. 16. Sample implementation of the selective-precharge matchline technique [43]. The first cell on the matchline is a NAND cell, while the other cells are NOR cells. Precharge occurs only in the case where there is a match in the first cell. If there is no match in the first cell, the precharge transistor is disconnected from the matchline, thus saving power.

Scheme from: "Content-Addressable Memory (CAM) Circuits and Architectures: A Tutorial and Survey", Kostas Pagiamtzis and Ali Sheikholeslami IEEE Journal of Solid-State Circuits, Vol. 41, NO. 3, March 2006

Scheme from: "A ternary content-addressable memory (TCAM) based on 4T static storage and including a Current-Race sensing scheme", Ali Sheikholeslami et al. IEEE Journal of Solid-State Circuits, Vol. 38, NO. 1, January 2003

CAM layer timing diagram



Power consumption rough estimates

We use the nominal simulation condition:

Transistor models : Typical

Power supply : 1.2 V

Temperature : 27 °C

Frequency : 100 MHz

To be verified with prototype measurement

These values do not take into account the standard cells part of the chip and the on chip power supply network distribution parasitic and other parasitic.

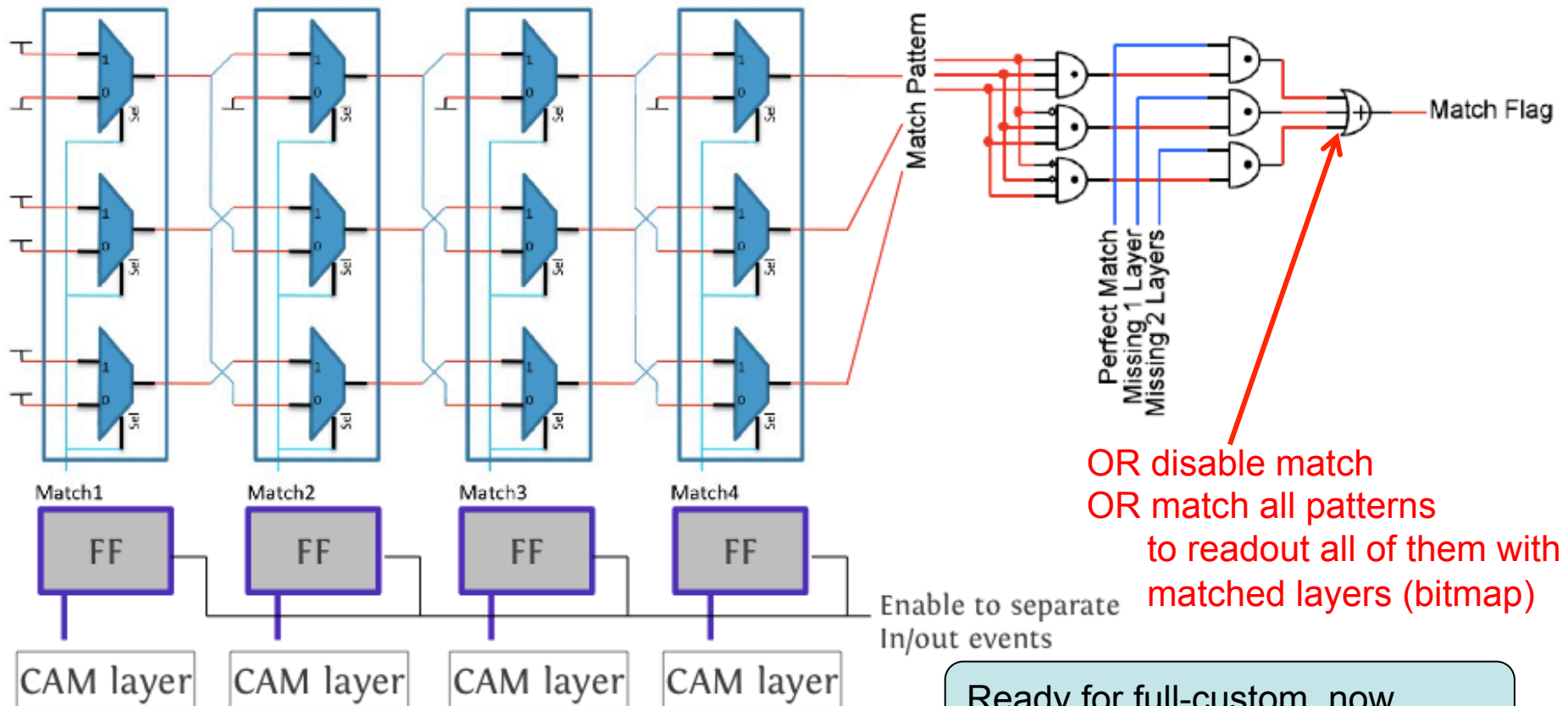
Memory state	Mean (mW)	Max (mW)	RMS (mW)
Write	23.04	557.57	28.12
Quiescent	21.89	22.09	21.88
Don't match	63.36	720.32	113.41
Match 1 out of 16 patterns	70.96	814.18	123.55
Match 1 out of 8 patterns	79.20	868.03	140.03
Match 1 out of 4 patterns	91.87	1045.44	172.34

Approximate consumption 80mW / 8kpattern / 100MHz \approx 100 μ W / kpattern / MHz
+ plus standard cell logic

AMCHIP04: MAJORITY LOGIC

Pattern match logic is made by identical logic for each layer:
 Receives in input 3 bits, if not matching shift down the output.

At the end of the chain the 3 bits are compared with the majority requirements:
 perfect match, 1 or 2 missing



OR disable match
 OR match all patterns
 to readout all of them with
 matched layers (bitmap)

Logic proposed by J. Hoff (Fermilab)

Ready for full-custom, now
 syntethized with standard cells

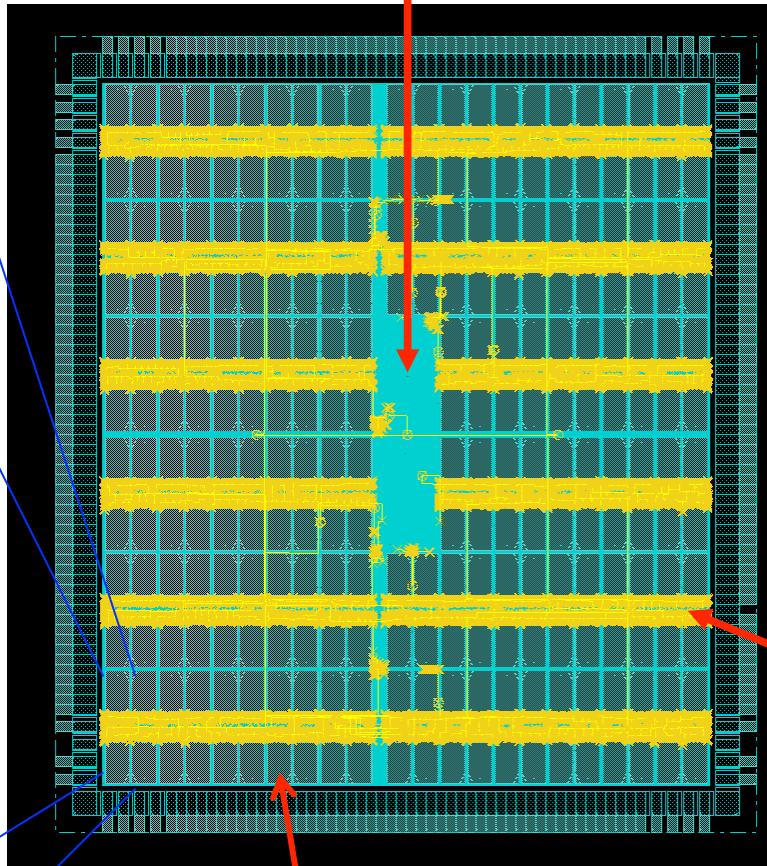
Layer matches from CAM layers are stored in a FF just before resetting the CAM layers.
 We can load an event in the CAM layers while we are reading the patterns found in the previous event.

Prototype Chip Layout

64 patterns



Control logic



The AMchip has an area of 14 mm²

CAM is organized as 22 column x 12 row matrix of full custom memory blocks

Each block is 64 x 4 layers

Between two rows of blocks there is the majority logic and the readout logic implemented with standard cells

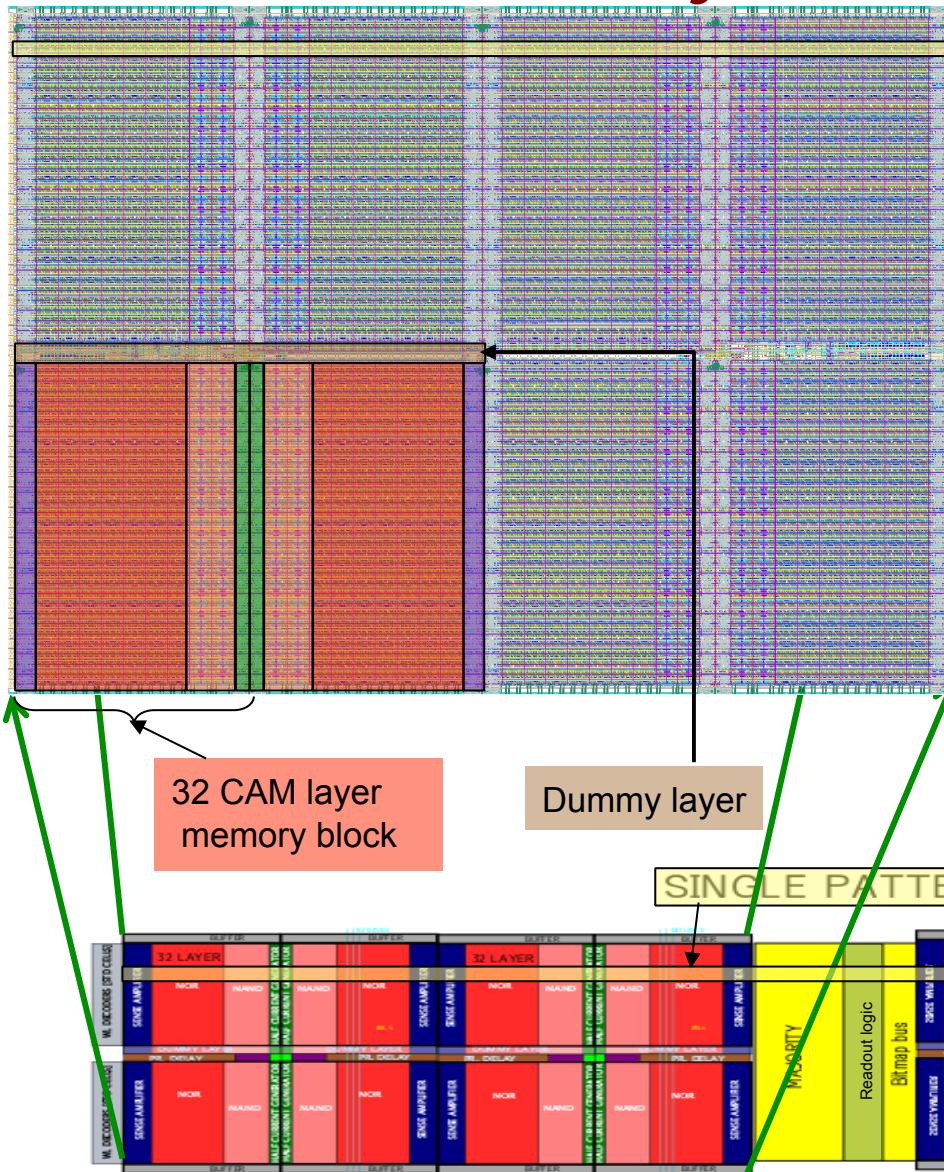
In the center there is the control logic implemented with standard cells

Majority logic and readout logic

size: 3510 μm \times 3985.0 μm

2x128 blocks: 64 half patterns each

Memory Block Layout



→ 4 Layers = 1/2 pattern

Full custom Layout of 64 x 4 CAM layers (half pattern):
w~226 μm h~123 μm

without including: majority logic,
readout logic, and control logic

Six metal layers are used to route
signals, power supply and ground.

Bit lines are routed vertically while
control lines and memory output are
routed horizontally

32 CAM layer
memory block

Dummy layer

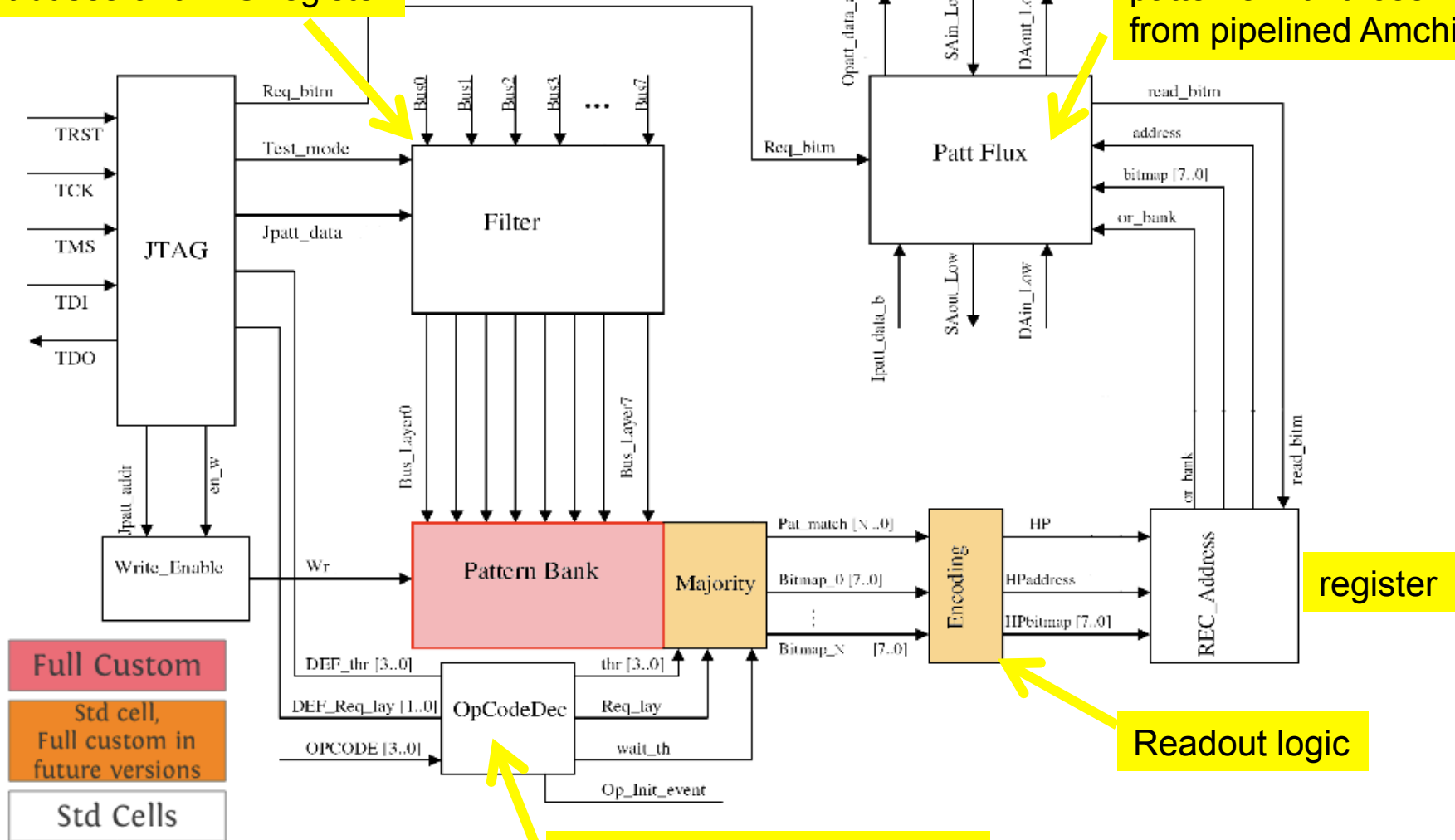
SINGLE PATTERN

64 patterns (vertically)

AMchip TOP level

Prepare data for match/write
Input buses or JTAG register

Multiplex internal
patterns with those
from pipelined Amchips



Full Custom

Std cell,
Full custom in
future versions

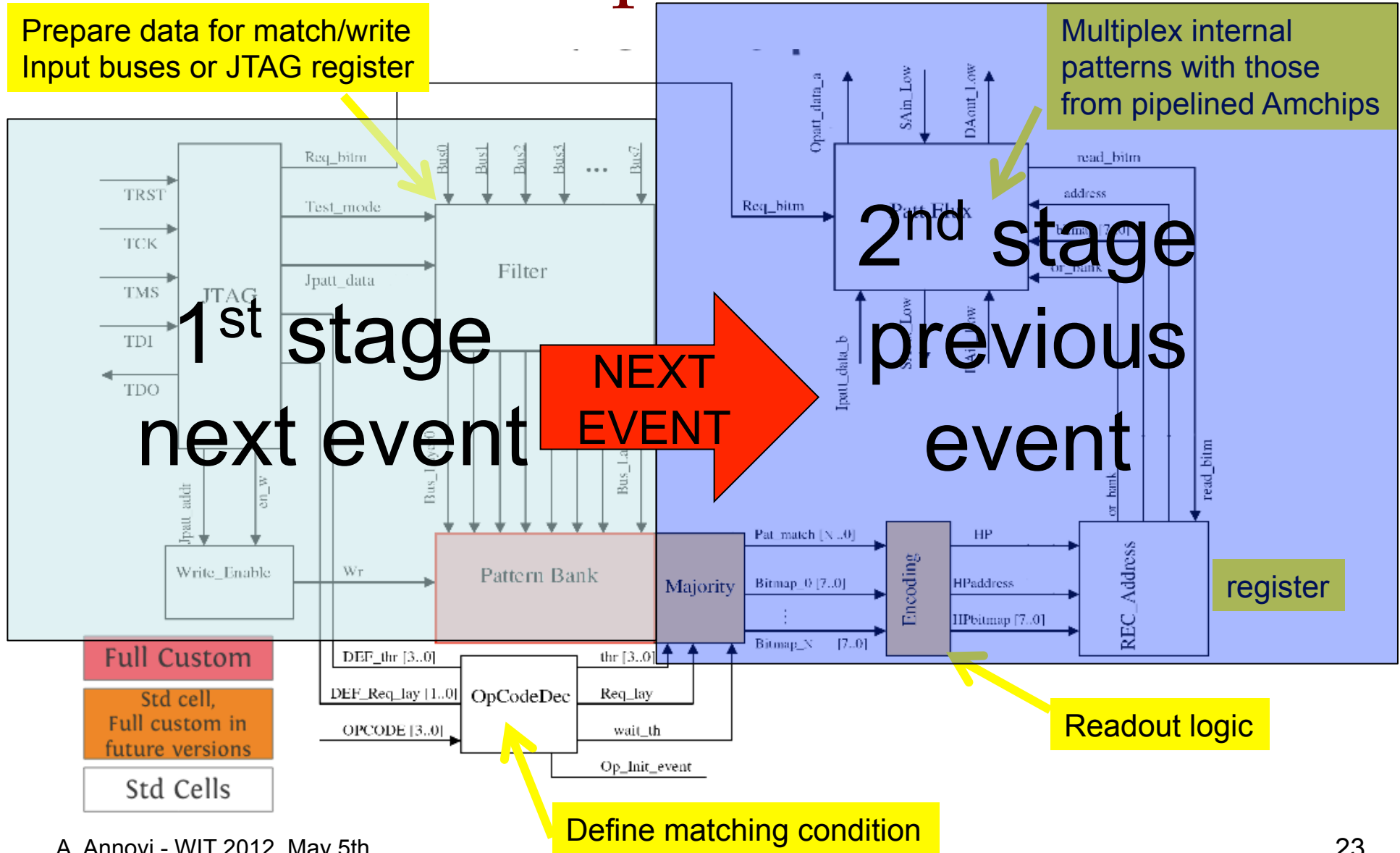
Std Cells

Define matching condition

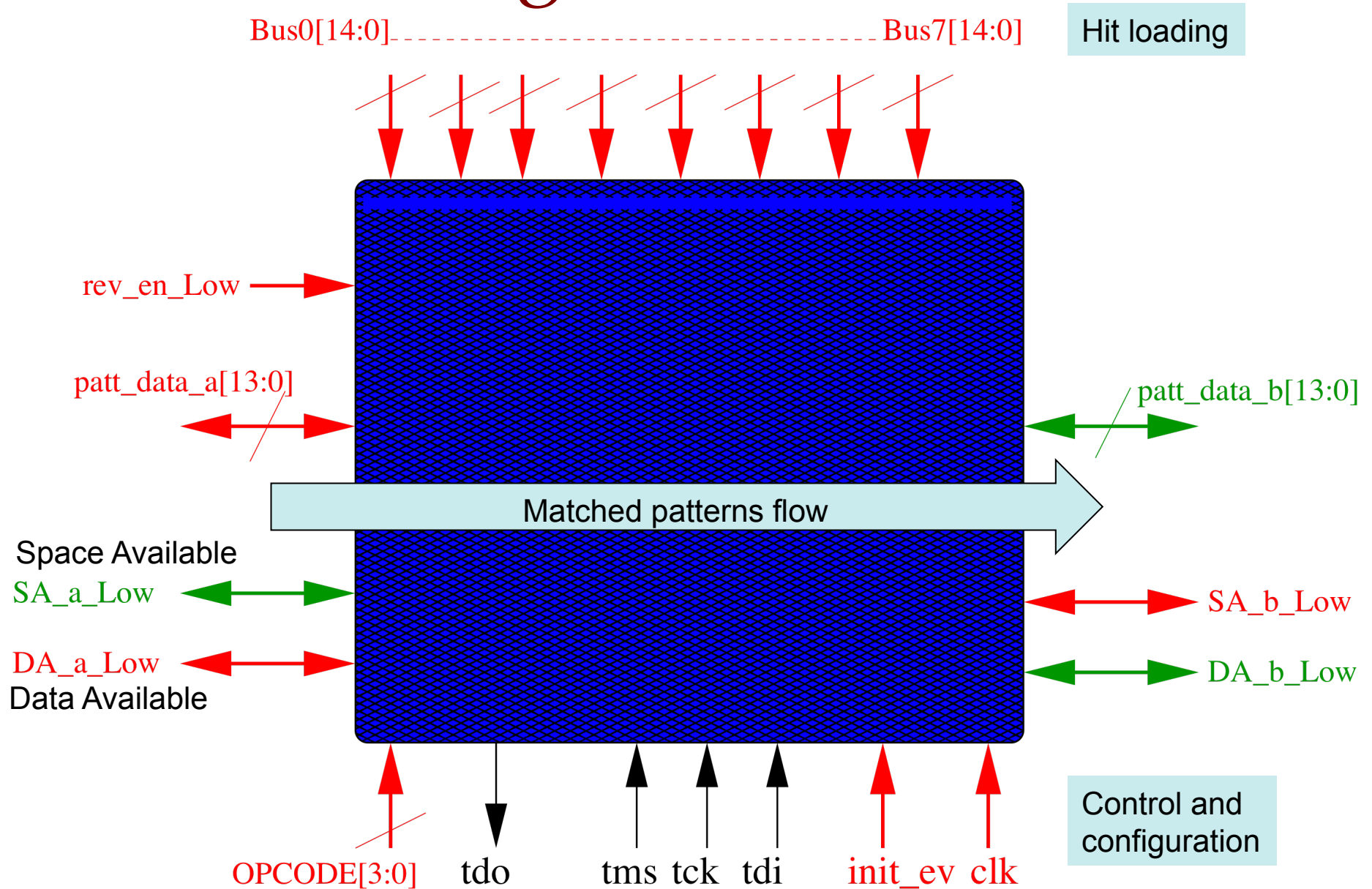
Readout logic

register

AMchip TOP level



Logical Pinout



OPCODES

Pattern matching condition changed dynamically with opocdes

code	mnemonic	delay	description
0	NOP	-	Do nothing.
1	SET_0MISS	2	Set the threshold to 0-miss.
2	SET_1MISS	2	Set the threshold to 1-miss.
3	SET_2MISS	2	Set the threshold to 2-miss.
5	INIT_EV	2	1st cycle: reset 2nd stage logic and set THR and required_layers to the default value ("DEF_" registers), copy matched layers to 2nd stage; 2nd cycle reset 1st stage logic
7	DEC_THR	2	Decrease THR register by 1.
10	FORCE_MATCH	2	Set THR to 0 and force all patterns to be readout along with their bitmap.
11	DIS_MATCH	2	Disable the matches.
14	TOGGLE_REQ	2	Enable/disable (toggle) the request of layer0 match
4, 6 8, 9, 12 13,15	reserved	-	

Other applications

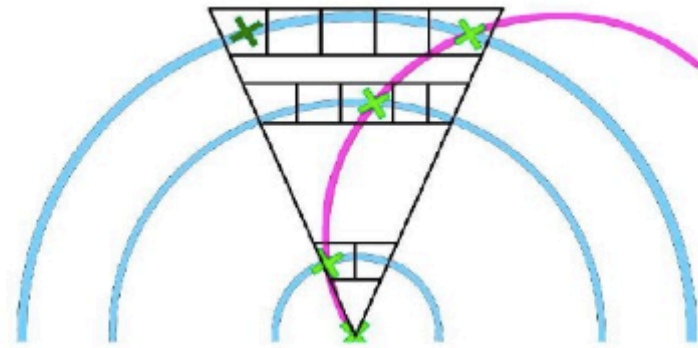
- CMS L1 track trigger
 - First study of Amchip used to trigger with doublets
 - See Thursday G. Boudoul's talk
- ATLAS L1 track trigger
 - Being evaluated
 - See Thursday R. Brenner's talk
- ATLAS Muon upgrade trigger (?)
 - Find micromega muon segments

Summary

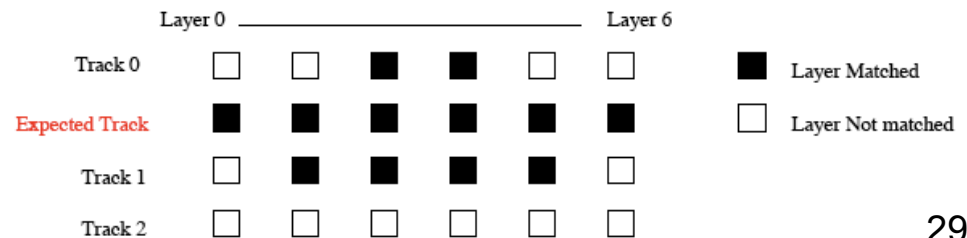
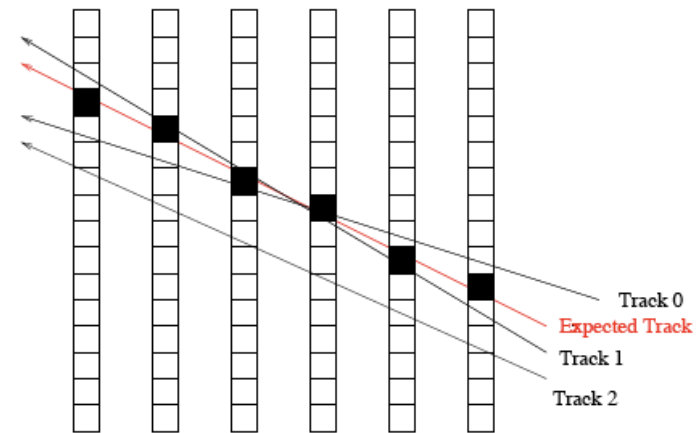
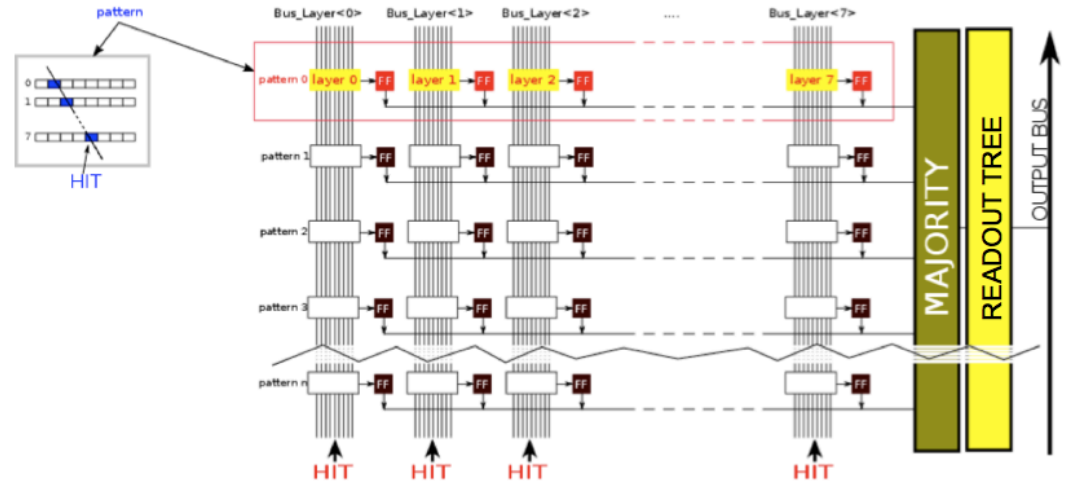
- Designed a new Associative Memory
- First application: ATLAS Fast-Tracker
- Special care to minimize power consumption & increase patt. density
- NEW: introduce powerful variable resolution pattern-matching !!!
 - any coincidence based trigger can profit
 - equivalent to a factor 3-5 extra patterns
- Prototype main goal: verify functionality of new features and full-custom cell
- Expecting delivery of AMchip04 prototype this month

BACKUP

AM working principle



	Lay0	Lay1	Lay2
Template:	01	10	000
Search data:	01	10	000 100



CAM cell configuration

- 18 CAM bits per layer: 4 NAND and 14 NOR
 - NOR pairs can make a ternary cell
- Default 12 bits + 3 ternary (minimum)
 - 15 bits per input bus (maximum)
 - (14:7) NOR, (6:3) 4 NAND, (2:0) 3 NOR-pairs
- 6 bits + 6 ternary (maximum)
 - Use only 12 bits per input bus
 - (11:10) NOR, (9:6) 4 NAND, (5:0) 6 NOR-pairs
- Ternary cells (NOR pairs) mapped to LSBs
- NAND cells are mapped to LSBs after the ternary cells, when they don't match small power consump.

Ternary CAM Cell with two NOR type cells

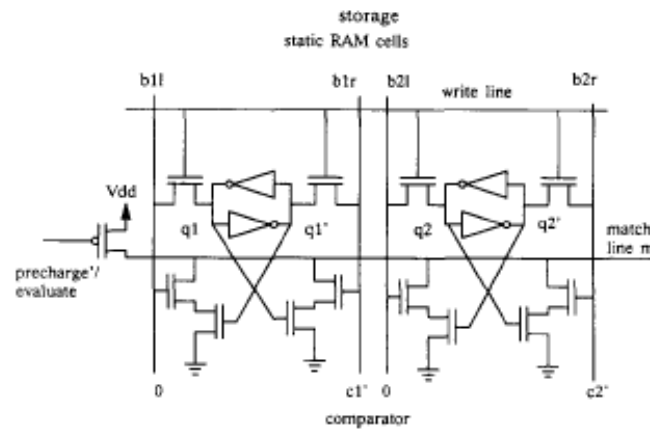


Fig. 8. Two adjacent static binary CAM cells.

storage scheme
stored values

q1q2
0 01
1 10
* 00
(a)

retrieval scheme
presented values

presented ternary value	encoded value in the bit lines of two binary static CAMs.	binary CAM equivalent operation
	c1c2 b1l b1r b2l b2r	l r
0	01 0 1 0 0	0 M*
1	10 0 0 0 1	M 0
*	11 0 0 0 0	M M

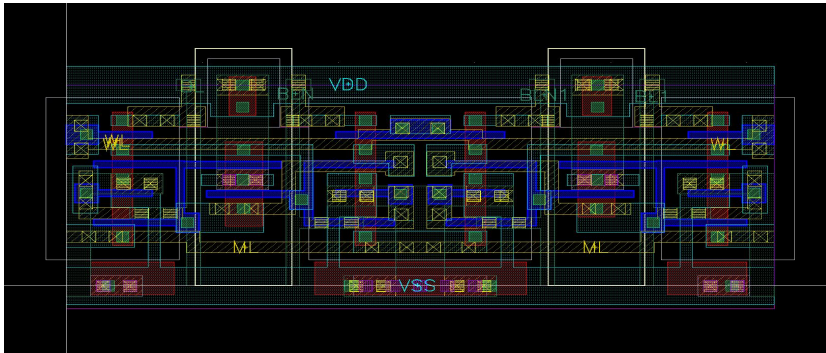
*M is the masking of a bit operation common in commercial binary CAMs.

(b)

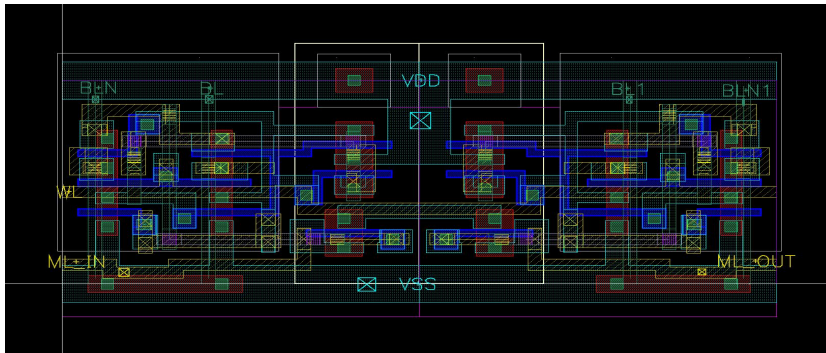
Fig. 9. Encoding and retrieval schemes for don't-care in two static binary CAM's cells with masking capability. (a) Encoding scheme. (b) Retrieval scheme.

Images from: "Encoding Don't Cares in Static and Dynamic Content-Addressable Memories", Sergio R. Ramirez-Chavez, IEEE Transactions on circuits and systems-II: Analog and Digital Signal Processing, Vol. 39 NO. 8, August 1992

NOR and NAND Cell Layout



Double NOR cell layout
10 transistors each cell
Dimensions: 5.5 X 1.8 μm^2



Double NAND cell layout
9 transistors each cell
Dimensions: 5.27 X 1.8 μm^2

Layout of memory cells is of type “wide” with uniform orientation of all the transistors in the cells to provide better reproducibility (W and L matching) and transistors V_{th} matching.