# High-Throughput Low-Energy Content-Addressable Memory Based on Self-Timed Overlapped Search Mechanism

Dr. Naoya Onizawa

Postdoctoral Fellow, McGill Univerisity

naoya.onizawa@mail.mcgill.ca

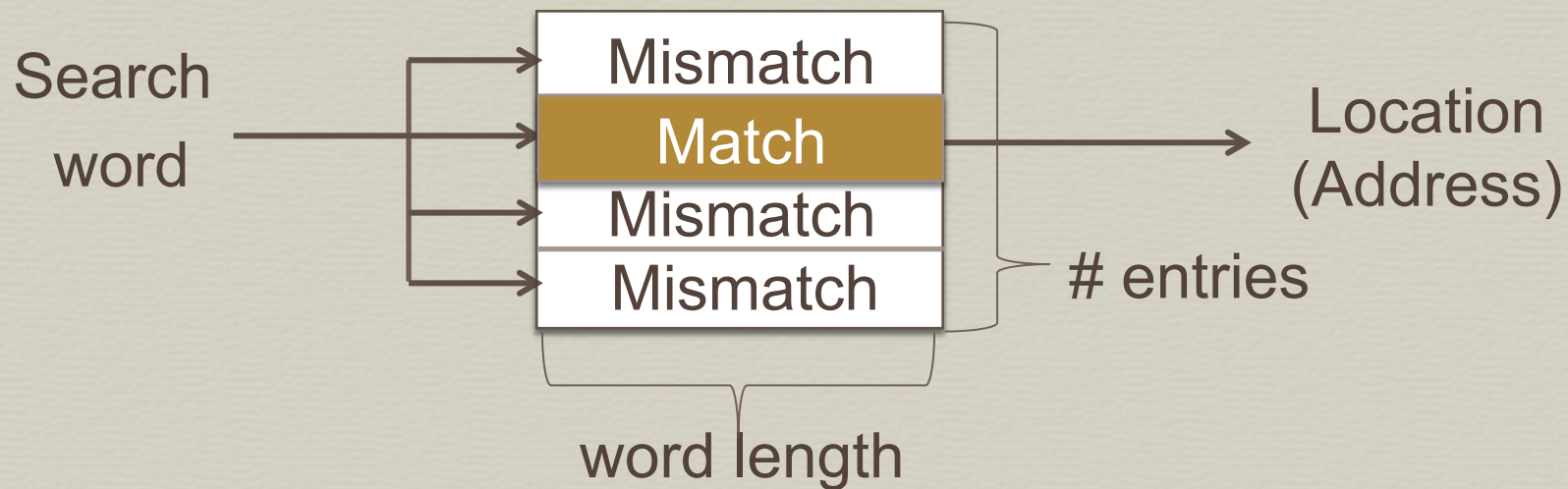Collaborators: S. Matsunaga, V. Gaudet, and T. Hanyu

# Content-Addressable Memory (CAM)

- Associative memory
- Parallel searching
- Applications
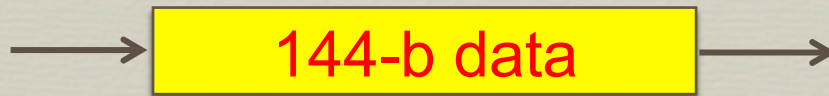  - Cache, Virus checking
  - Packet forwarding (40G, 100Gbps)

Search word → | Mismatch / **Match** / Mismatch / Mismatch | → Location (Address)

\# entries

word length

# Hardware-Implementation Issue

Speed restriction

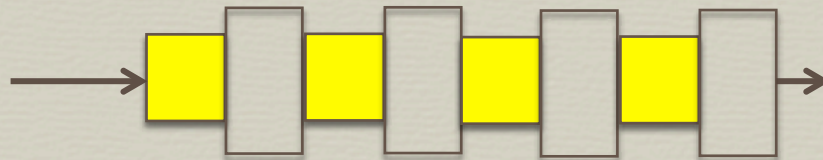● Packet length - 32bit (IPv4), 128,144bit (IPv6)

Search word → | 144-b data | → ➢ Large matching delay

Pipelined approach  K. Pagiamtzis, et al (JSSC'04 vol.39-9)

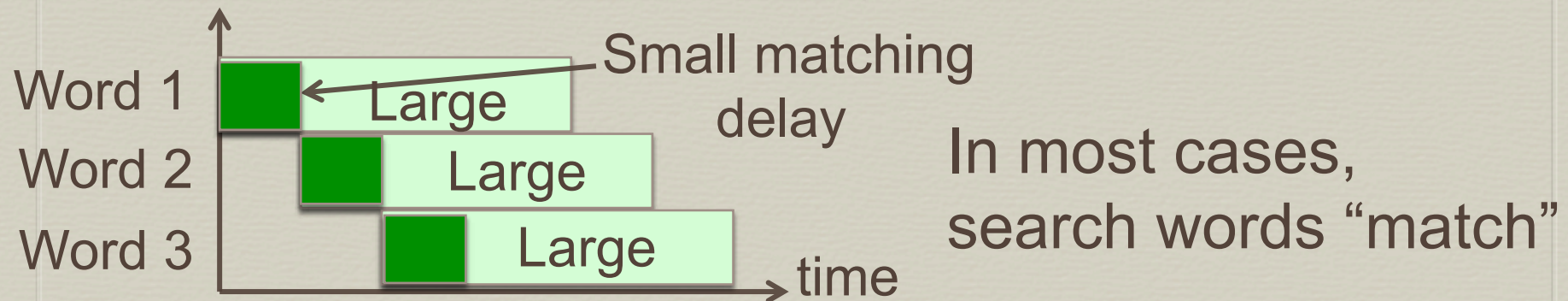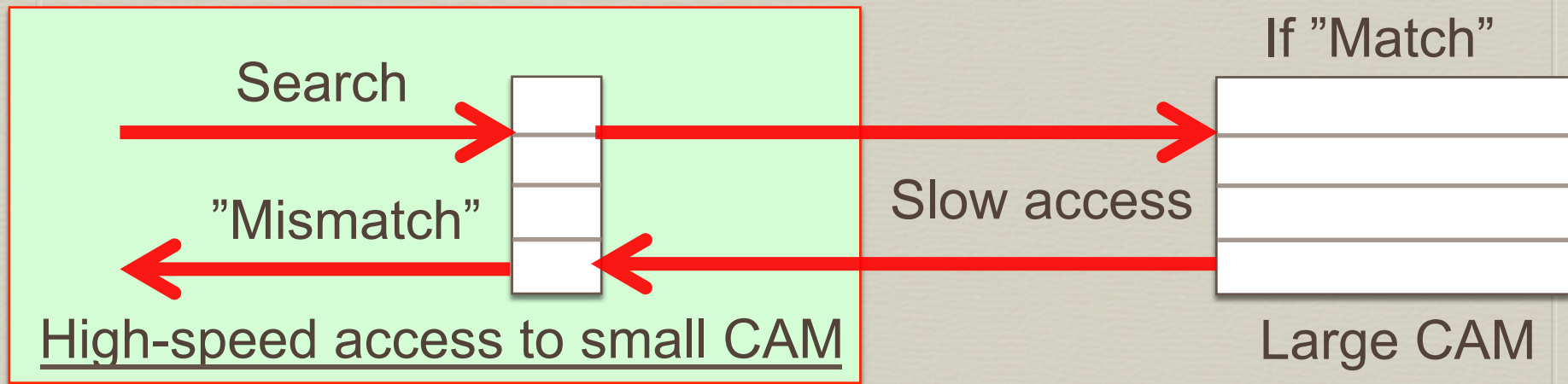Search word → ➢ Large area and power dissipation

Goal: High-throughput low-overhead CAM

12/05/07

# Concept

- Operate as comparable to small CAM



Search

"Mismatch"

High-speed access to small CAM

If "Match"

Slow access

Large CAM

Word 1
Word 2
Word 3

Large
Large
Large

Small matching delay

time

In most cases, search words "match"

Hide large matching delay to improve throughput

# Approach

- Assign search words to unused blocks

After searching first few bits, most blocks are mismatched

➤ If unused blocks are found, it doesn't need to wait to search new words until the current search is complete.

Pre-computation block

(Find unused blocks before sending)

| Unused |
| In use |
| Unused |
| Unused |

Search new words

after searching few bits

Partitioning

5.57x higher throughput at 8% cost of area

# Table of Contents
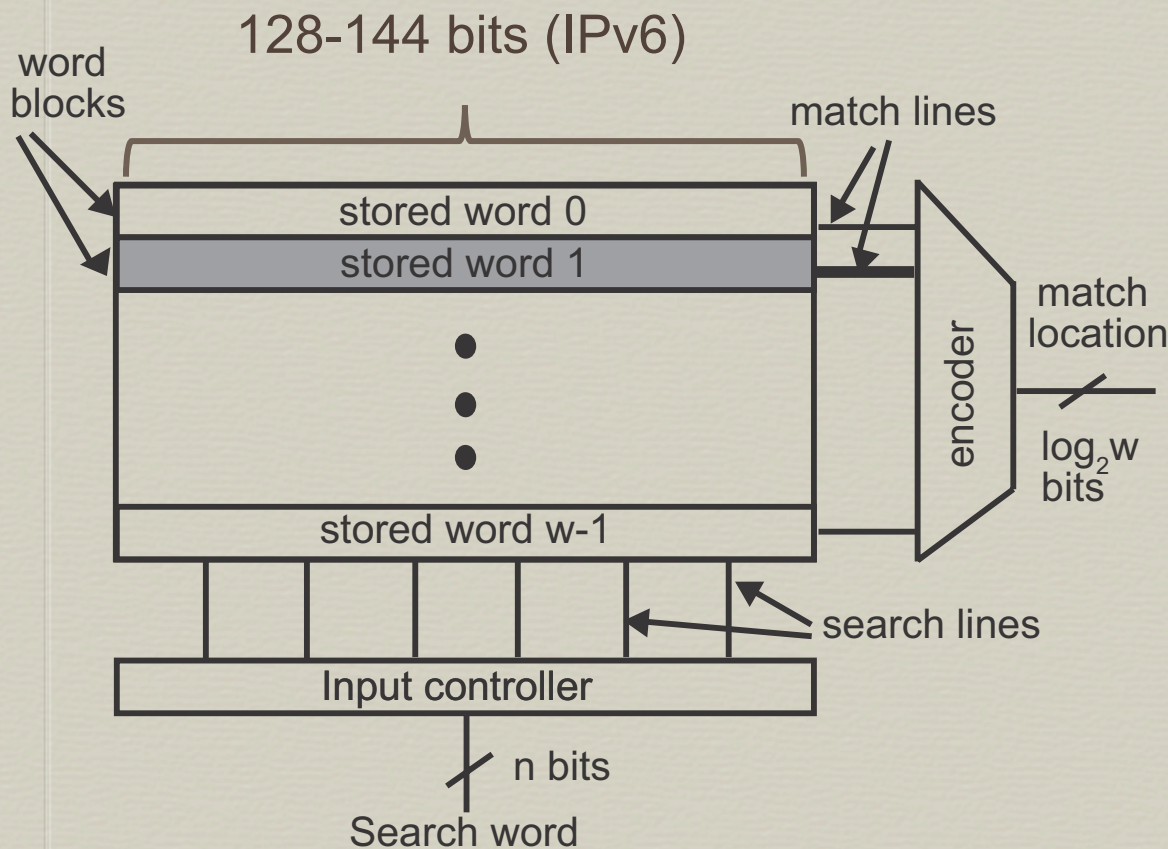
- Introduction to content-addressable memory
- Overlapped search mechanism
  - Word overlapped search
  - Phase overlapped processing
- Hardware implementation
- Evaluation
- Conclusion and future prospect

12/05/07

# CAM

128-144 bits (IPv6)

word blocks

match lines

stored word 0

stored word 1

encoder

match location

$\log_2 w$ bits

stored word w-1

search lines

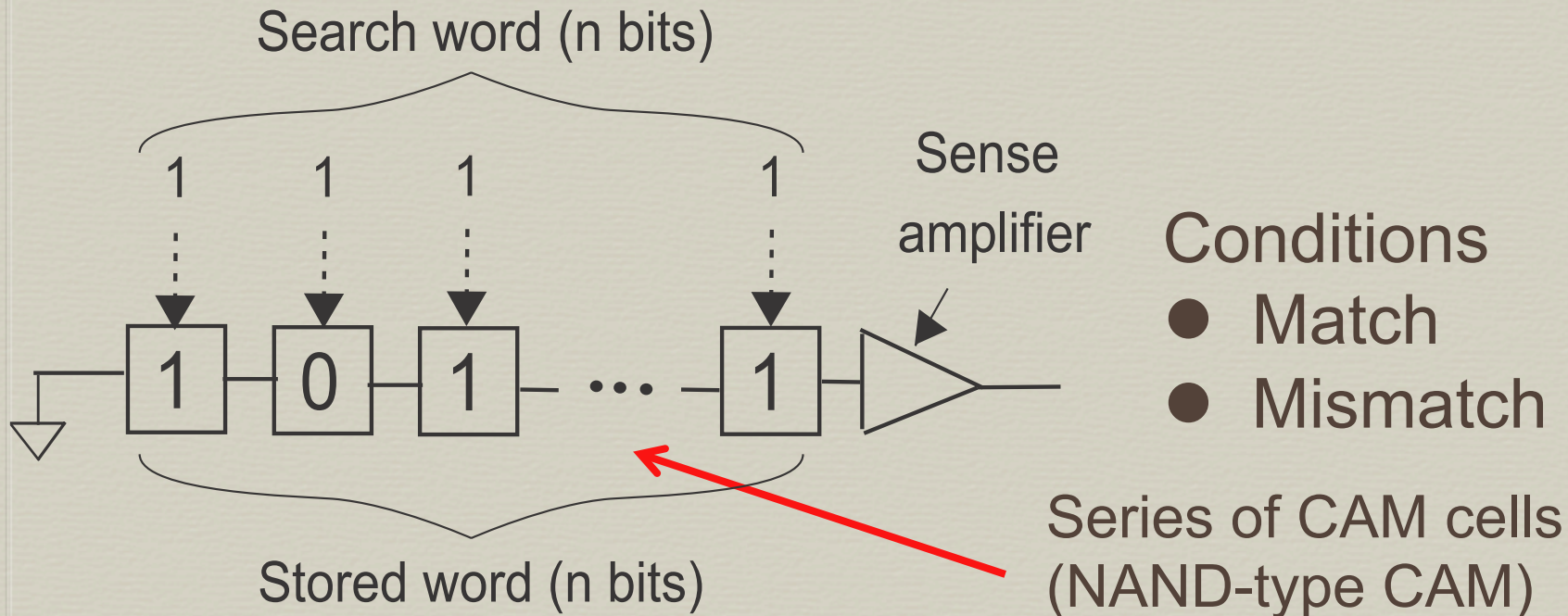Input controller

n bits

Search word

## Operation

1. Search all words in parallel

2. Find a matched word block

3. Output a matched location (address)

## Word-parallel search in single cycle

# CAM Word Block

Search word (n bits)

1    1    1         1

Sense amplifier

1    0    1  ...  1

Conditions
- Match
- Mismatch

Stored word (n bits)

Series of CAM cells (NAND-type CAM)

Throughput determined by word length in conventional CAM

Long word length degrades throughput

# CAM Characteristics

Matching probability of word blocks after k-bit search is

$$p_{matched} = \left(\frac{1}{2}\right)^k$$

Most word blocks are not used after k-bit search
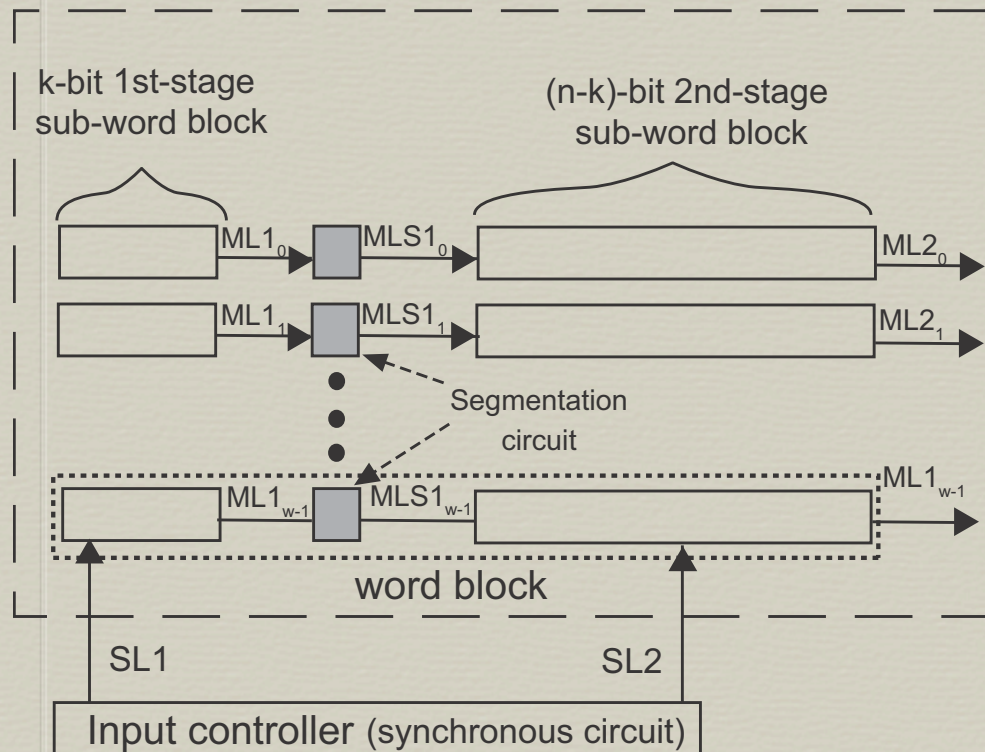(We set k=8 in the hardware implementation)

➢ Use unused blocks to improve throughput

Most word blocks are unused (mismatched).

# Word Overlapped Search (WOS)

## CAM architecture based on segmentation method

CAM block (self-timed circuit)

k-bit 1st-stage sub-word block

(n-k)-bit 2nd-stage sub-word block

$ML1_0$  $MLS1_0$  $ML2_0$

$ML1_1$  $MLS1_1$  $ML2_1$

Segmentation circuit

$ML1_{w-1}$  $MLS1_{w-1}$  $ML1_{w-1}$

word block

SL1  SL2
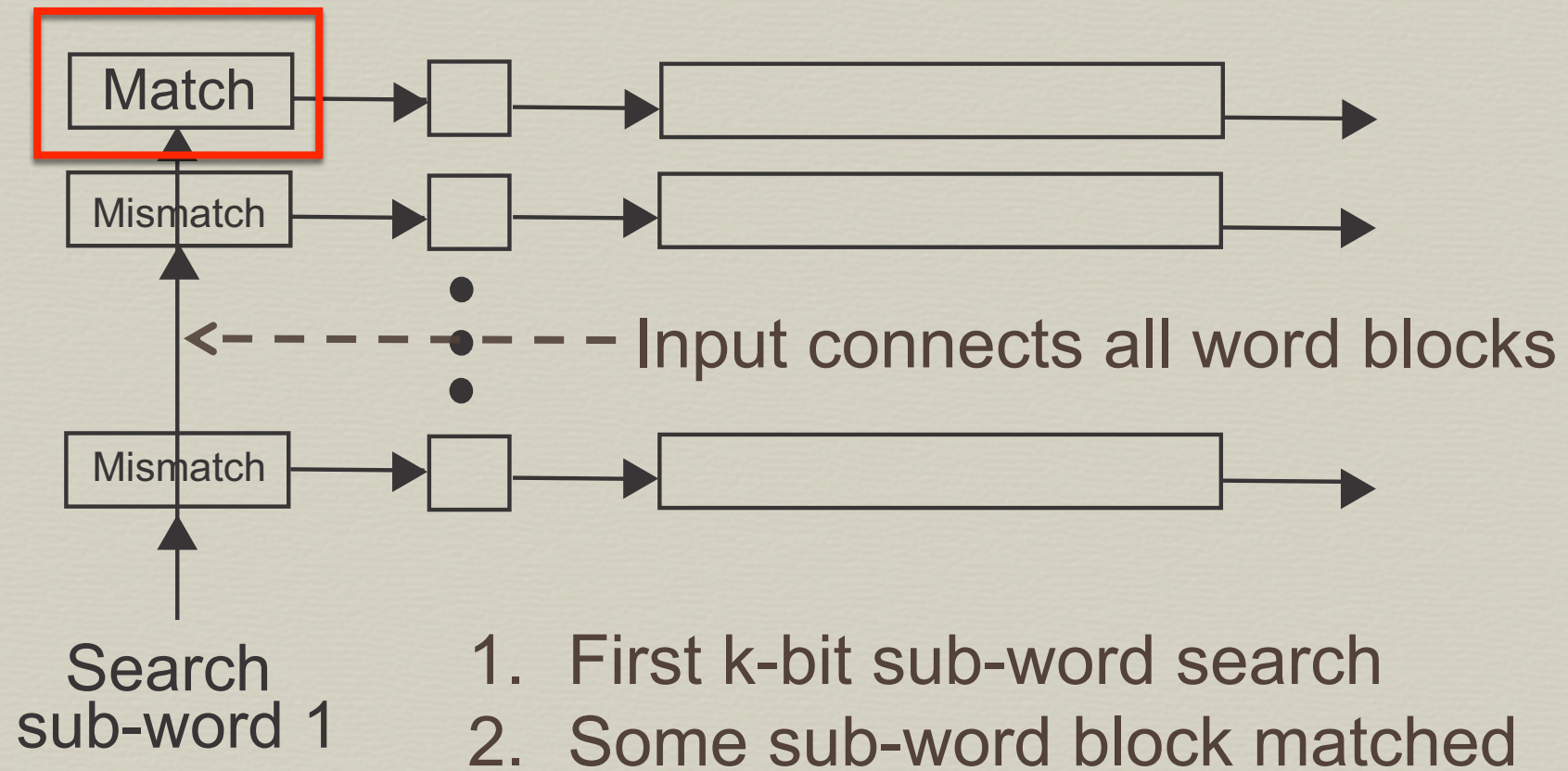
Input controller (synchronous circuit)

1. Partition word block to:
   a) small k-bit block and
   b) large (n-k)-bit block
   by segmentation block

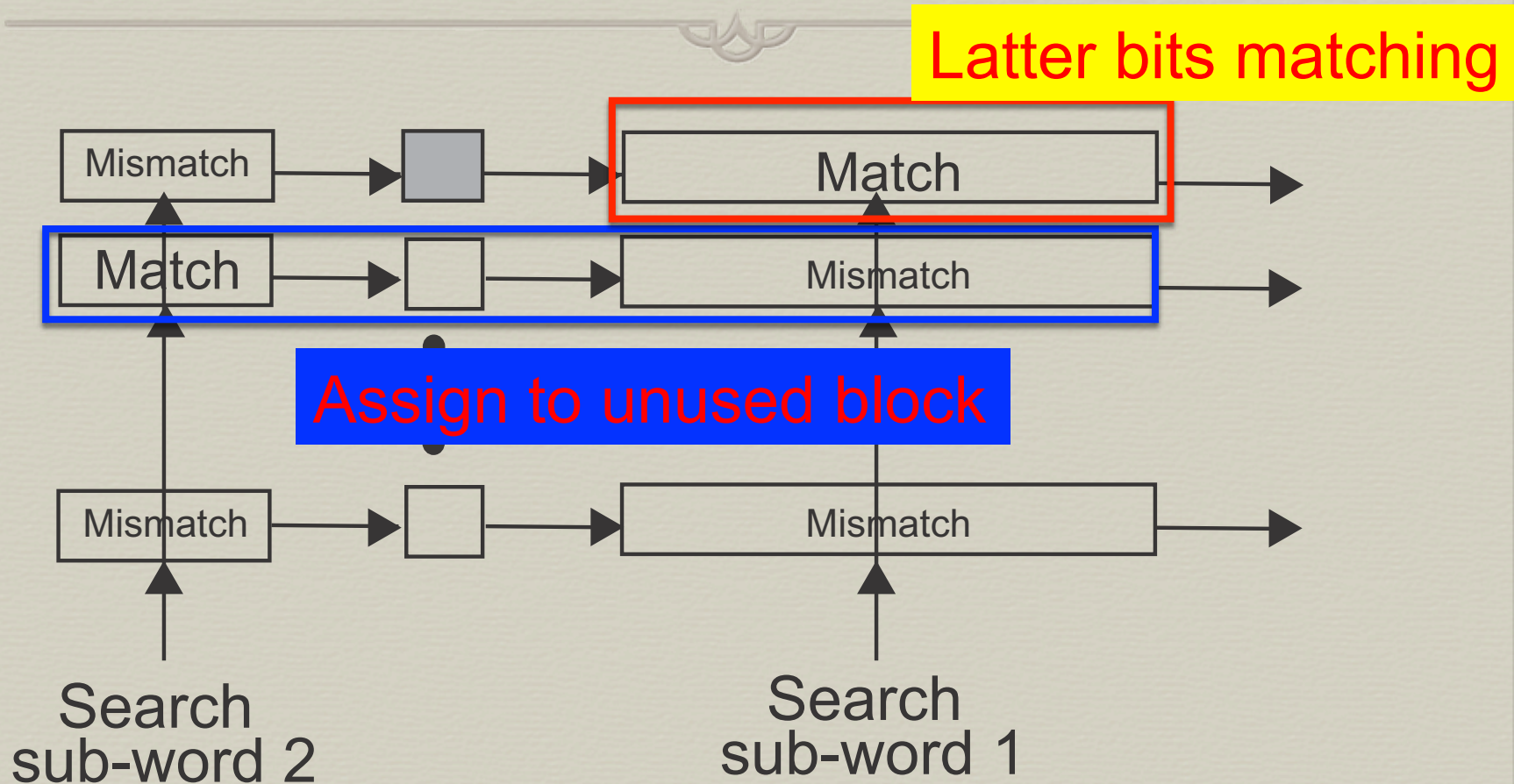2. Segmentation block stores Its k-bit matched result

3. Latter block operates when the first block matches

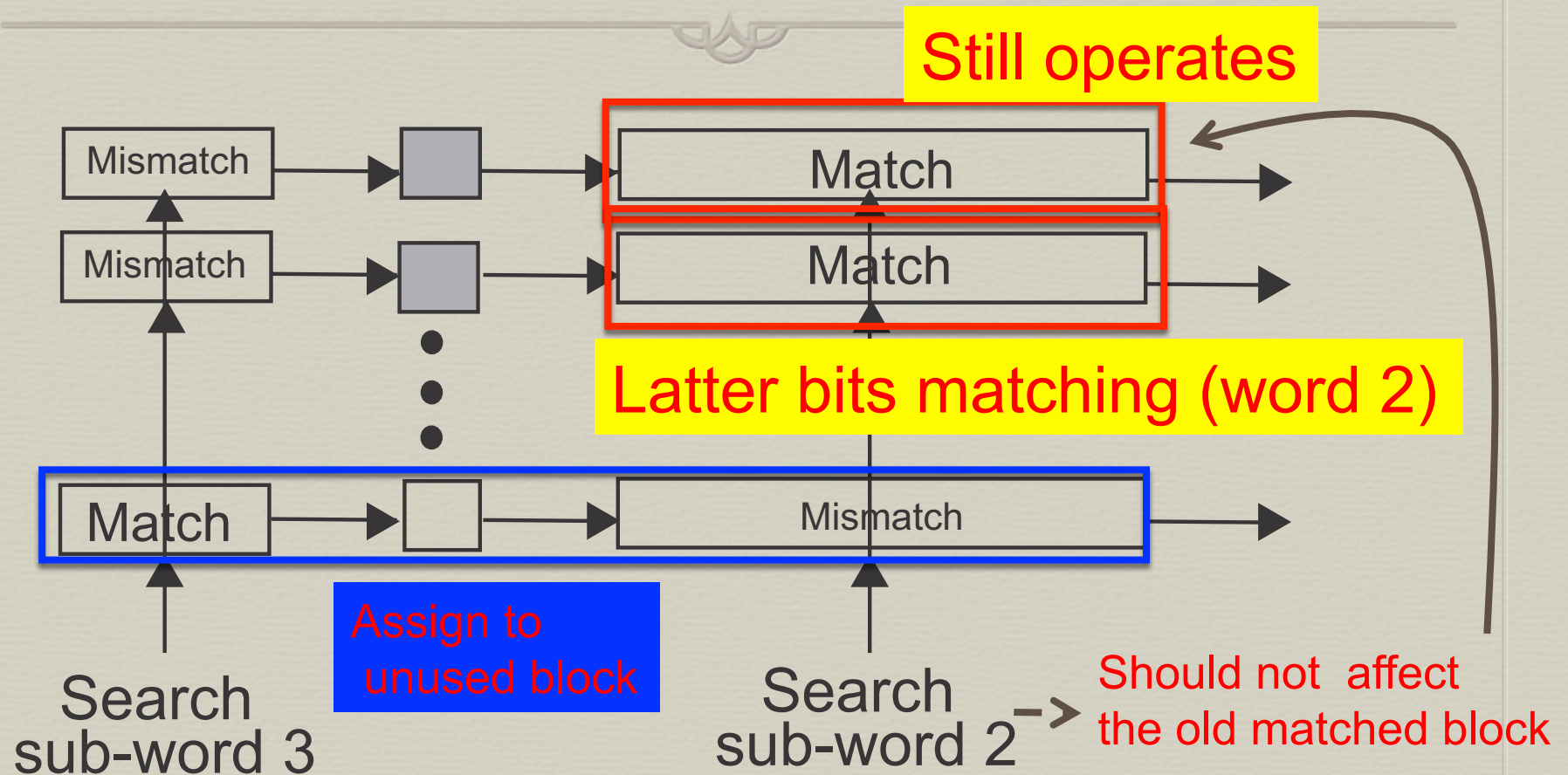## Word block partitioned by segmentation block

# WOS operation



Match

Mismatch

Input connects all word blocks

Mismatch

Search
sub-word 1

1. First k-bit sub-word search
2. Some sub-word block matched

# WOS operation

Latter bits matching

| Mismatch | → | | → | Match | → |
| Match | → | | → | Mismatch | → |
| Mismatch | → | | → | Mismatch | → |

Assign to unused block

Search sub-word 2

Search sub-word 1

After k-bit search, new search starts

# WOS operation

Still operates

Mismatch → ☐ → Match

Mismatch → ☐ → Match

Latter bits matching (word 2)

Match → ☐ → Mismatch

Assign to unused block

Search sub-word 3

Search sub-word 2 –>

Should not affect the old matched block

## How to assign search words to unused blocks?

# Categorize Word Blocks

Same stored
 k-bit data

Group A

| 00000000 |

| 00000000 |

| 00000001 |

Group B

Categorize based on the first k-bit stored word

# Pre-Computation



SL1    SL2

ctrl

Search line registers

Search line registers

Comparator

Search line registers

*enable*

*mode*

Mode controller

k bits

(n-k) bits
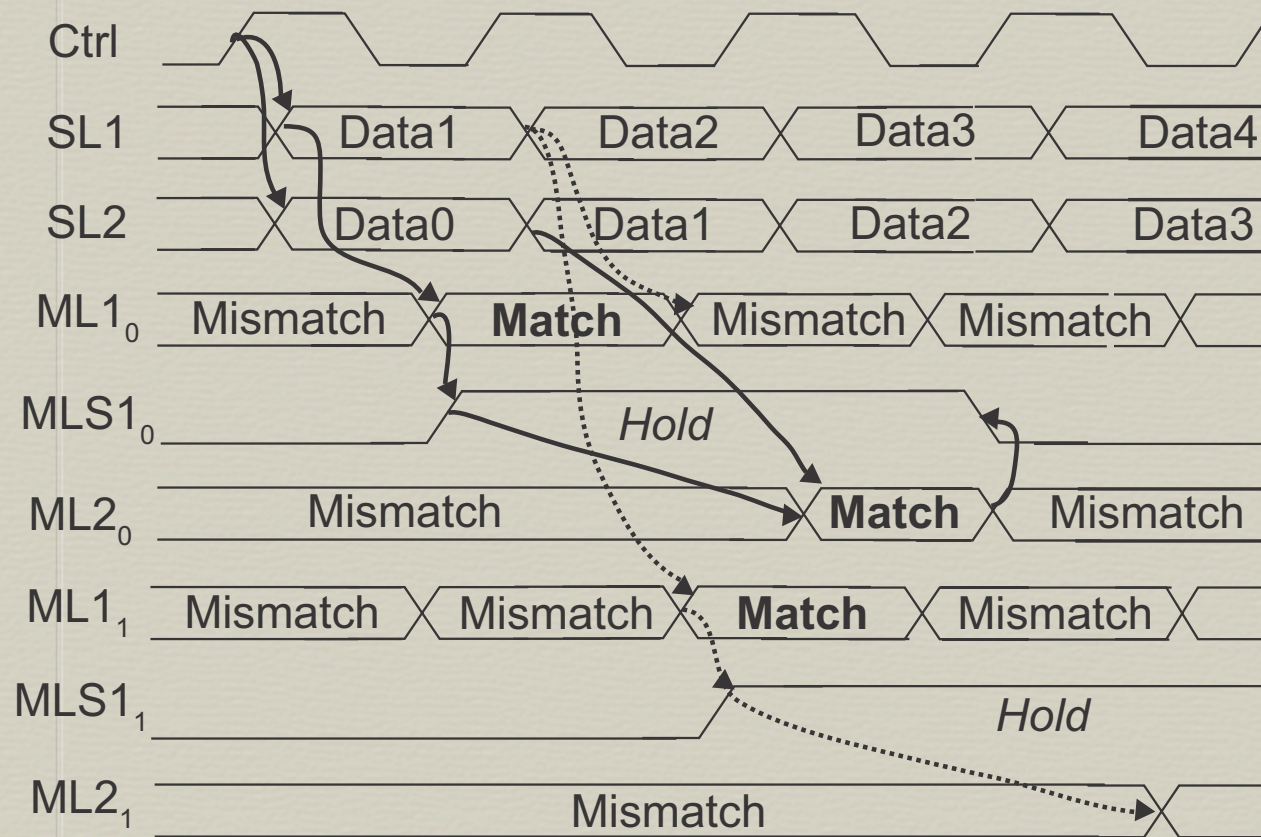
Search data

Input controller (m=1)

Compare "m" consecutive k-bit search words

If they are different, they are in different groups (Category 1: fast mode)

Otherwise, they are in the same group (Category 2: slow mode)

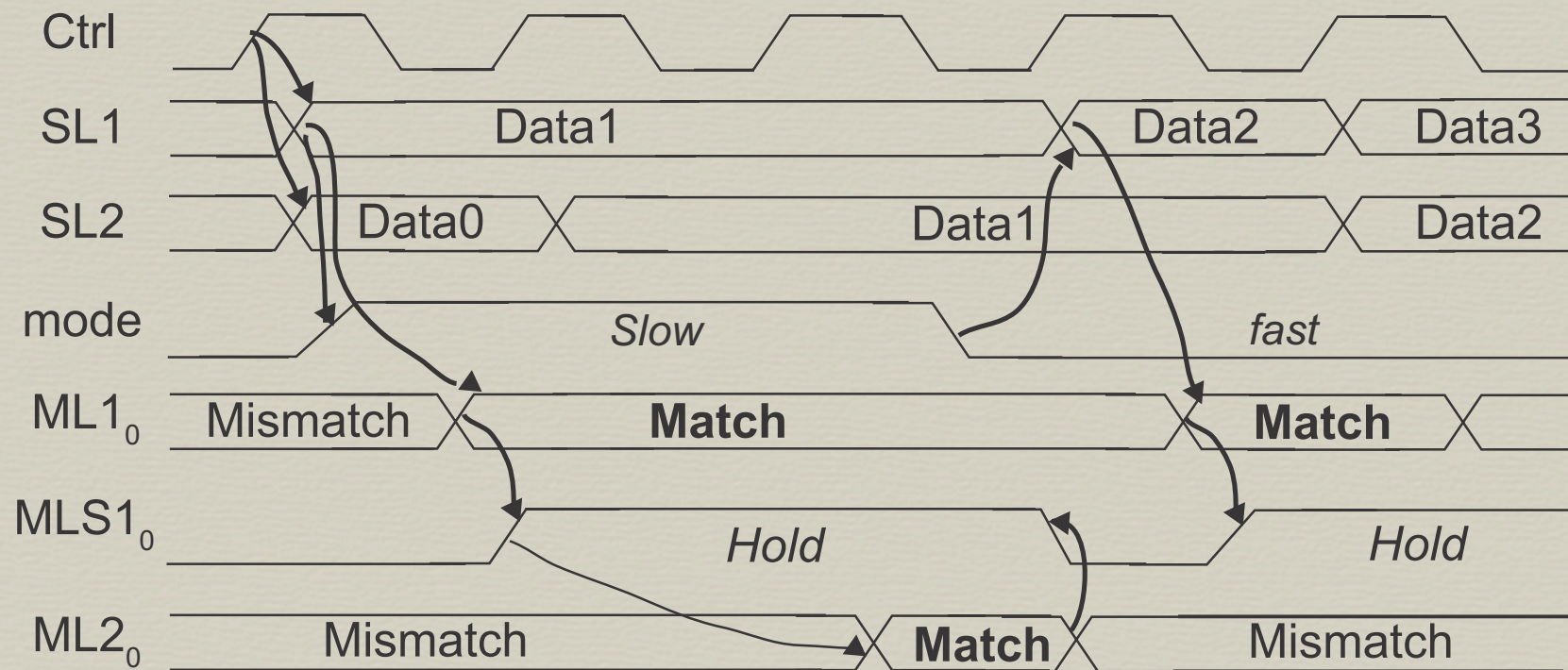Categorize search words using comparator

# Timing Diagram (fast mode)



| Ctrl | | | | |
|------|------|------|------|------|
| SL1 | Data1 | Data2 | Data3 | Data4 |
| SL2 | Data0 | Data1 | Data2 | Data3 |
| $ML1_0$ | Mismatch | **Match** | Mismatch | Mismatch |
| $MLS1_0$ | | *Hold* | | |
| $ML2_0$ | Mismatch | | **Match** | Mismatch |
| $ML1_1$ | Mismatch | Mismatch | **Match** | Mismatch |
| $MLS1_1$ | | | *Hold* | |
| $ML2_1$ | Mismatch | | | |

(a)

Send search words based on short delay $T_{tst}$

Consecutive words are assigned to unused blocks

High-speed searching based on $T_{1st}$

# Timing Diagram (slow mode)



Two consecutive words use the same word block

Wait until the current search is complete

# Average Search Delay

- Category 1 – fast mode

Send search words based on the first k-bit delay ($T_{1st}$)

- Category 2 – slow mode

Send a new word after the current n-bit search is complete ($T_{slow}$)

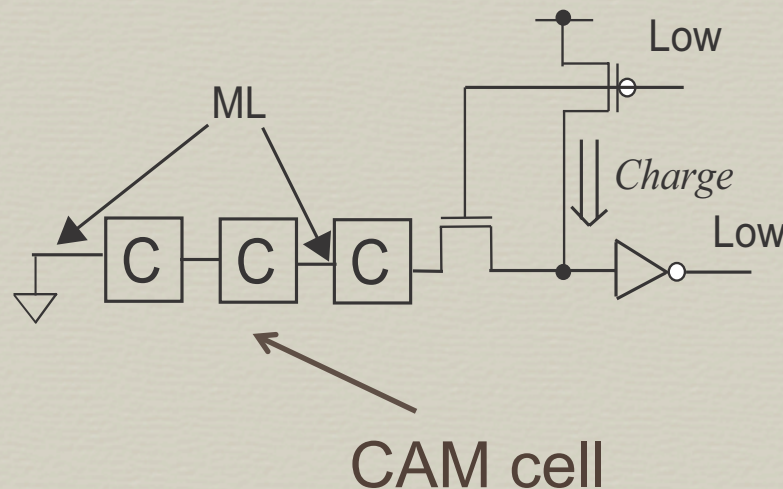$$T_{sa} = T_{1st}\left(1 - m\left(\frac{1}{2}\right)^k\right) + T_{slow}\left(m\left(\frac{1}{2}\right)^k\right)$$

# Table of Contents

- Introduction to content-addressable memory
- Overlapped search mechanism
  - Word overlapped search
  - Phase overlapped processing
- Hardware implementation
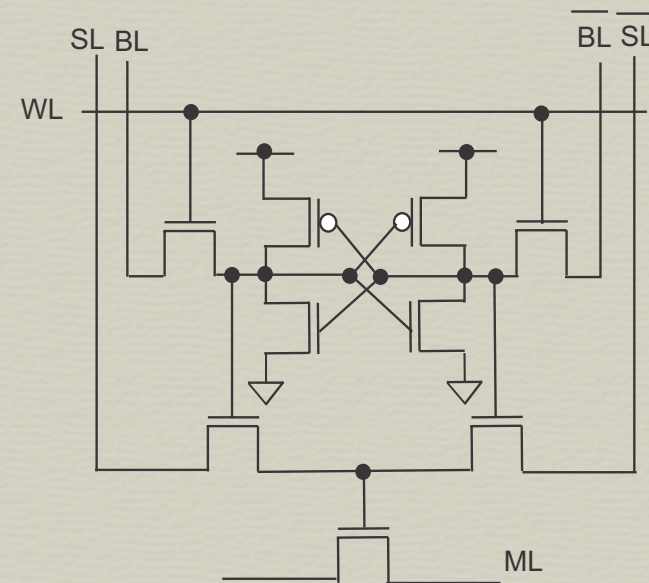- Evaluation
- Conclusion and future prospect

# Word Circuit (precharge)

## NAND-type word circuit



CAM cell

## NAND-type CAM cell



- Dynamic logic
- Series of pass transistors

- Match "ON",
- Mismatch "OFF"

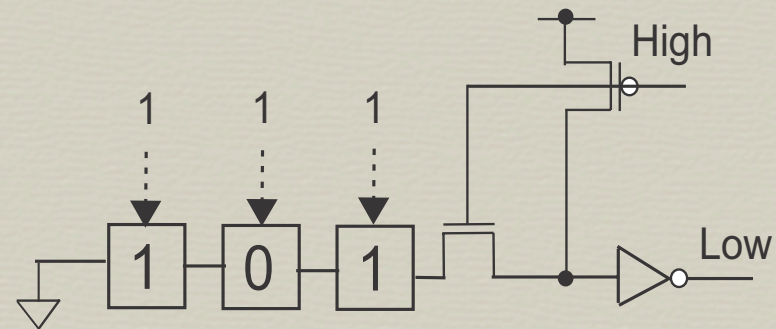**Charge capacitance on match line**

# Word Circuit (evaluate)

## Match operation

Search word

High

1    0    1

High

1    0    1

⟸ *Discharge*

Discharging capacitance
on match line
➢ Output goes high

## Mismatch operation

High

1    1    1

Low

1    0    1

Not discharging
➢ Output remains low

Match line remains high in mismatched case

12/05/07

# Synchronous Control (conventional)



*Evaluate*

Search data     Clk

1   0   1     *High*

$ML_0$ *Low*
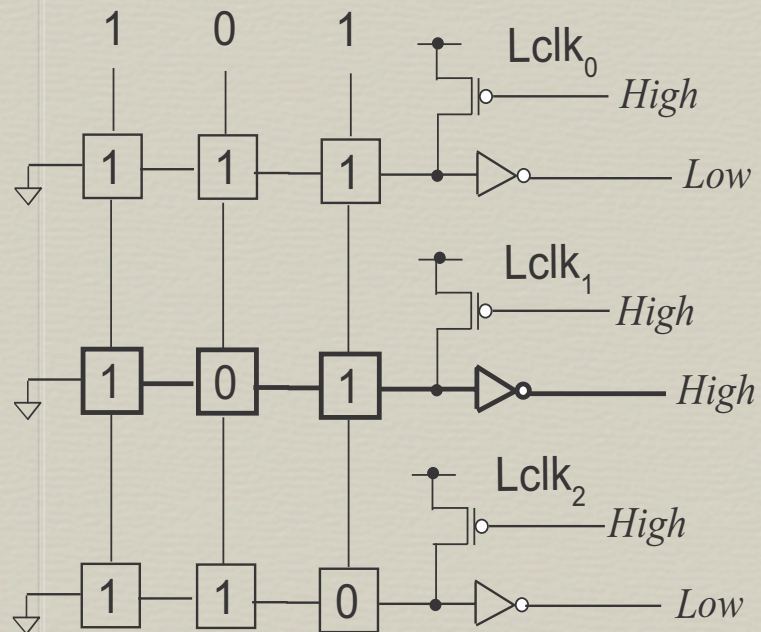
$ML_1$ *High*

$ML_2$ *Low*

*Precharge*

1   0   1     *Low*

*Low*

*Low*

*Low*

All word circuits are controlled by a global clock signal

## 2 phases are required every search

# Phase Overlapped Processing (POP)

1    0    1

| 1 | 1 | 1 |

Lclk$_0$ — *High*

*Low*

| 1 | 0 | 1 |

Lclk$_1$ — *High*

*High*

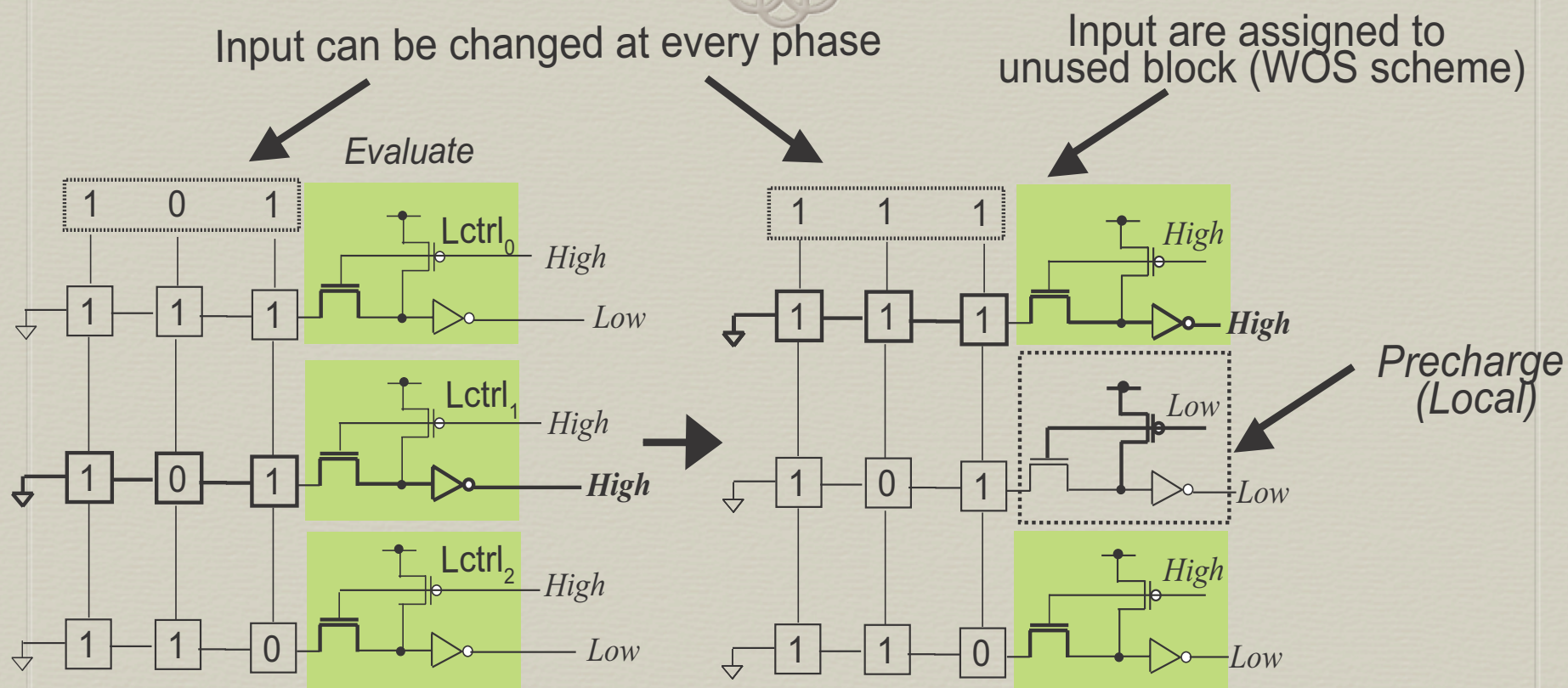| 1 | 1 | 0 |

Lclk$_2$ — *High*

*Low*

Each circuit is independently controlled using local control signals

- Matched word circuit
  Move on to precharge phase
- Mismatched word circuit
  Stay in evaluate phase

➔ Lowering switching activity of pre-charging signals

## Mismatched blocks always process new word

# WOS based POP



Input can be changed at every phase

Input are assigned to unused block (WOS scheme)

Precharge (Local)

Unused block can process without waiting precharge phase

## Searching words requires just 1 phase

12/05/07

# Throughput Ratio

Conventional $T_{CS} = 2T_{SS} = 2(T_{reg} + T_{1st} + T_{2nd})$

Proposed
$$T_{CA} = T_{SA} = T_{1st}\left(1 - m\left(\frac{1}{2}\right)^k\right) + T_{slow}\left(m\left(\frac{1}{2}\right)^k\right)$$
$$\cong T_{1st}$$

$$Throughput \ \ ratio = \frac{T_{CS}}{T_{CA}} = \frac{2(T_{reg} + T_{1st} + T_{2nd})}{T_{1st}}$$

$T_{SS}$   Synchronous search delay (evaluate phase)
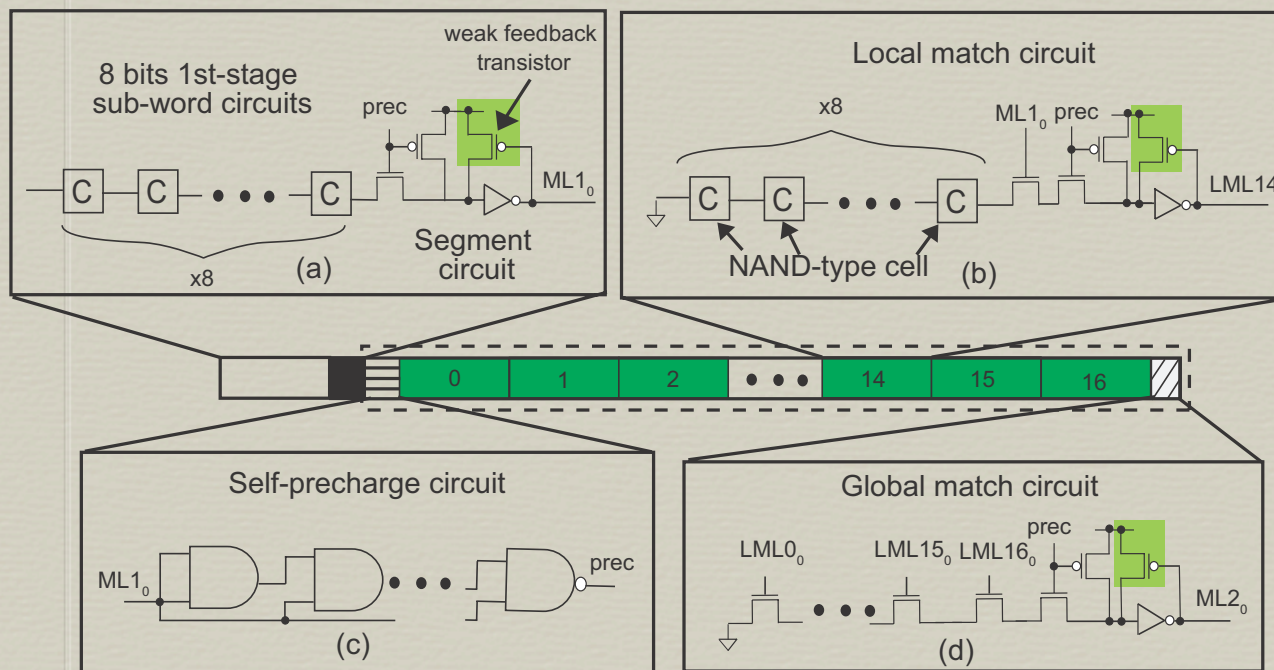
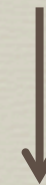$T_{SA}$   Asynchronous search delay

12/05/07

# Table of Contents

➢ Introduction to content-addressable memory
➢ Overlapped search mechanism
  ➢ Word overlapped search
  ➢ Phase overlapped processing
➢ Hardware implementation
➢ Evaluation
➢ Conclusion and future prospect

# Circuit Implementation

- 144-bit CAM word block with self-precharge circuit
- Self-precharge circuit pre-charges after 2nd stage is complete.
- Hierarchical 2nd stage block (17 local and 1 global match circuit)
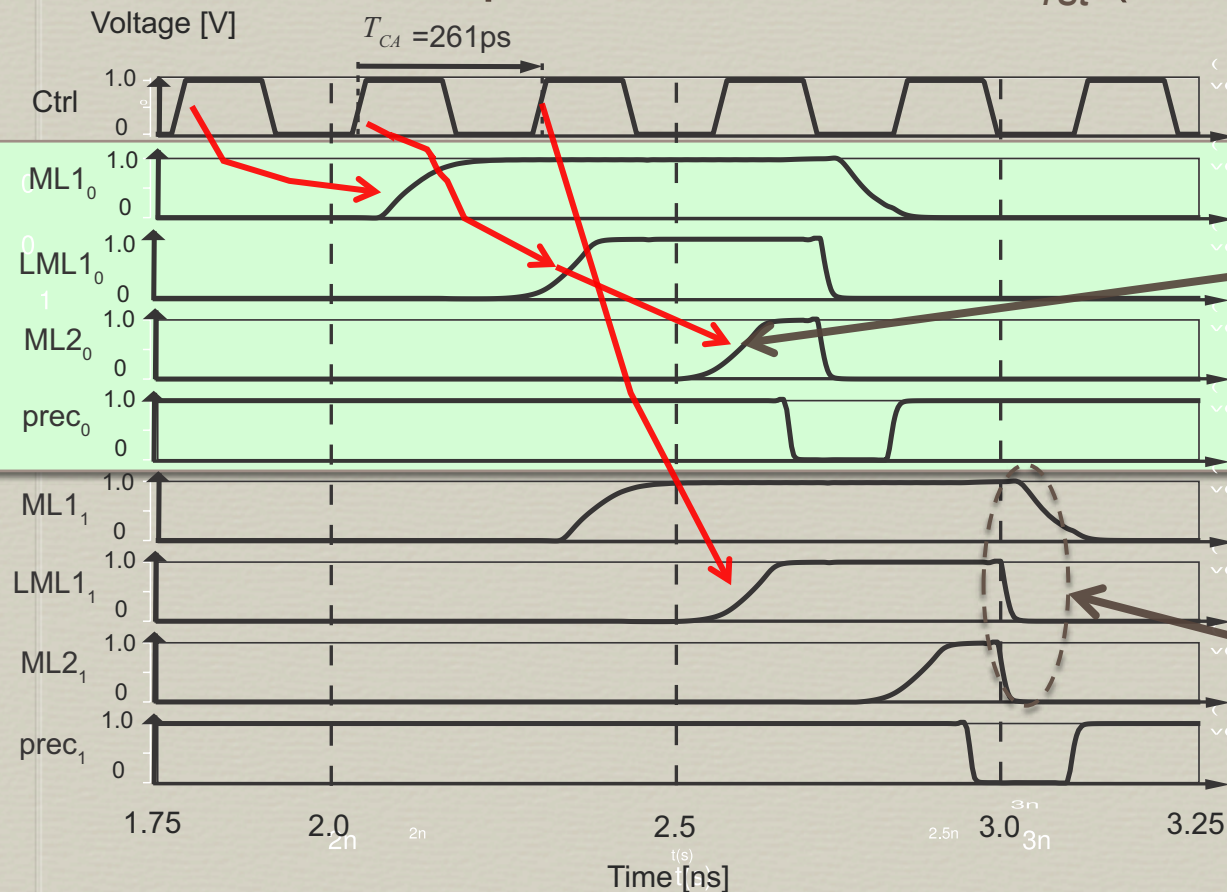


Store local matched result

It isn't affected by input changing

## Self-precharge circuit controls its word circuit

# Simulated Waveforms

CAM operates based on $T_{1st}$ (259ps)

Voltage [V]

$T_{CA}$ =261ps

Ctrl

$ML1_0$

$LML1_0$

$ML2_0$

$prec_0$

$ML1_1$

$LML1_1$

$ML2_1$

$prec_1$
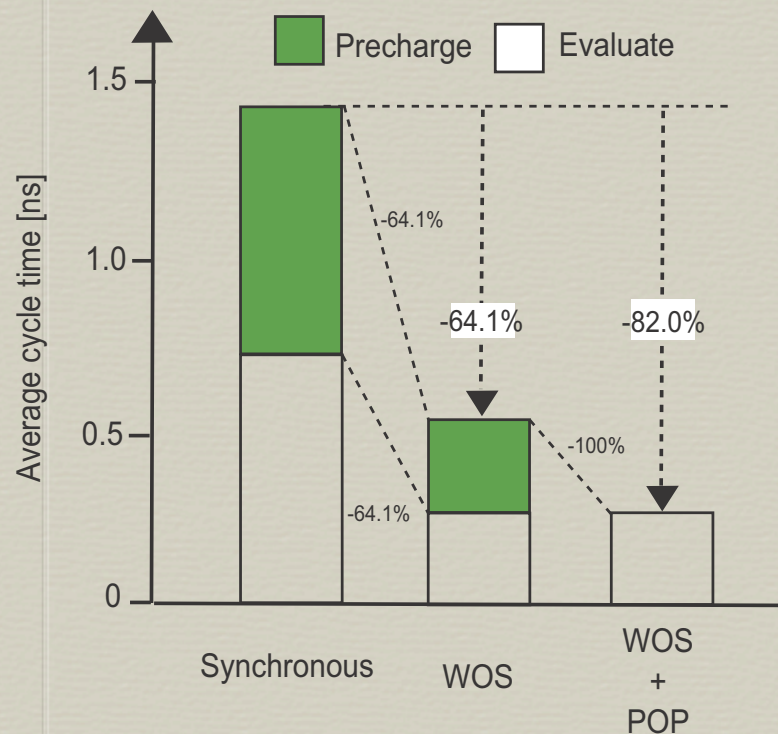
1.75   2.0   2.5   3.0   3.25

Time [ns]

Global match circuit uses only local matched result.

After search is complete, self pre-charging is locally done.

HSPICE simulation under a 90nm CMOS technology

# Performance Comparison
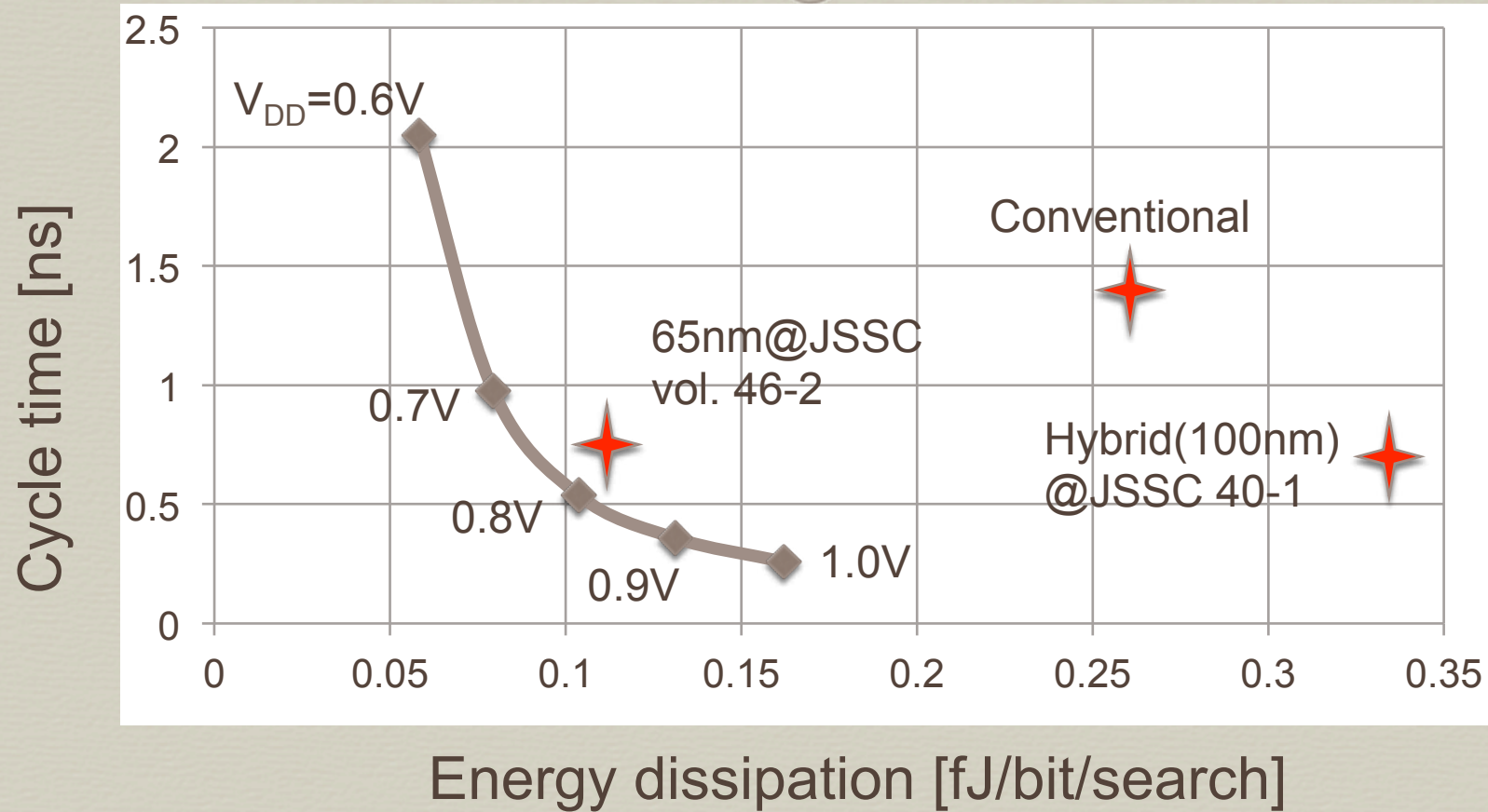
## 256-word 144-bit binary CAM@90nm CMOS



| | | Conventional | Proposed |
|---|---|---|---|
| Cycle delay [ps] | | 1454 | 261 |
| Energy [fJ/bit/ search] | Match | 0.0003 | 0.0006 |
| | Search | 0.160 | 0.160 |
| | Control | 0.103 | 0.001 |
| | Total | 0.263 | 0.162 |
| Area [Trs.] | | 372K | 408K |

Independent control reduces switching activity of pre-charging

## 5.57x throughput and 38% energy saving

# Performance Comparison



Better energy-delay product

# Conclusion

High-throughput low-energy CAM

- Word overlapped search
  - Use unused word block
  - Assign based on pre-computation
- Phase overlapped search
  - Independent control of each word block
  - Search without waiting for precharge
  - ➢ 5.57x throughput, 38% energy saving, 8% cost of area

12/05/07

# Future Prospects

- Circuit design considerations
  - Number of partitions
  - Timing robustness
- Extend to Ternary CAM (TCAM)
  - Redesign input controller
- Application specific design
  - Cache (TLB), virus checker

12/05/07