# Real-Time Executions of Program Codes in NAND Flash Memory

2009. 10. 20

Seoul National University
Chang-Gun Lee

# Increasing Market of Flash Memory

Mobile embedded devices → shock resistance
Data and Codes are dramatically increasing → high volume
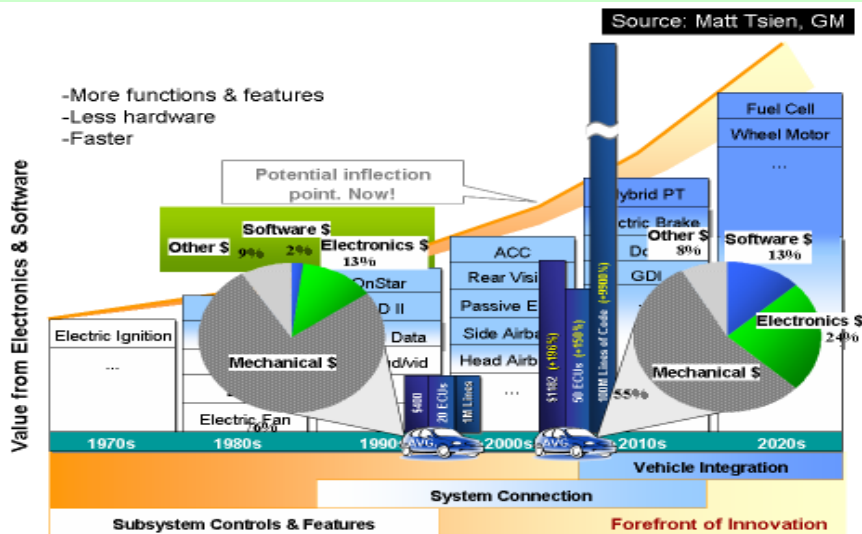Flash Memory: Non-volatile, shock resistance, high volume

Traditional
Embedded
Systems

Soft Real-time
Embedded Systems

Hard Real-time
Embedded Systems

NAND: Data Storage
NOR: Code Storage

# Why NAND is attractive for codes?
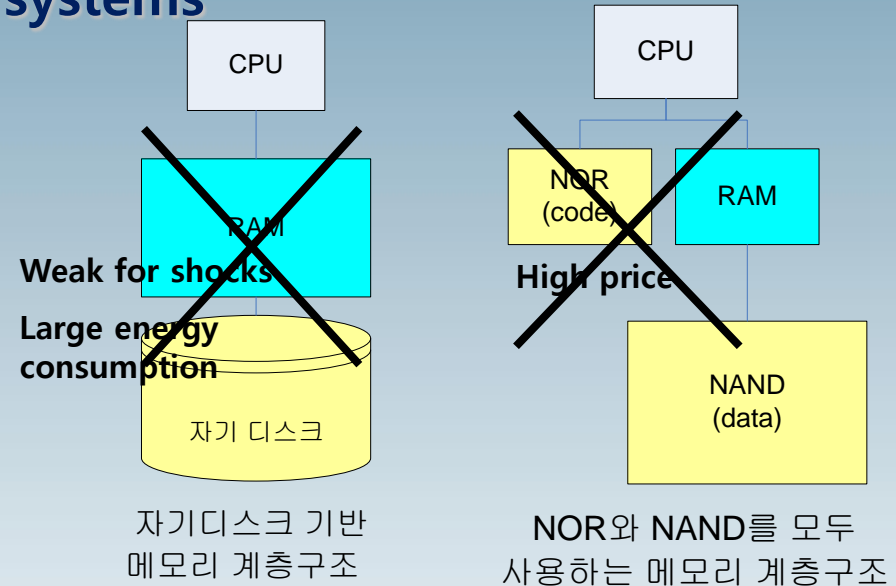
## Increasing SW Complexity → Huge program codes

### Soft Real-Time Embedded Systems (e.g., Multimedia Smart Phones)
- more than 5M source code lines in a smart phone
- cf. 5M source code lines in banking system

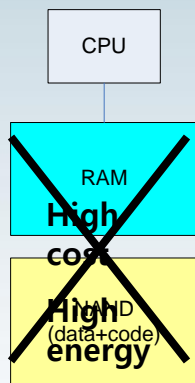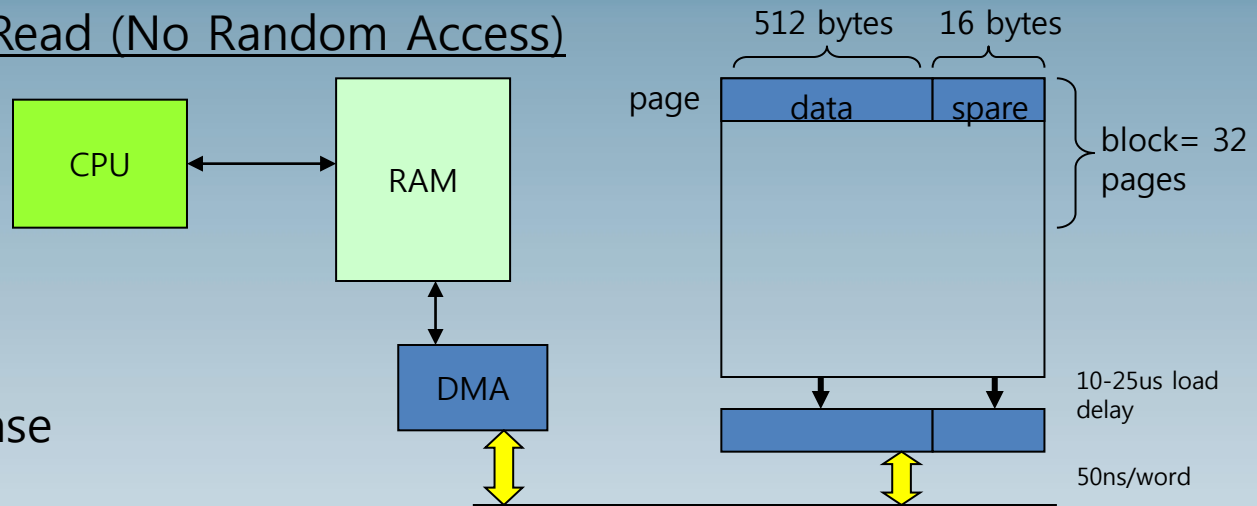### Hard Real-Time Embedded Systems (e.g., Automotive)



Source: Matt Tsien, GM

### Memory-hierarchy for embedded systems



CPU

RAM

**Weak for shocks**

**Large energy consumption**

자기 디스크

자기디스크 기반 메모리 계층구조

CPU

NOR (code)   RAM

**High price**

NAND (data)

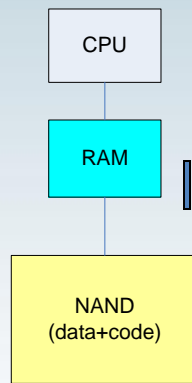NOR와 NAND를 모두 사용하는 메모리 계층구조

*NAND-based low-cost, large-size Code execution technology*

3

# What is a big challenge of NAND for codes?

- Page based sequential Read (No Random Access)
    - Read: 130us/page
- Page based write
    - 300us/page
- Block based Erase
    - 2 ms/block
- No overwrite before erase

CPU ↔ RAM

RAM ↕ DMA

512 bytes | 16 bytes

page | data | spare

block= 32 pages

10-25us load delay

50ns/word

CPU
RAM
High cost
NAND (data+code)
High energy

NAND 기반
Shadowing

CPU
RAM
NAND (data+code)

NAND 기반
가상메모리

**How to guarantee Program's real-time execution with smallest RAM?**

- RT-PLRU
  - Soft real-time
  - Single task
- mRT-PLRU
  - Extension to multiple tasks
- HRT-PLRU
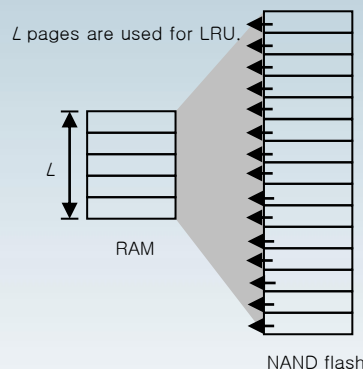  - Extension to hard real-time

# RT-PLRU: Soft real-time single task

- ## Two Important Goals
  - Developer-transparency
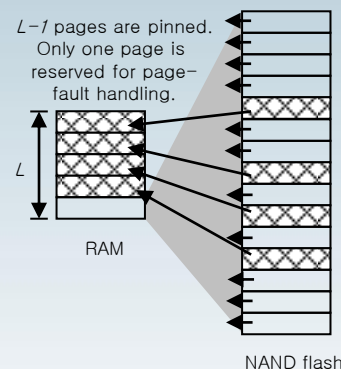  - Probabilistic guarantee of real-time with minimum DRAM

- ## Solution approach
  - Kernel-level auto-discovery of apps. temporal intension
  - Kernel-level auto-tracing of page reference sequences
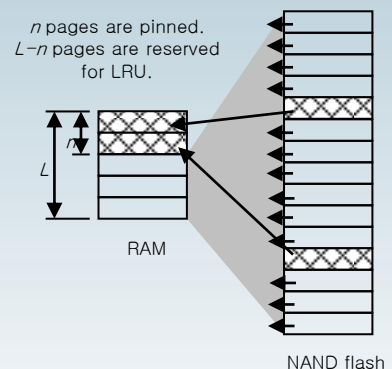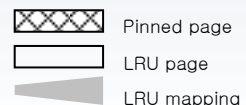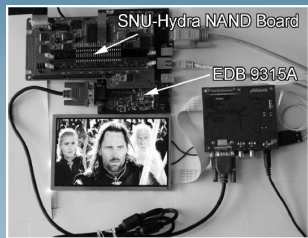  - Kernel-level auto-configuration (optimal) of pinning and LRU (**RT-PLRU**)

$L$ pages are used for LRU.

$L-1$ pages are pinned. Only one page is reserved for page-fault handling.

$n$ pages are pinned. $L-n$ pages are reserved for LRU.

RAM

RAM

RAM

$L$

$L$

$n$

$L$

NAND flash

NAND flash

NAND flash

(a) LRU only

(b) Pinning only

(c) Pinning + LRU

Pinned page

LRU page

LRU mapping
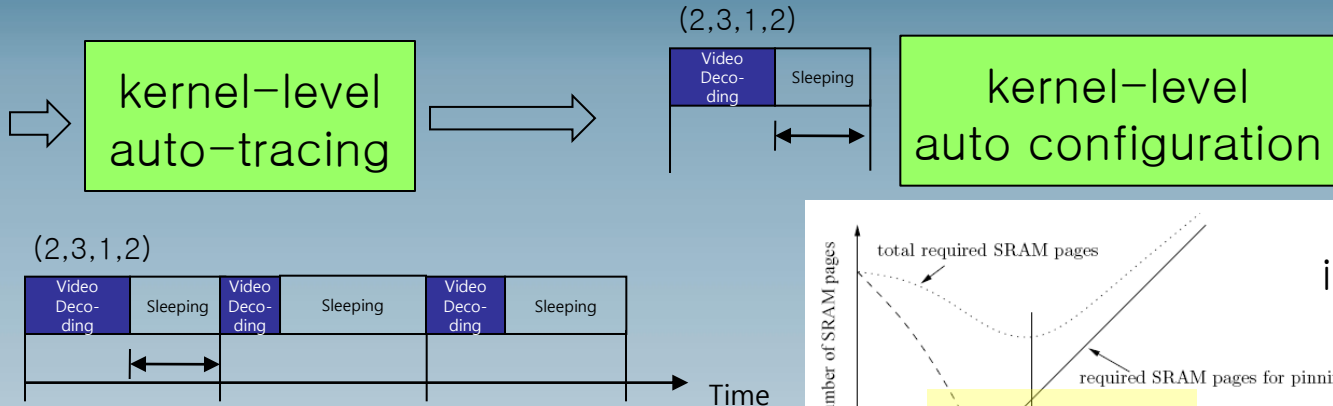
# Overall Design Flow (RT-PLRU)



prototype with
sample movie

kernel-level
auto-tracing

(2,3,1,2)

kernel-level
auto configuration

single
instance

RT-PLRU

probabilistic
extension for
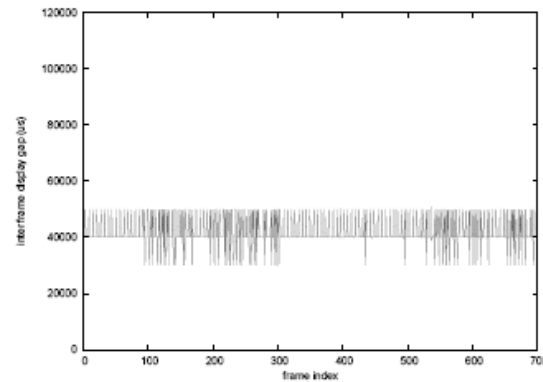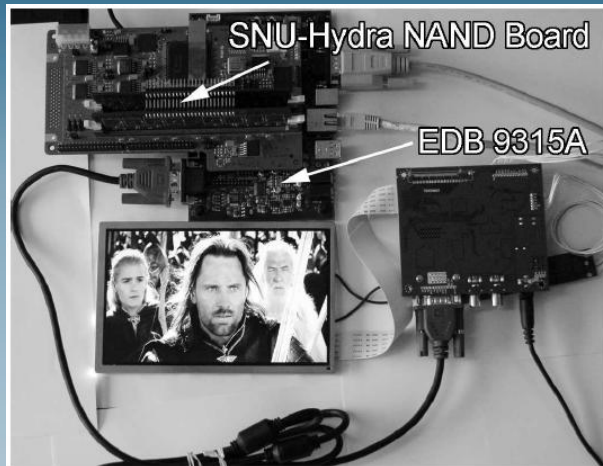multiple instances

production with
RT-PLRU

# Comparison of required DRAM sizes

# Implementations



SNU-Hydra NAND Board

EDB 9315A
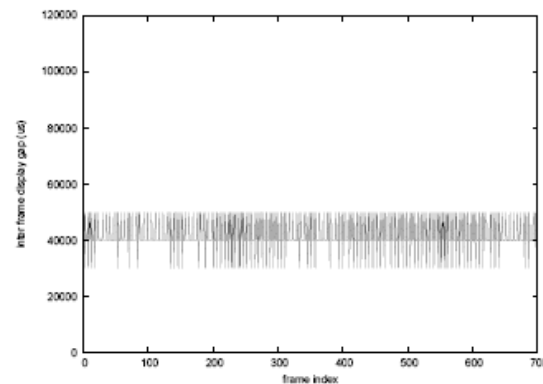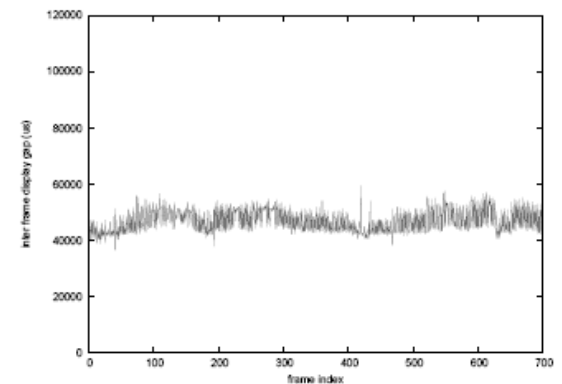


(a) shadowing (The Lord Of The Rings 1)
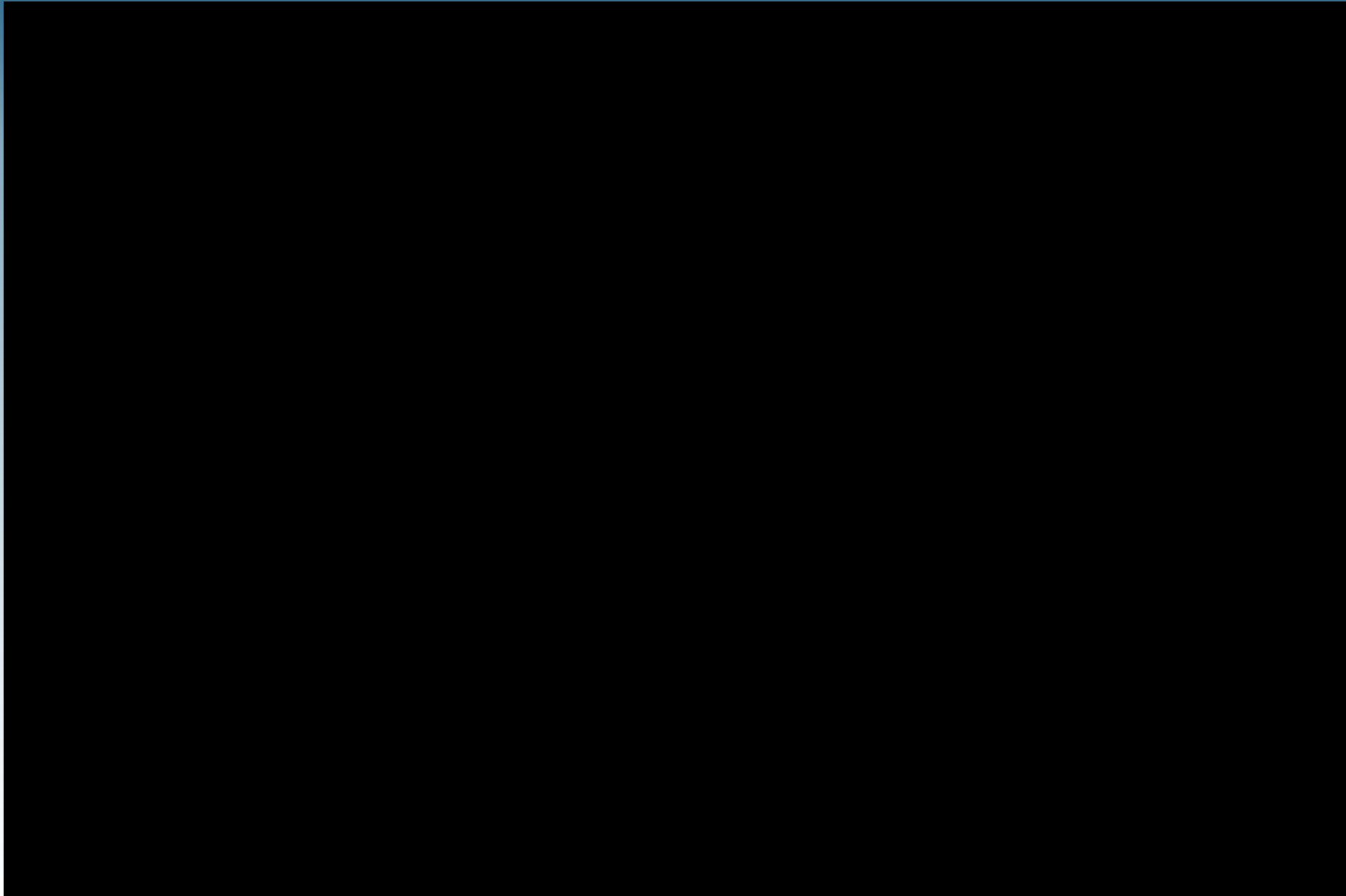
(b) RT-PLRU (The Lord Of The Rings 1)

(c) shadowing (Starwars Ep2)

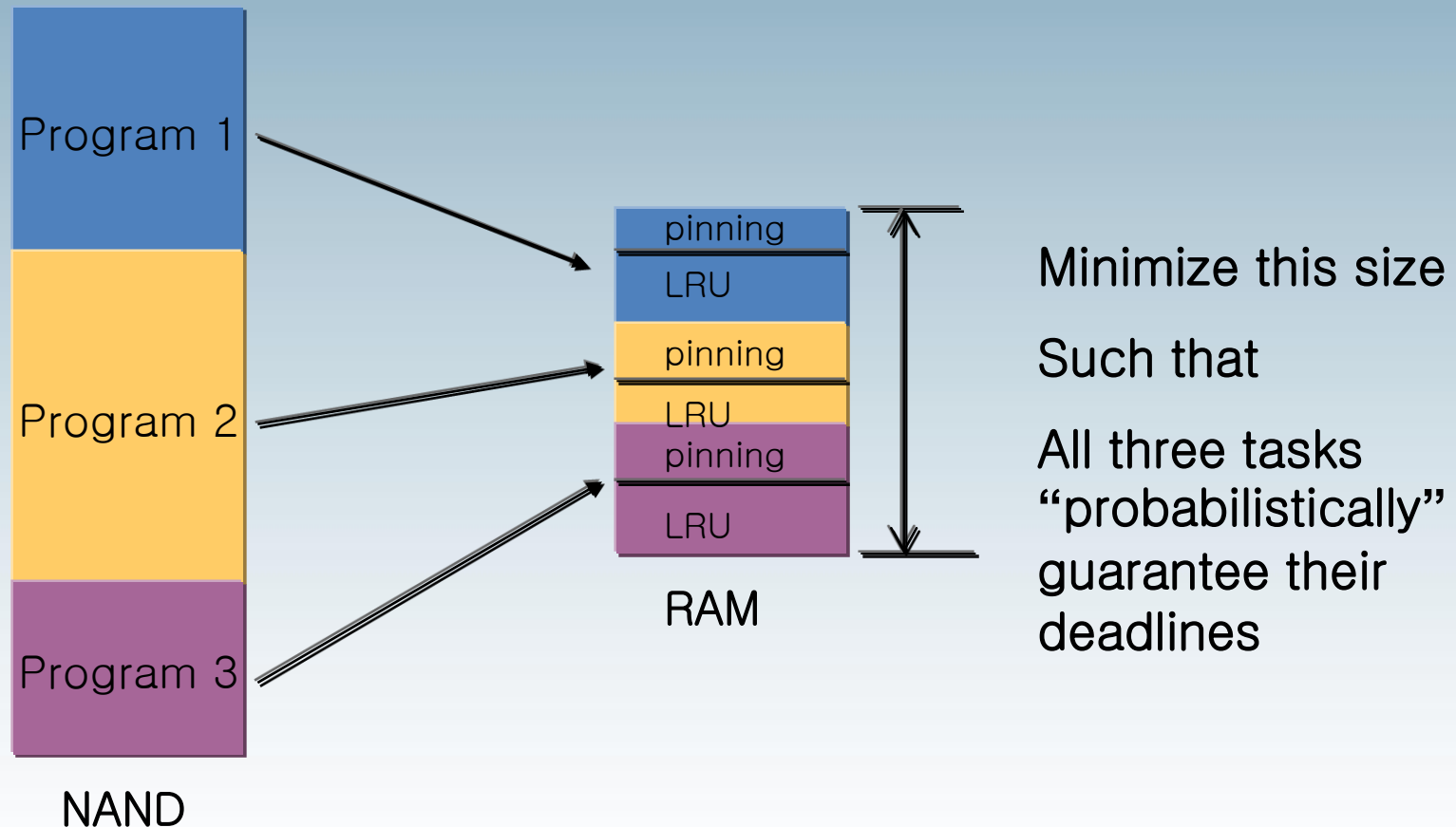(d) RT-PLRU (Starwars Ep2)

# Demo

- RT-PLRU
  - Soft real-time
  - Single task
- mRT-PLRU
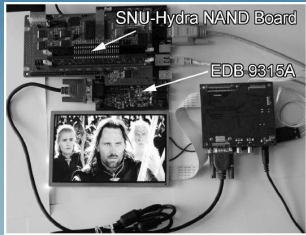  - Extension to multiple tasks
- HRT-PLRU
  - Extension to hard real-time

# mRT-PLRU:
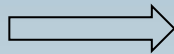# Soft real-time Multiple tasks

- Problems to answer



Program 1

Program 2

Program 3

NAND

pinning
LRU
pinning
LRU
pinning
LRU

RAM

Minimize this size

Such that

All three tasks "probabilistically" guarantee their deadlines

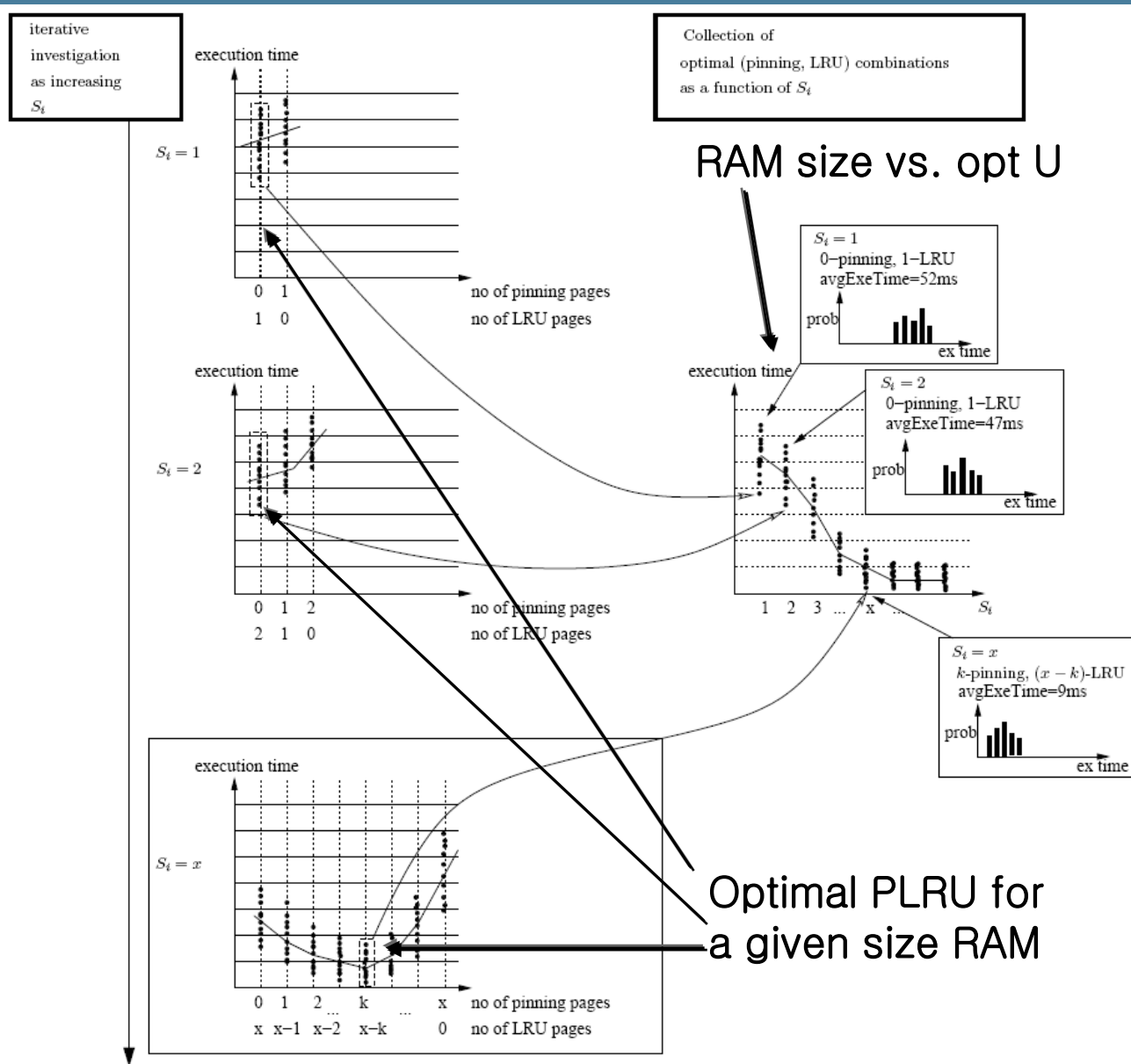# Step 1: Per-task analysis



SNU-Hydra NAND Board

EDB 9315A

prototype with
sample content

kernel-level
auto-tracing

(2,3,1,2)

| Video Deco-ding | Sleeping | Video Deco-ding | Sleeping | Video Deco-ding | Sleeping |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

iterative investigation as increasing $S_i$

Collection of optimal (pinning, LRU) combinations as a function of $S_i$

execution time

$S_i = 1$

0    1    no of pinning pages
1    0    no of LRU pages

execution time

$S_i = 2$

0    1    2    no of pinning pages
2    1    0    no of LRU pages

execution time

$S_i = x$

0    1    2 ...  k  ...  x    no of pinning pages
x  x−1  x−2    x−k    0    no of LRU pages

RAM size vs. opt U

$S_i = 1$
0−pinning, 1−LRU
avgExeTime=52ms
prob        ex time

$S_i = 2$
0−pinning, 1−LRU
avgExeTime=47ms
prob        ex time

execution time

1   2   3 ... x ...   $S_i$

$S_i = x$
$k$-pinning, $(x-k)$-LRU
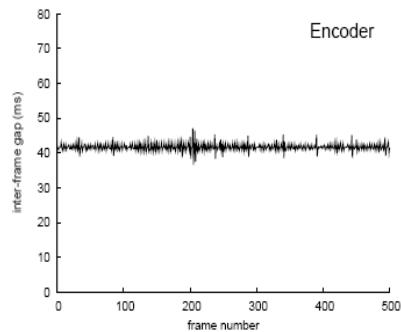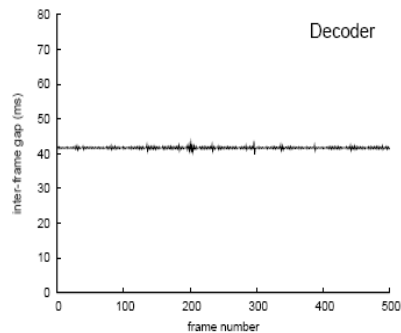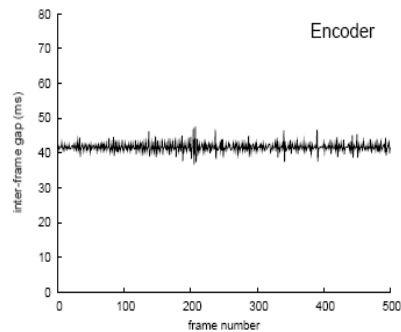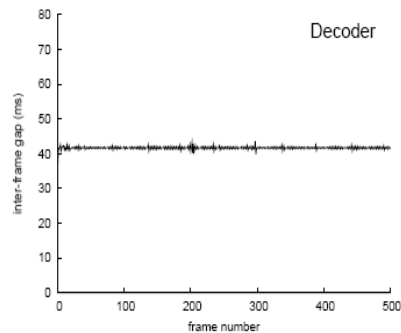avgExeTime=9ms
prob        ex time

Optimal PLRU for
a given size RAM
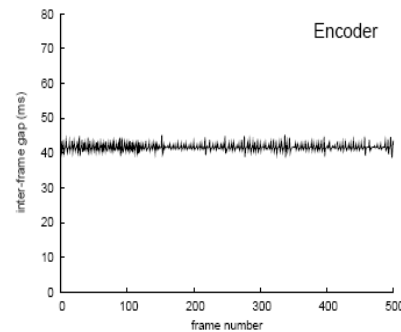
# Step 2: Convex optimization
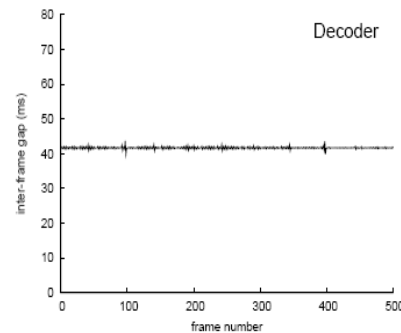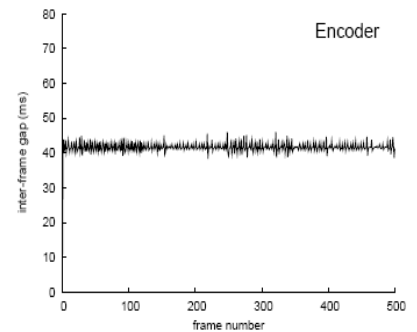
# How much RAM saved?
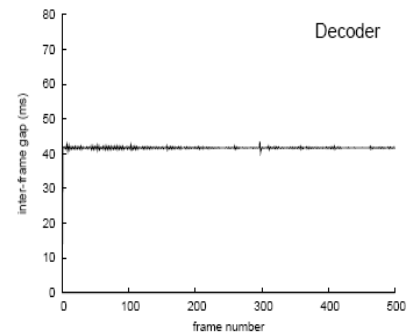
# Really work?



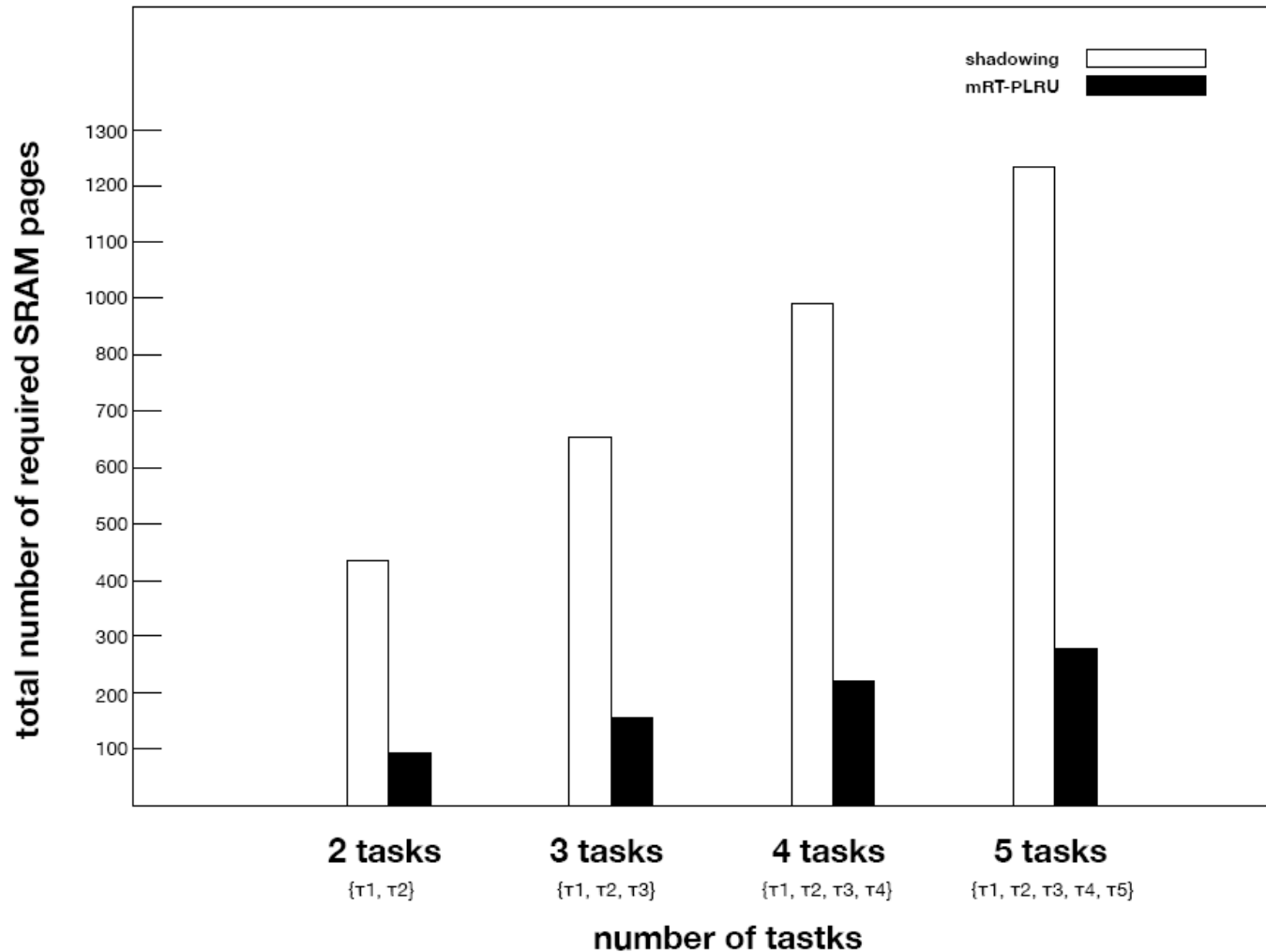(a) shadowing ("Content A")   (b) mRT-PLRU ("Content A")   (c) shadowing ("Content B")   (d) mRT-PLRU ("Content B")
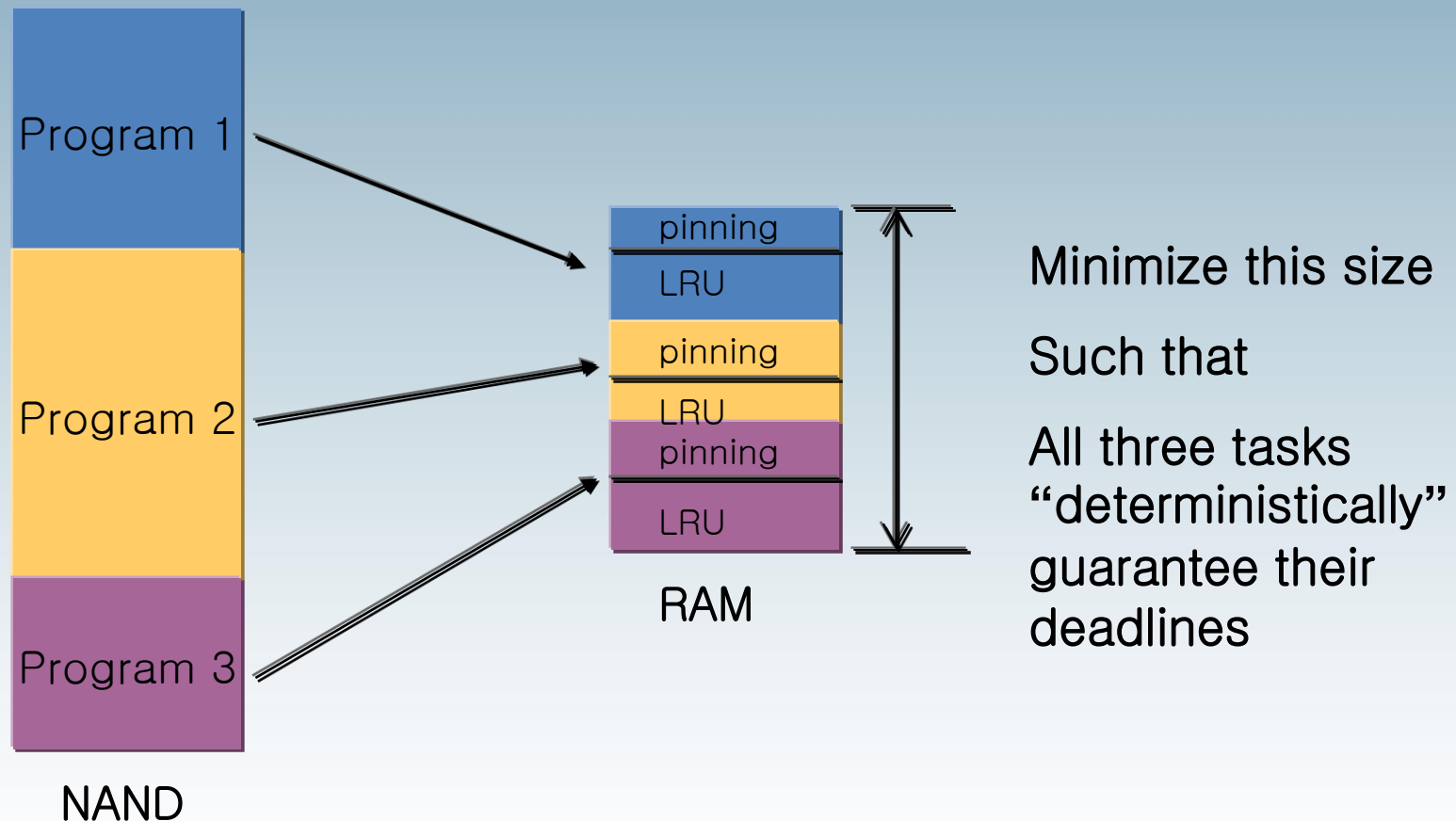
# More than two tasks?

- RT-PLRU
  - Soft real-time
  - Single task
- mRT-PLRU
  - Extension to multiple tasks
- HRT-PLRU
  - Extension to hard real-time
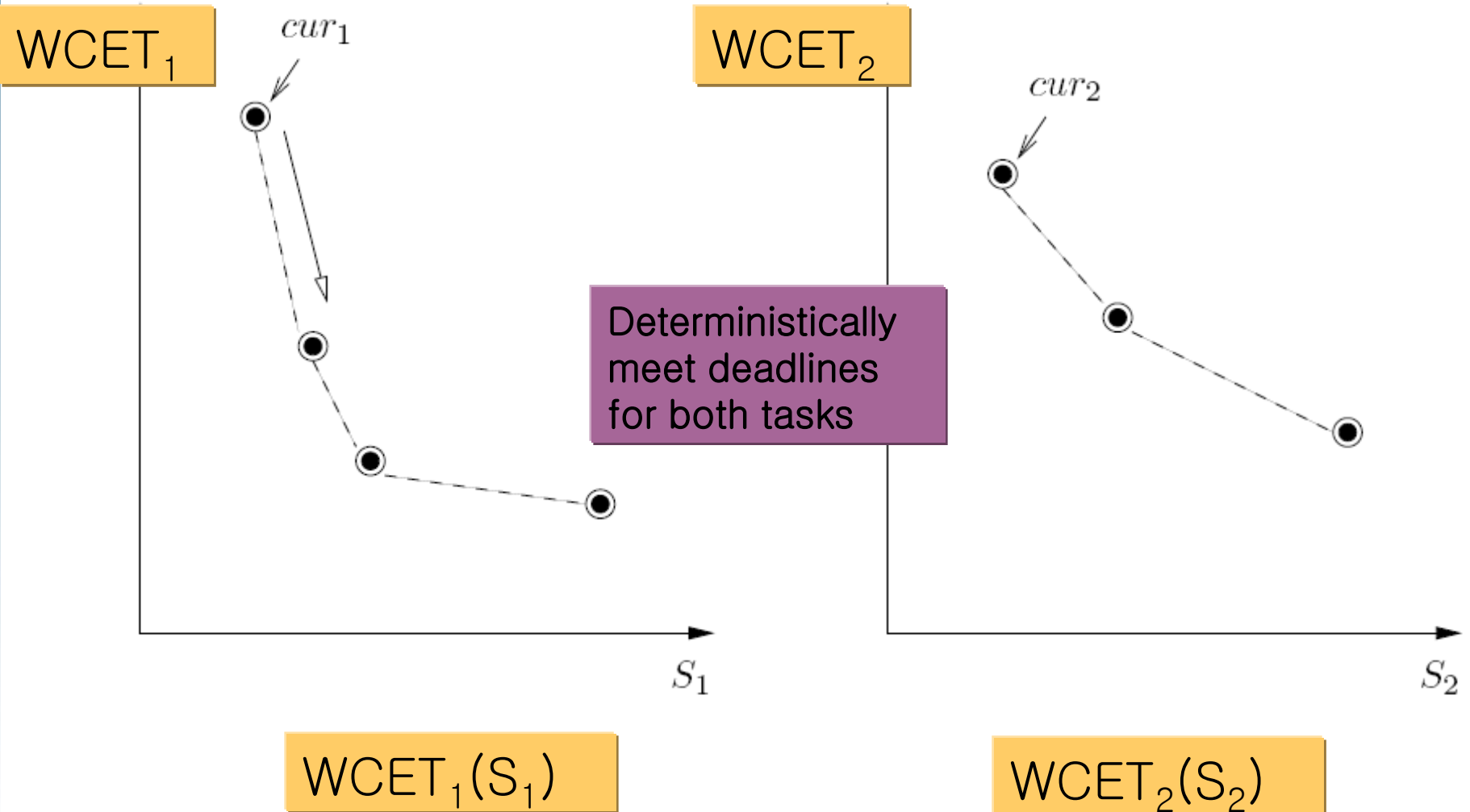
# HRT-PLRU:
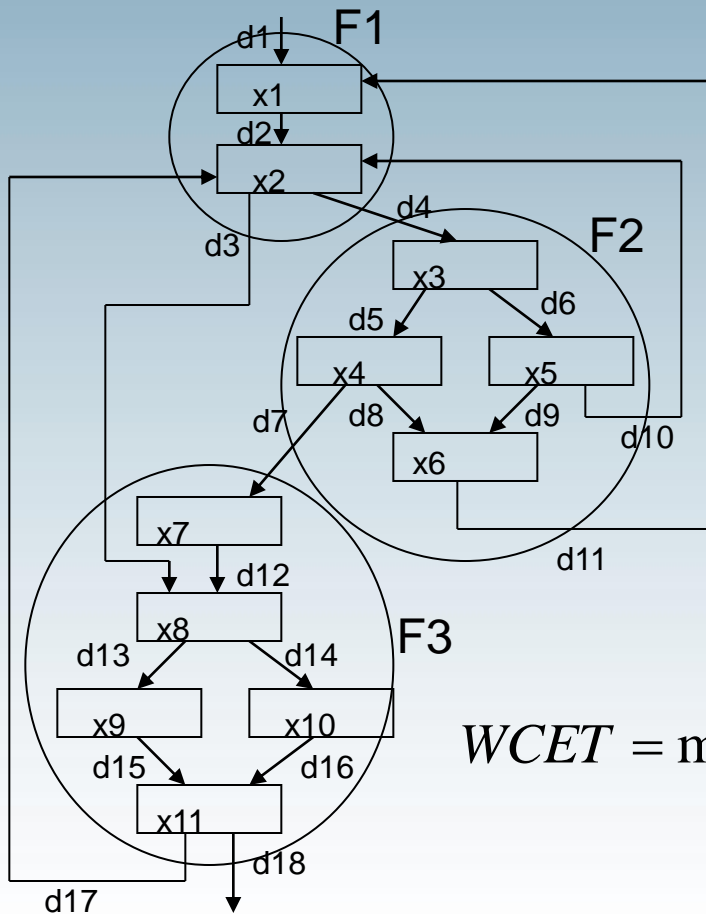# Hard real-time Multiple tasks

- Problems to answer



Program 1

Program 2

Program 3

NAND

pinning
LRU
pinning
LRU
pinning
LRU

RAM

Minimize this size

Such that

All three tasks "deterministically" guarantee their deadlines

# Per-task analysis and Convex optimization



WCET$_1$

$cur_1$

WCET$_2$

$cur_2$

Deterministically meet deadlines for both tasks

$S_1$

$S_2$

WCET$_1$(S$_1$)

WCET$_2$(S$_2$)

# Step 1: Per-task analysis
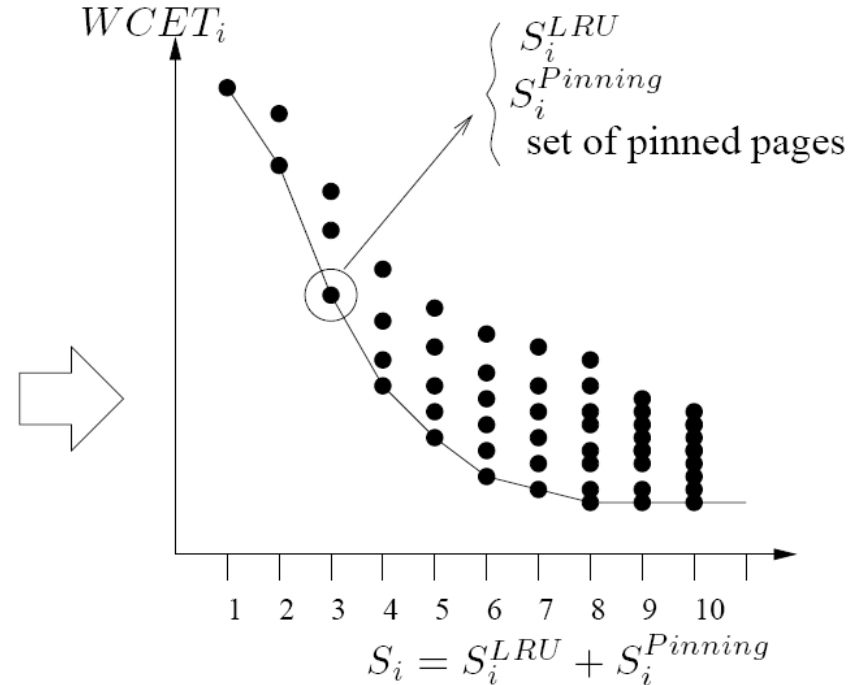
- WCET for a PLRU combination



ILP can solve this!

$$WCET = \max\left(\left(\sum_{i=1}^{11} e_i x_i\right) + \left(\sum_{j \in PageTransition} d_j^{miss}\right) \cdot PageFaultDelay\right)$$

# Step 1: Per-task analysis

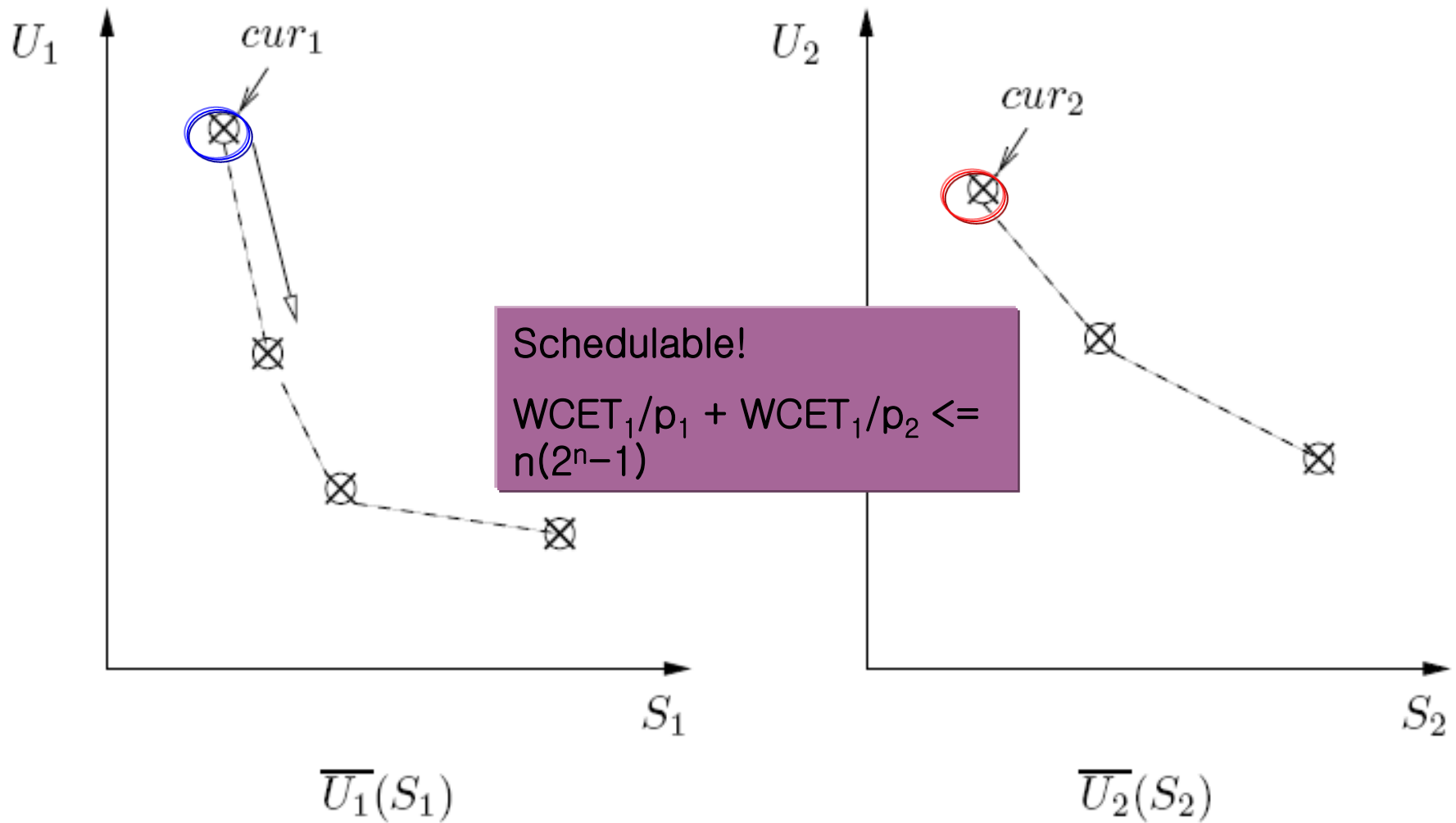- RAM size vs. opt PLRU in terms of WCET



(a) $WCET_i(S_i^{LRU}, S_i^{Pinning})$ table

(b) RAM size $S_i$ vs. WCET relation

# Step 2: Convex optimization



Schedulable!

$WCET_1/p_1 + WCET_1/p_2 <= n(2^n-1)$

# How much RAM saved?

# Conclusion

- RT-PLRU for
  - Soft real-time single task → RT-PLRU
  - Soft real-time multiple tasks → mRT-PLRU
  - Hard real-time multiple tasks → HRT-PLRU
- It provides a potential to use NAND for code executions of real-time applications
- More study needed for practical applications
  - Trade-off between RAM cost and energy consumption
  - System bus conflict problems
  - etc.