# EECS 507: Introduction to Embedded Systems Research
## Memory Hierarchies in Embedded Systems

Robert Dick

University of Michigan

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

## Announcements

W. Strunk Jr. and E. B. White, *The Elements of Style*. Macmillan
Publishing Co., Inc., 2000.

Application specific.

Real-time.

Design-time information.

# Memory hierarchies in embedded systems

Highly varied.

MMU-less.

MP-less.

Marker-based debugging.

OS controlled swap relatively rare.

Overlays sometimes used.

Special hardware memory bus controllers common.

Common to find SRAM, DRAM, EEPROM, and FLASH in same system.

Low-power sleep mode might clear RAM.

Streaming systems that avoid transition through RAM.

Scratchpad memory.

Cache locking for real-time systems.

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Problem background

RAM quantity limits application functionality

RAM price dropping but usage growing faster

Secure Internet access, email, music, and games

## How much RAM?

- Functionality
- Cost
- Power consumption
- Size

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Ideal hardware–software design process

### Ideal case

Hardware and software engineers collaborate on system-level design from start to finish

### We teach the advantages of this in our classes

It doesn't always happen

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Real hardware–software design process



HW engineers        SW engineers

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Options

## Option 1: Add more memory

Implications: Hardware redesign, miss shipping target, get fired

## Option 2: Rip out memory-hungry application features

Implications: Lose market to competitors, fail to recoup design and production costs, get fired

## Option 3: Make it seem as if memory increased

- Do not change hardware
- Do not change applications
- Do not decrease performance
- Do not increase power

Nobody knew how to do this

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Options

## Option 1: Add more memory

Implications: Hardware redesign, miss shipping target, get fired

## Option 2: Rip out memory-hungry application features

Implications: Lose market to competitors, fail to recoup design and production costs, get fired

## Option 3: Make it seem as if memory increased

- Do not change hardware
- Do not change applications
- Do not decrease performance
- Do not increase power

Nobody knew how to do this

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

## Goals

Allow application RAM requirements to overrun initial estimates even after hardware design

Reduce physical RAM, negligible performance and energy cost

Improve functionality or performance with same physical RAM

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Hardware RAM compression: Tremaine 2001, Benini 2002, Moore 2003

- Hardware (de)compression unit between cache and RAM
- Hardware redesign and application-specific compression hardware
- Past work claimed application-specific hardware essential to keep power and performance overhead low

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Conventional view: Hardware required

"Compression speed becomes the primary cost measure for assessing the quality of a compression scheme; this automatically rules out any software-based solution, and makes hardware-assisted approaches the only viable option in this context [cache or memory compression] . . ."

– A famous IEEE Fellow, TVLSI 2002

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Code compression: Lekatsas 2000, Xu 2004

- Store code compressed, decompress during execution
- Compress off-line, decompress on-line
- For RAM, less important than on-line data compression

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression for swap performance

## Compressed caching

- Douglis 1993, Russinovich 1996, Wilson 1999, Kjelso 1999
- Add compressed software cache to VM

## Swap compression

- RamDoubler, Cortez 2000, Roy 2001, Chihaia 2005
- Compress swapped-out pages and store them in software cache

## Both techniques

- Target: general purpose system with disks
- Goal: improve system performance
- Interface to backing store (disk)

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

## Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# CRAMES subproblems

Page selection

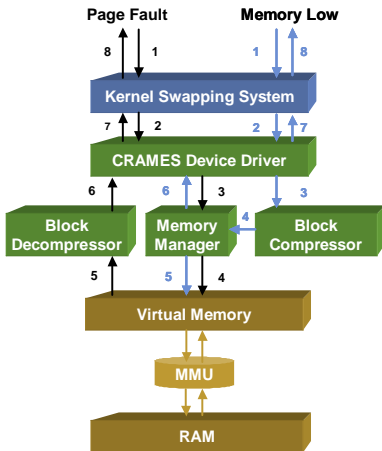Scheduling compression and decompression

Organizing compressed and uncompressed regions
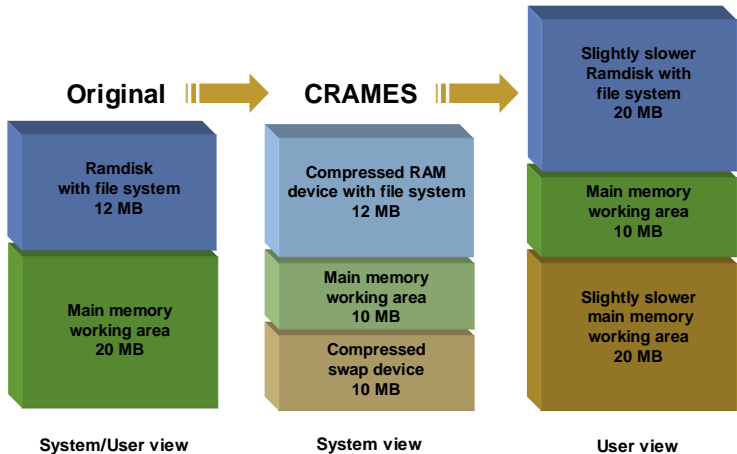
Dynamically adjust compressed region size

## Compression scheme

- High performance
- Energy efficient
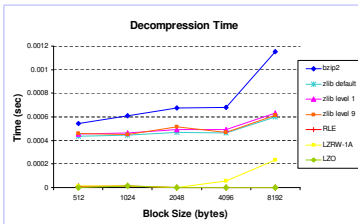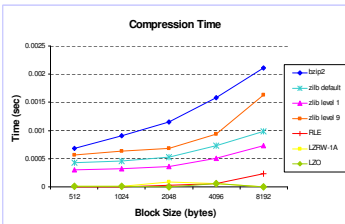- Good compression ratio
- Low memory requirement

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# System configuration



**Original** ▶  **CRAMES** ▶

Ramdisk
with file system
12 MB

Main memory
working area
20 MB

Compressed RAM
device with file system
12 MB

Main memory
working area
10 MB

Compressed
swap device
10 MB

Slightly slower
Ramdisk with
file system
20 MB

Main memory
working area
10 MB

Slightly slower
main memory
working area
20 MB

**System/User view**       **System view**       **User view**

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm



Compression Ratio

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm



**Working Memory**

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compression algorithm

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Memory management

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compressed page allocation



Allocation/Free Time (sec)

Memory Usage (TotalBytes/rm–64KB)

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Compressed page allocation



Allocation/Free Time (sec)    Memory Usage (TotalBytes/rm–64KB)

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Original problem solved, raise the goal

Applications that used to run still run fast

Applications requiring twice as much memory run

Can we go farther?

What happens if we remove half the RAM and run the same applications?

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Original problem solved, raise the goal

Applications that used to run still run fast

Applications requiring twice as much memory run

Can we go farther?

What happens if we remove half the RAM and run the same applications?

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Weak link: compression algorithm

LZO average performance penalty 9.5% when RAM reduced to 40%

Developed simulation environment to permit profiling

Compression and decompression were taking most time

Needed a better compression algorithm

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Weak link: compression algorithm

LZO average performance penalty 9.5% when RAM reduced to 40%

Developed simulation environment to permit profiling

Compression and decompression were taking most time

Needed a better compression algorithm

Introduction, motivation, and past work
CRAMES design
**Compression algorithm design**
Experimental evaluation
Commercialization
Action items

# Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Pattern-based partial dictionary match coding

Introduction, motivation, and past work
CRAMES design
**Compression algorithm design**
Experimental evaluation
Commercialization
Action items

# Data regularity

Introduction, motivation, and past work
CRAMES design
**Compression algorithm design**
Experimental evaluation
Commercialization
Action items

# Data regularity

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Pattern-based partial match

## Algorithm

- Consider each 32-bit word as an input
- Allow partial dictionary match
- Detect frequent patterns selected based on statistical analysis

## Optimizations

- Two-way associative LRU 16-entry hash-mapped dictionary
- Optimized coding scheme
- Early termination for incompressible data
- Fine-grained operation parallelization

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
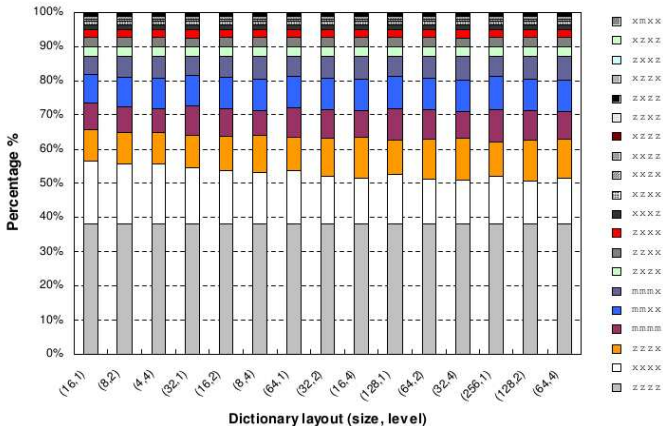Commercialization
Action items

# Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Experimental setup



## Sharp Zaurus SL-5600 PDA

- Intel XScale PXA250
- 32 MB flash memory
- 32 MB RAM
- Embedix (Linux 2.4.18 kernel)
- Qt/Qtopia PDA edition

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

## Benchmarks

PBPM twice as fast as LZO

### Benchmark applications

Results for

- ADPCM: Speech compression application from MediaBench
- JPEG: Image encoding application from MediaBench
- MPEG2: Video CODEC application from MediaBench
- Straight-forward matrix multiplication
    - Intentionally difficult for CRAMES

Also tested on

- 10 GUI applications that came with Qtopia
- Next-generation cellphone prototype

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Results

Base case: 20 MB RAM, no compression

Reduced RAM from 20 MB to 8 MB

### Without CRAMES
- All suffered significant performance penalties
- Matrix multiplication cannot execute

### With CRAMES
- LZO average case 9.5% overhead, worst case 29%
- PBPM average case 2.1% overhead, worst case 9.2%

Also works on arbitrary in-RAM filesystems

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
**Commercialization**
Action items

## Status

CRAMES used in NEC cellphones, first in N904i, June 2007

NEC was selling 10 million cellphones per year

Technology won Computerworld Horizon Award in 2007

What is the impact of this technology, essentially?

Introduction, motivation, and past work
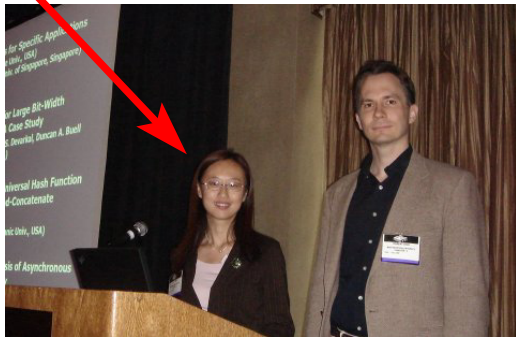CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# What did this require?

DRAM process shrink?

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# What did this require?

~~DRAM process shrink?~~
No! One good Ph.D. student. One good student can have a big impact.

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Other compression work

## CMP cache compression and data migration algorithms

- Move to many-core increases pressure on cache
- New technique that unifies cache compression and data migration
- Evaluated with SPEC OMP, NASA Parallel, SPEC2000, DIS
- 34% average throughput relative to private L2 caches

## Efficient cache compression hardware

- Most cache compression research assumes efficient hardware
- However, fast cache compression hardware design is challenging
- Developed PBPM-based cache compression hardware unit
- 58% system-wide compression ratio, 6.6 ns latency in 65 nm
- DCC'08

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

# Outline

1. Introduction, motivation, and past work

2. CRAMES design

3. Compression algorithm design

4. Experimental evaluation

5. Commercialization

6. Action items

Introduction, motivation, and past work
CRAMES design
Compression algorithm design
Experimental evaluation
Commercialization
Action items

## Assignments

12 Sep: R. Banakar, S. Steinke, B.-S. Lee, M. Balakrishnan, and P. Marwedel, "Scratchpad memory: A design alternative for cache on-chip memory in embedded systems," in *Proc. Int. Wkshp. Hardware/Software Co-Design*, May 2002, pp. 73–78.

17 Sep: C. L. Liu and J. W. Layland, "Scheduling algorithms for multiprogramming in a hard-real-time environment," *J. of the ACM*, vol. 20, no. 1, pp. 46–61, Jan. 1973.