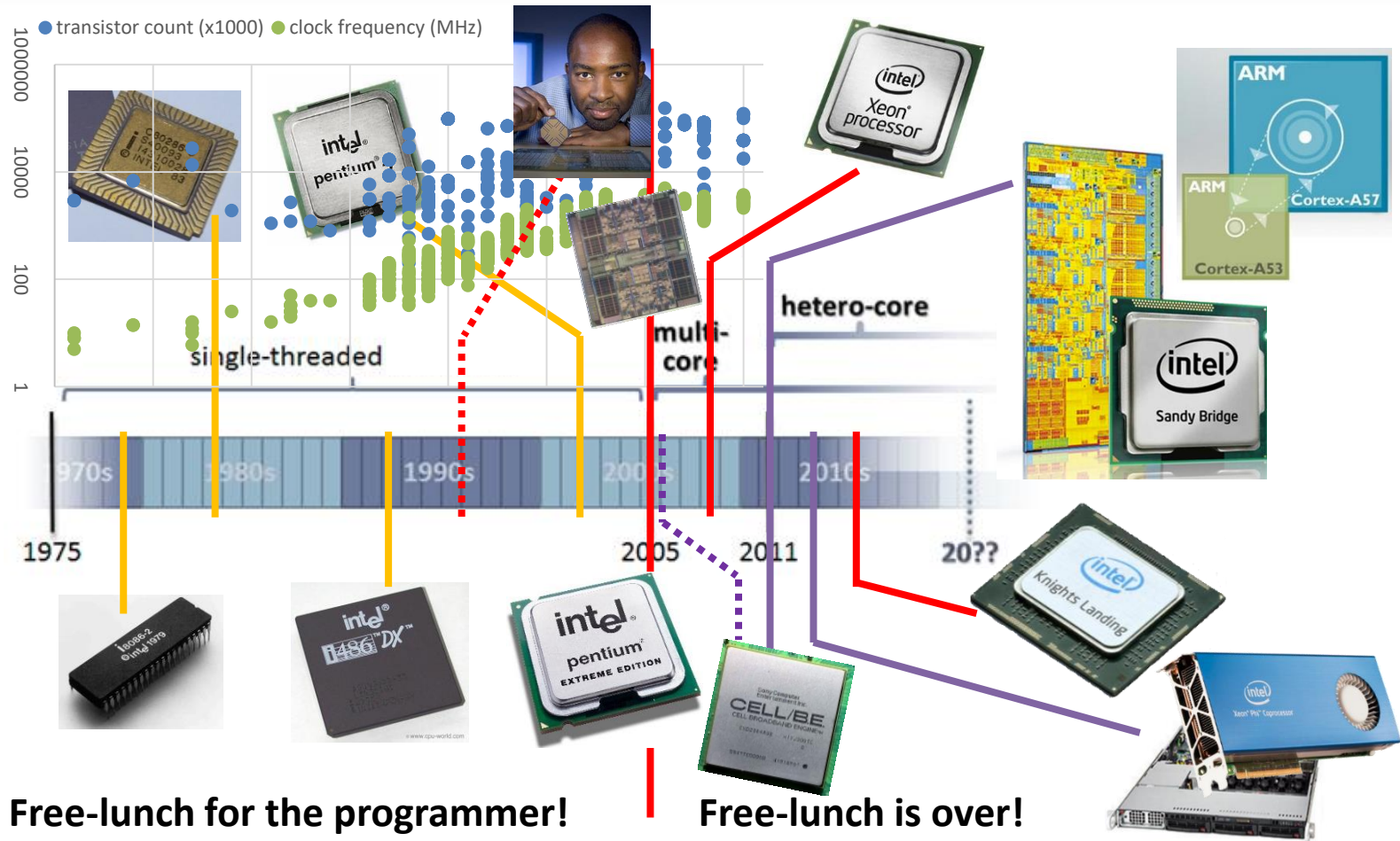# Popcorn Linux: System Software for Emerging Heterogeneous Platforms

Rob Lyerly (rlyerly@vt.edu)

**Systems Software Research Group at Virginia Tech**
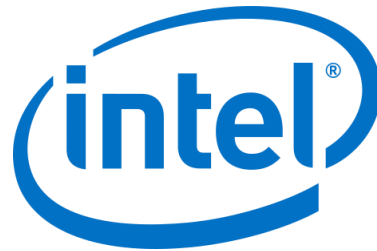ssrg.ece.vt.edu

ece

The BRADLEY DEPARTMENT of ELECTRICAL and COMPUTER ENGINEERING

**VIRGINIA TECH**

# Introduction



**Free-lunch for the programmer!**   **Free-lunch is over!**

Each image is copyright of the respective company or manufacturer.  Images used here for educational purposes.
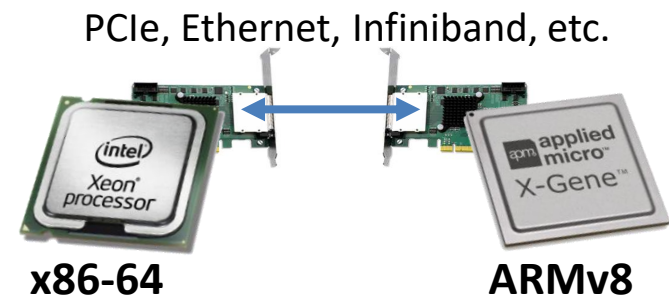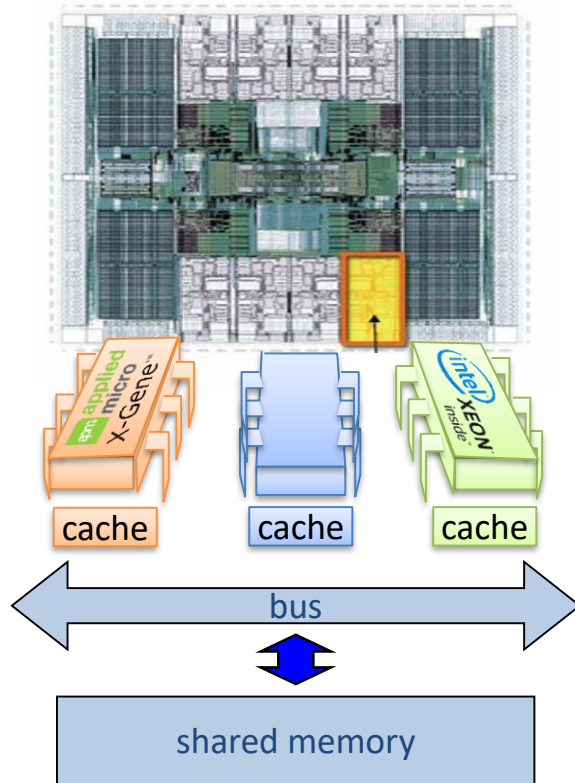
# Introduction

- Not only heterogeneous at the chip level – datacenters incorporating heterogeneous ISAs

Each image is copyright of the respective company or manufacturer.  Images used here for educational purposes.

# Introduction

- On the horizon – fully general purpose, OS-capable, heterogeneous-ISA chip multiprocessors & rack-scale systems



cache · cache · cache

bus

shared memory



PCIe, Ethernet, Infiniband, etc.



**x86-64** · **ARMv8**

# Why Heterogeneous ISAs?

- Rack-level – vastly different performance/efficiency designs

Intel Xeon E5-1650v2
- x86-64, 64-bit
- 6 cores/12 threads per socket
- 3.5GHz base, 3.9GHz turbo boost
- L1: 32KB I$, 32KB D$ per core
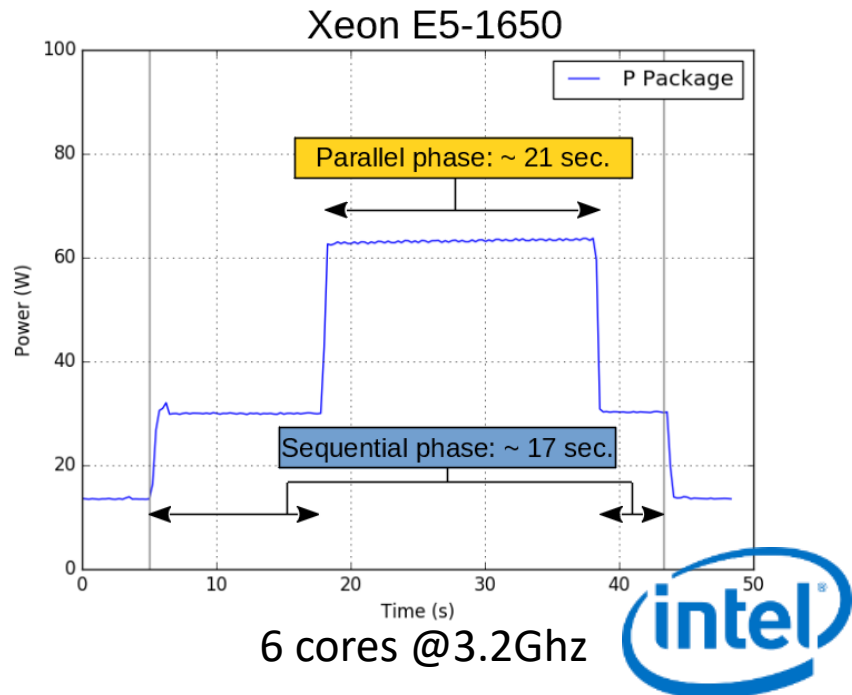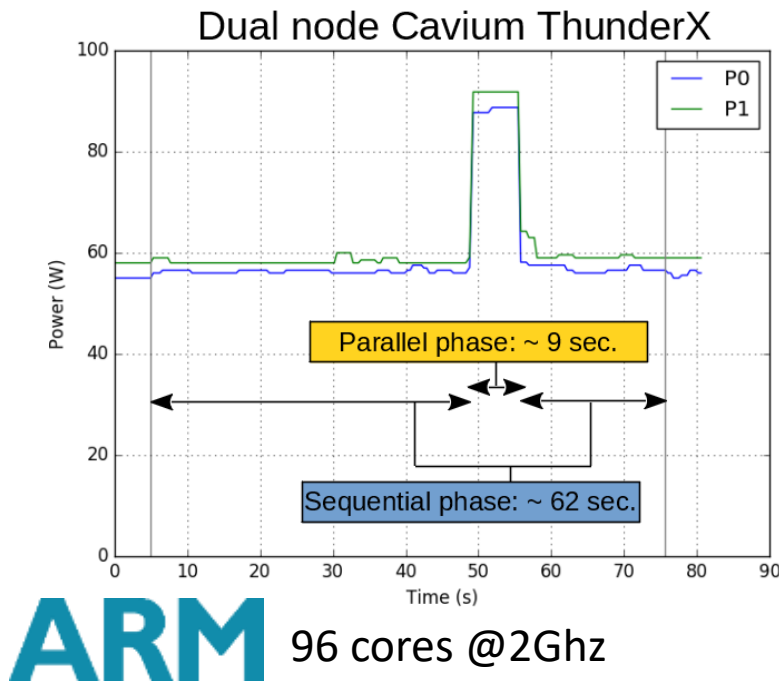- L2: 256KB per-core
- L3: 12MB (shared)
- 130W TDP

Cavium ThunderX
- ARMv8, 64-bit
- 48 cores per socket
- 2GHz per core
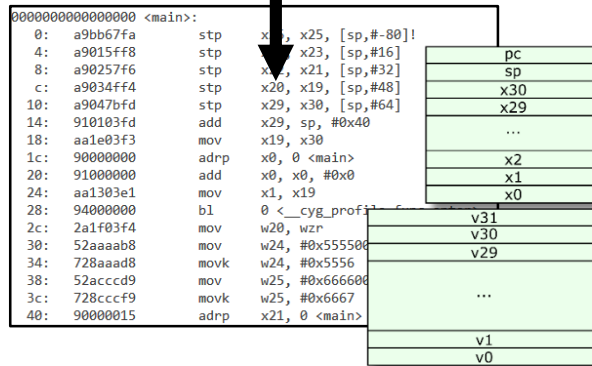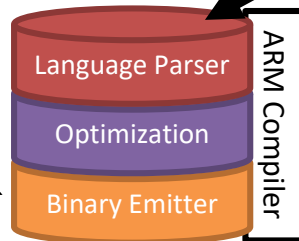- L1: 78KB I$, 32KB D$ per core
- L2: 16MB (shared)
- 120W TDP

# Why Heterogeneous ISAs?

- Smart consolidation, load balancing for performance & energy gains
- ISA/Machine affinity:
  - **Program phases** illustration: PARSEC *blackscholes* native input
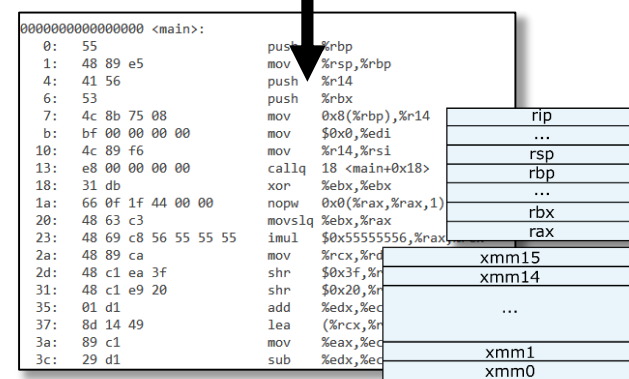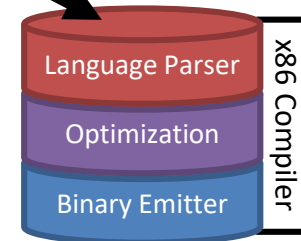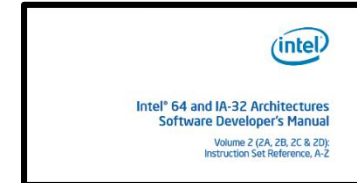


96 cores @2Ghz

6 cores @3.2Ghz

# Why Heterogeneous ISAs?

# Why Heterogeneous ISAs?

- Chip-level – micro-architectural heterogeneity is already here!

ARM big.LITTLE, e.g.,



iPhone 8/X          Galaxy S8
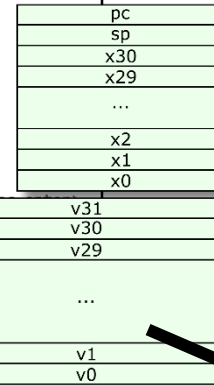
```
0000000000000000 <main>:
   0:   a9bb67fa    stp     x26, x25, [sp,#-80]!
   4:   a9015ff8    stp     x24, x23, [sp,#16]
   8:   a90257f6    stp     x22, x21, [sp,#32]
   c:   a9034ff4    stp     x20, x19, [sp,#48]
  10:   a9047bfd    stp     x29, x30, [sp,#64]
  14:   910103fd    add     x29, sp, #0x40
  18:   aa1e03f3    mov     x19, x30
  1c:   90000000    adrp    x0, 0 <main>
  20:   91000000    add     x0, x0, #0x0
  24:   aa1303e1    mov     x1, x19
  28:   94000000    bl      0 <__cyg_profile_f
  2c:   2a1f03f4    mov     w20, wzr
  30:   52aaaab8    mov     w24, #0x555500
  34:   728aaad8    movk    w24, #0x5556
  38:   52acccd9    mov     w25, #0x666600
  3c:   728cccf9    movk    w25, #0x6667
  40:   90000015    adrp    x21, 0 <main>
```

| pc  |
| --- |
| sp  |
| x30 |
| x29 |
| ... |
| x2  |
| x1  |
| x0  |
| v31 |
| v30 |
| v29 |
| ... |
| v1  |
| v0  |

Power

Big core

Little core

Compute capacity (Performance)

Little Cores
1.0-1.5 GHz

Cache (little cores)

Big Cores
1.8-2.5 GHz

Cache (big cores)

# Why Heterogeneous ISAs?

- The case for heterogeneous-ISA multicores

**Phase 1**



**Phase 2**



**Performance of bzip2's two different phases for different peak power budgets**

■ **Homogeneous**
- Alpha

■ **Single-ISA**
- big Alpha
- medium Alpha
- little Alpha

■ **Heterogeneous-ISA**
- ARM's thumb
- x86_64
- Alpha

**Single thread performance for different area budgets**



**Core Area (mm²)**

*Smaller the core area, greater # of cores can be placed on the silicon*

"Harnessing ISA Diversity: Design of a Heterogeneous-ISA Chip Multiprocessor," A. Venkat and D. M. Tullsen, ISCA 2014.

# Heterogeneous-ISA Execution

- Big questions for programming heterogeneous-ISA systems
  1. What is the programming language/model?
  2. How do I abstract away ISA differences, e.g., code, data layout?
  3. How is memory accessed/shared across distinct memory regions?

# Heterogeneous-ISA Execution

## Message Passing Interface (MPI)

+ **High performance**
− **Complex code development/refactoring**
− **Hardcoded application partitions**



## ISA Virtualization

- Managed languages, e.g., Java
  - **Rewrite application from scratch**
  - **Performance overheads**
- Dynamic binary translation, e.g., QEMU
  - + **Run unmodified binaries**
  - − **Order of magnitude slowdown**

# Heterogeneous-ISA Execution

- ## Using **shared memory** gives sanity back to the developer!

  - Well understood programming model – POSIX shared memory (30 years), OpenMP (20 years)

  - One common memory region – no data marshaling!

  - Higher programmability versus offloading!

| Benchmark | CG | EP | FT | IS | MG |
|---|---|---|---|---|---|
| OpenMP LOC | 1150 | 297 | 1106 | 1108 | 1481 |
| **MPI modified** | **98%** | **44%** | **98%** | **46%** | **97%** |

OpenMP and MPI version of NASA NPB

  - High performance – no language or system VMs necessary!

  - Flexible platform-wide resource management!

# Popcorn Linux

- System software stack for **migrating** compiled applications between heterogeneous-ISA servers
    - **Replicated-kernel OS** for thread and data migration
    - **Compiler** for creating a mostly-common virtual address space, generating metadata about ISA-specific execution state
    - **Runtime** for transforming ISA-specific execution state

- Allow developers to write **shared memory compiled applications** and leverage heterogeneity
    - Legacy code works too!

# Popcorn Linux: Operating System

- Multiple kernels provides *single system image* (SSI) allowing threads to migrate freely between nodes

# Popcorn Linux: Operating System

- ## Thread migration & heterogeneous continuations
  - Threads invoke migration via syscall
  - Kernels cooperate to migrate user-space thread contexts between ISAs
  - Kernel maps user-space PC, SP and FBP registers between ISAs

- ## On-demand page migration
  - Migrate memory pages between kernels when accessed by application
    - Intercept & redirect the page fault handler
    - Kept coherent using MSI-like protocol
  - Memory region aliasing for ISA-specific sections (e.g., `.text`)

Systems Software Research Group

VIRGINIA TECH

# Popcorn Linux: Operating System

# Popcorn Linux: Compiler

- ## Compiler toolchain builds *multi-ISA binaries*

  - Create mostly-common virtual address space (data, code, heap)
    - Pointers are valid across all ISAs
  - Dynamically transform thread execution state (stack, registers) between ISA-specific formats at migration time
  - Instrument generated code with migration points

# Popcorn Linux: Compiler

- ## Built on top of clang/LLVM
  - clang/LLVM 3.7.1, GNU gold 2.27, musl-libc 1.1.18
  - Custom address space alignment, post-processing tools
  - State transformation/migration libraries

# Popcorn Linux: Compiler

- Insert **migration points** into code
  - Can only transform stack at *equivalence points*
    - Direct mapping of execution state between ISA-specific formats
  - Scheduler cannot migrate threads at arbitrary points, must signal threads to initiate migration process

# Popcorn Linux: Runtime

- Transform registers & stack between ISA-specific formats

- Runtime transforms state before migration
  - Attaches to a thread's registers/stack
  - Reads compiler metadata describing function activation layouts
  - Rewrites stack in its entirety from source to destination ISA format

- After transformation, runtime invokes migration
  - Passes destination ISA's register state and stack to OS's thread migration service

Systems Software Research Group

VIRGINIA TECH™

# Popcorn Linux: Runtime

Source

Destination

Function: foo
Call site: 193
Call frame size: 32 bytes
Return address: 0x412820

Function: foo
Call site: 193
Call frame size: 40 bytes
Return address: 0x412700

foo() call frame

1

Function: bar
Call site: 37
Call frame size: 16 bytes
Return address: 0x410204

Function: bar
Call site: 37
Call frame size: 32 bytes
Return address: 0x410198

bar() call frame

2

Function: baz
Call site: 10
Call frame size: 32 bytes
Return address: 0x410548

Function: baz
Call site: 10
Call frame size: 48 bytes
Return address: 0x410532

baz() call frame

3

Runtime

Top of Stack

# Popcorn Linux: Runtime

# Evaluation

- ## APM X-Gene 1
  - 8 cores @ 2.4GHz
  - 8MB LLC, 32GB RAM
  - 40nm process, 50W TDP
    - Measured via on-board sensor
    - Estimated power consumption scaled to 22nm using McPAT

- ## Intel Xeon E5-1650v2
  - 6 cores @ 3.5GHz (3.9GHz turbo)
    - Hyperthreading disabled
  - 12MB LLC, 16GB RAM
  - 22nm process, 130W TDP
    - Measured via RAPL

PCIe Gen 3

Dolphin PXH810
PCIe point-to-point connection, 64Gbps

# Evaluation

- ## Load balance across architectures
  - Periodic workload – each set consists of 5 waves of up to 14 jobs
    - Uniformly sampled from NAS Parallel Benchmarks (NPB), class A, B & C
    - Waves arrive every 60-240 seconds
  - Comparison against 2 x Intel Xeon E5-1650v2 w/o migration

# Results

# Summary

- Initial working prototype on 2 node ARM64 & x86-64

- Load balancing across high-performance/energy-efficient processors enabled energy savings versus homogeneous setup

# Ongoing Research: Rack-Scale

- Centralized page management for implementation feasibility
  - Track the ownership for each page at the origin
  - Remote threads request missing pages to the origin
  - The origin handles the requests accordingly

- Model the simple RMI cache coherent protocol
  - Read-Modify Invalidate
  - Read replicates a page
  - Write gets exclusive access to the page
    - Invalidates the page from other nodes

# Ongoing Research: Rack-Scale

# Ongoing Research: Rack-Scale

- Reduce number of page migrations
  - Co-locate threads with needed data – graph partitioning problem
  - Invert mechanism – migrate threads to data
  - Change work ↔ data mapping



Node A

Node B

# New Project: Secure Popcorn



**NEWS**

**Yahoo exe[cutive] breach, inv[estigation]**

Yahoo's security properly investig[ation]

**SALTED**
By Steve Raga[n]

**NEWS**

**Raising [a]**

## The RSA hack FAQ

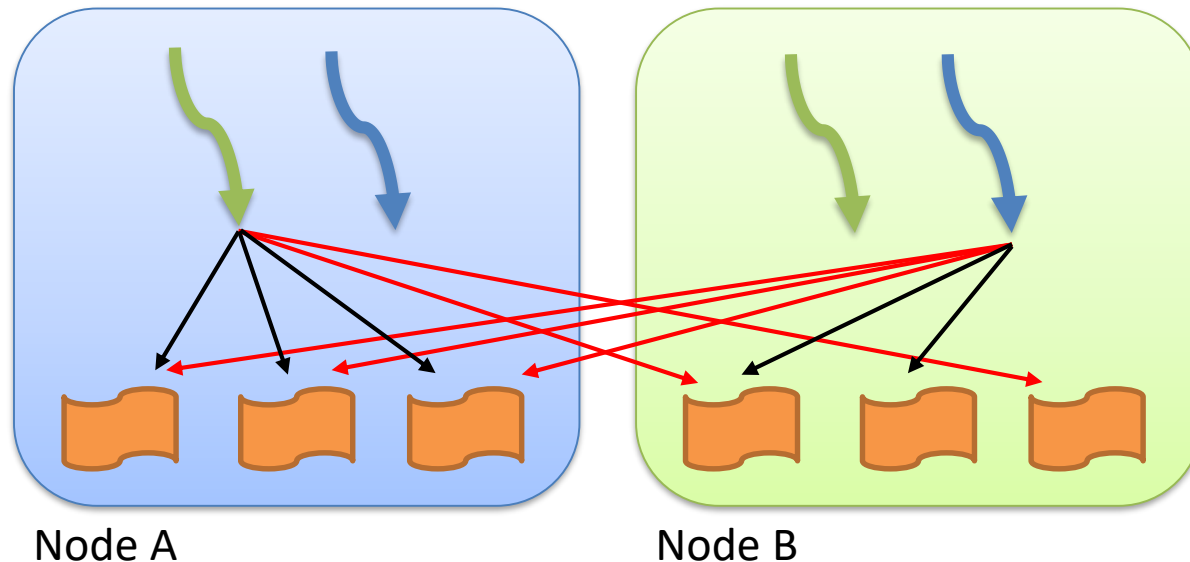In the aftermath of RSA saying that its SecureID two-factor authentication tokens may have been compromised in a data breach of the company's network, here are some key questions and answers about the situation.

**By Tim Greene**
Executive Editor, CSO | MAR 18, 2011 8:00 AM PT

**Impact: Possibly 40 million employee records stolen**

**Tw[o]**
**go[...]**

Recent reports show declining grades for government agencies' efforts to improve cybersecurity. Experts weigh in on what needs to be done.

**By Maria Korolov**
Contributing Writer, CSO | JUN 20, 2017 3:00 AM PT

**Impact: Personal information of 22 million current and former federal employees**

## 90% of today's attacks utilize ROP

Systems Software Research Group

30

Virginia Tech

# New Project: Secure Popcorn

Read-only text section

Gadget

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
....
ret
....
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
....
....
xor %eax,%eax
ret
....
....
pop %ebx
ret
```

Stack

| |
|---|
| 0x20d1b0 |
| 0x17049d |
| 0x10ad |
| 0x80436a |
| 0xbfff8a40 |
| 0x74636570 |
| |

Caller frame ends here

Systems Software Research Group

VIRGINIA TECH™

# New Project: Secure Popcorn

Read-only text section

lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
….
ret
….
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
….
….
xor %eax,%eax
ret
….
….
pop %ebx
ret

Gadget

**Exploit buffer overflow**

Stack

0x80d165
0x870d9d
0xbff076
0x804264
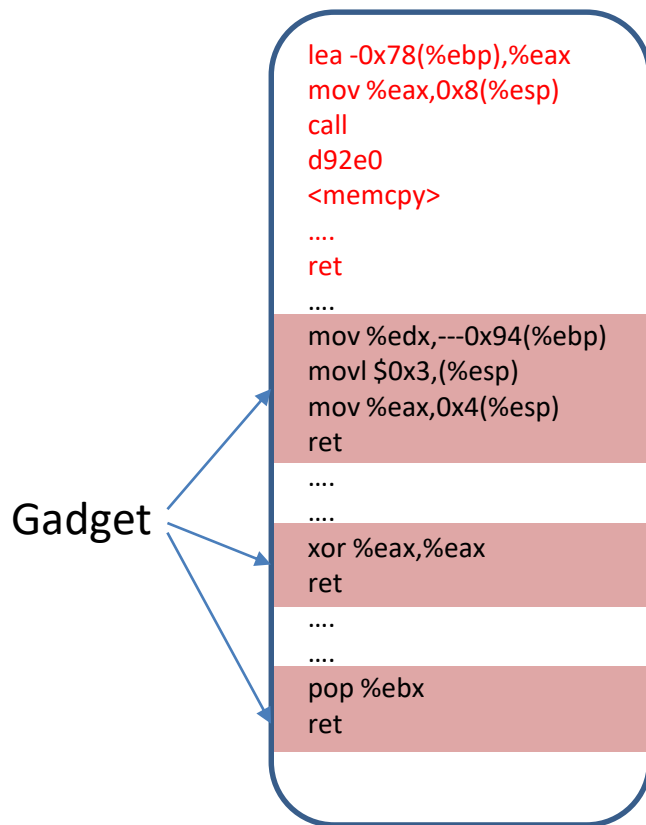0xbff840
0x74638570

Caller frame
ends here

# New Project: Secure Popcorn

Read-only text section

**Return to gadget 1**

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
....
ret
....
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
....
....
xor %eax,%eax
ret
....
....
pop %ebx
ret
```

Stack

| |
|---|
| 0x870f65 |
| 0x87098d |
| 0xbfff8076 |
| **0x870234** |
| 0x432a123 |
| 0x65708ad6 |
| |

Dynamic execution stream

pop %ebx

Systems Software Research Group

VIRGINIA TECH™

# New Project: Secure Popcorn

Read-only text section

**Return to gadget 2**

Stack

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
....
ret
....
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
....
....
xor %eax,%eax
ret
....
....
pop %ebx
ret
```

| |
|---|
| 0x870f65 |
| **0x87098d** |
| 0xbfff8076 |
| **0x870234** |
| 0x432a123 |
| 0x65708ad6 |

Dynamic execution stream

```
pop %ebx

xor %eax,%eax
```

Systems Software Research Group

VIRGINIA TECH™

# New Project: Secure Popcorn

Read-only text section

**Return to gadget 3**

Stack

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
....
ret
....
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
....
....
xor %eax,%eax
ret
....
....
pop %ebx
ret
```

| |
|---|
| |
| **0x870f65** |
| 0x87098d |
| 0xbfff8076 |
| **0x870234** |
| 0x432a123 |
| 0x65708ad6 |
| |

Dynamic execution stream

```
pop %ebx

xor %eax,%eax

mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
```

# New Project: Secure Popcorn

Read-only text section

**Return to gadget 3**

Stack

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
….
ret
….
mov %edx,---0x94(%ebp)
movl $0x3,(%esp)
mov %eax,0x4(%esp)
ret
….
….
xor %eax,%eax
```

| |
|---|
| **0x870f65** |
| 0x87098d |
| 0xbfff8076 |
| **0x870234** |
| 0x432a123 |
| 0x65708ad6 |

Dynamic execution stream

pop %ebx

xor %eax,%eax

**ROP is Turing-complete given sufficiently large binary**

Systems Software Research Group

VIRGINIA TECH™

# New Project: Secure Popcorn

**black hat USA 2015**

**REGISTER NOW**

AUGUST 1 – 6, 2015
MANDALAY BAY | LAS VEGAS, NV

## Exploiting the DRAM rowhammer bug to gain kernel privileges

Row hammer attack

Mark Seaborn and Thomas Dullien

# New Project: Secure Popcorn
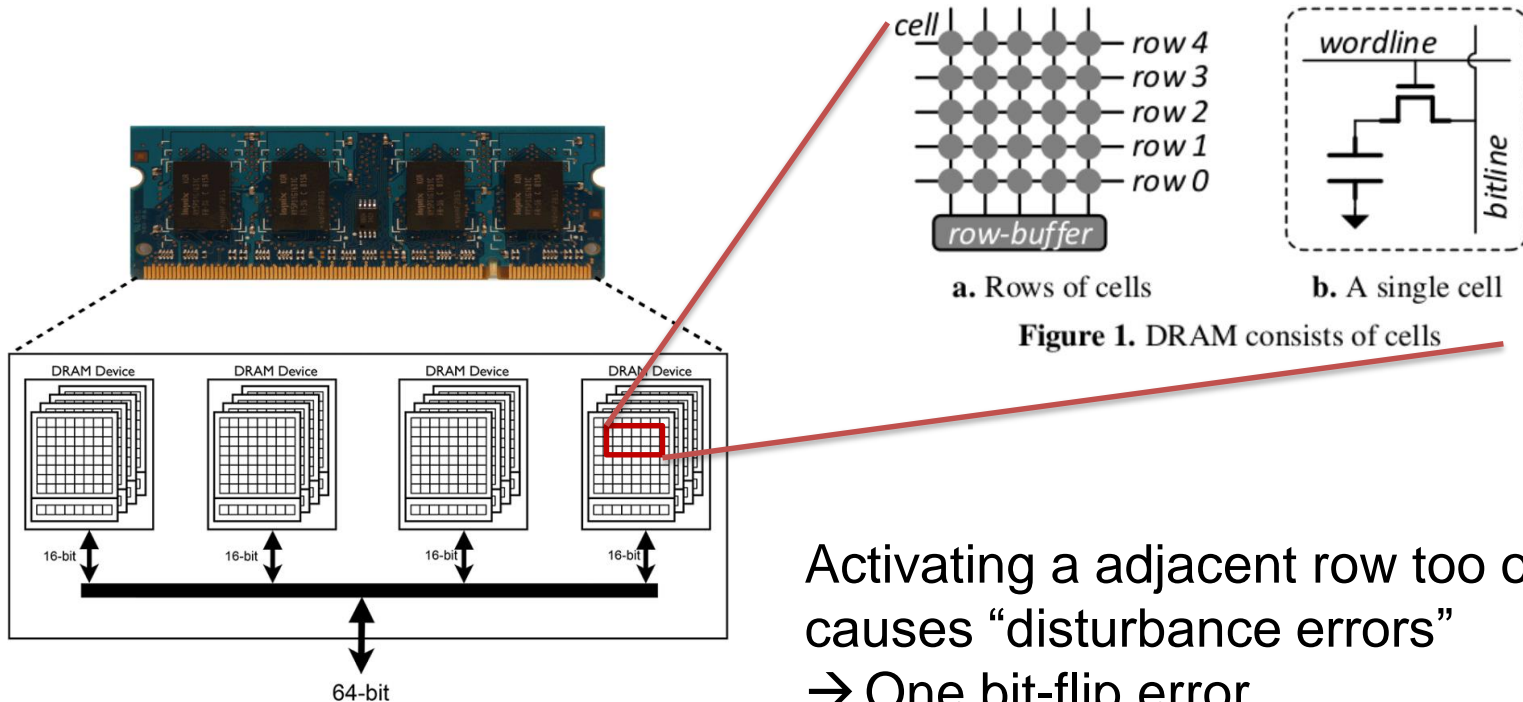


**a.** Rows of cells  **b.** A single cell

**Figure 1.** DRAM consists of cells

Activating a adjacent row too often causes "disturbance errors"
→ One bit-flip error
→ NaCl exploit, kernel exploit

# New Project: Secure Popcorn

- **Runtime ISA migration**
  - Move application threads to physically different machine with different ISA
  - Mitigate memory crosstalk attack by running on a physical different machine

- **Runtime ISA randomization**
  - Upon migration, randomize code and program status
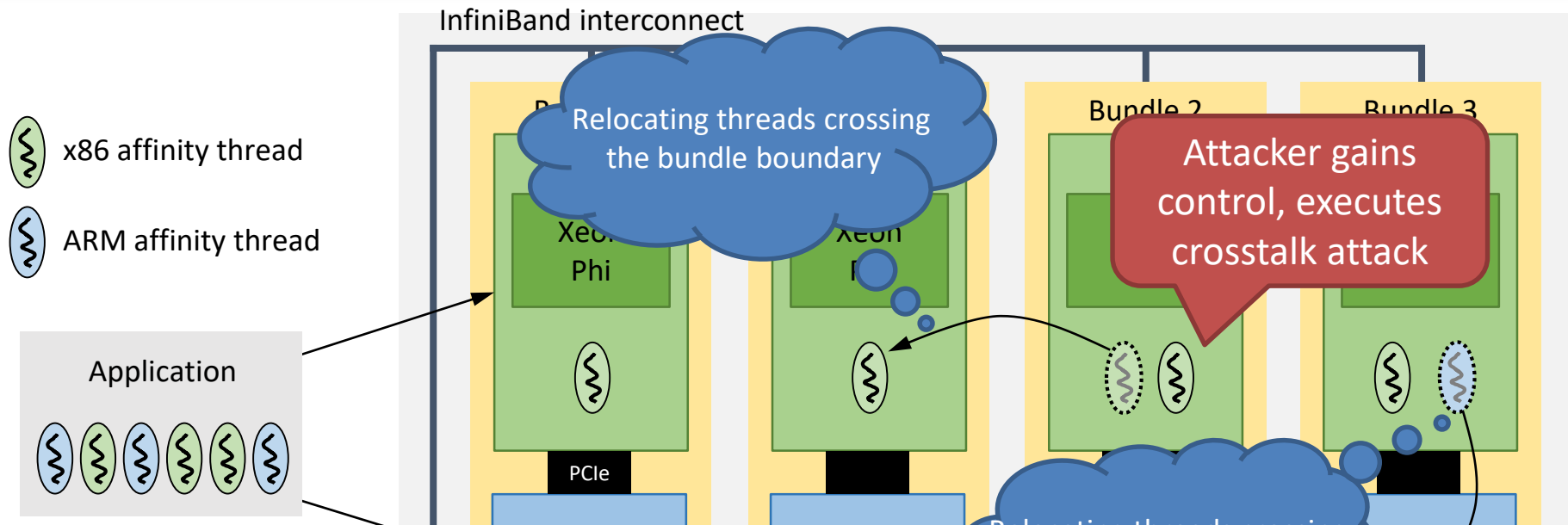  - Thwart attacker's knowledge by changing program code and data

- **Swift continuous code re-randomization**
  - Continuously re-randomize program code with very low overhead
  - Introduce real-time deadline (one re-randomization cycle) to attackers

- **Runtime integrity check (CFI & DFI)**
  - Check integrity of code and data at ISA boundary
  - Make low-overhead integrity check possible

# New Project: Secure Popcorn



Unleash program execution from a machine
→ Mitigate memory/cache crosstalk attack

# New Project: Secure Popcorn

Read-only text section

```
lea -0x78(%ebp),%eax
mov %eax,0x8(%esp)
call
d92e0
<memcpy>
….
ret
….
ldr s1, [x9]
orr x9, xzr, #0x4
add x8, x8, x9
str s0, [x8]
….
….
xor %eax,%eax
ret
….
….
pop %ebx
ret
```

Migrate to new ISA when intrusion detected or with some periodic frequency
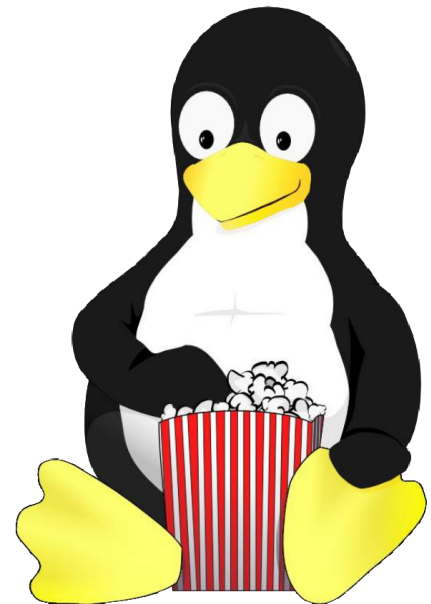
0x870f65
0x8709ad
0xbfff80

Gadget returns to garbage!

Systems Software Research Group

VIRGINIA TECH™

# Conclusion

- Heterogeneity is permeating all corners of computer architecture

- Popcorn Linux gives developers a better way to program heterogeneous-ISA systems – shared memory!
  - A compiler which builds multi-ISA binaries
  - An OS which enables cross-ISA thread and data migration
  - A state transformation runtime which converts ISA-specific data

- Allows developers to transparently take advantage of heterogeneity
  - Performance
  - Power/energy efficiency
  - Security

# More Information

- Popcorn Linux is open source and available online at http://popcornlinux.org

Systems Software Research Group

VIRGINIA TECH™

# Questions?