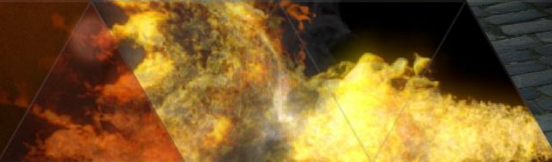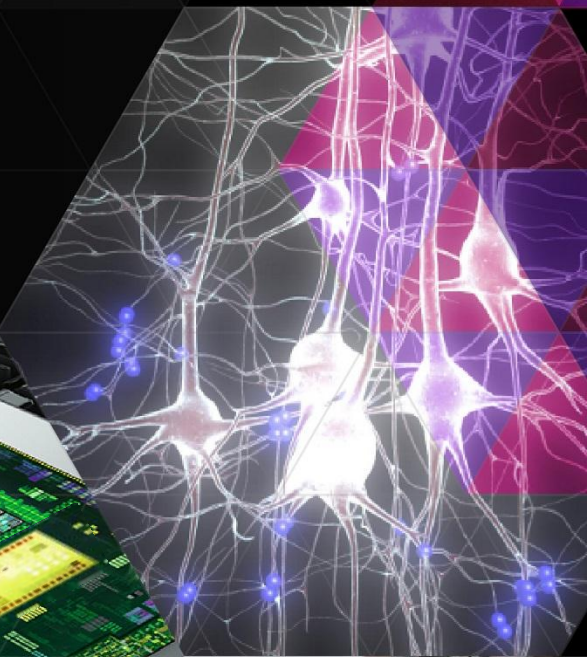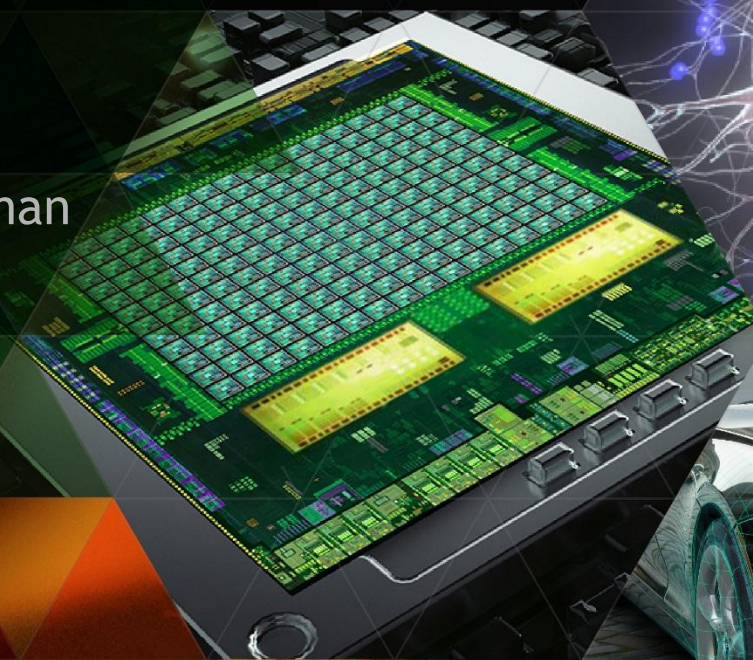# HOT CHIPS 2014
# NVIDIA'S DENVER PROCESSOR

Darrell Boggs, CPU Architecture

Co-authors:  Gary Brown, Bill Rozas,

Nathan Tuck, K S Venkatraman

# TEGRA K1
## with Dual Denver CPUs

*The First 64-bit Android Kepler-Class Chip*

NVIDIA.

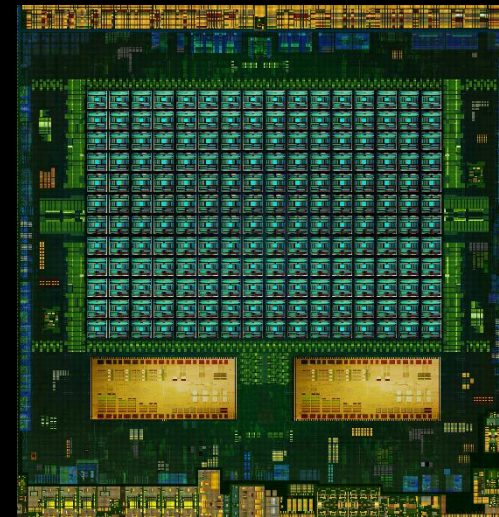# One Chip — Two Versions

# TEGRA K1

## 192-core Kepler-Class Chip

Pin Compatible

| Quad A15 CPUs | Dual Denver CPUs |
|---|---|
| 32-bit | 64-bit |
| 3-way Superscalar | 7-way Superscalar |
| Up to 2.3GHz | Up to 2.5GHz |
| 32K+32K L1$ | 128K+64K L1$ |

# DENVER VALUE PROPOSITION

**Highest performance and very energy-efficient ARMv8 processor**

▷ Greater dynamic sharing with GPU
▷ Extended battery life
▷ Low latency power-state transition
▷ Best web browsing experience

**Designed to bring PC-class performance to the ARM ecosystem**

▷ Content creation
▷ Gaming
▷ Enterprise applications

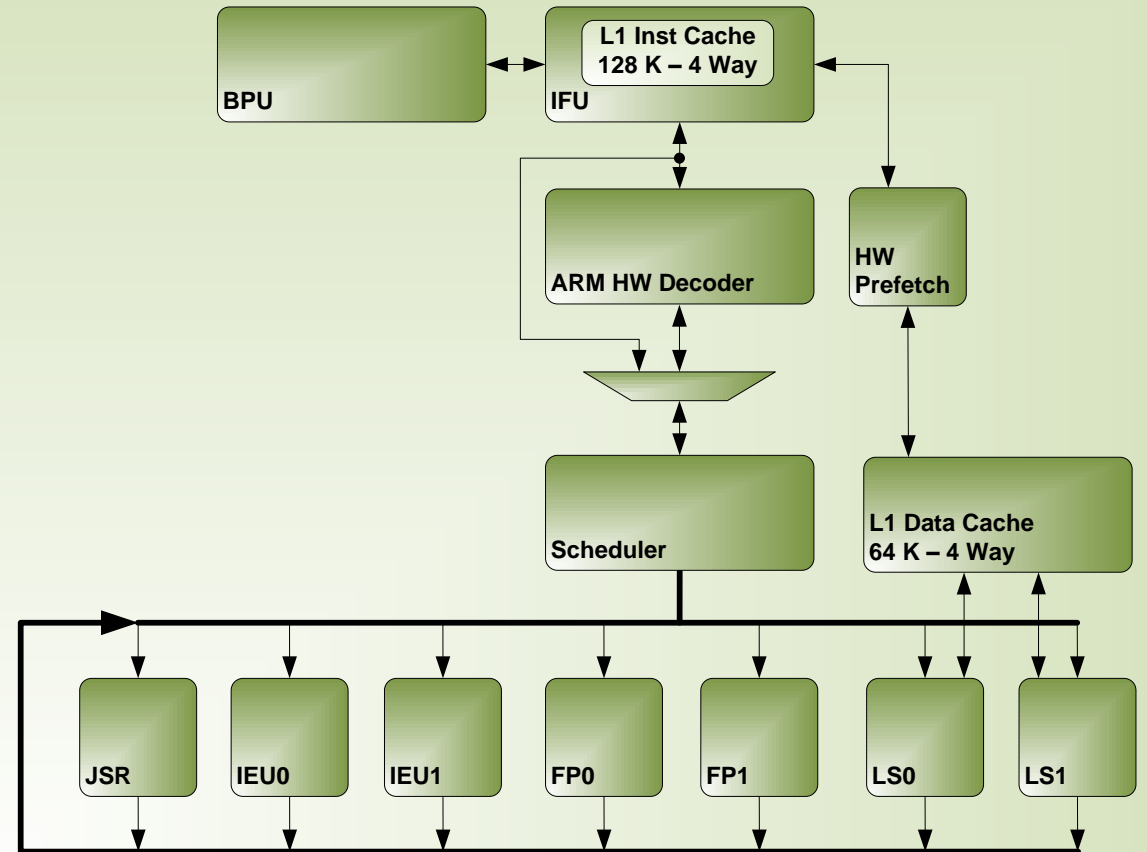4 ⬢ nVIDIA.

# DENVER CPU

## Highest Perf ARMv8 CPU
- 7-wide superscalar
- Aggressive HW prefetcher

## Dynamic Code Optimization
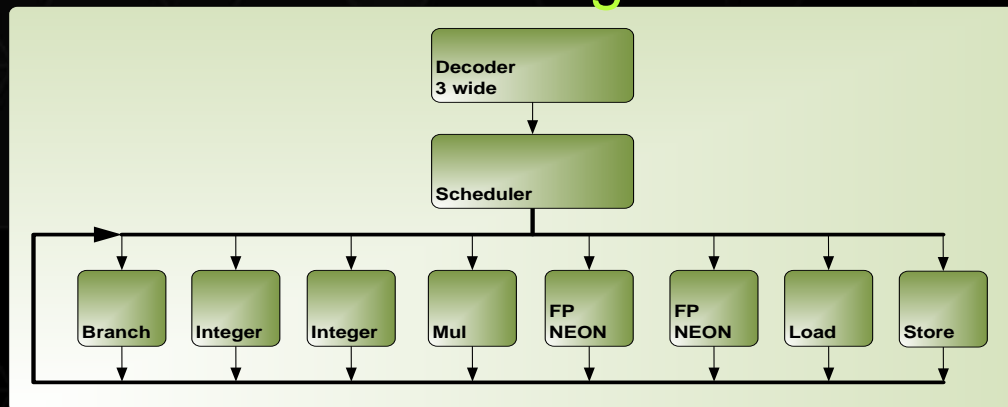- Optimize once, use many times
- OOO execution without the power



Denver Core

BPU

L1 Inst Cache
128 K – 4 Way

IFU

ARM HW Decoder

HW Prefetch

Scheduler

L1 Data Cache
64 K – 4 Way

JSR    IEU0    IEU1    FP0    FP1    LS0    LS1

# TEGRA K1 SUPERSCALAR ARCHITECTURE

## Cortex-A15  Tegra K1-32

Decoder 3 wide
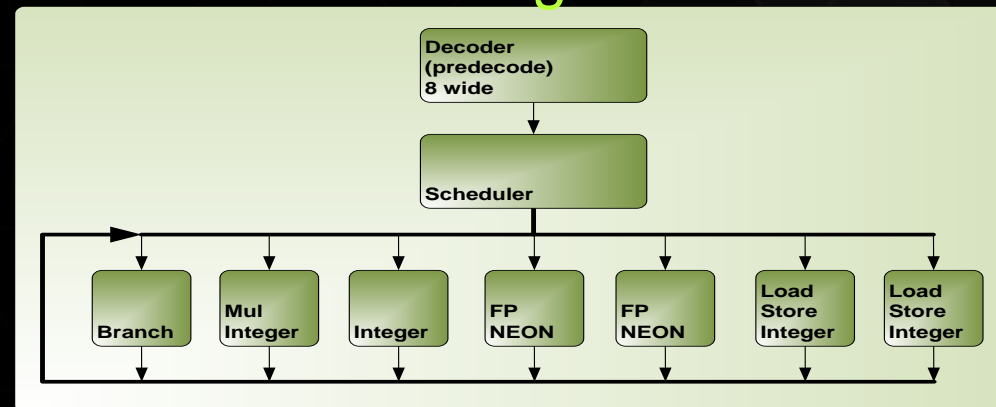
Scheduler

Branch | Integer | Integer | Mul | FP NEON | FP NEON | Load | Store

- Branch: 1
- Integer: 2
- Multiply: 1
- Floating Point/Neon: 2 x 64-bit
- LD/ST: 1 LD and 1 ST

Peak IPC 3

## Denver  Tegra K1-64

Decoder (predecode) 8 wide

Scheduler

Branch | Mul Integer | Integer | FP NEON | FP NEON | Load Store Integer | Load Store Integer
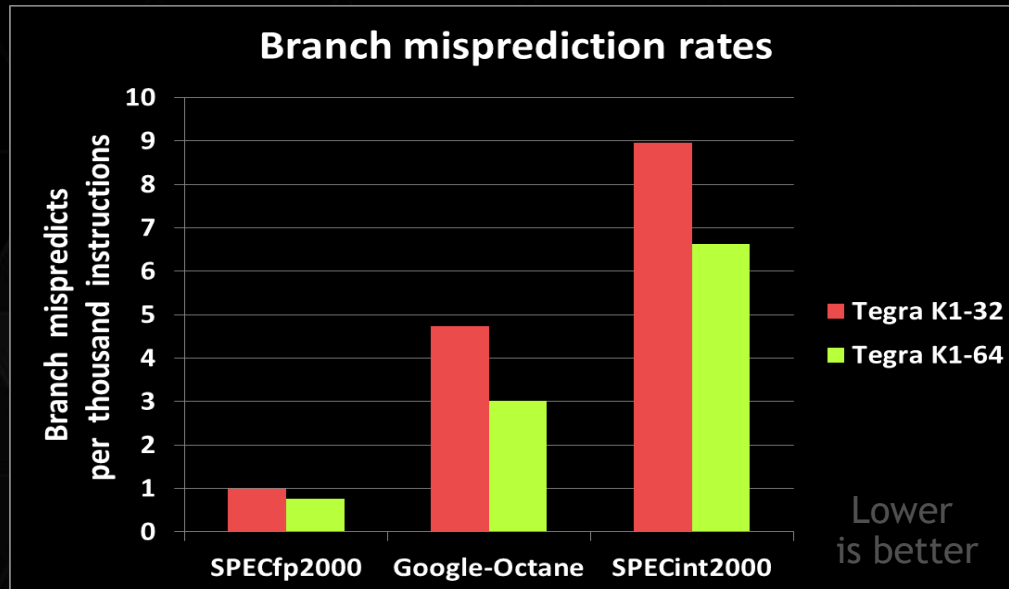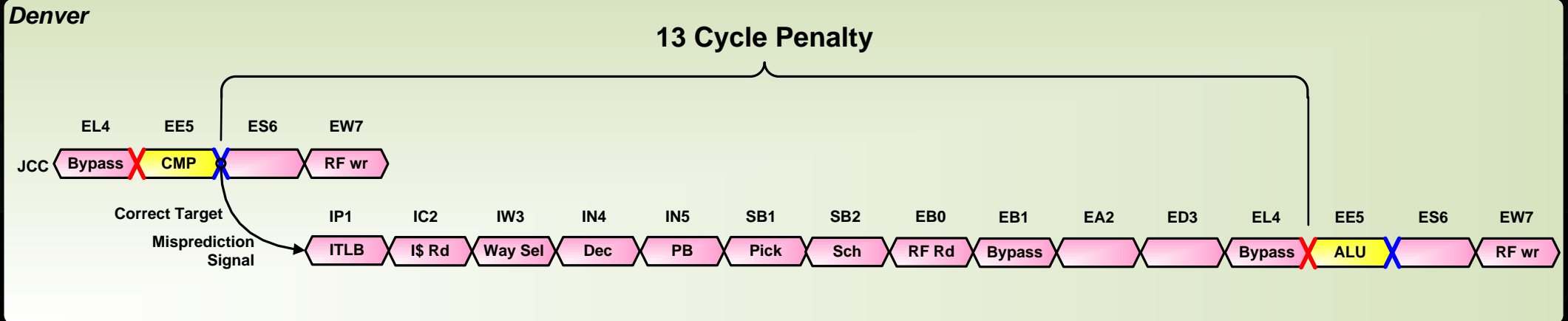
- Branch: 1
- Integer: 2 (+ Mul) + 2
- Floating Point/Neon: 2 x 128-bit
- LD/ST: 2 LD and/or ST

Peak IPC 7+

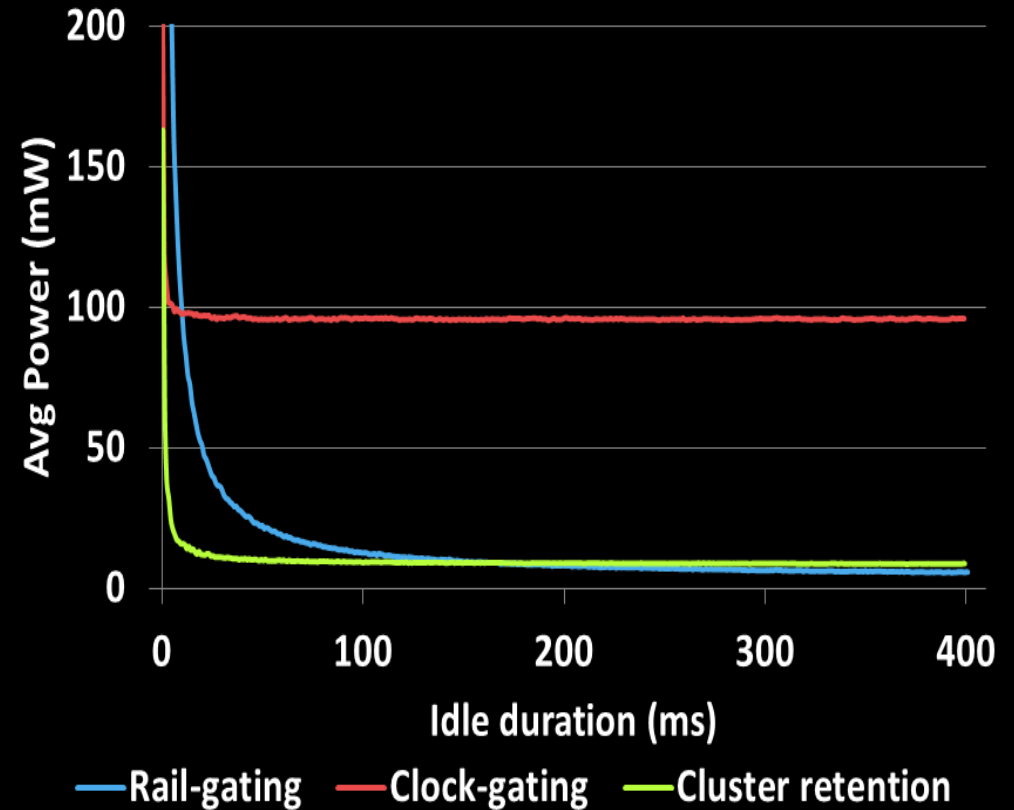NVIDIA.

# Pipeline Microarchitecture - Mispredict Penalty



**Denver**

**13 Cycle Penalty**

| EL4 | EE5 | ES6 | EW7 |
| --- | --- | --- | --- |

JCC: Bypass — CMP — RF wr

Correct Target
Misprediction Signal

| IP1 | IC2 | IW3 | IN4 | IN5 | SB1 | SB2 | EB0 | EB1 | EA2 | ED3 | EL4 | EE5 | ES6 | EW7 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ITLB | I$ Rd | Way Sel | Dec | PB | Pick | Sch | RF Rd | Bypass | | | Bypass | ALU | | RF wr |

## Branch misprediction rates



Branch mispredicts per thousand instructions

- Tegra K1-32
- Tegra K1-64

SPECfp2000   Google-Octane   SPECint2000

Lower is better

- **Tegra K1-32**
  - 15 cycle mispredict

- **Tegra K1-64**
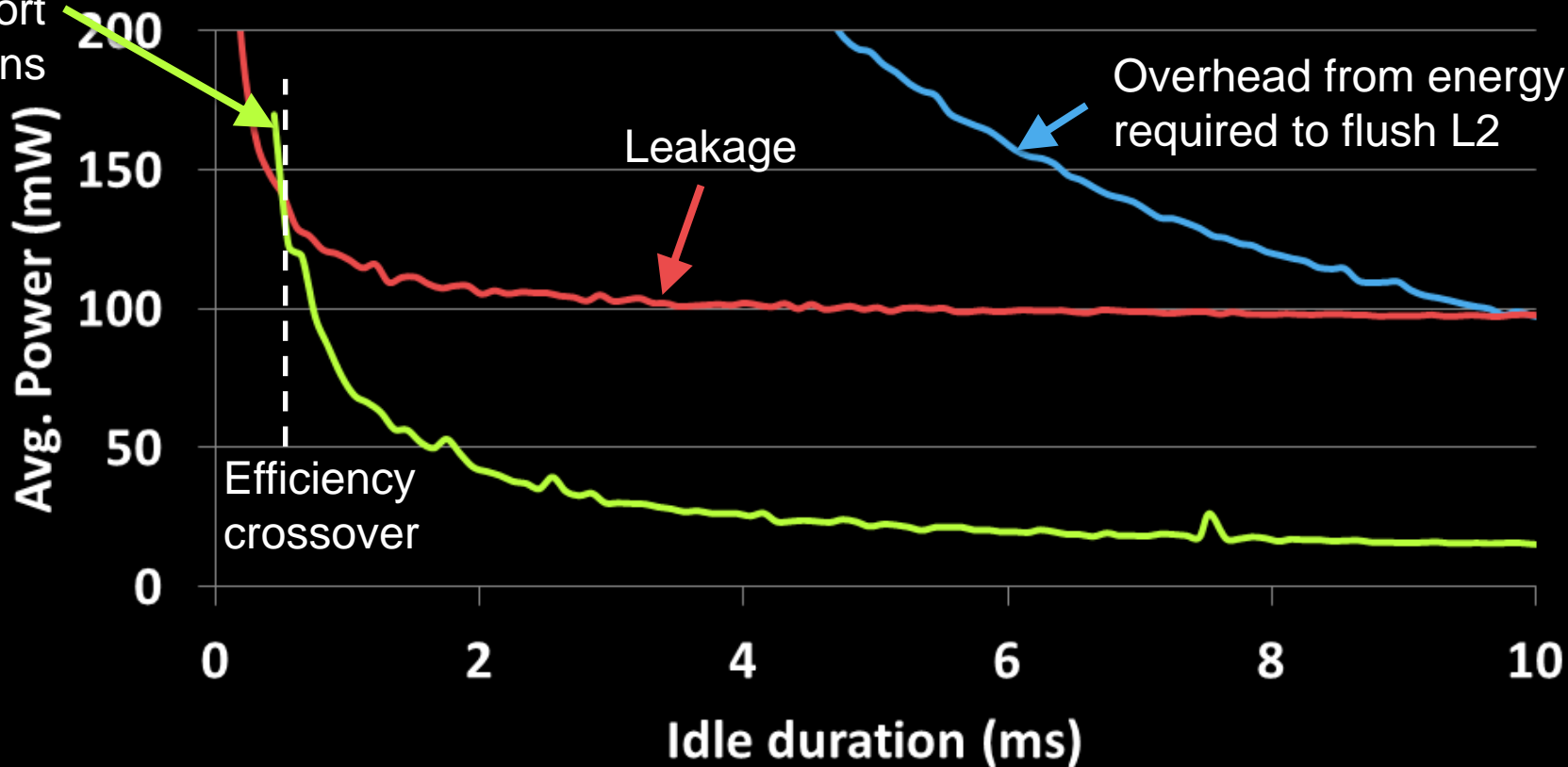  - 13 cycle mispredict

**Higher ILP and efficiency**

⬢ nVIDIA.

# CORE CLUSTER RETENTION STATE

- New power management state: CC4

- Allows cache and architectural state retention

- Allows voltage to be reduced below Vmin to a retention voltage

- Fast entry and exit latencies enable wider range of use



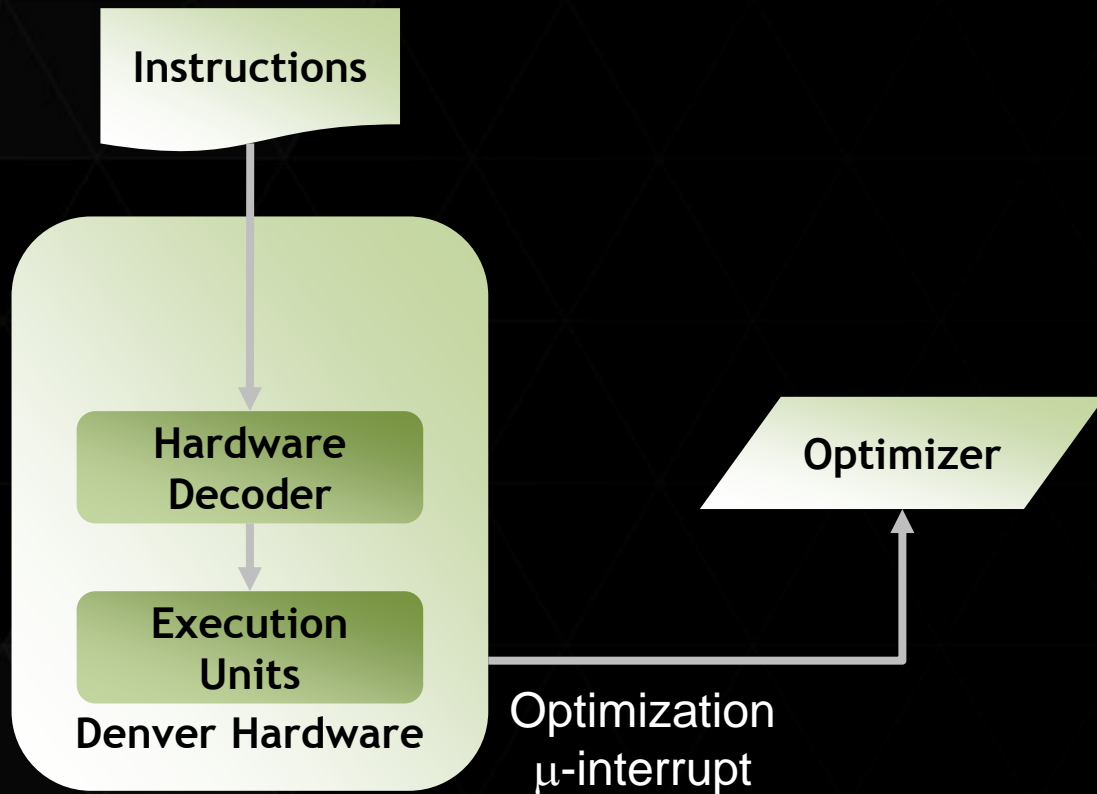NVIDIA.

# DENVER IDLE POWER IMPROVES WITH RETENTION



Power penalty if entered for short durations

Overhead from energy required to flush L2

Leakage

Efficiency crossover

Avg. Power (mW)

200

150

100

50

0

0    2    4    6    8    10

Idle duration (ms)

—— Rail-gating    —— Clock-gating    —— Cluster retention

# DYNAMIC CODE OPTIMIZATION
## OPTIMIZE ONCE, USE MANY TIMES

Instructions

Hardware Decoder

Optimizer

Execution Units

Denver Hardware

Optimization μ-interrupt

NVIDIA.

# DYNAMIC CODE OPTIMIZATION
## OPTIMIZE ONCE, USE MANY TIMES

**Instructions**

**Hardware Decoder**

**Execution Units**

**Denver Hardware**

Dynamic Profile Information

**Optimizer**

▸ Unrolls Loops
▸ Renames registers
▸ Reorders Loads and Stores
▸ Improves control flow
▸ Removes unused computation
▸ Hoists redundant computation
▸ Sinks uncommonly executed computation
▸ Improves scheduling

**Optimized μcode**

**Optimization Cache**

# DYNAMIC CODE OPTIMIZATION
## OPTIMIZE ONCE, USE MANY TIMES

Instructions

Optimization Lookup

Hardware Decoder

Execution Units

Denver Hardware

Optimized code execution from optimization cache

Optimized μcode

Optimization Cache

NVIDIA.

# DYNAMIC CODE OPTIMIZATION
## OPTIMIZE ONCE, USE MANY TIMES

Instructions

Optimization Lookup

Hardware Decoder

Execution Units

**Denver Hardware**

Dynamic Profile Information

Optimizer

Chaining

Optimized μcode

Optimized μcode

Optimized μcode

Optimized μcode

**Optimization Cache**

⬢ nVIDIA.

# DENVER PERFORMANCE



_Chart — Performance Relative to Tegra K1 32, by benchmark (DMIPS, SPECInt 2K, SPECFP 2K, AnTuTu 4, Geekbench 3 Single-Core, Google Octane v2.0, 16MB Memcpy (GB/s), 16MB Memset (GB/s), 16MB Memread (GB/s)). Legend: Baytrail (Celeron N2910), Krait-400 (8974-AA), iPhone 5s (A7 Cyclone), Haswell (Celeron 2955U), Tegra K1 64._

# DCO: AN EXAMPLE SPECINT – CRAFTY EXECUTION

Full benchmark run

Dynamic Code Optimizer

% Exec Types

Profile of execution

Ticks

HW Decoder execution

Optimized ucode execution

# DCO: AN EXAMPLE SPECINT - CRAFTY EXECUTION
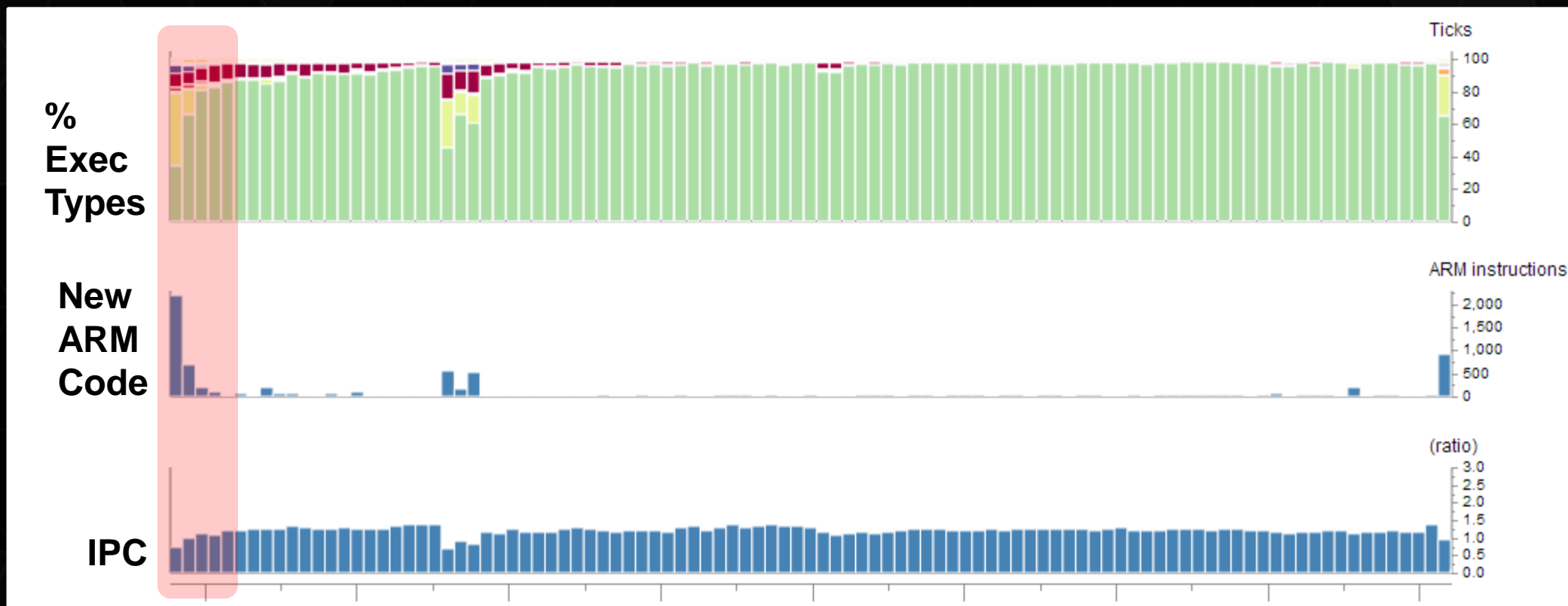
# DCO: AN EXAMPLE SPECINT - CRAFTY EXECUTION

# DCO: AN EXAMPLE SPECINT – CRAFTY EXECUTION

Full benchmark run



NVIDIA.

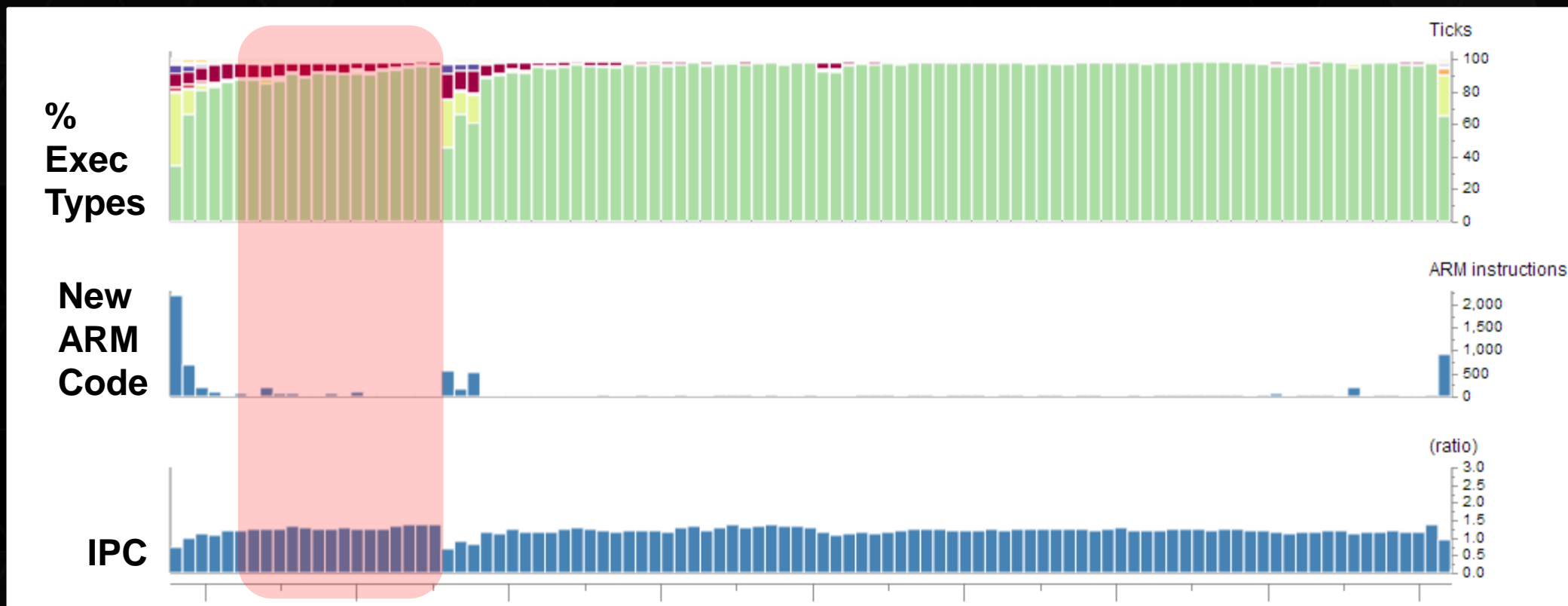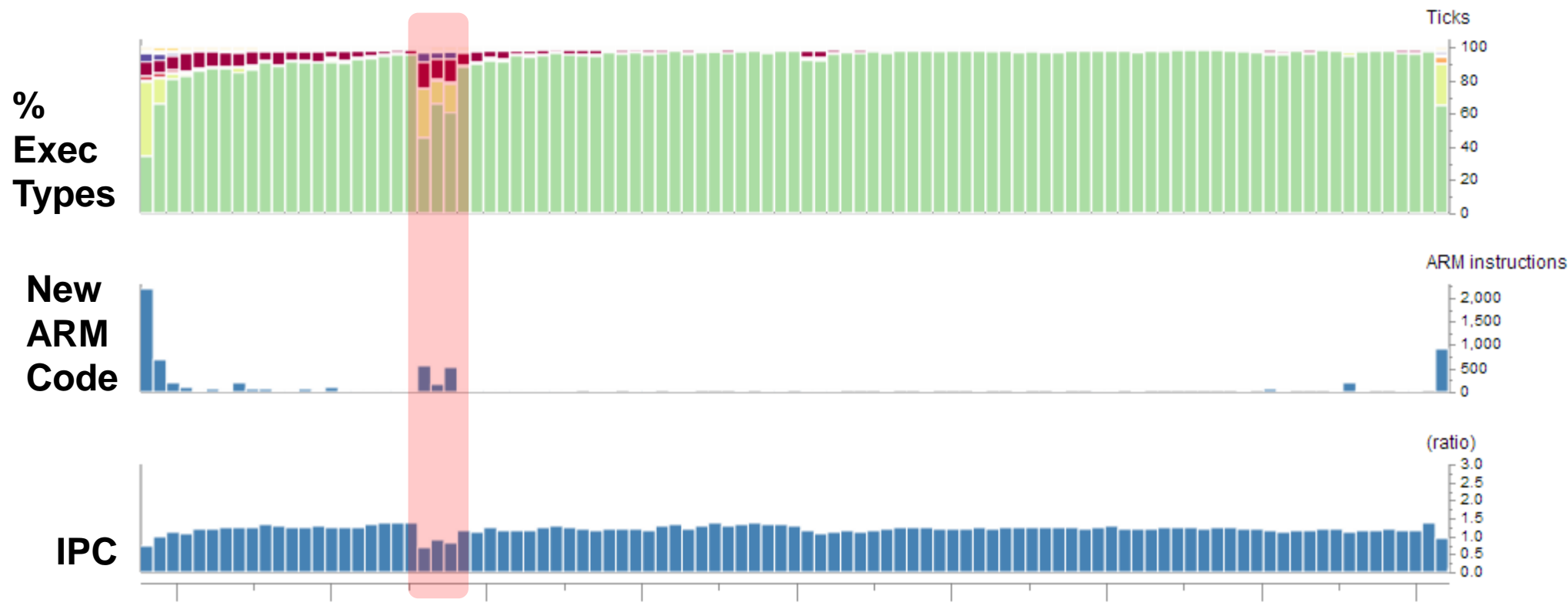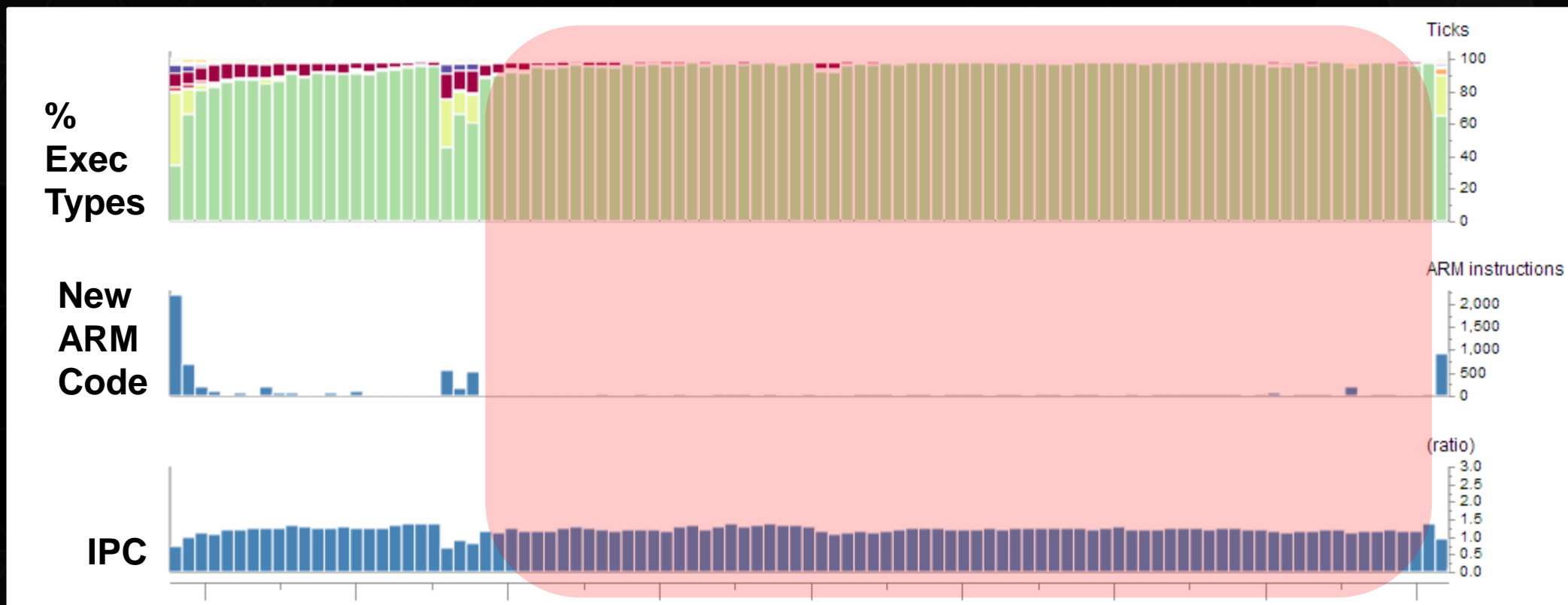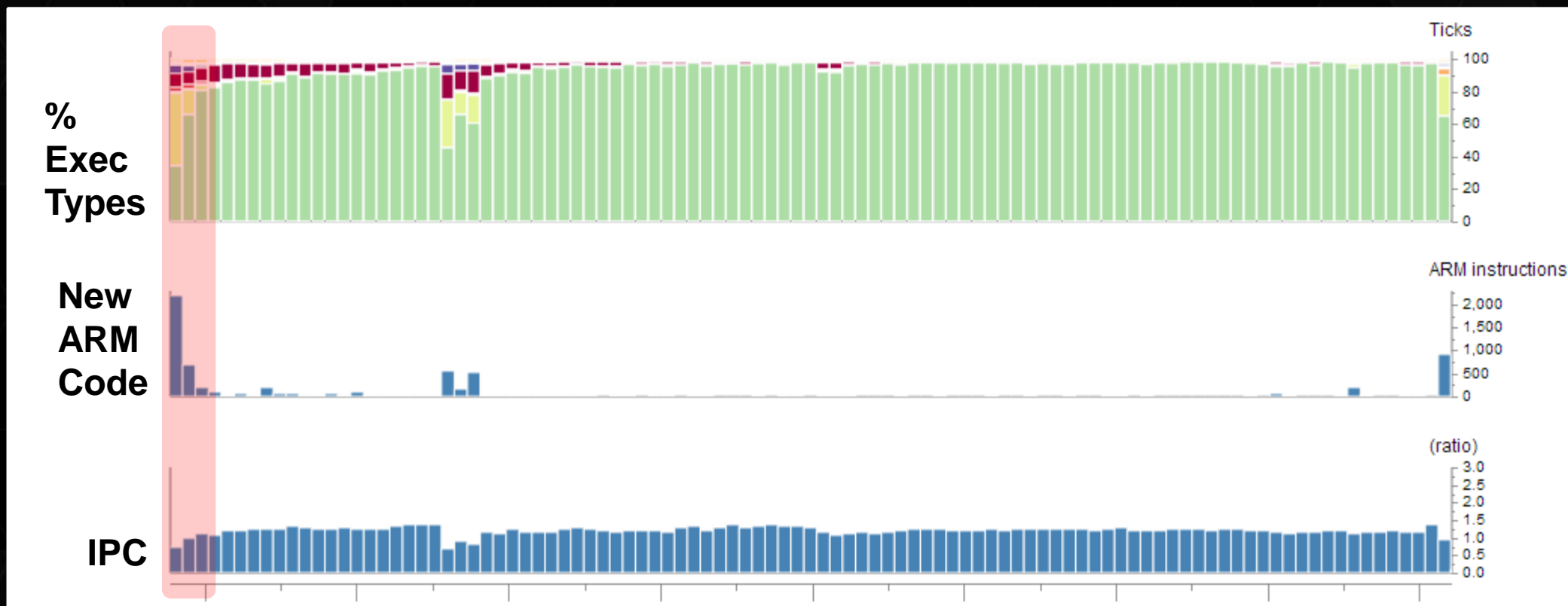# DCO: AN EXAMPLE SPECINT - CRAFTY EXECUTION

# DCO: AN EXAMPLE SPECINT - CRAFTY EXECUTION

# DCO: AN EXAMPLE SPECINT - CRAFTY EXECUTION



Full benchmark run

% Exec Types

Ticks

New ARM Code

ARM instructions

IPC

(ratio)

# DCO: AN EXAMPLE SPECINT – CRAFTY EXECUTION

3% of benchmark run

% Exec Types

New ARM Code

IPC

# CONCLUSION

▸ Dynamic Code Optimization is the architecture of the future

   ▹ Breaks the out-of-order window physical limitation

   ▹ Opens synergy between HW and SW that current architectures lack

   ▹ Improves efficiency by optimizing once and using many times

▸ Delivering PC-class performance to mobile form factors

   ▹ Enables PC-class gaming experience

   ▹ Enables true enterprise applications

   ▹ Enables content creation

<span>NVIDIA.</span>

# ACKNOWLEDGMENT

▷ We would like to thank the CPU team in NVIDIA for all the creativity, hard work, and dedication to bring this vision to a reality.

NVIDIA.