

Safe and Secure Real-Time Computing Infrastructure for Intelligent Cyber-Physical Systems

Heechul Yun

Associate Professor, University of Kansas

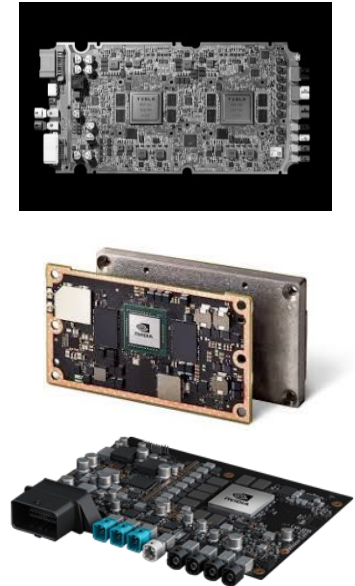
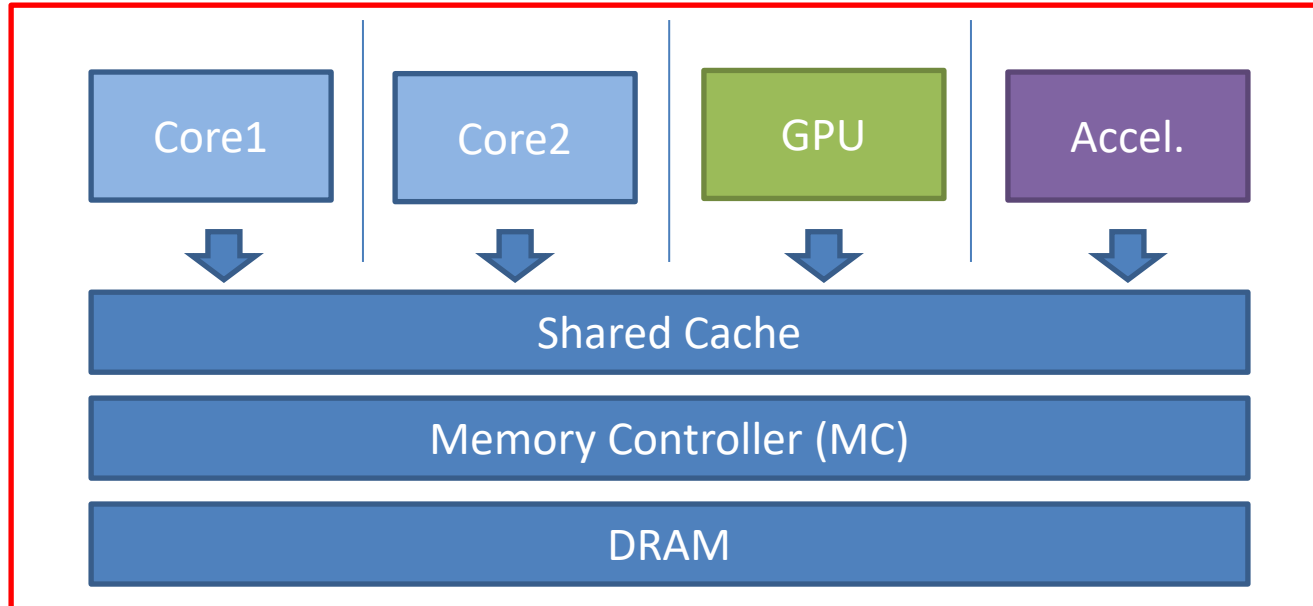
March, 2021

Modern Cyber-Physical Systems

- Cyber Physical Systems (CPS)
 - Cyber (Computer) + Physical (Plant)
- Real-time
 - Control physical process in real-time
- Safety-critical
 - Can harm people/things
- **Intelligent**
 - Can function autonomously

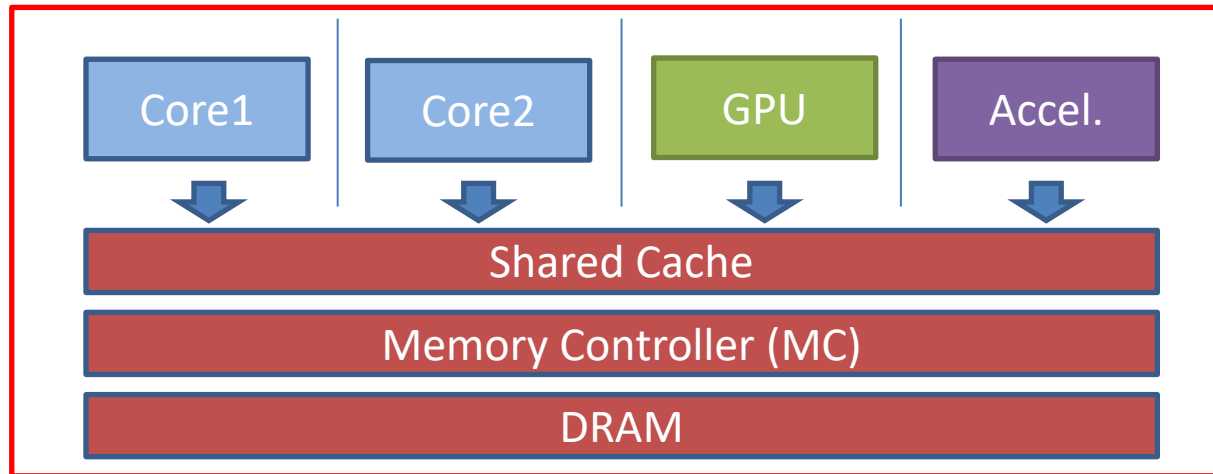


Modern System-on-Chip (SoC)



- Integrate multiple cores, GPU, accelerators
- Good performance, cost, size, weight, power
- **But ...**

Time Predictability and Safety



- Many important hardware resources are *shared*
- Difficult to guarantee predictable timing
- A **safety problem** in CPS

Timing Channels and Security



Meltdown

Meltdown breaks the most fundamental isolation between user applications and the operating system. This attack allows a program to access the memory, and thus also the secrets, of other programs and the operating system.



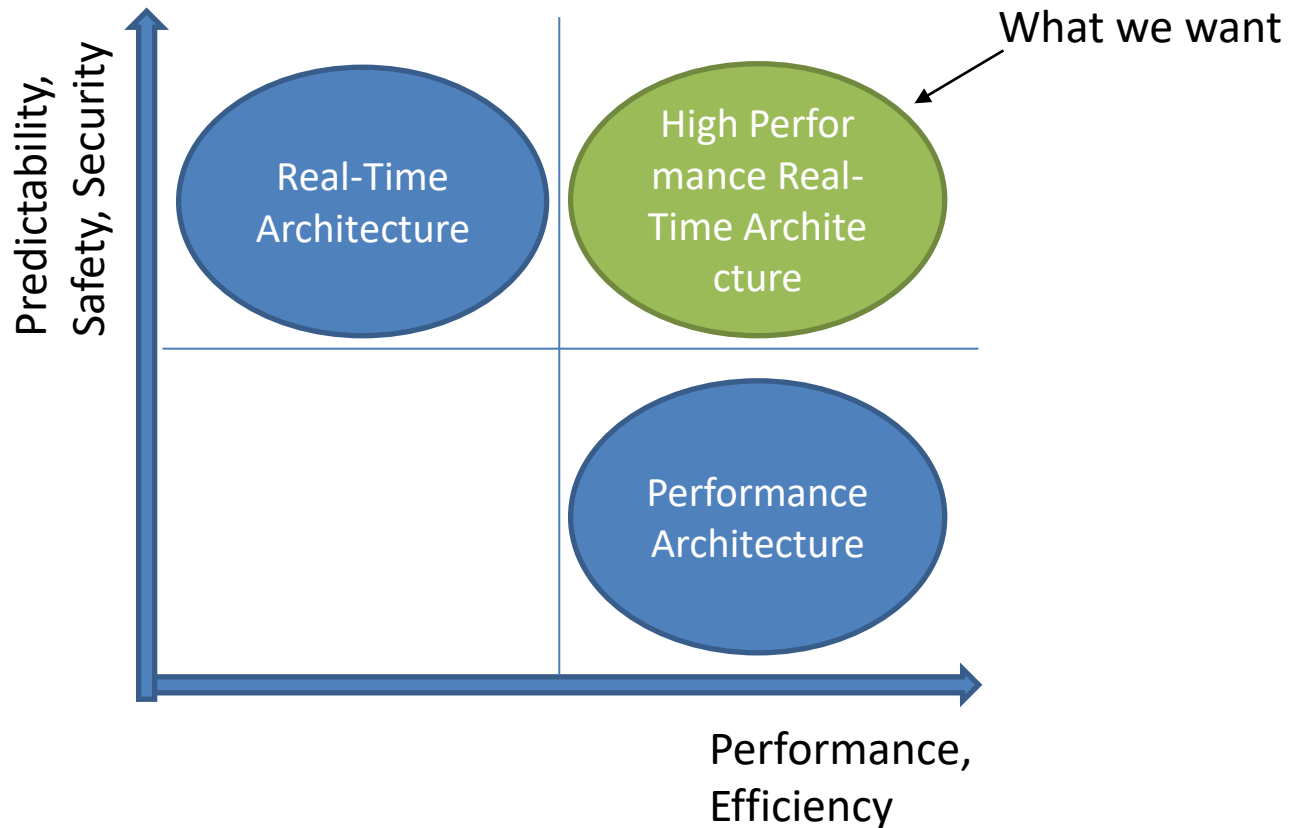
Spectre

Spectre breaks the isolation between different applications. It allows an attacker to trick error-free programs, which follow best practices, into leaking their secrets. In fact, the safety checks of said best practices actually increase the attack surface and may make applications more susceptible to Spectre

<https://meltdownattack.com/>

- Measurable timing differences in accessing *shared* hardware resources can leak secret: a **security problem**

Computing Infrastructure for CPS



- Predictable, secure, and high-performance computing

Our Research

- Build **safe and secure real-time computing infrastructure** for the next generation of intelligent Cyber Physical Systems (CPS).

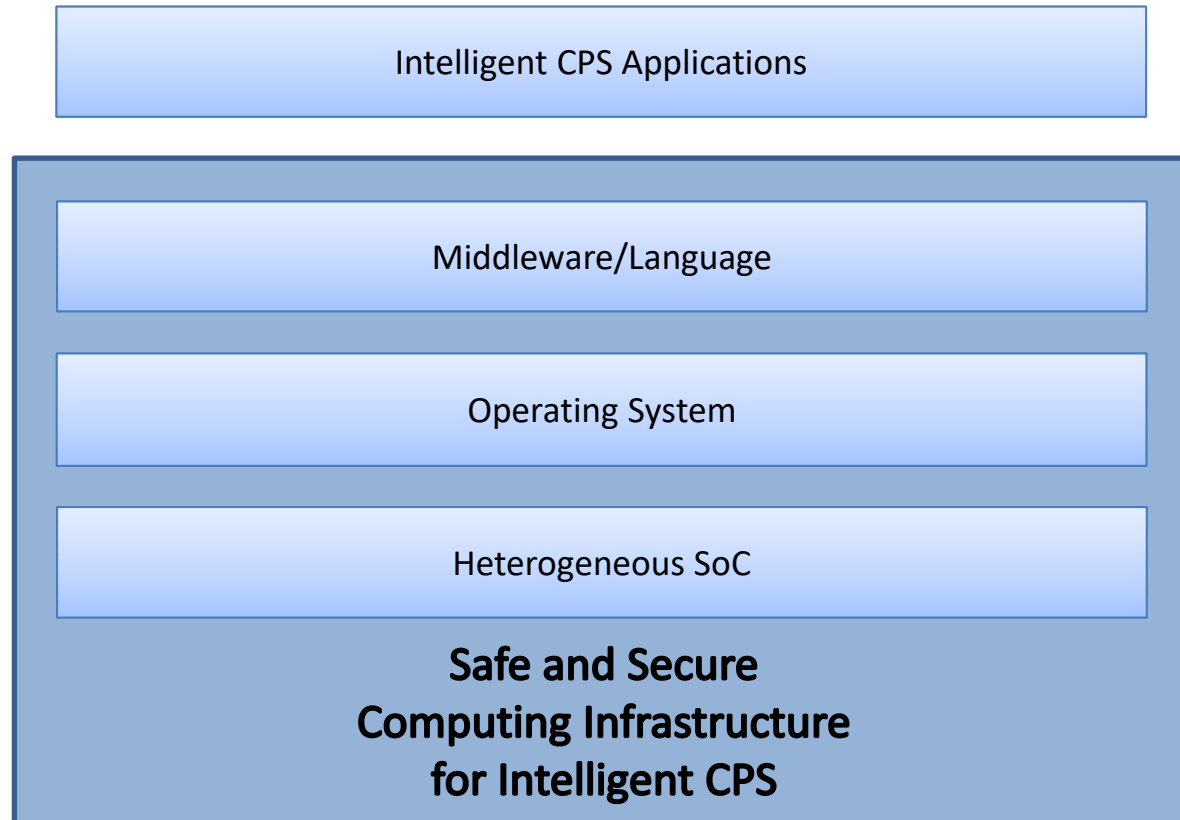
Research Sponsors



Equipment Sponsors

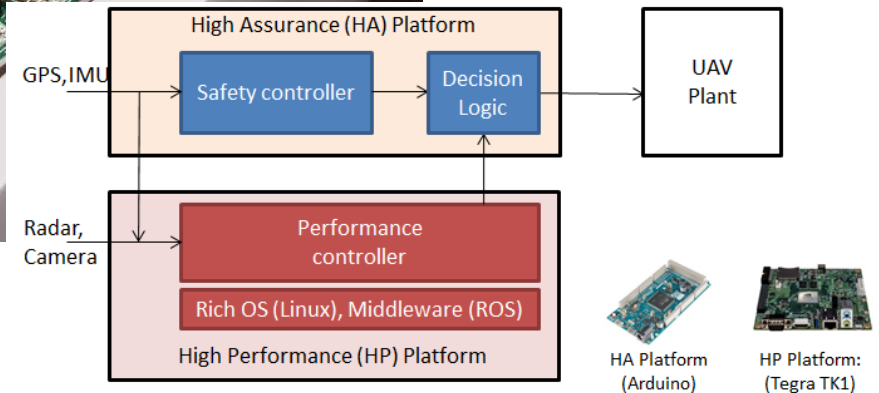
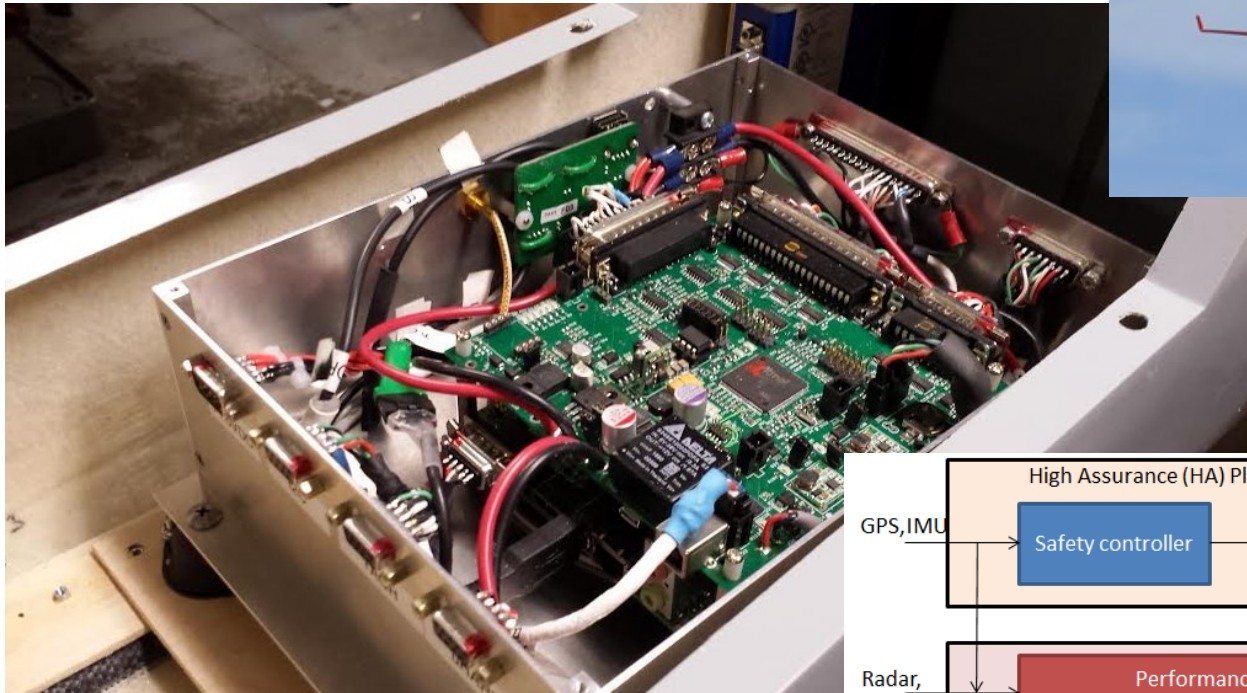


Our Approach



- Holistic, cross-layer approach for safe and secure CPS

KU Fixed-wing UAV



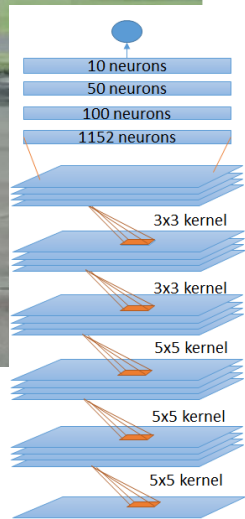
Safe and Intelligent UAV architecture

DeepPicar

- End-to-end deep learning: *pixels* to *steering*
- Using **identical** DNN with NVIDIA's DAVE-2



DAVE-2



Same DNN

output: steering angle
fc4: fully-connected layer
fc3: fully-connected layer
fc2: fully-connected layer
fc1: fully-connected layer

conv5: 64@1x18
convolutional layer

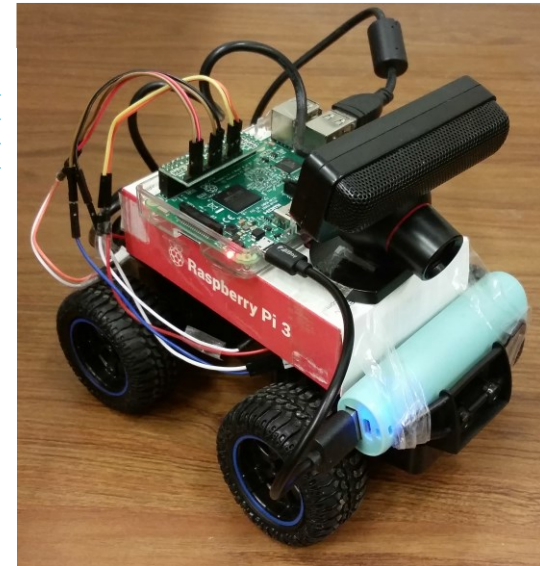
conv4: 64@3x20
convolutional layer

conv3: 48@5x22
convolutional layer

conv2: 36@14x47
convolutional layer

conv1: 24@31x98
convolutional layer

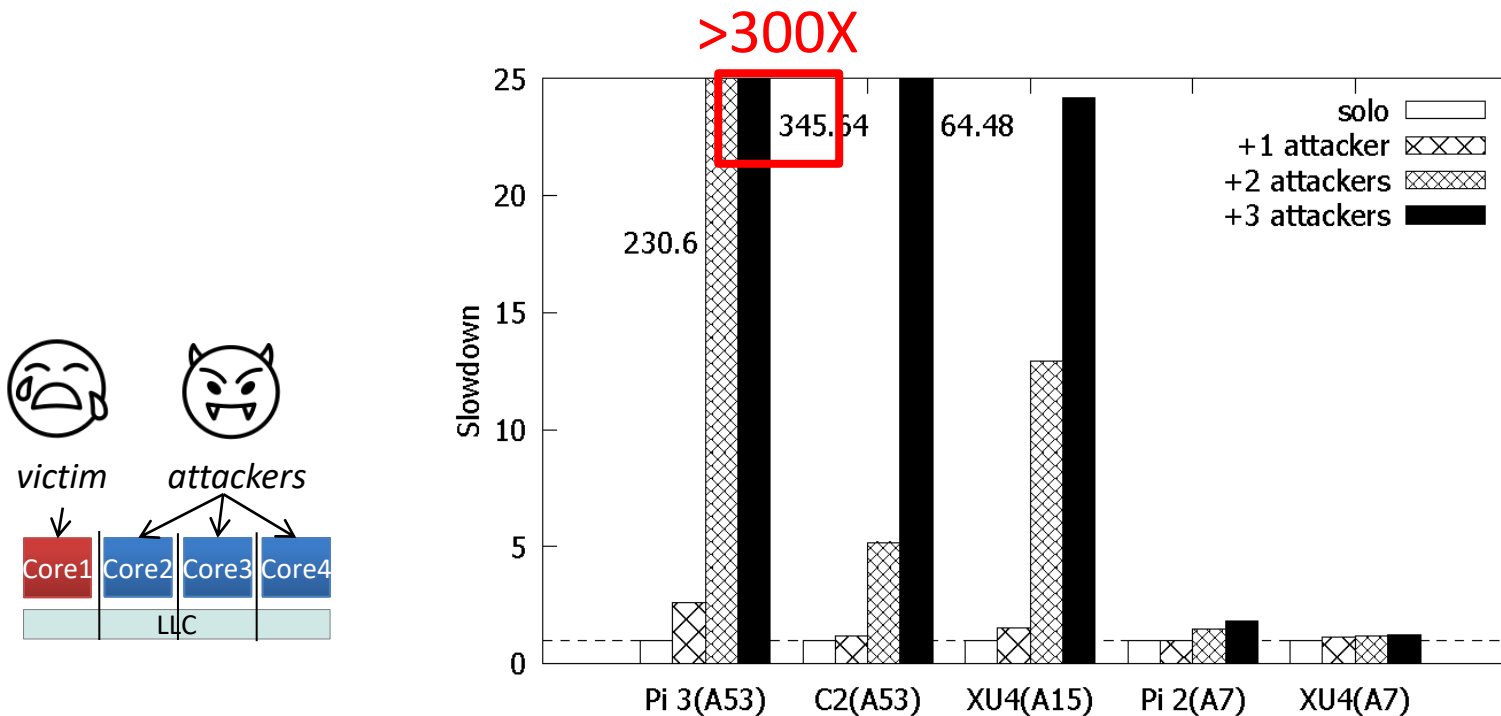
input: 200x66 RGB pixels



DeepPicar

DNN based real-time control in embedded multicore CPU

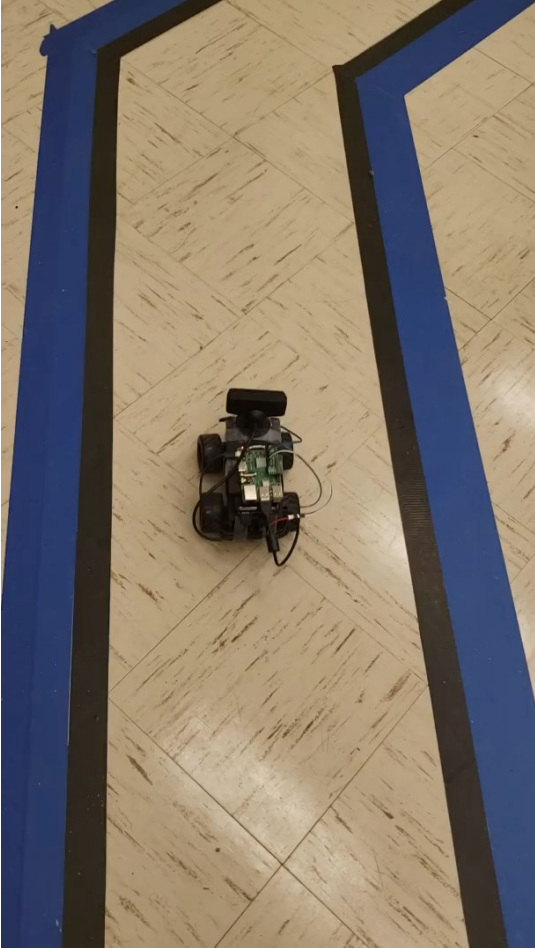
Cache Denial-of-Service Attacks



- Observed worst-case: >300X (times) slowdown
 - On popular in-order multicore processors

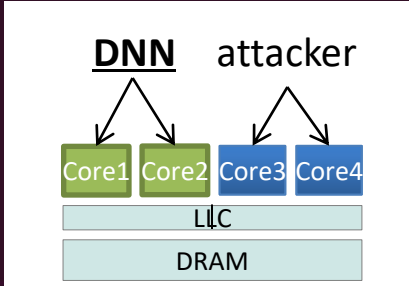
Found serious DoS attack vulnerability in COTS multicore

Effect of Cache DoS Attack



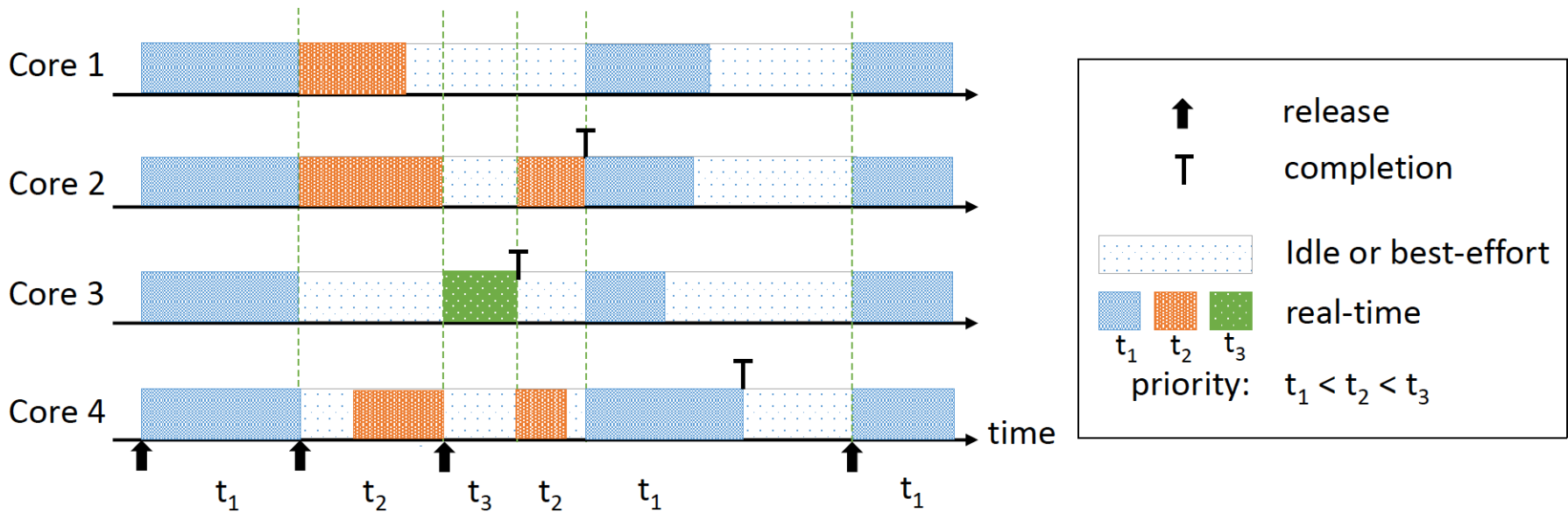
```
pi@raspberrypi:~/Documents/DeepPicar-v2 $ ./drive.sh
DNN is on
Initilize camera.
start camera thread
camera init completed.
Load TF

pi@raspberrypi:~/Documents/DeepPicar-v2 $ ./attack.sh
```



<https://youtu.be/Jm6KSDqlqiU>

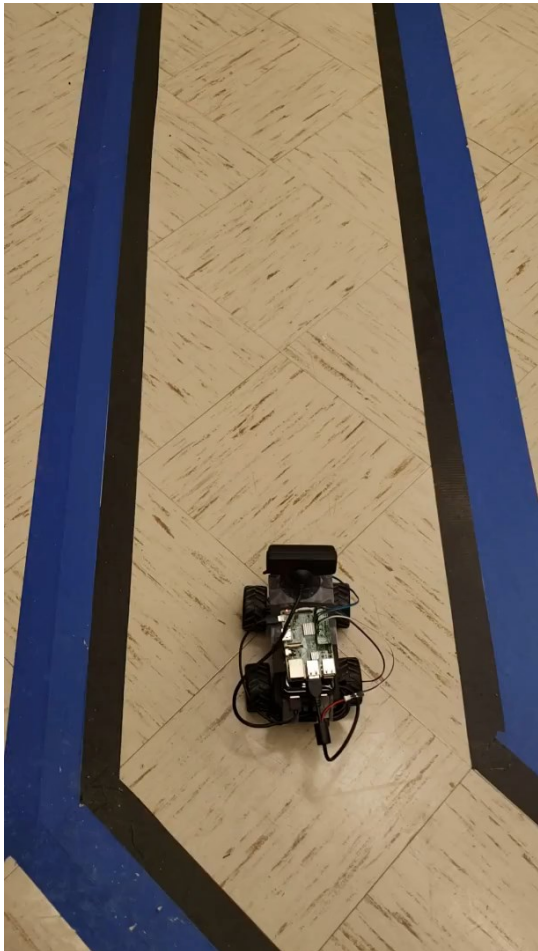
RT-Gang



- **One parallel real-time task---a gang---at a time**
 - Eliminate inter-task interference by construction
- Schedule best-effort tasks during slacks **w/ throttling**

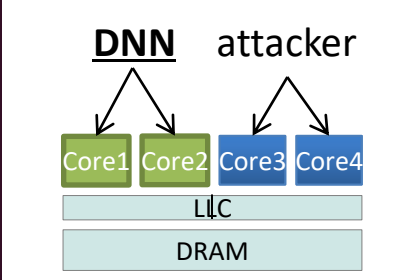
OS solution for predictable real-time computing on COTS

Effect of RT-Gang



```
pi@raspberrypi:~/Documents/DeepPicar-v2 $ ./drive.sh
DNN is on
Initilize camera.
start camera thread
camera init completed.
Load TF

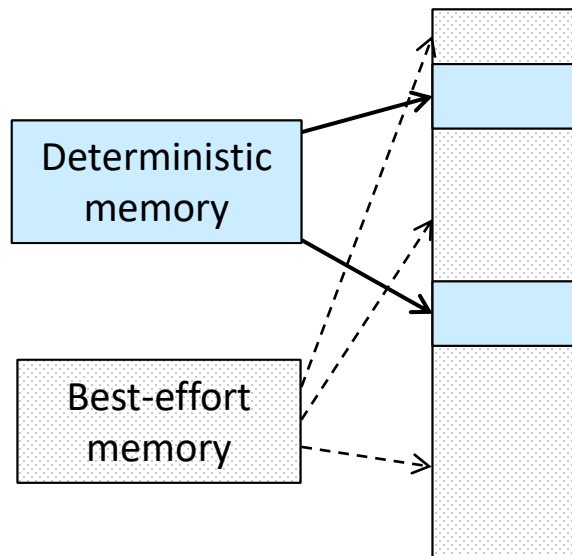
pi@raspberrypi:~/Documents/DeepPicar-v2 $ ./attack.sh
```



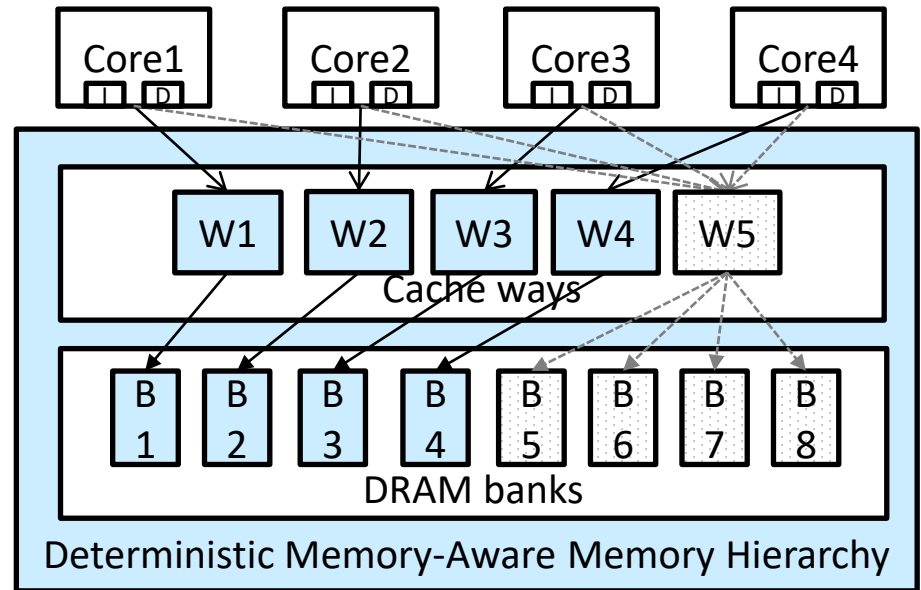
The diagram illustrates the hardware architecture. It shows four cores: Core1 and Core2 are green and labeled 'DNN'; Core3 and Core4 are blue and labeled 'attacker'. All four cores are connected to a shared L2 cache (L2C), which is connected to DRAM.

Deterministic Memory

- Declare all or part of address space as deterministic memory
- DM-aware **end-to-end resource management**



Application view (logical)

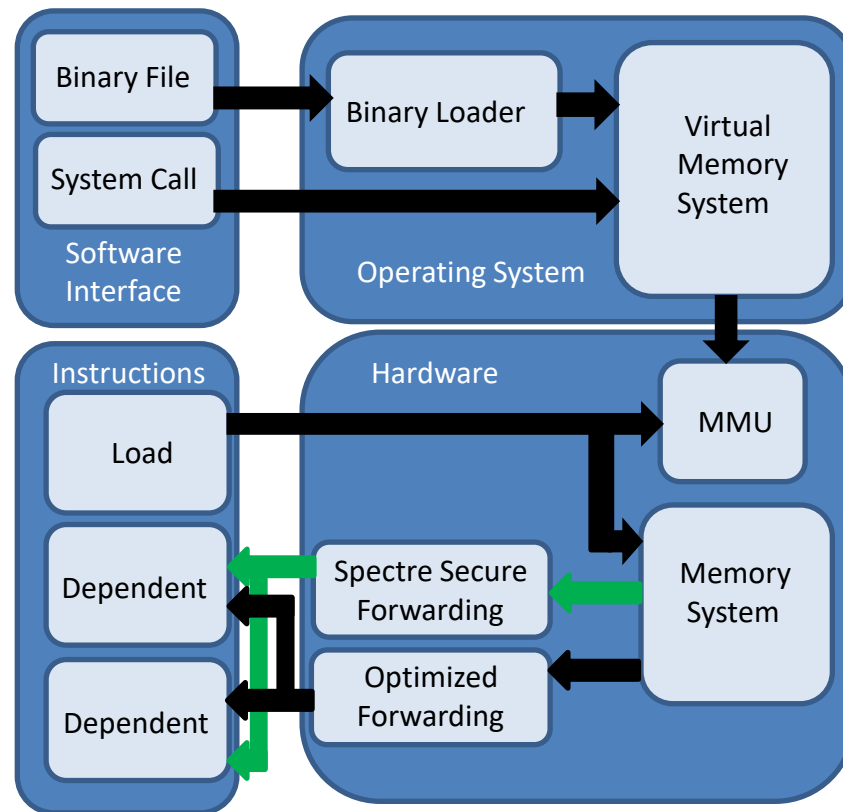


System-level view (physical)

Data-centric cross-layer approach for **real-time**

SpectreGuard

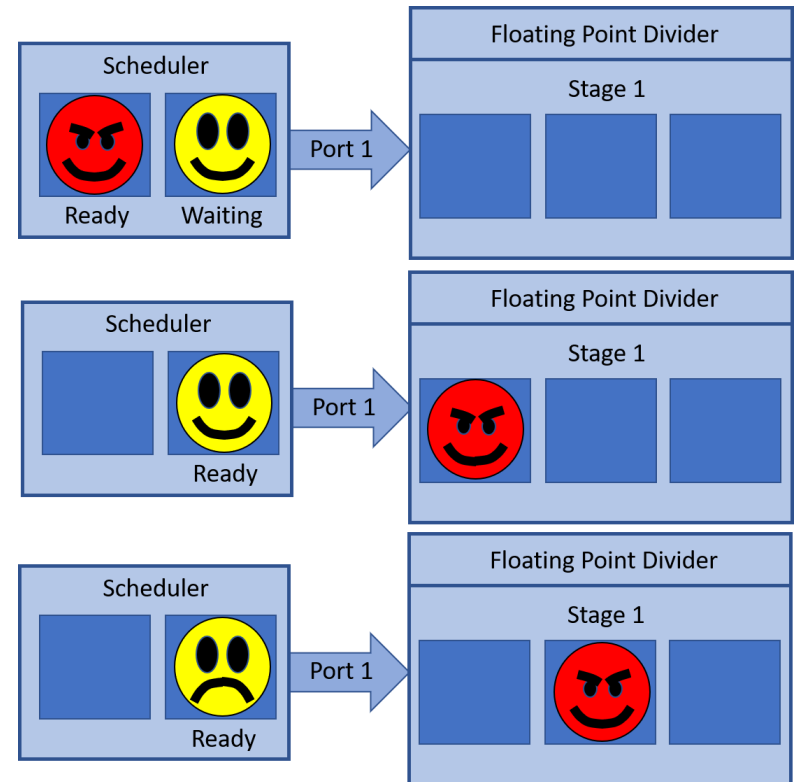
- Step 1: Software tells OS what **data** is secret
- Step 2: OS updates the page table entries
- Step 3: Load of the secret data is identified by MMU
- Step 4: secret data forwarding is **delayed** until safe



Data-centric cross-layer approach for security

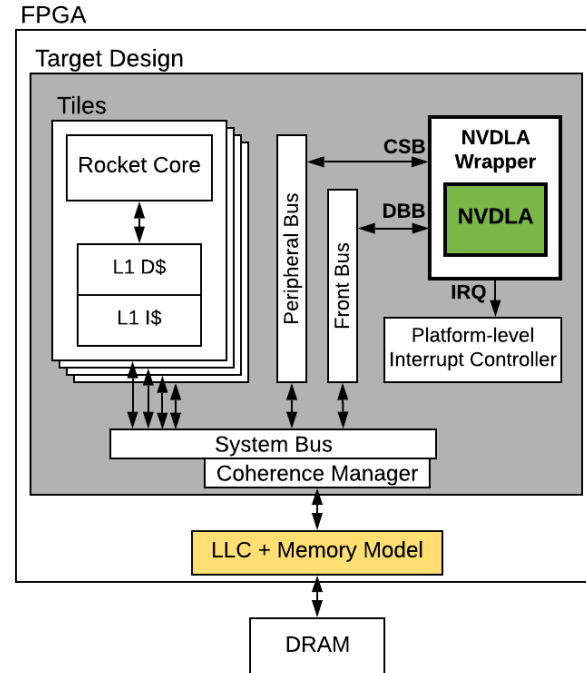
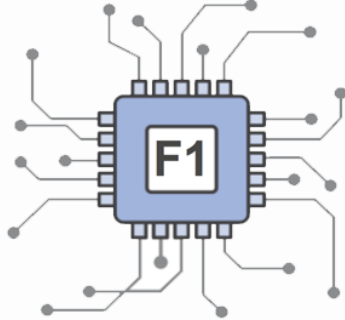
SpectreRewind

- New covert channel for Spectre-like speculative execution attacks
- Exploit contention on non-pipelined functional units
- Leak secret to past instructions, bypassing mitigation techniques
- High performance, low noise



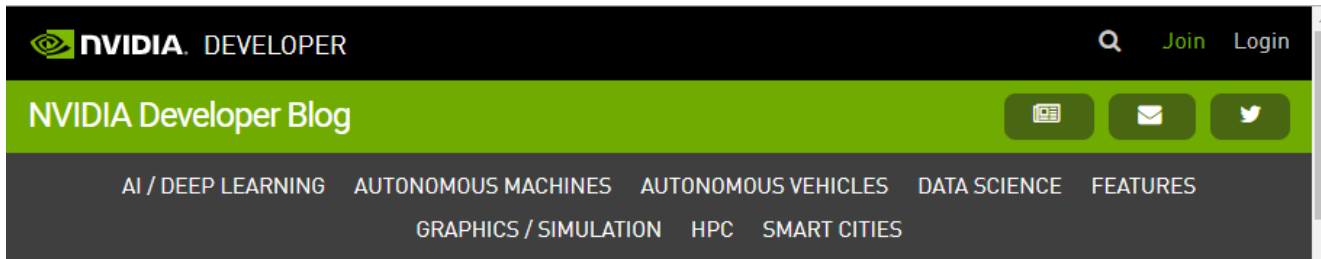
Discovered new speculative contention-based attacks

RISC-V + NVDLA SoC Platform



- Full-featured quad-core SoC with hardware DNN accelerator on Amazon FPGA cloud
 - Run Linux, YOLO v3 object detection

Open-source hardware for research



AI / DEEP LEARNING

NVDLA Deep Learning Inference Compiler is Now Open Source

By Rekha Mukund, Prashant Gaikwad and Mitch Harwell | September 11, 2019

Tags: compilers, Inference, Jetson, Machine Learning and AI, Object Detection

Designing new custom hardware accelerators for deep learning requires a balance of performance and efficiency with a narrow focus on the target application.

Two years ago, NVIDIA opened the source code for the NVIDIA Deep Learning Inference Compiler (NVDLA) to help advance the adoption of efficient deep learning hardware. The open-source release of NVDLA's optimized compiler is a key milestone in our point with the complete source for the NVDLA software and hardware design.

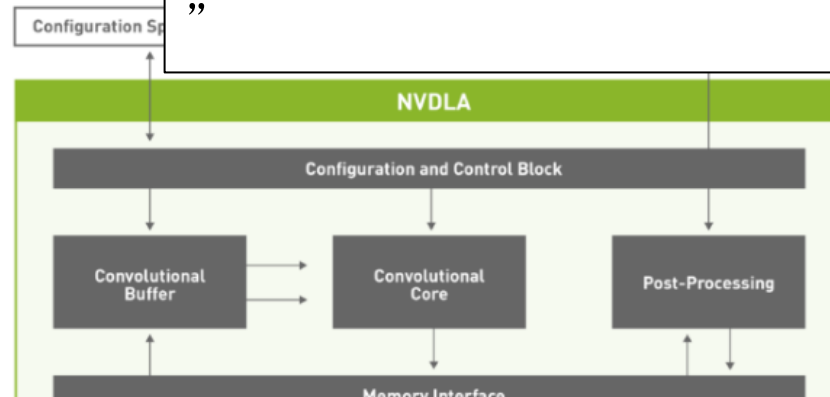
In this blog we'll explain the role that the NVIDIA Deep Learning Inference Compiler (NVDLA) software and hardware design play in the purpose-built hardware accelerator for deep learning inference.

“ One of the best ways to get started is to dive right in with object detection using YOLOv3 on NVDLA with RISC-V and FireSim in the cloud.

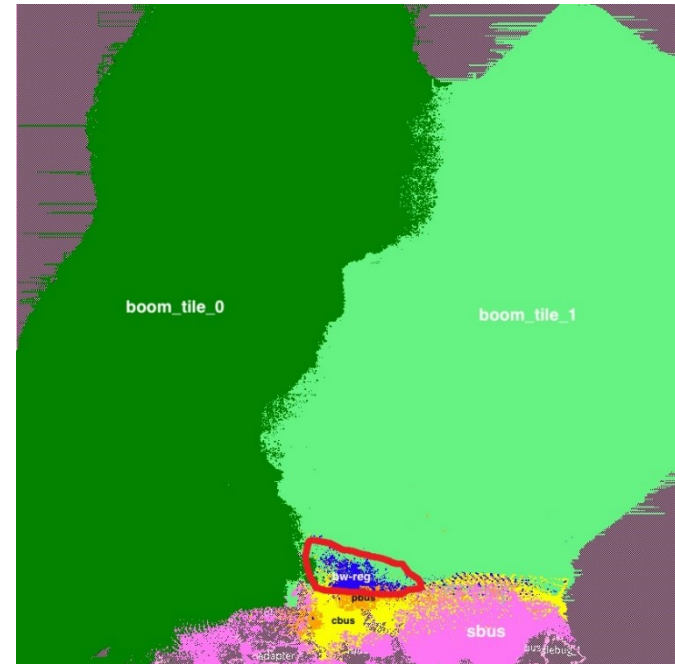
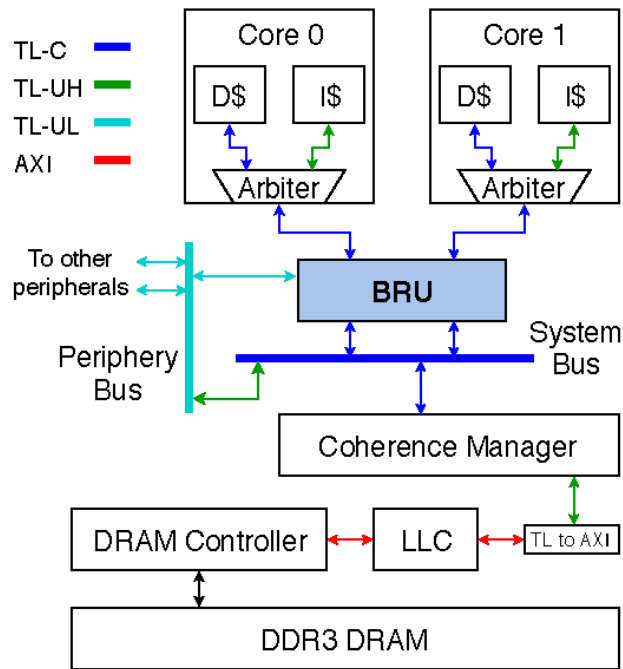
... git clone <https://github.com/CSL-KU/firesim-nvdl> ”

★ Star

108



Bandwidth Regulation Unit (BRU)



- Regulate per-core/group memory bandwidth
 - Drop-in addition to existing processor design

Real hardware-based resource control mechanism

Conclusion

- Intelligent cyber-physical systems need safe and secure computing infrastructure.
- Today's computing infrastructure is inadequate for safe and secure CPS
- Our research develops fundamental computing infrastructure technologies to enable safe and secure computing for intelligent CPS

Recent Publications

1. [DATE'21] Waqar Ali, Rodolfo Pellizzoni, Heechul Yun. Virtual Gang based Scheduling of Parallel Real-Time Tasks (to appear)
2. [RTSS'20] Shengzhong Liu, Shuochao Yao, Xinzhe Fu, Rohan Tabish, Simon Yu, Ayoosh Bansal, Heechul Yun, Lui Sha and Tarek Abdelzaher. On Removing Algorithmic Priority Inversion from Mission-critical Machine Inference Pipelines.
Best Paper Award
3. [ASHES'20] Jacob Fustos, Michael Garrett Bechtel, Heechul Yun. SpectreRewind: Leaking Secrets to Past Instructions.
4. [EMSOFT'20] Homa Aghilinasab, Waqar Ali, Heechul Yun, Rodolfo Pellizzoni. Dynamic Memory Bandwidth Allocation for Real-Time GPU-Based SoC Platforms.
5. [RTAS'20] Farzad Farshchi, Qijing Huang, and Heechul Yun. BRU: Bandwidth Regulation Unit for Real-Time Multicore Processors.
6. [DAC'19] Jacob Fustos, Farzad Farshchi, and Heechul Yun. SpectreGuard: An Efficient Data-centric Defense Mechanism against Spectre Attacks.
7. [ECRTS'19] Renato Mancuso, Heechul Yun, Isabelle Puaut. Impact of DM-LRU on WCET: a Static Analysis Approach.
8. [RTAS'19-2] Waqar Ali and Heechul Yun. RT-Gang: Real-Time Gang Scheduling Framework for Safety-Critical Systems.
9. [RTAS'19-1] Michael Bechtel and Heechul Yun. Denial-of-Service Attacks on Shared Cache in Multicore: Analysis and Prevention.
Outstanding Paper Award
10. [EMC2'19] Farzad Farshchi, Qijing Huang, and Heechul Yun. Integrating NVIDIA Deep Learning Accelerator (NVDLA) with RISC-V SoC on FireSim.
11. [RTCSA'18] Michael Bechtel, Elise McElhiney, Minje Kim, Heechul Yun. DeepPicar: A Low-cost Deep Neural Network-based Autonomous Car.
12. [ECRTS'18-2] Waqar Ali, Heechul Yun. Protecting Real-Time GPU Applications on Integrated CPU-GPU SoC Platforms.
13. [ECRTS'18-1] Farzad Farshchi, Prathap Kumar Valsan, Renato Mancuso, Heechul Yun. Deterministic Memory Abstraction and Supporting Multicore System Architecture.
14. [RTSJ'17] Prathap Valsan, Heechul Yun, Farzad Farshchi. Addressing Isolation Challenges of Non-blocking Caches for Multicore Real-Time Systems.
15. [TC'17] Heechul Yun, Waqar Ali, Santosh Gondi, Siddhartha Biswas. BWLOCK: A Dynamic Memory Access Control Framework for Soft Real-Time Applications on Multicore Platforms.
16. [RTCSA'16] Prasanth Vivekanandan, Gonzalo Garcia, Heechul Yun, Shawn Keshmiri. A Simplex Architecture for Intelligent and Safe Unmanned Aerial Vehicles.
Best Student Paper Nomination
17. [RTAS'16] Prathap Valsan, Heechul Yun, Farzad Farshchi . Taming Non-blocking Caches to Improve Isolation in Multicore Real-Time Systems.
Best Paper Award
18. [TC'16] Heechul Yun, Gang Yao, Rodolfo Pellizzoni, Marco Caccamo, and Lui Sha. Memory Bandwidth Management for Efficient Performance Isolation in Multi-core Platforms.
Editor's Pick of the Year

Thank You!

Acknowledgement:

This research has been supported in part by NSA Science of Security initiative contract #H98230-18-D-0009; NSF CNS 1718880, 1815959; NASA SBIR

