

Evaluation of High Performance Compute Architectures

Thorsten Wilmer, Visteon ADAS System Architect

TechAD Berlin 2019



Agenda

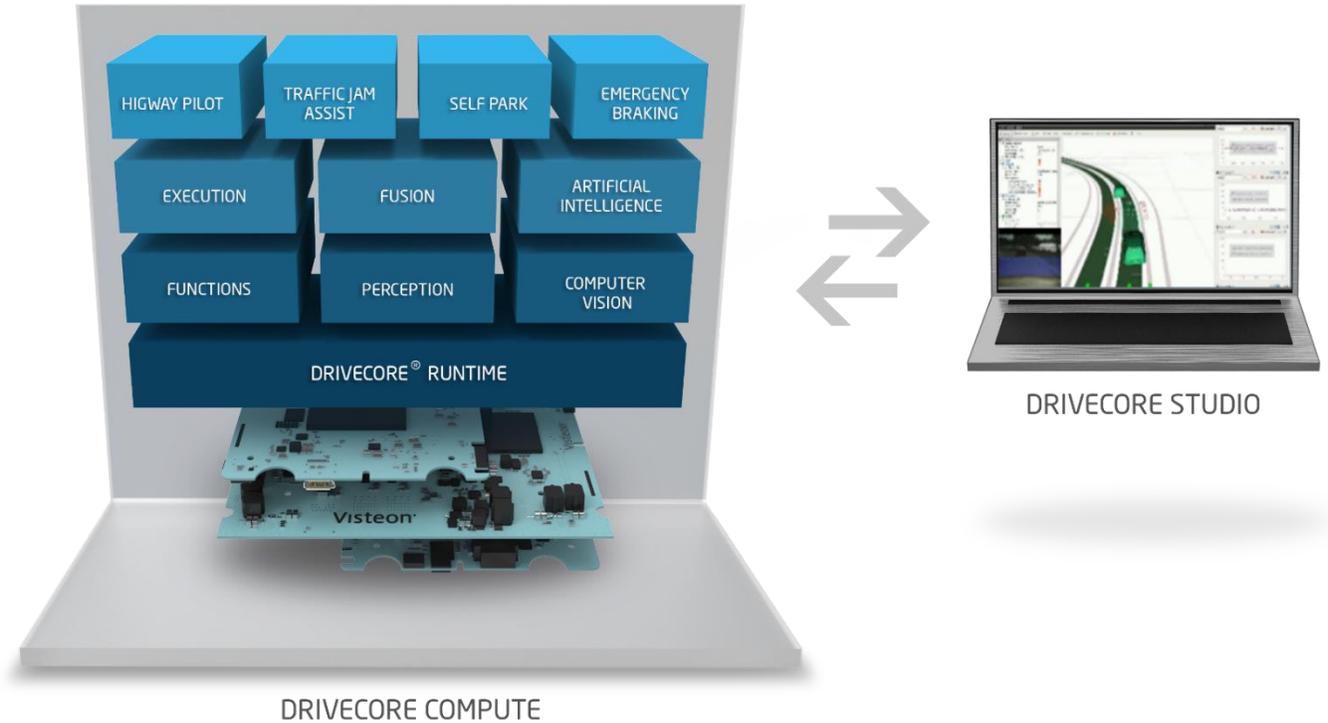
- Overview
- Comparisons
 - Costs
 - Performance Benmarking
 - Communication Latencies
- Conclusion

DriveCore™

Visteon's Autonomous Driving Platform

- DriveCore™ Studio: OPEN for Common Development
- DriveCore™ Compute: Modular Design with SoC Flexibility
- DriveCore™ Runtime: Safe Communication and Easy Integration in Vehicle

Visteon®



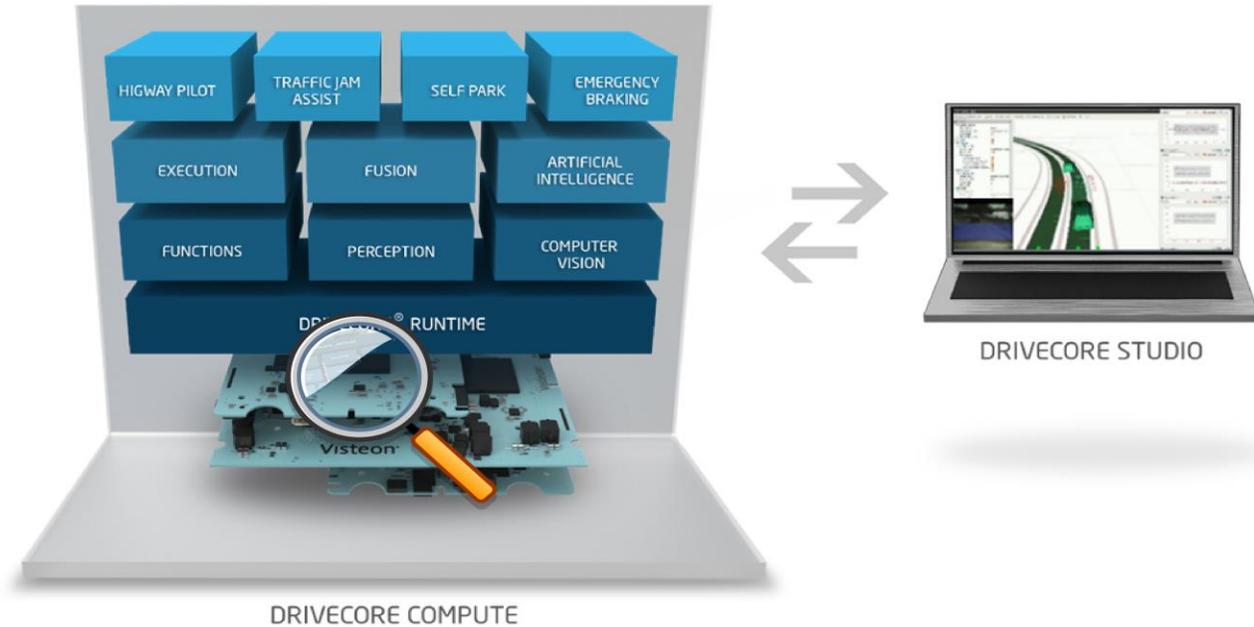
drivecore™
Compute | Runtime | Studio

DriveCore™

Visteon's Autonomous Driving Platform

- DriveCore™ Studio: OPEN for Common Development
- DriveCore™ Compute: Modular Design with SoC Flexibility
- DriveCore™ Runtime: Safe Communication and Easy Integration in Vehicle

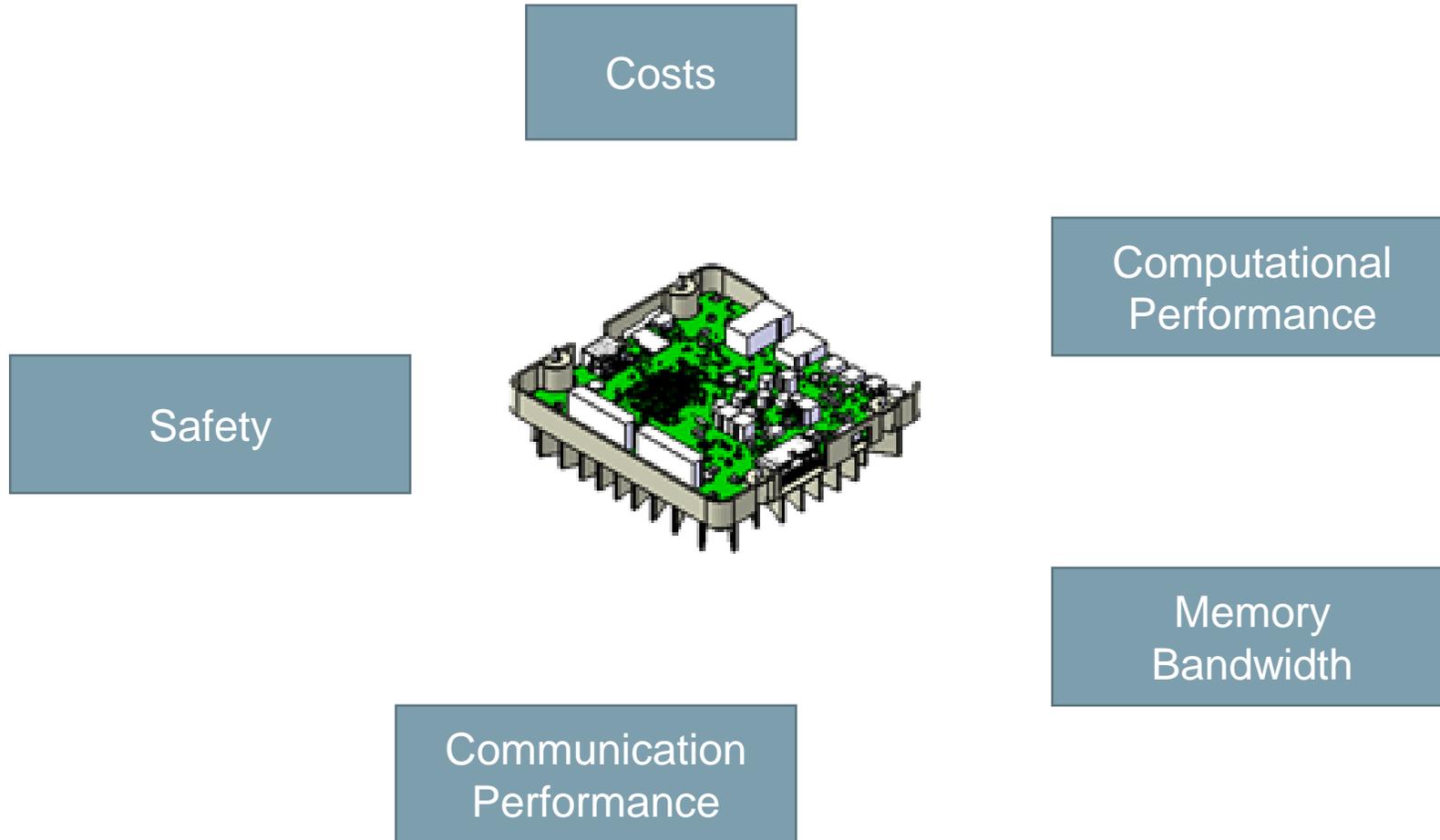
Visteon®



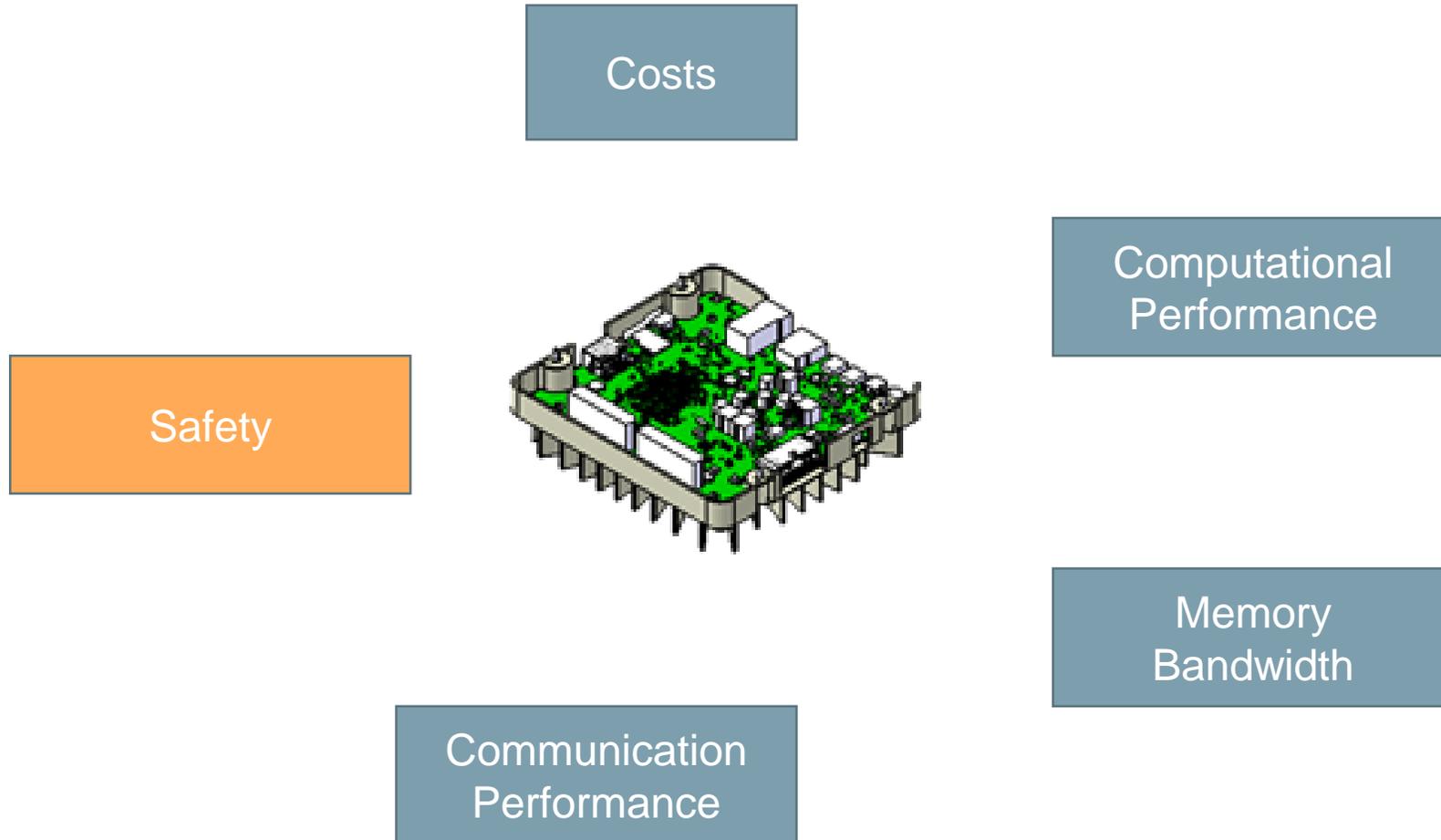
drivecore™

Compute | Runtime | Studio

Performance Optimization

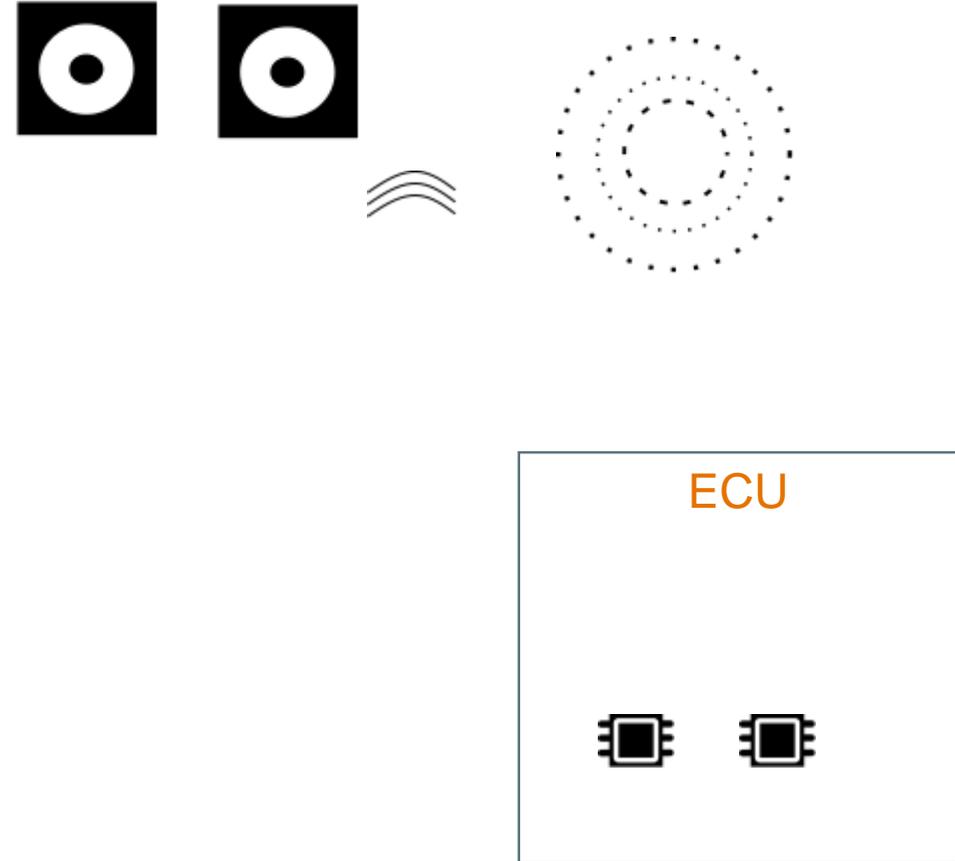


Performance Optimization



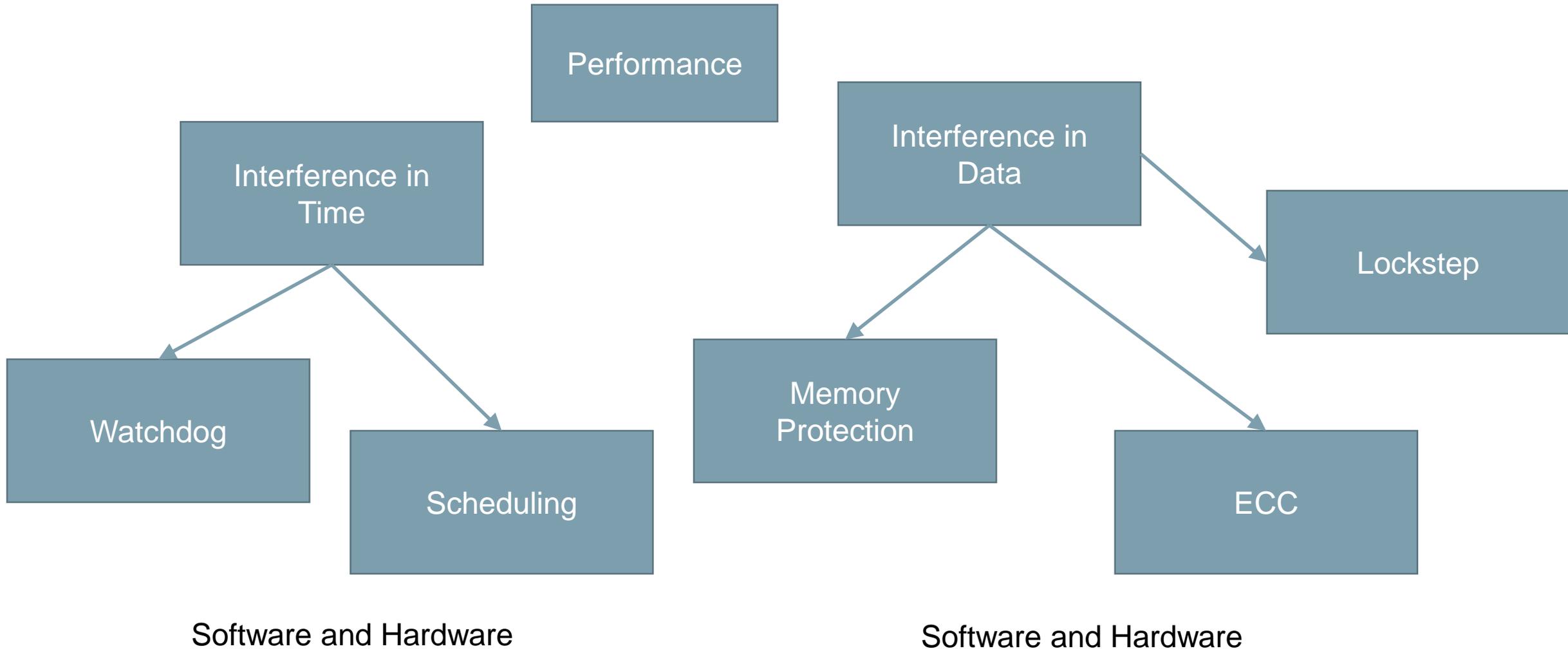
Redundancies On All Levels

- Two front cameras
- Lidar, radar and camera
- Redundant ECUs
- Redundant SoCs
- Redundant virtual partitions
- Redundant SWC
- ECC and repetition of messages
- Repetition of the execution
- Plausibility checks
- HW ECC on busses
- HW lock step
- Triple voting gates



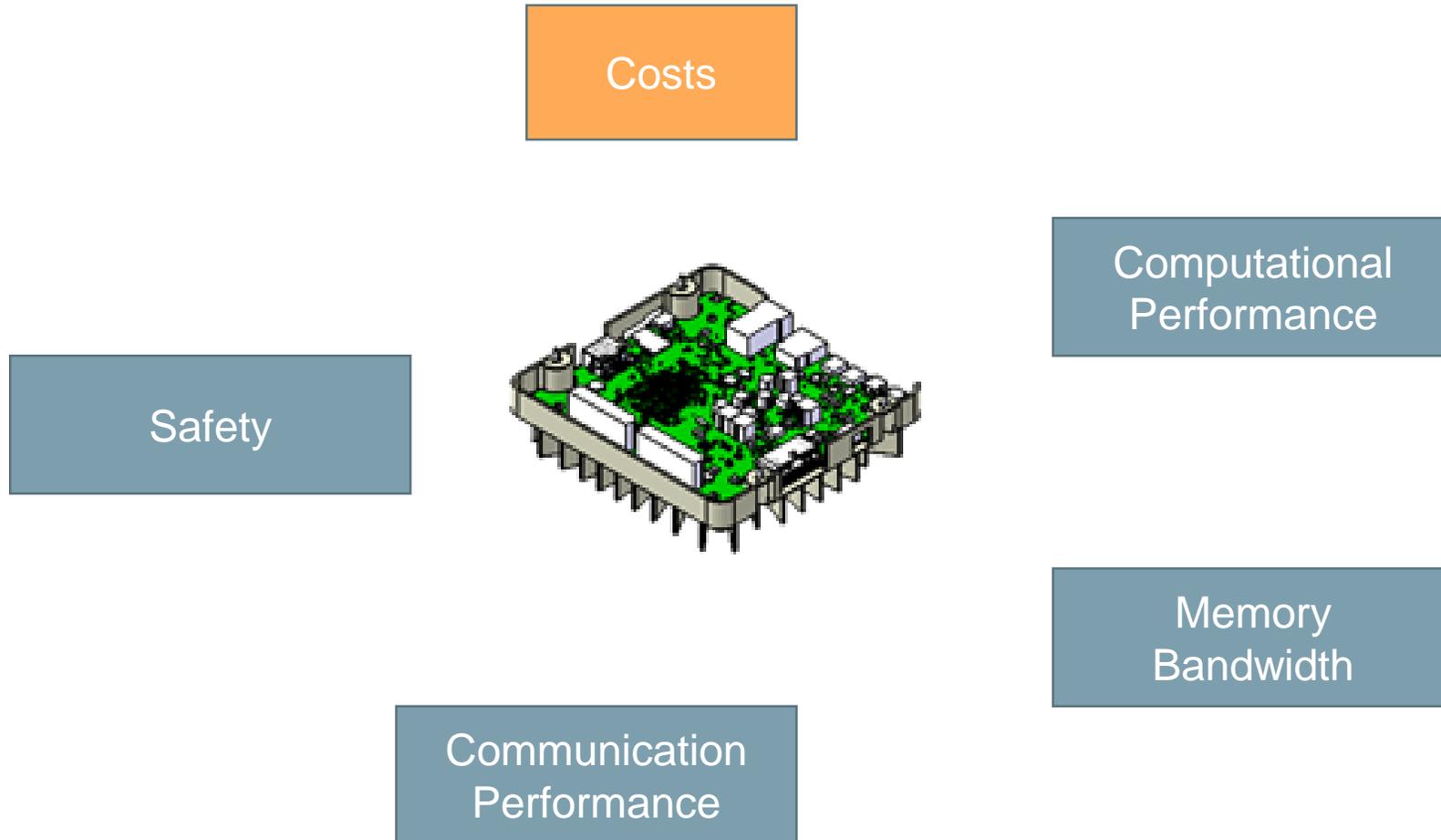
Resilient System Design

Functional Safety



Not every SoC supports all Safety features

Performance Optimization



Costs

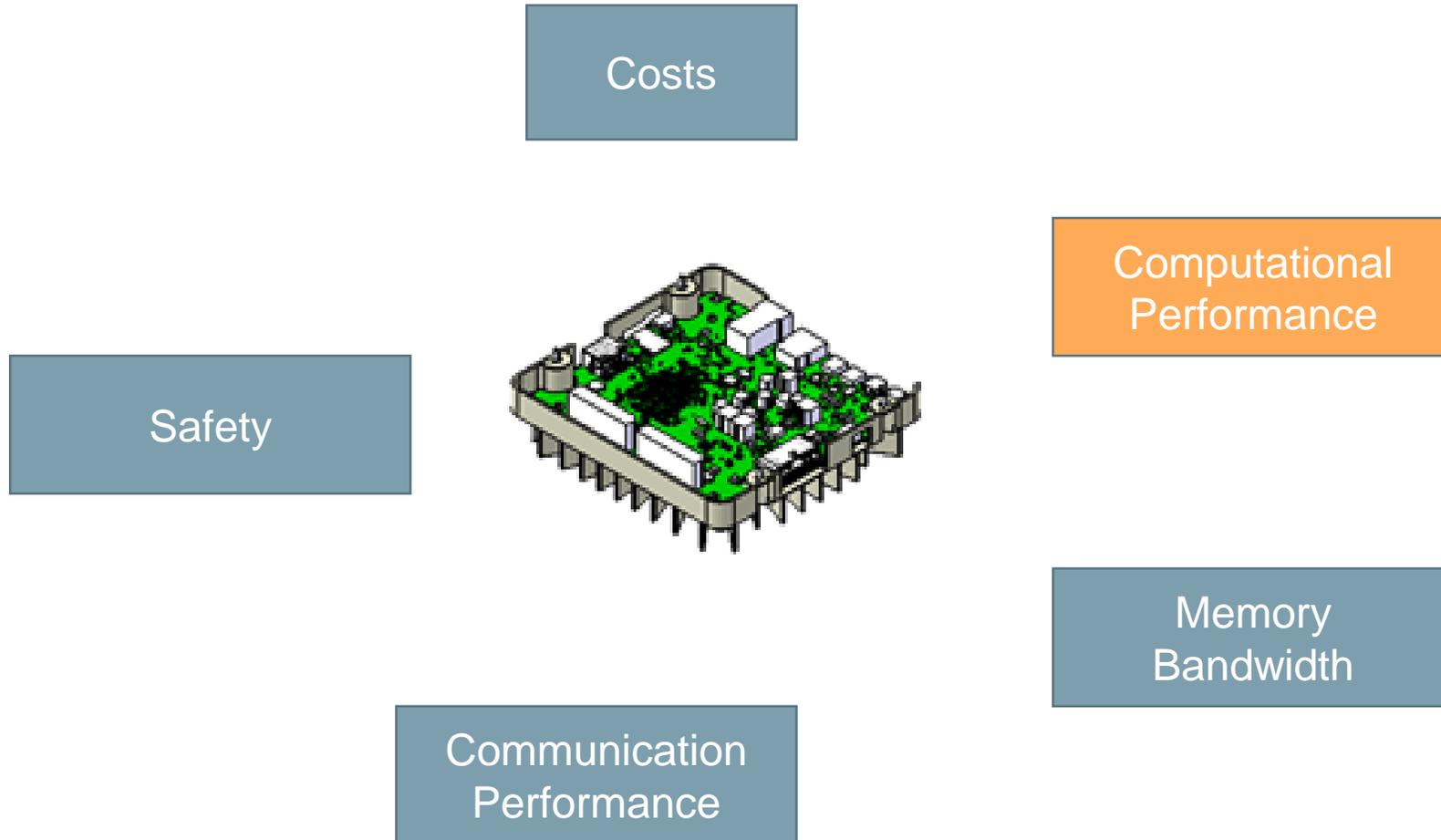
- Scalable SoC
- Scalable System
- Unused features cost
- Searching cheapest Solution which solves the Problem

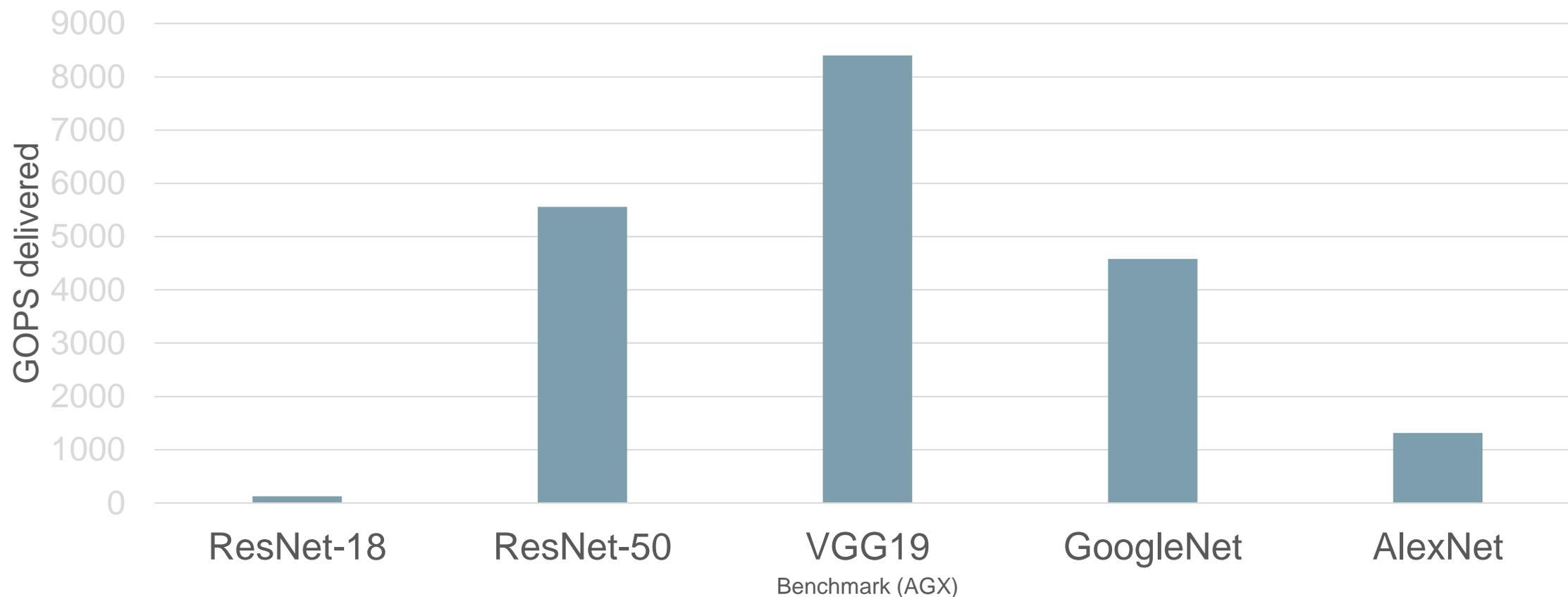


© Visteon 2019

Choice to select different SoC

Performance Optimization



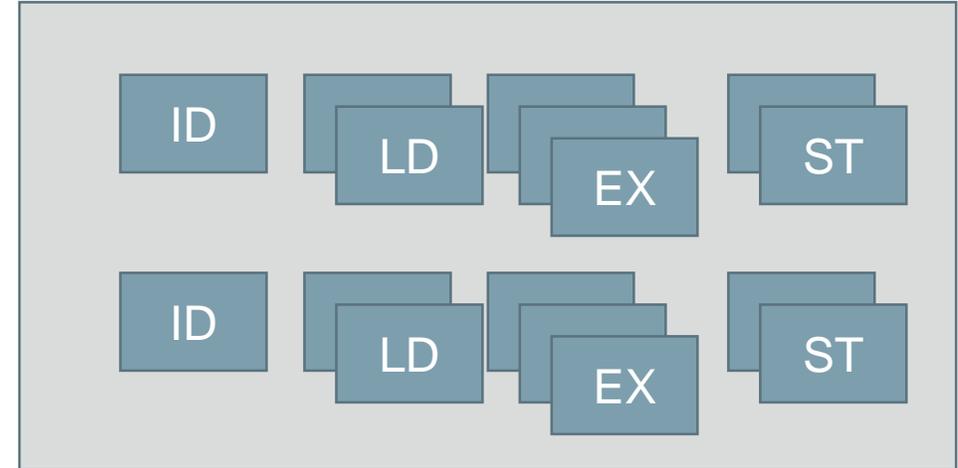


Source:
<https://developer.nvidia.com/embedded/jetson-agx-xavier-dl-inference-benchmarks>
<https://github.com/albanie/convnet-burden>

Different Benchmarks deliver different performance

Computational Performance

- Operations per Clock
- Pipeline Depth, penalty of Branches
- Register Set Size
- Parallelizm
 - Pipelining
 - Parallel Step Execution
 - Multi Threading
 - Multi Cores



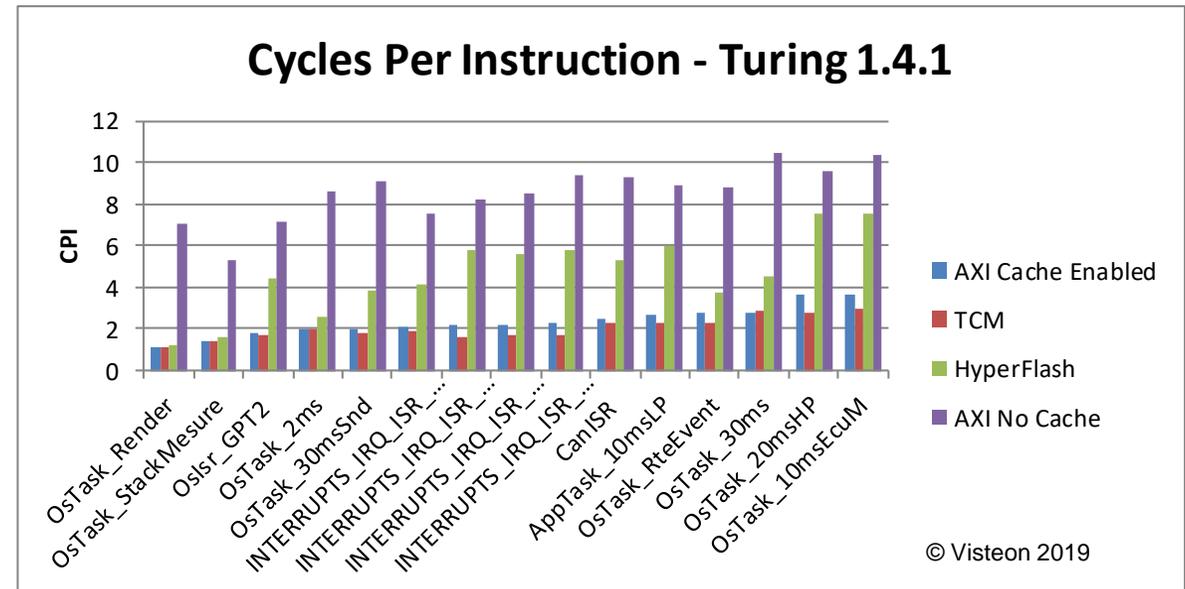
Computational Performance

Reality

- Not utilized Instruction slots
- Not enough parallelizm
- Cache misses (size, associativity)
- Memory Bandwidth limitation
- Branch Miss prediction
- Random Access Patterns
- Inefficient imlementaion

Theroy

- Clock Cycle each Execution unit is performing one usefull Operation

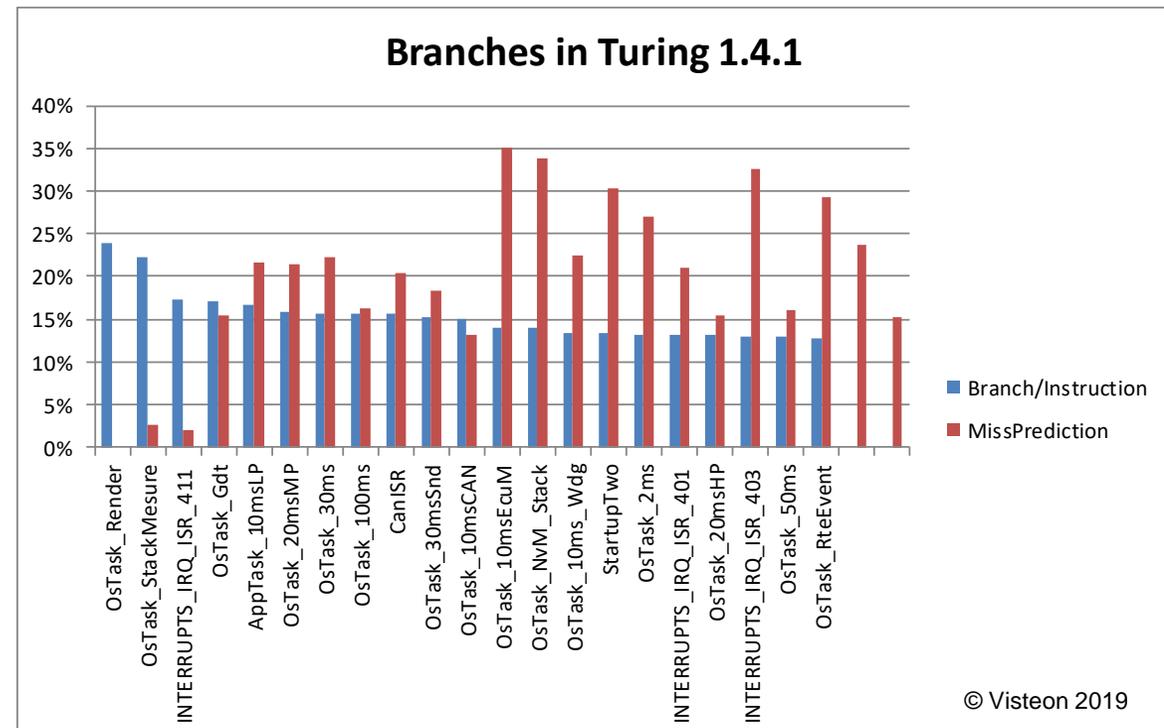


Reality

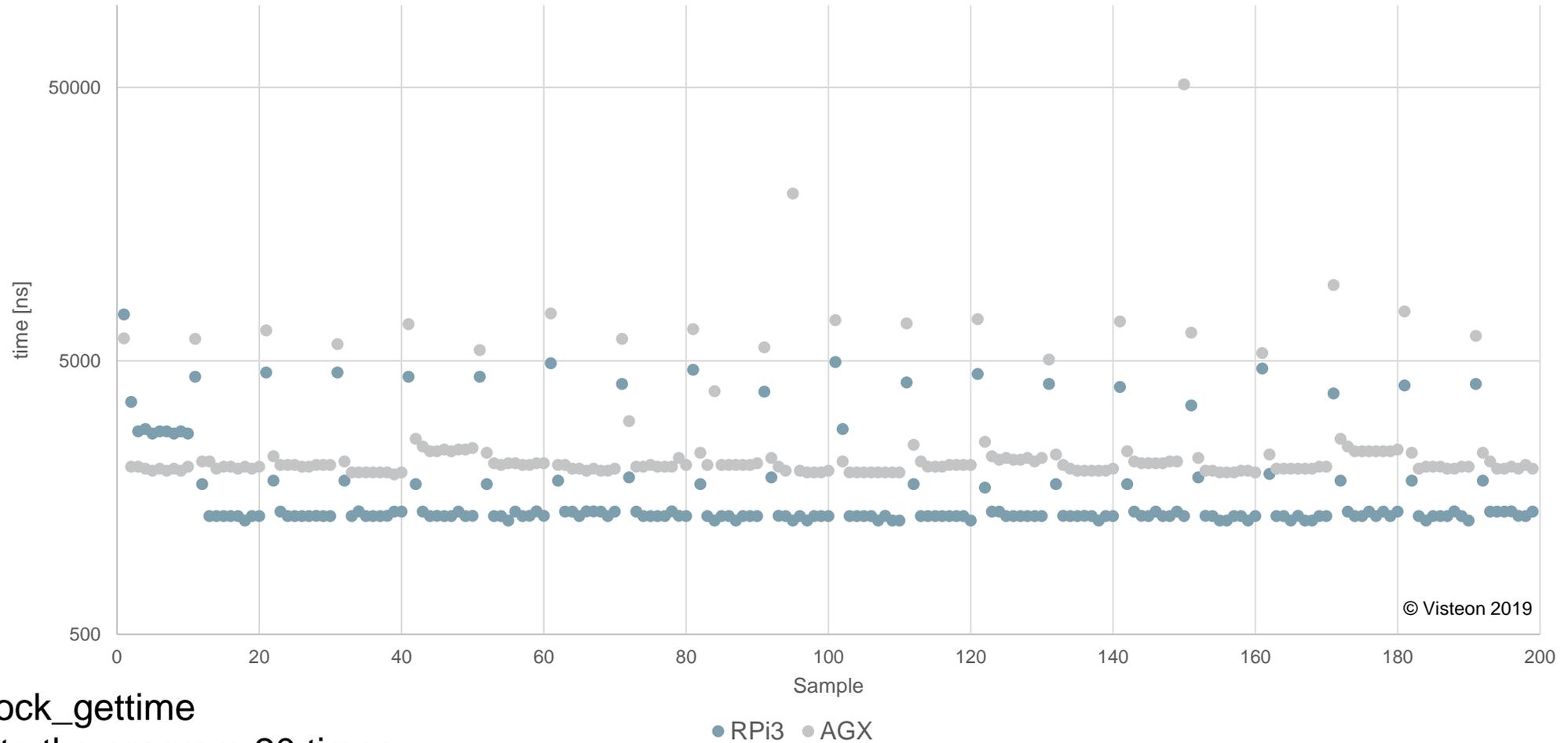
- Not utilized Instruction slots
- Not enough parallelizm
- Cache misses (size, associativity)
- Memory Bandwidth limitation
- Branch Miss prediction
- Random Access Patterns
- Inefficient imlementaion

Theroy

- Clock Cycle each Execution unit is performing one usefull Operation



Measurement overhead

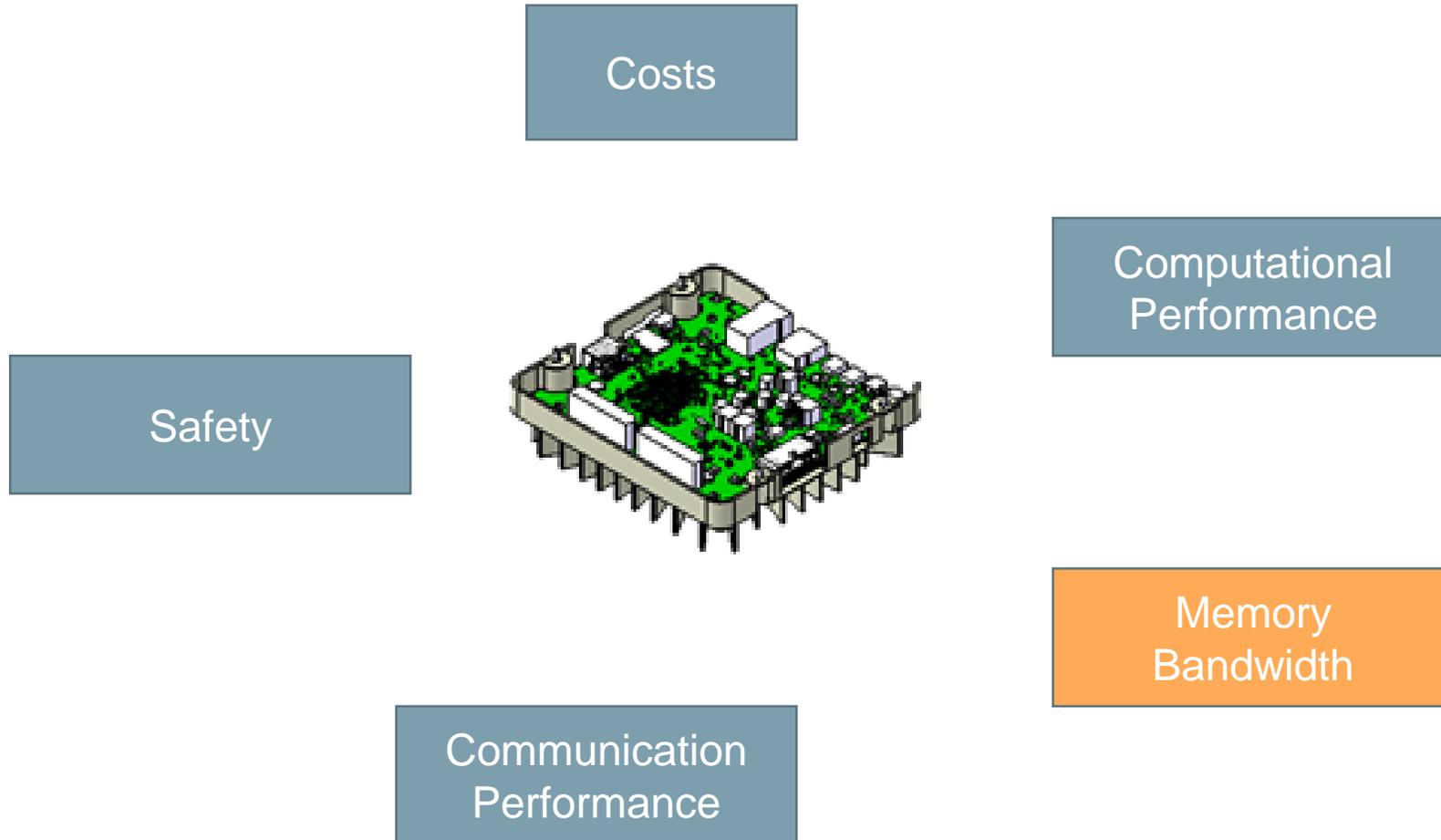


- 10x clock_gettime
- Execute the program 20 times

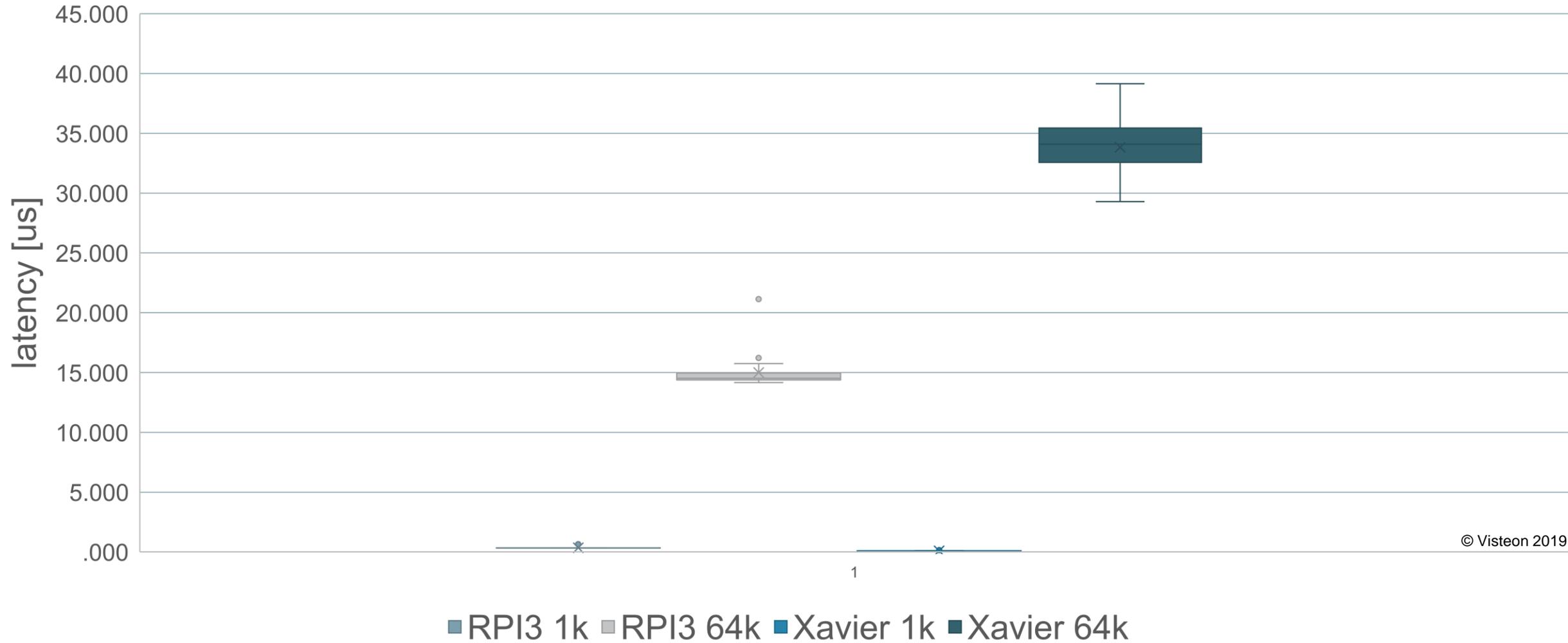
Cache effects

Use performance counters to measure

Performance Optimization



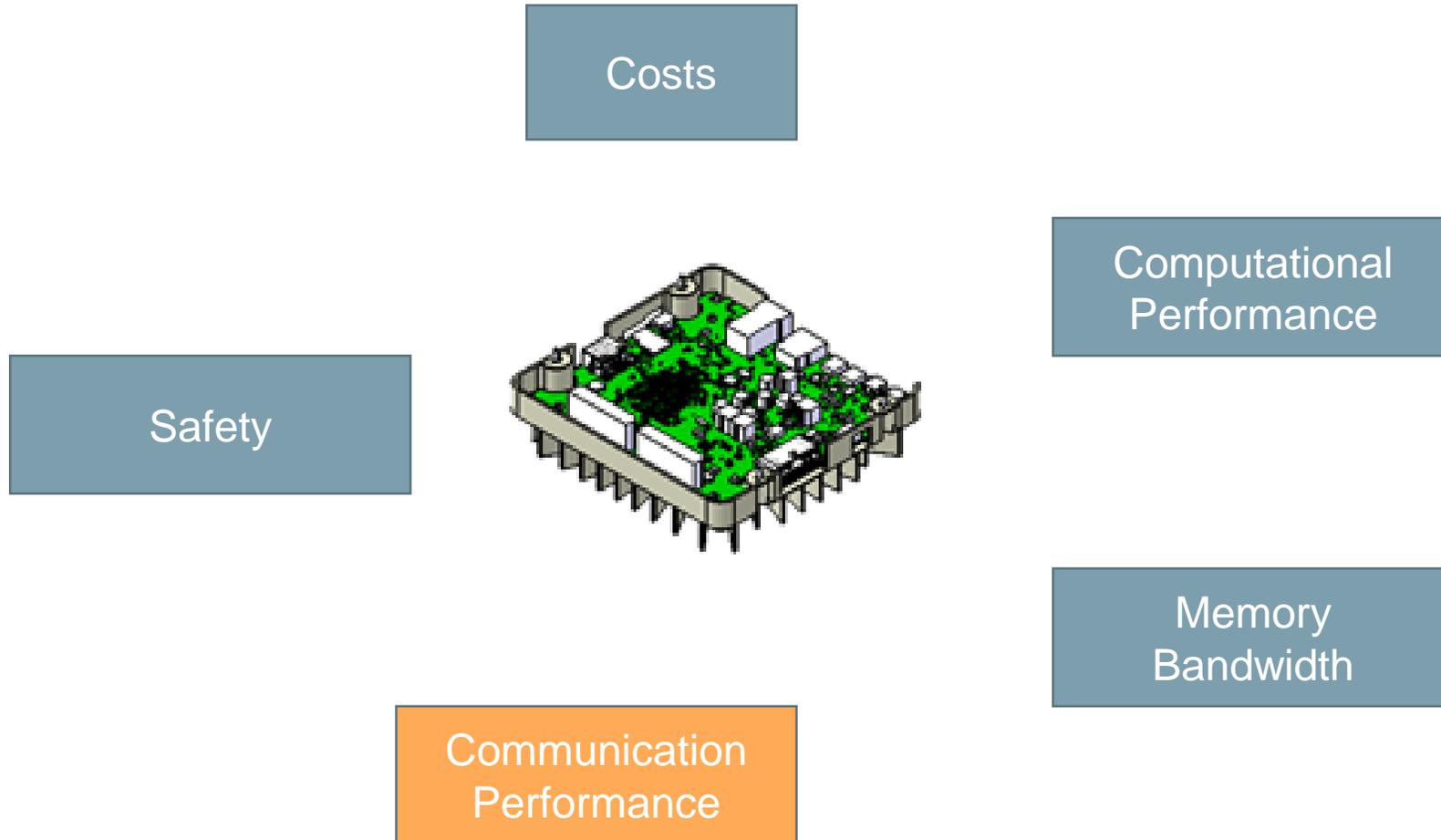
Memory Latency of Messages



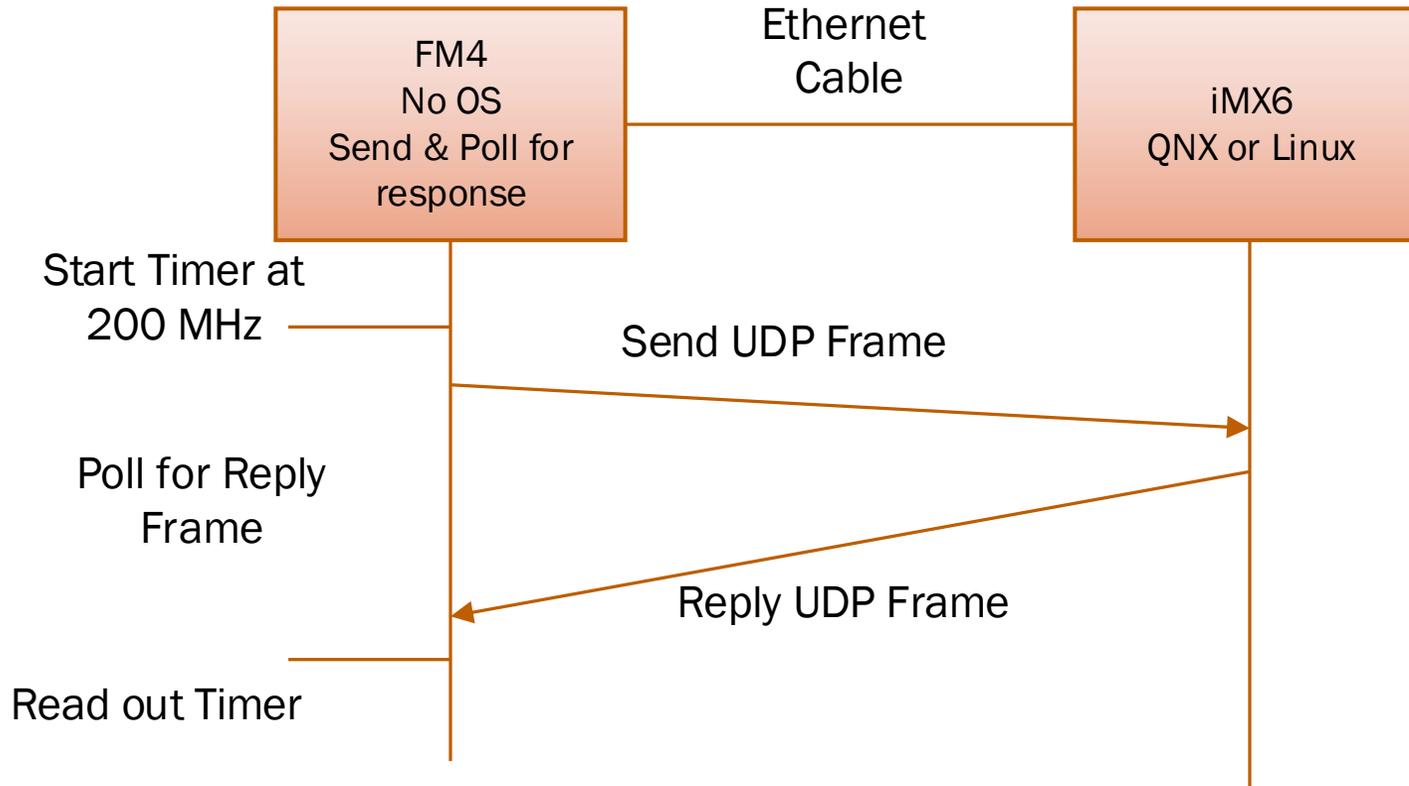
© Visteon 2019

Latency depends on the size of a message

Performance Optimization

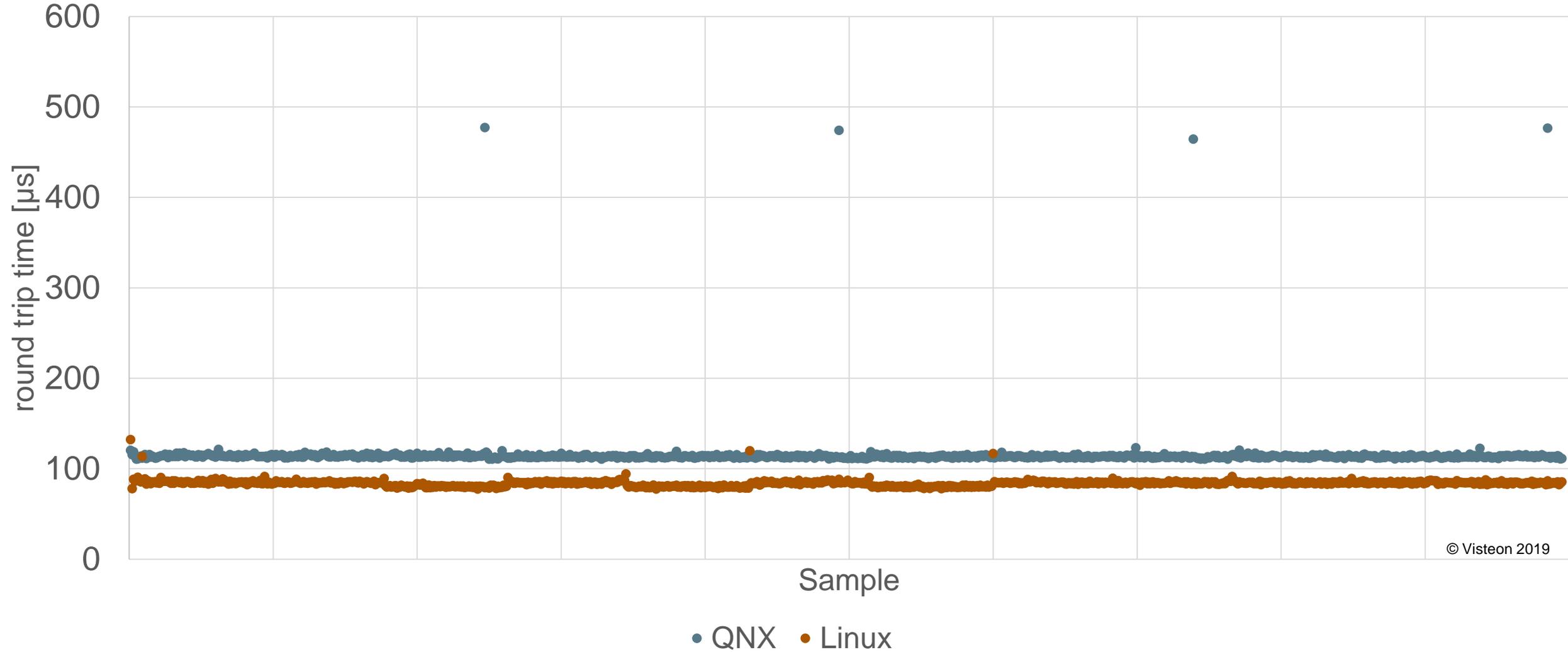


Measure UDP Round Trip Time



Prefer bare metal HW to measure timings

UDP Round Trip Time



Without optimization Linux has better UDP Latency

Conclusion

- It is not possible to find “the” best SoC
- Key is how Software uses Hardware
- Latency is a concern

Visteon®

