

On Monitoring and Resiliency for Machine-Learning-Based Autonomous Systems

Michael Paulitsch
2021-06-25

[IFIP WG10.4](#) Talk



intel®

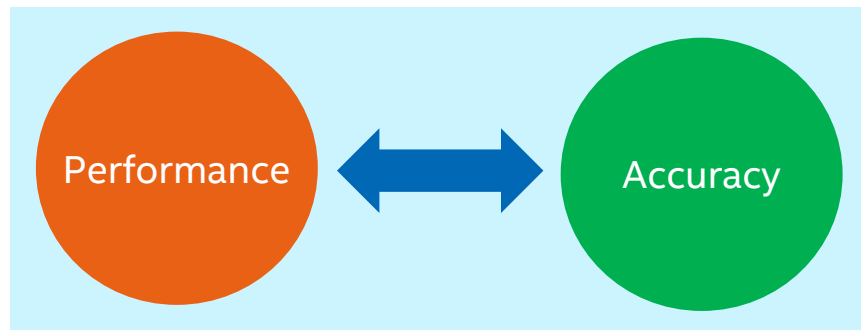
intel® labs

Overview of Two Related Topics

1. Resiliency: dependable AI/ML considering platform faults
2. Monitoring: safe perception

Accelerators

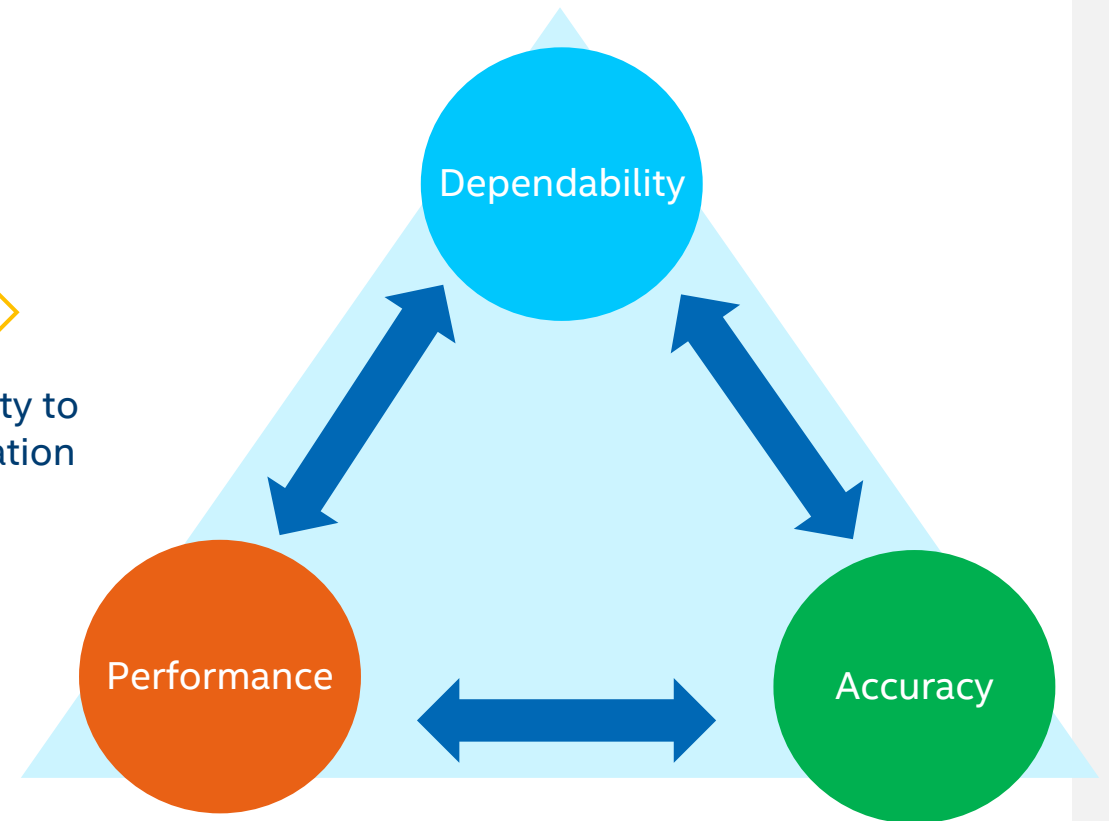
Dependability as Integral Part of Machine Learning System Optimization



State-of-the-Art System Optimization with Performance vs. Accuracy Trade-Off (e.g., float32 => bfloat16)

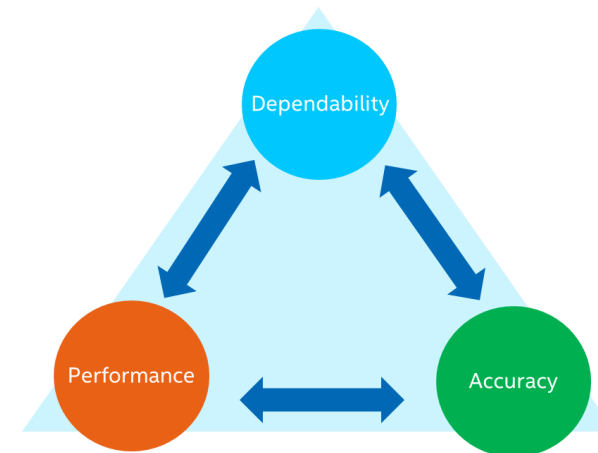


Add Dependability to System Optimization



Research Questions?

- **Dependability:**
 - What are the faults and failure modes to consider?
 - Is dependability really a problem in machine learning?
 - What is the target? 10^{e-9} range for hazards? Failures due to silent data errors?
 - Can we trust AI/ML at all? Diversity arguments? Safety monitoring?
- **Cost:** overhead runtime and development
 - Can it be nearly free?
 - How to automate?
- **Balance:** the quest for win-win situation
 - Is there an optimal balance?
 - Monitoring versus generalizability?



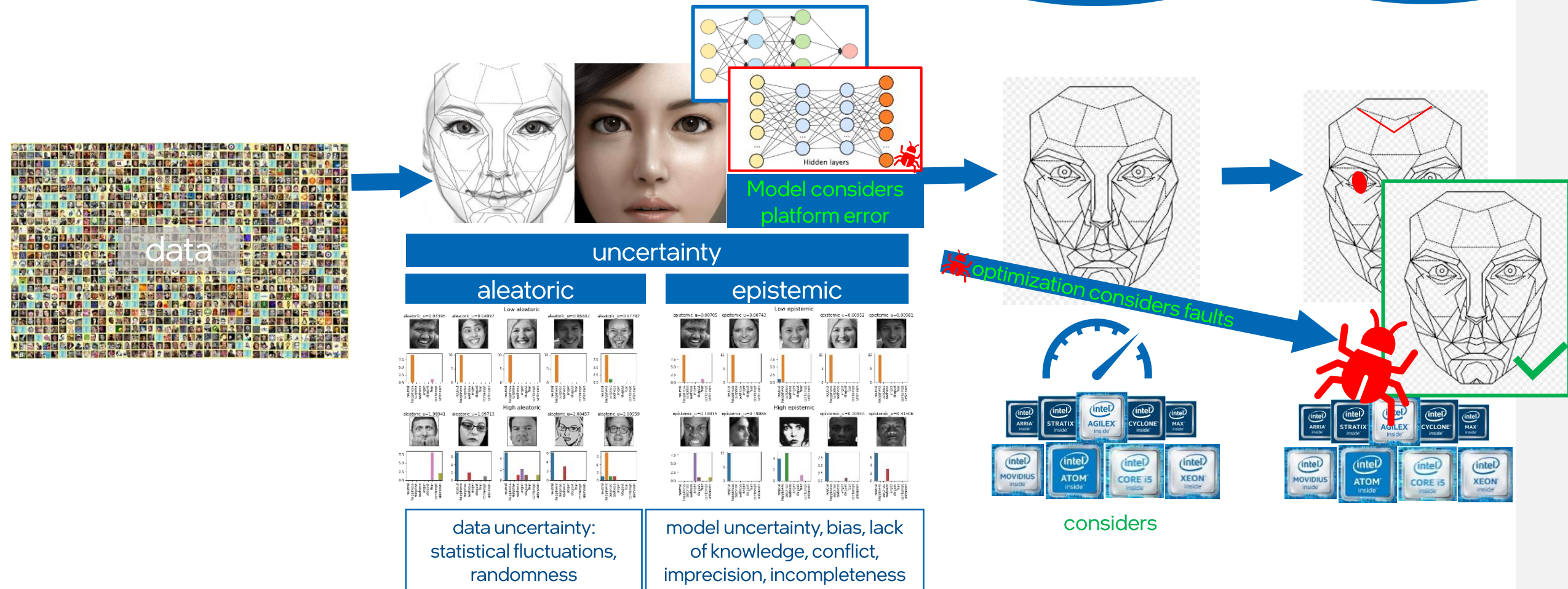
From Data to DNNs (CNNs) to Execution

data

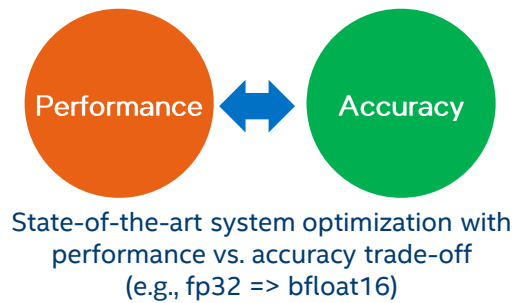
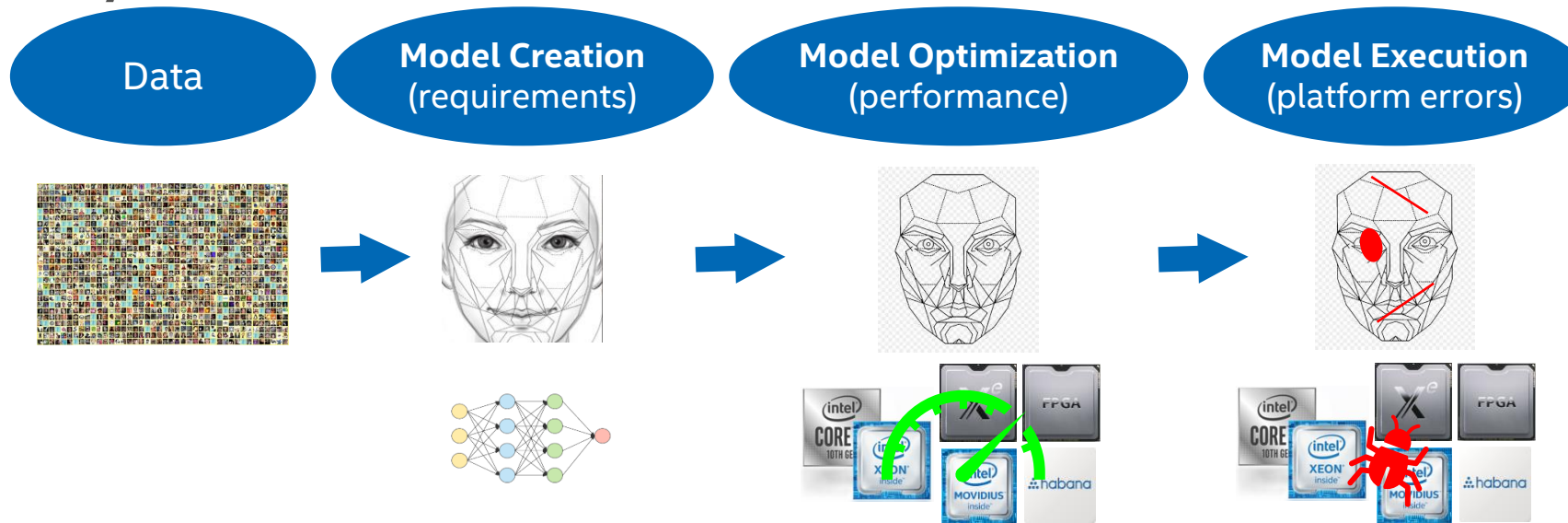
"near" ideal model
(requirements/ spec)

model
implementation
(performance optim.)

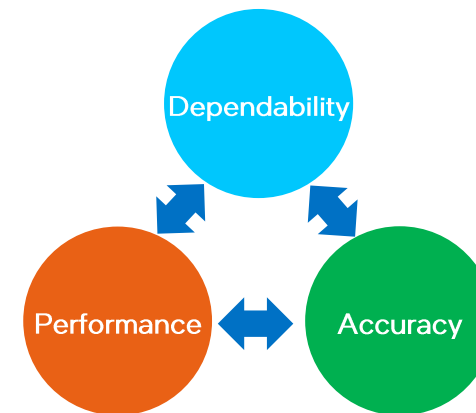
model impl.
(errors due to
platform)



Summary

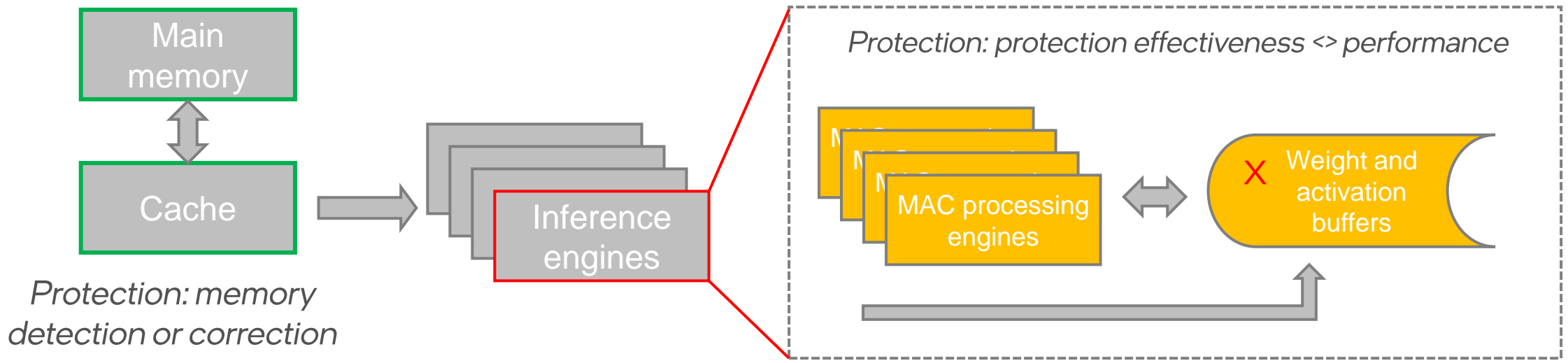
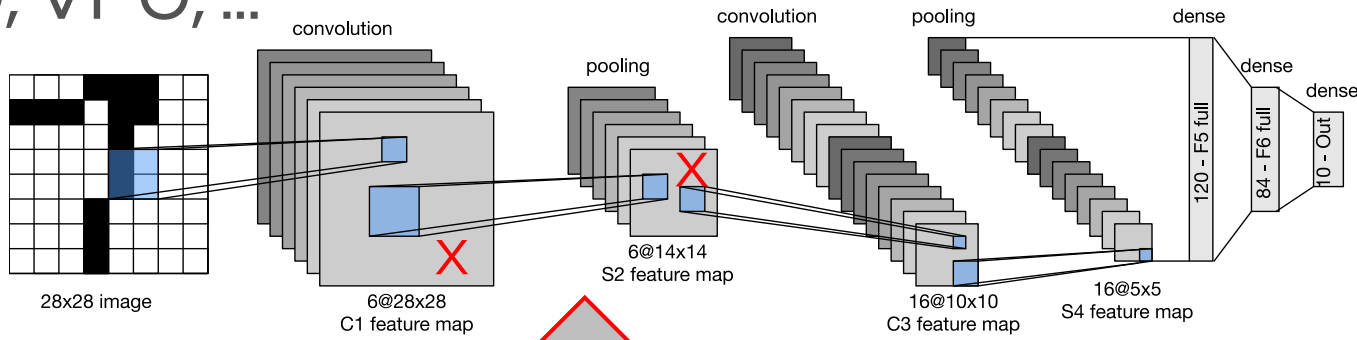


add dependability to system optimization



Common Elements in Intel CNN AI/ML accelerator HW

CPU extensions, GPGPU, VPU, ...



Protection: memory detection or correction

Accelerators in Practice Tiger Lake and Ice Lake

- Tiger Lake with Intel X^e accelerator
- Ice Lake: Advanced Vector Extensions (AVX)

3rd Gen Intel® Xeon® Scalable Platform

Feature	2nd Gen Intel® Xeon® Scalable Processor (Cascade Lake)	3rd Gen Intel® Xeon® Scalable Processors (Ice Lake)	Notes
Cores per Socket	4-28	8-40	New Sunny Cove architecture
L1/L2/L3 cache per core	32KB/1MB/1.375MB	48KB/1.25MB/1.5MB	Larger caches to enable fast access to data
Memory Channels and DIMM Speed	6 Up to 2933	8 Up to 3200	Huge boost in memory bandwidth & support for Intel® Optane™ PMem 200
Processor Interconnect: UPI links, speed	2 or 3, 10.4 GT/s	2 or 3, 11.2 GT/s	Improved bandwidth between processors
PCIe lanes per socket	PCIe 3.0, 48 Lanes (x16, x8, x4)	PCIe 4.0, 64 lanes (x16, x8, x4)	2x bandwidth and more PCIe lanes to support new Gen 4 SSD, Ethernet and other adjacencies
Workload Acceleration Instructions	AVX-512 VNNI DDIO	AVX-512, VNNI, DDIO vAES, VPCLMULQDQ, VPMADD52, VBMI, PFR, Crypto, SHA extensions, TME, SGX	Enable new capabilities and speedup performance
Platform Adjacencies		Intel® Optane™ PMem 200 series, Intel® Optane™ P5800X SSD, Intel DC P5510 SSD, Intel E810-C ethernet	

Designed to Move Faster, Store More, Process Everything

Performance made flexible.

37

source: 3rd Gen Intel Xeon Processor on [intc.com](https://www.intel.com)

Introducing 11th Gen Intel® Core™ Processor

New Willow Cove Cores

Up to 4 Cores / 8 Threads
Up to 4.8GHz

New Converged Chassis Fabric

High Bandwidth / Low Latency
IP and Core Scalable

New Memory Controller

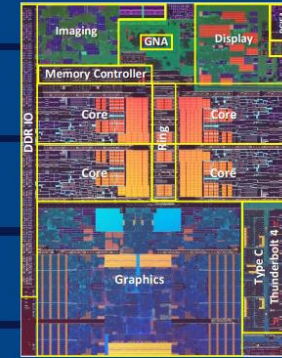
LP4/x-4266 4x32b up to 32GB
DDR4-3200 2x64b up to 64GB

1st Integrated Thunderbolt™ 4

Full 4x DP/USB/PCIe mux on-die
Up to 40Gbps bi-directional per port

1st Integrated PCIe Gen 4 (CPU)

Low Latency, High Bandwidth
SSD or Discrete Graphics Direct CPU Attach



New Iris® X^e Graphics

Up to 96EU – Up to 2x Higher Performance
Intel® Deep Learning Boost: DP4A for AI

New 2x MEDIA Encoders

Up to 4K60 10b 4:4:4
Up to 8K30 10b 4:2:0

New 4 x Display Pipes

Up to 1 x 8K60 or 4 x 4K60
DP1.4 HBR3, BT.2020

New Image Processing Unit (IPU6)

Video up to 4K90 resolutions (initially 4K30)
Still image up to 42 megapixels (initially 27MP)

New GNA 2.0

Enhanced Power Management
Autonomous DVFS

For more complete information about performance and benchmark results, visit www.intel.com/11thgen (configuration details in section 3).

intel

Low Power, High Performance Intel® Iris® X^e Graphics

New X^e-LP microarchitecture

Up to 96 EUs
Up to 1.35 GHz

New high-efficiency thread control with software score boarding

New 8-wide vector units with support for Intel® DL Boost: DP4a

New L1 data cache Up to 3.8 MB L3



2X bandwidth to the memory fabric

Up to 48 texels/clock
Up to 24 pixels/clock

End-to-End Compression

Variable Rate Shading

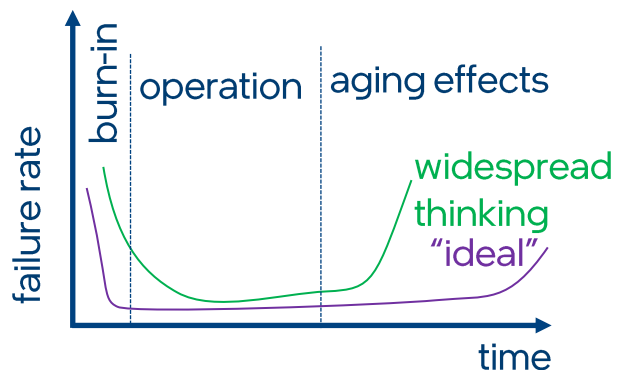
For more complete information about performance and benchmark results, visit www.intel.com/11thgen (configuration details in section 4).

intel

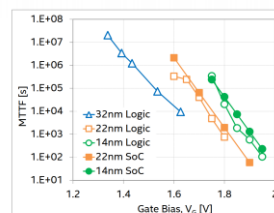
source: 11th Gen Intel Core Processor on [intc.com](https://www.intel.com)

Fault / Error / Failures

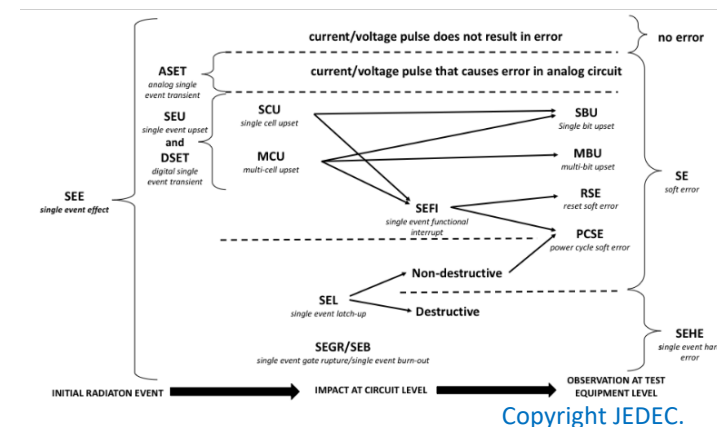
Hardware Faults – Error - Failures



Core: NMOS Dielectric Reliability



- SoC: same reliability as corresponding Logic node
- Strong reliability gains on both Tri-gate technologies



PERMANENT

INTERMITTENT

TRANSIENT

STUCK-AT ERROR

OPEN CIRCUIT ERROR

BRIDGING ERROR

SINGLE EVENT HARD ERROR (LATCHUP)

COMPLEX ERRORS

TRANSIENT OR HARD ERROR

SINGLE EVENT TRANSIENT, UPSET

MULTIPLE CELL UPSET

MULTIPLE BIT UPSET

production (reliability, DPM)

design errors

aging

external - radiation

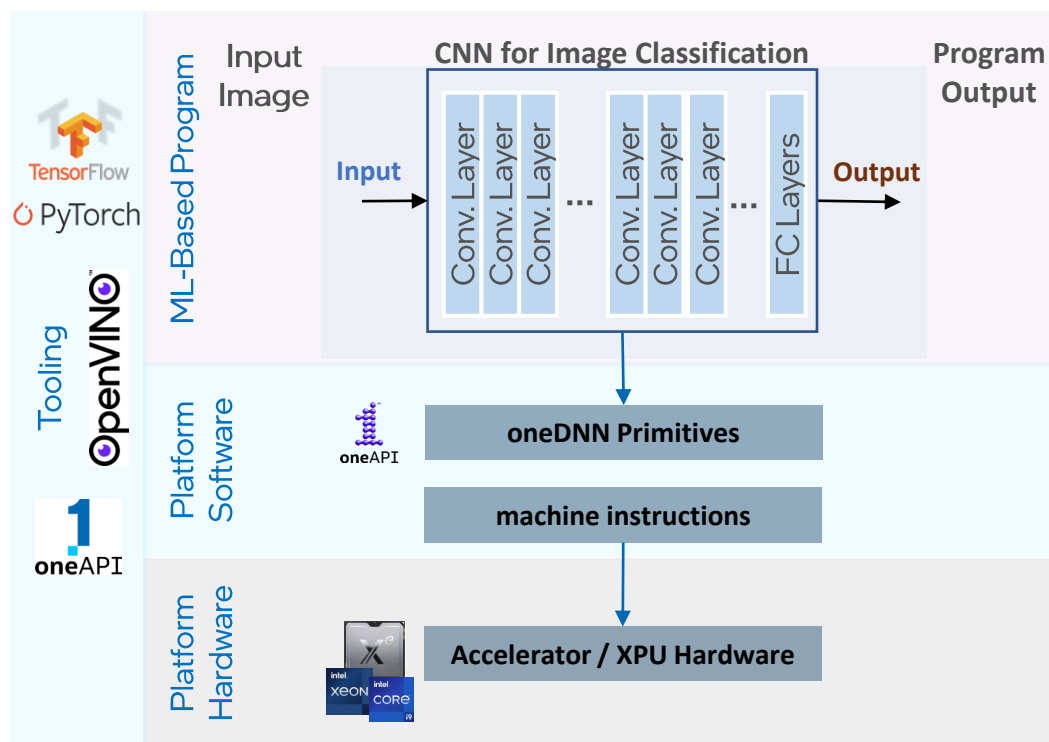
Dependable AI / ML - Resiliency



High-Performance Computing



Safety-critical Applications



System View

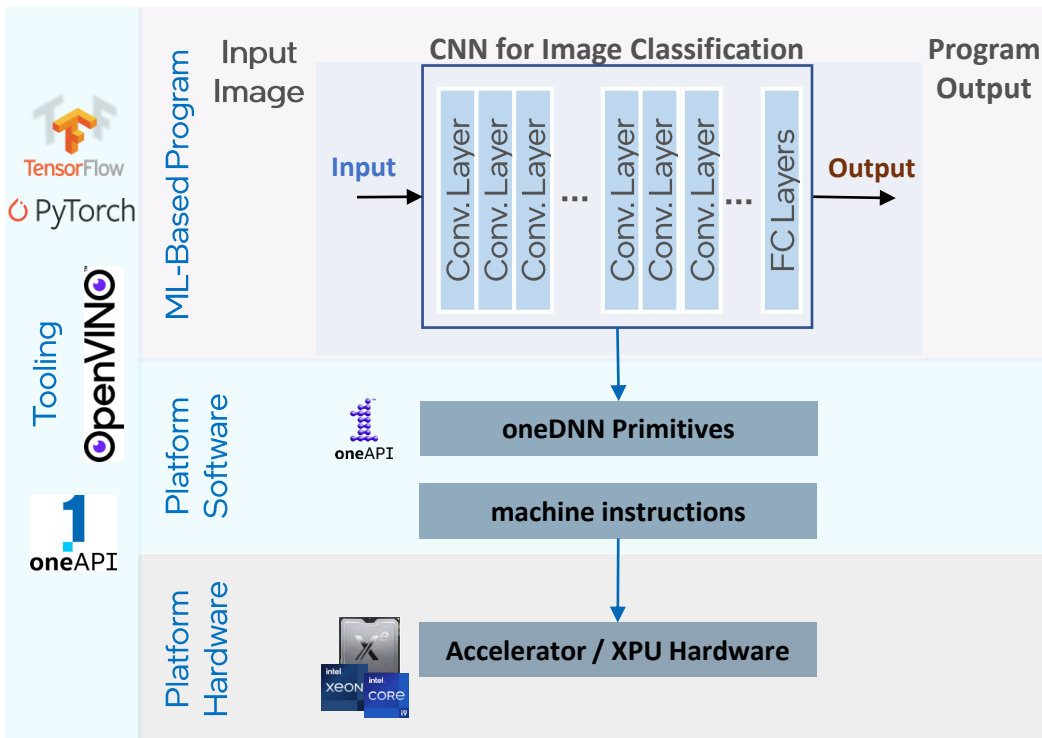
Dependable AI / ML - Resiliency



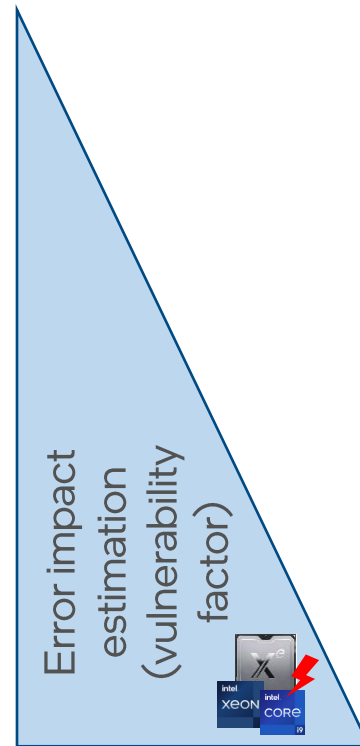
High-Performance Computing



Safety-critical Applications



System View



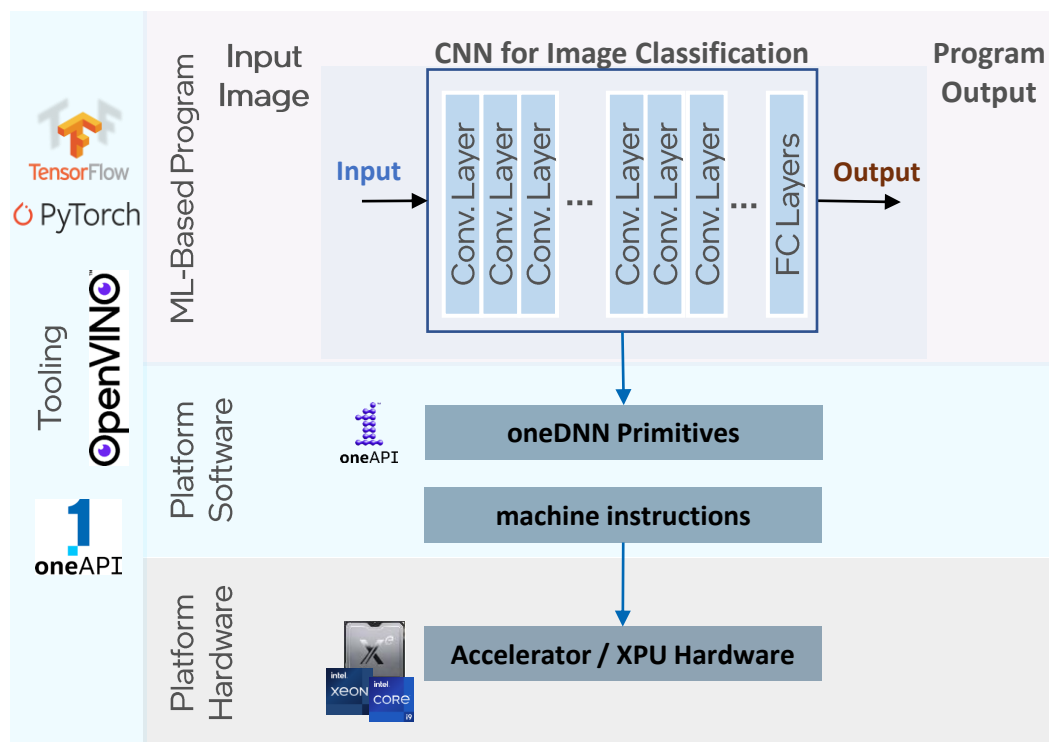
Dependability Threats

Dependable AI / ML - Resiliency

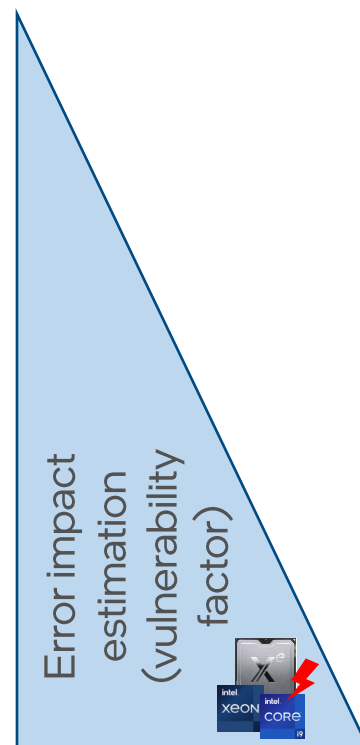


High-Performance Computing Safety-critical Applications

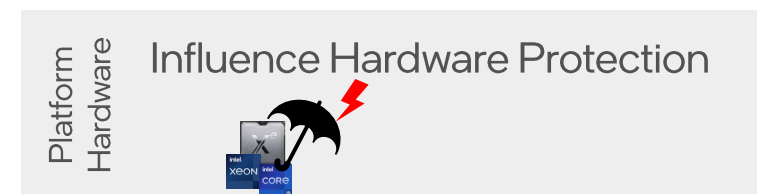
dependability / reliability target



System View



Dependability Threats



Error Impact Estimation Influences Hardware Protection

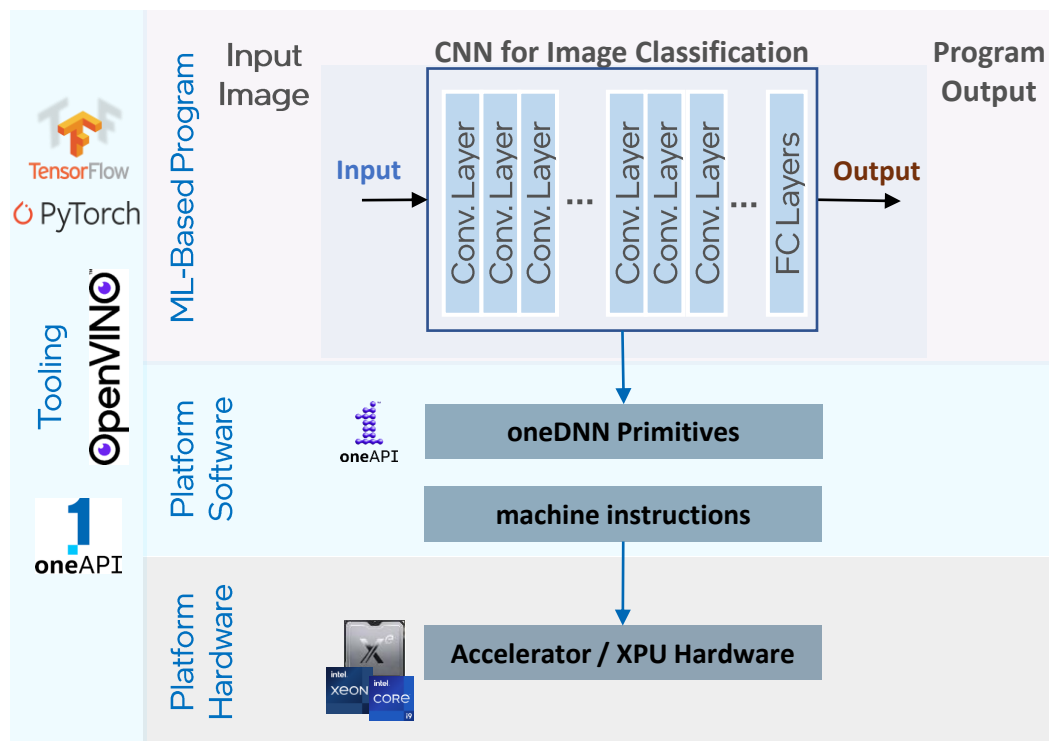
Dependable AI / ML - Resiliency



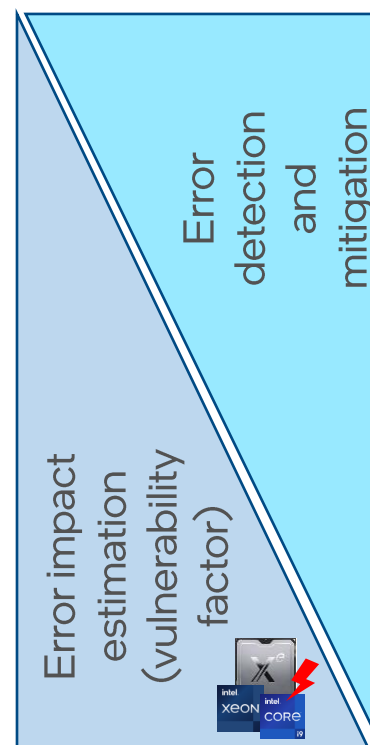
High-Performance Computing Safety-critical Applications

dependability / reliability target

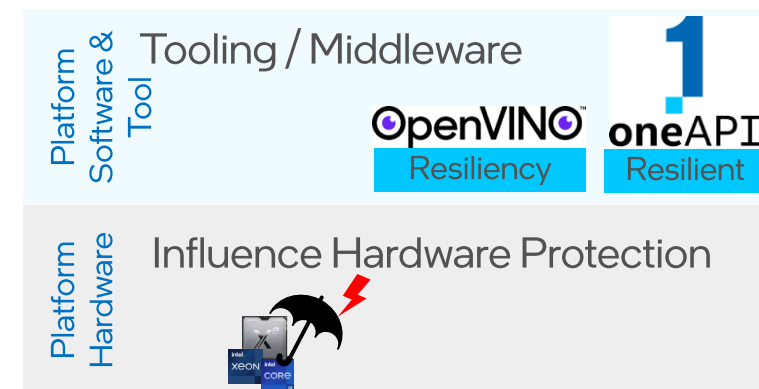
Dependability Means



System View



Dependability Threats



Details on Dependability Means

Research Question: Mitigation Mechanisms at SW & Tooling Level

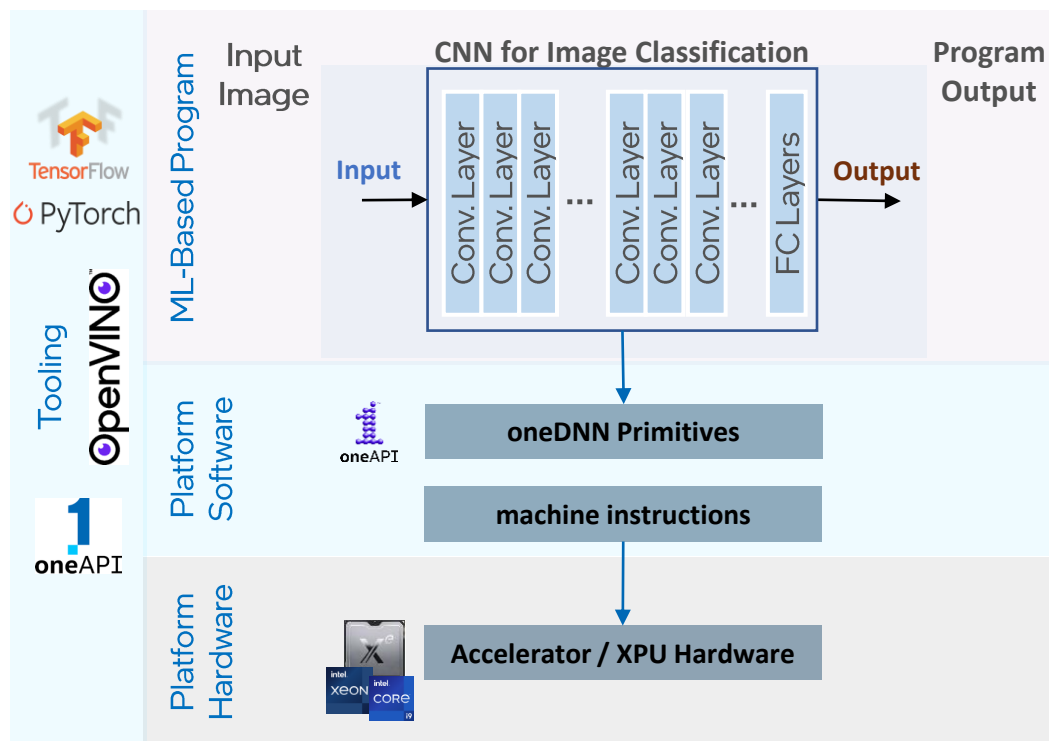
Dependable AI / ML - Resiliency



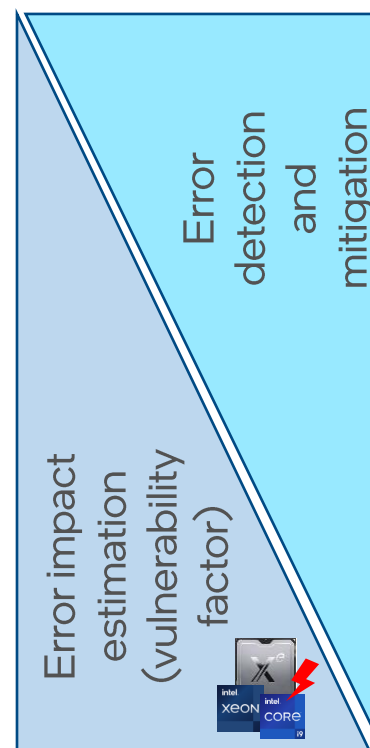
High-Performance Computing Safety-critical Applications

dependability / reliability target

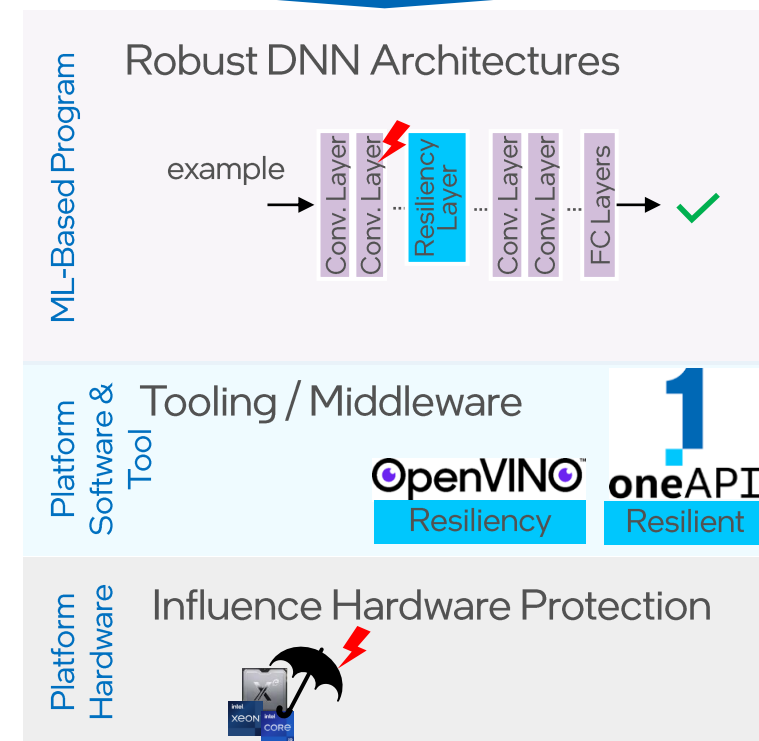
Dependability Means



System View



Dependability Threats



Details on Dependability Means

Research Question: Resilient Networks

Range Supervision

Example for a simple mitigation approach (joint work with Univ. of British Columbia – Prof. Karthik Pattabiraman)

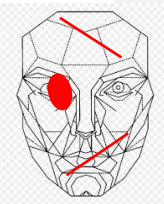
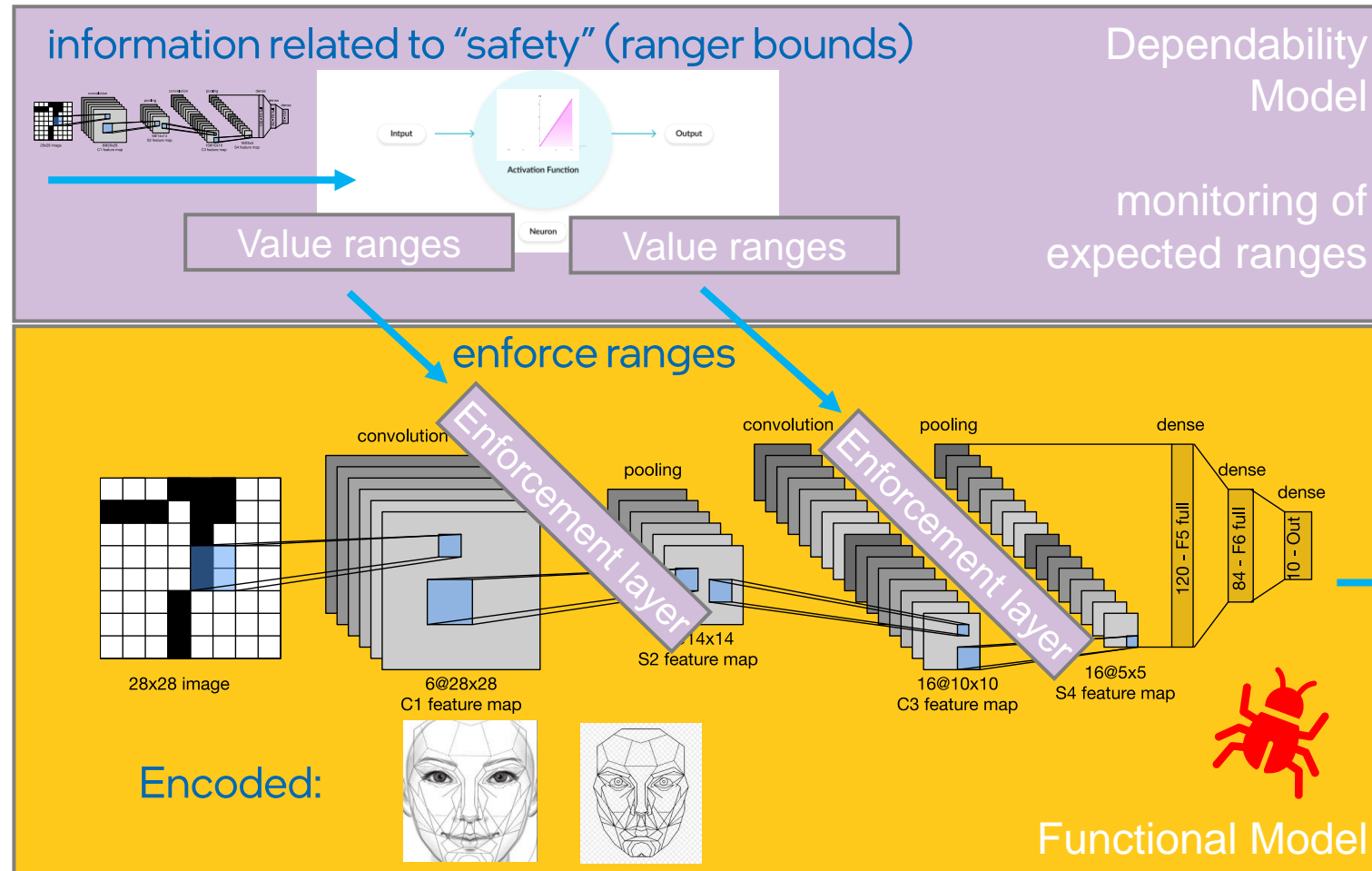
"Independent" Additional Info?

Data

Used to extract activation value ranges

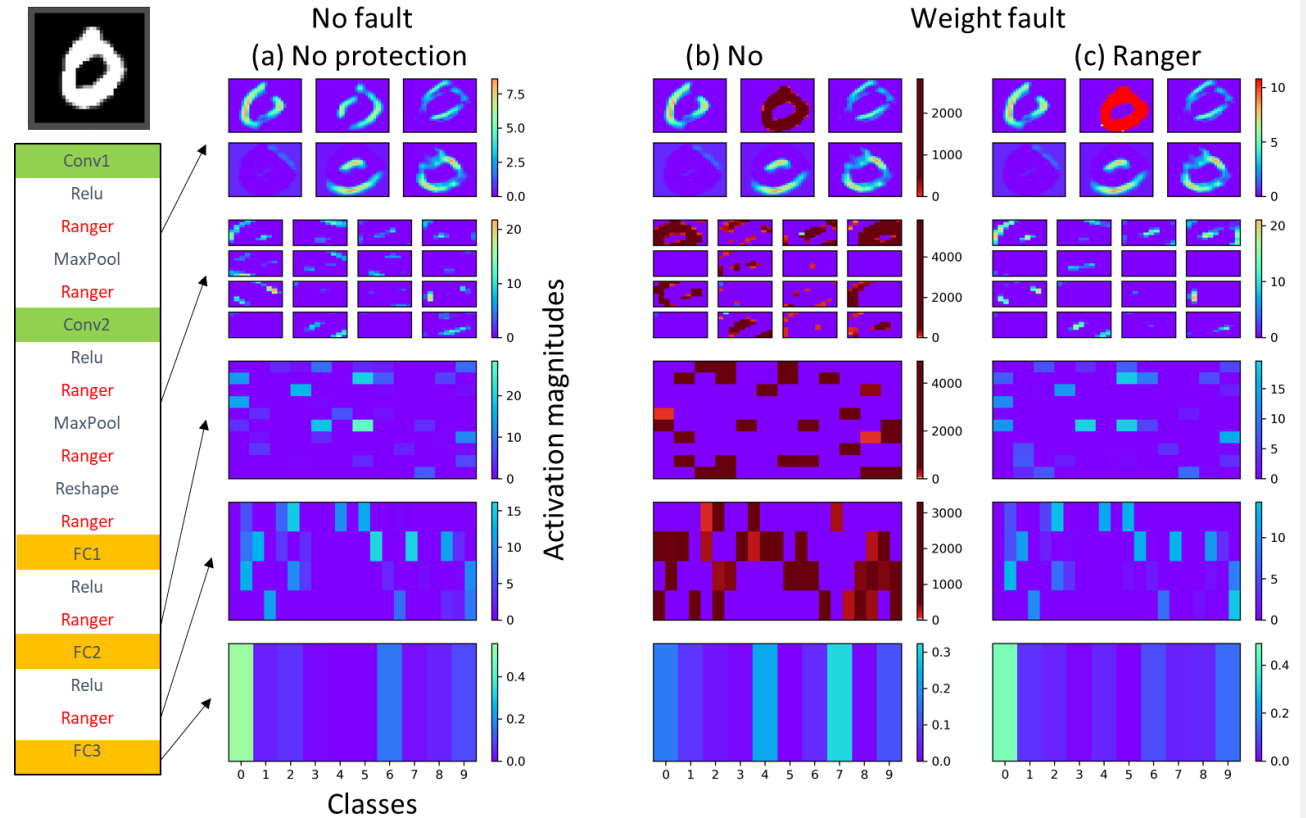
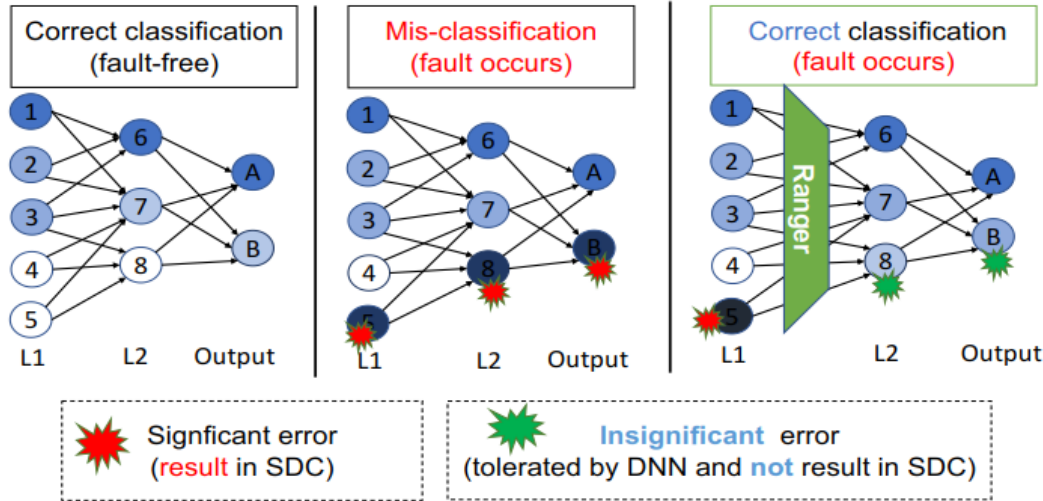


Used to train



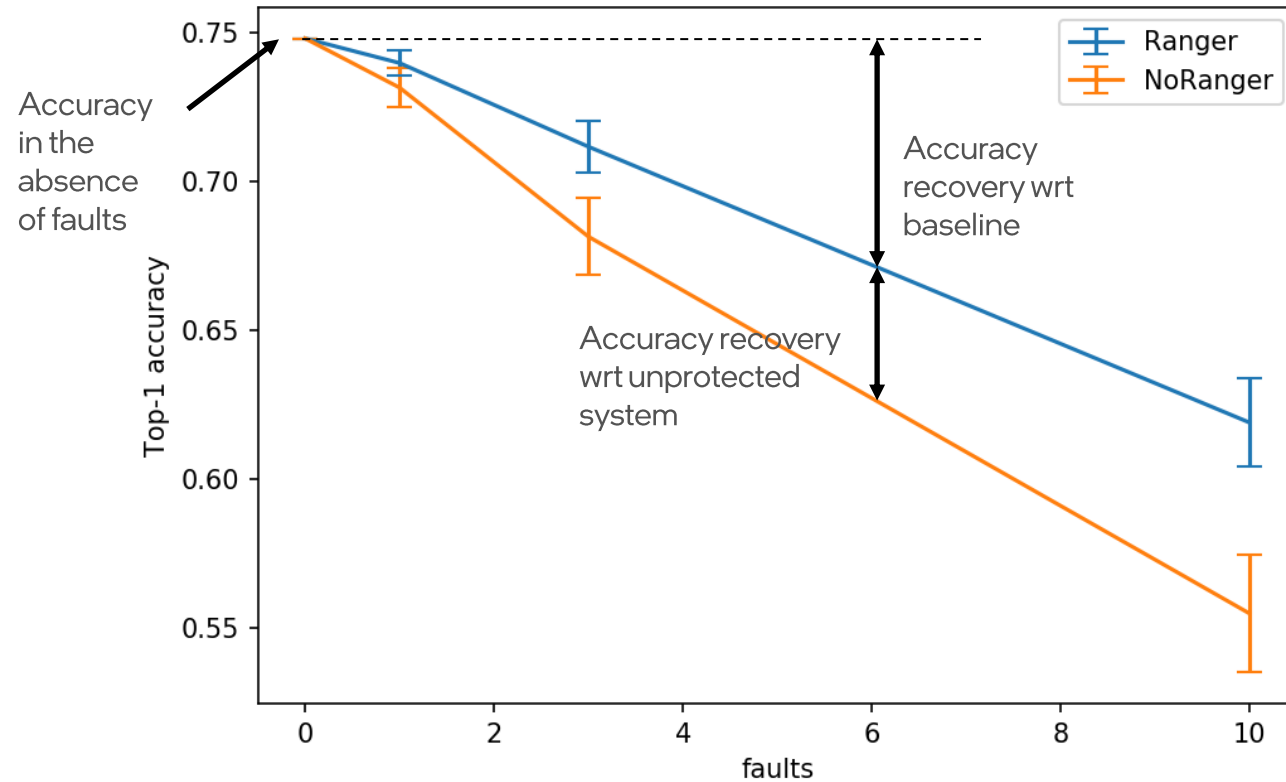
Ranger – Intuition (I)

Chen et al 2020 (“Ranger”)
 Li et al, 2017
 Hong et al, 2019
 Hoang et al, 2019 (“ClipAct”)



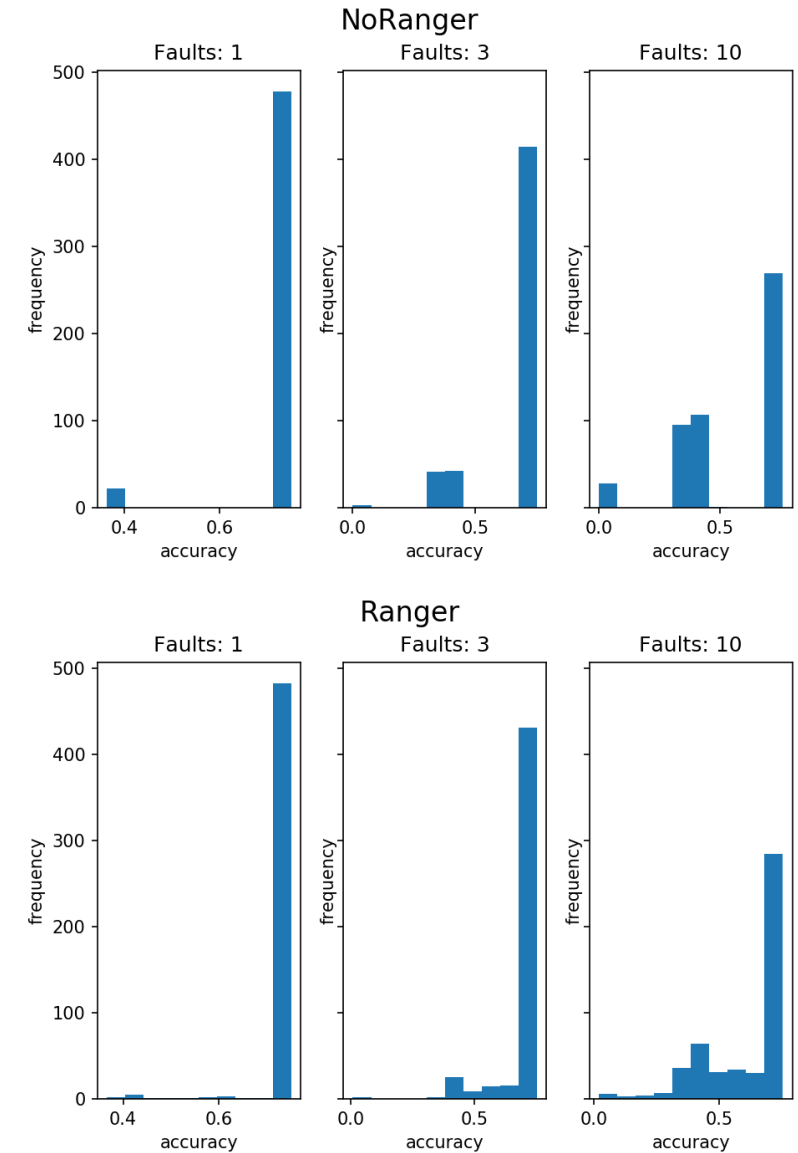
- Insert customizable protection layers (“Ranger layers”) for activation range restriction
- Bound extraction from an independent dataset (e.g. training data)
- No retraining of parameters needed

Ranger – Intuition (II)

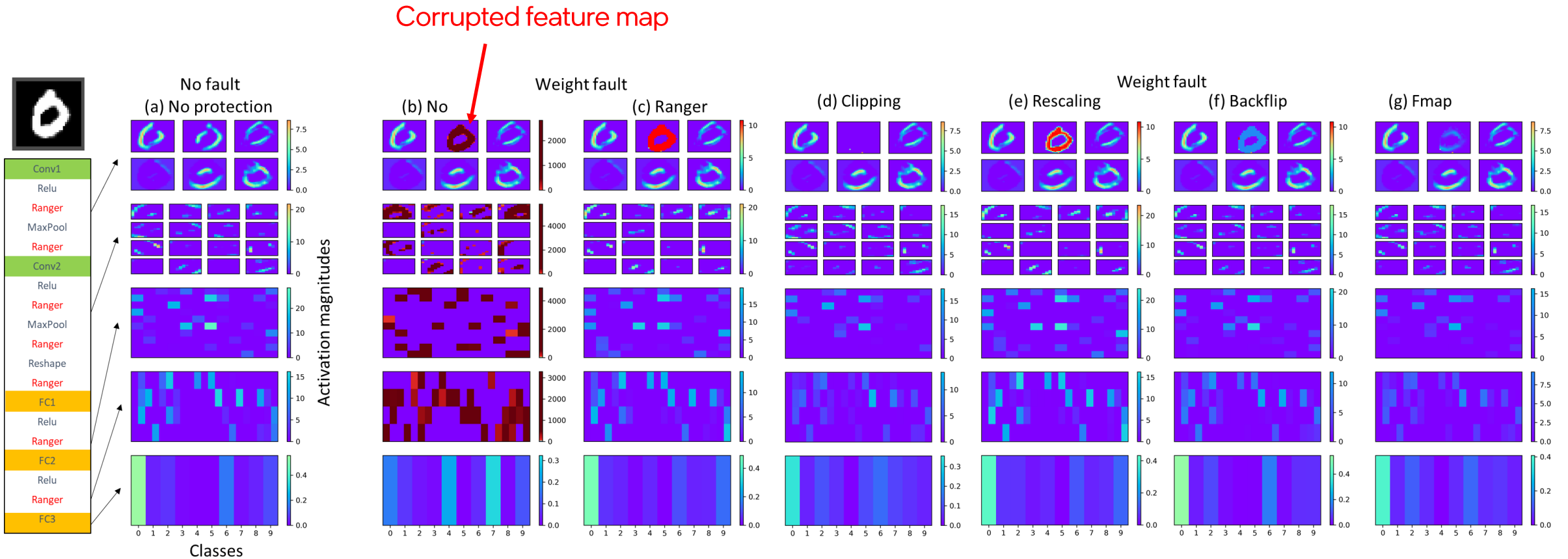


- Large fault injection space: 1 weight fault in conv layers means fault rate of $\sim 6.7e-8$.
- Ranger mitigates the detrimental effect of faults per epoch by eliminating “outliers”, and shift bulk towards maximum accuracy

Distribution of accuracy results in 500 epochs:



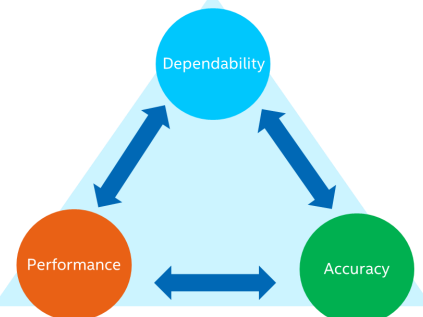
Range restriction alternatives



Goal is to restore the topology of feature maps after a soft error.

Some Ranger Findings

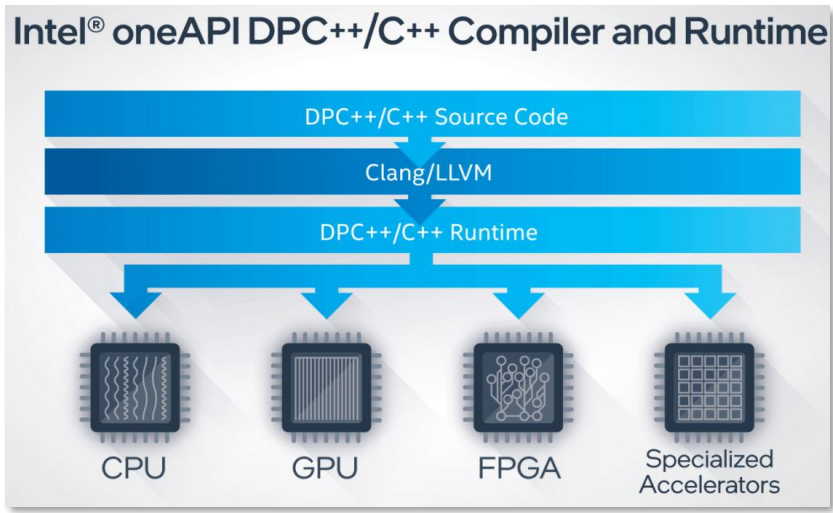
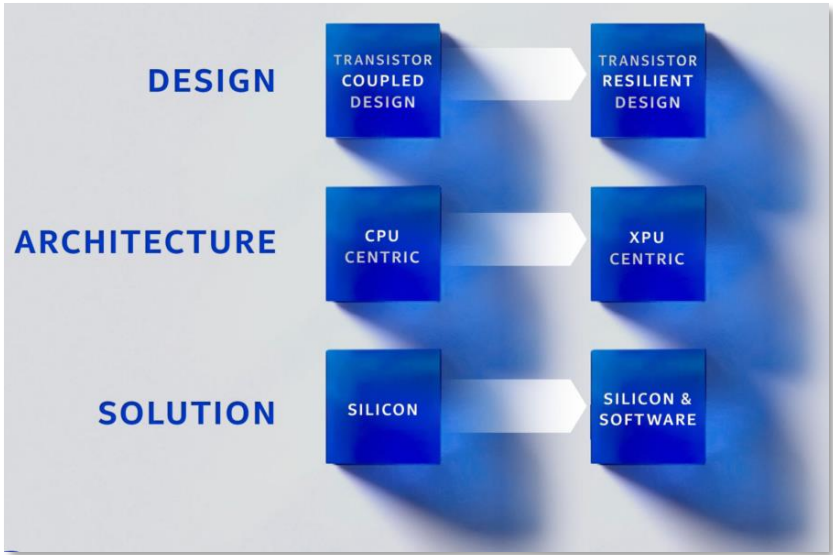
- If **bounds are extracted** from appropriate data, there is almost no reduction of the **baseline accuracy**.
- Impact depends on **data representation**: For FP32 and the given setup, almost all misclassifications happen due to flip in the MSB.
- Error detection: Strong correlation between out-of-bound events and misclassification. Range supervision provides **very high recall** (>0.99), **precision can be lower** (>0.82) due to false positives.
- Error mitigation: **Very high**, especially **Clipper/Backflip**. SDC rate is reduced by up to $\sim 50x$, to $<0.5\%$ in the studied setups.
- Some details @ Geissler et al. Towards a Safety Case for Hardware Fault Tolerance in Convolutional Neural Networks Using Activation Range Supervision, AISafety WS 2021
- Overhead: practically for free in practical scenarios



Summary & Outlook (Resiliency)

- **Dependability:** some faults are real (see also recent publications [1, 2])
- **Cost:** automation and low overhead key to acceptance
- **Software and tools** play a larger role these days
 - Open source and open languages
 - Libraries oneDNN / oneAPI
 - Cross-architecture languages, compilers, and tools
 - DPC++ = ISO C++ and Khronos SYCL™ and community extensions
 - **OpenVINO™** tool
 - Researchers can engage in open source push

Special thanks goes to DRL team and Karthik Pattabiraman



[1] <https://sigops.org/s/conferences/hotos/2021/papers/hotos21-s01-hochschild.pdf> [2] <https://arxiv.org/abs/2102.11245>

Monitoring for Safe Perception

Application and System Level

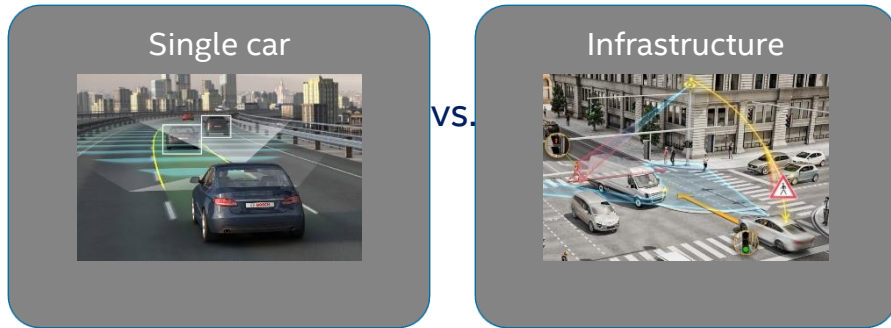
Safe Perception - Monitoring

- Perception in complex environment with **multi-level approaches** to improve safety:
 - Application-level context (view angles, infrastructure involvement, ...)
 - System-level context (diverse space and object representation, sensors, monitors, ...)

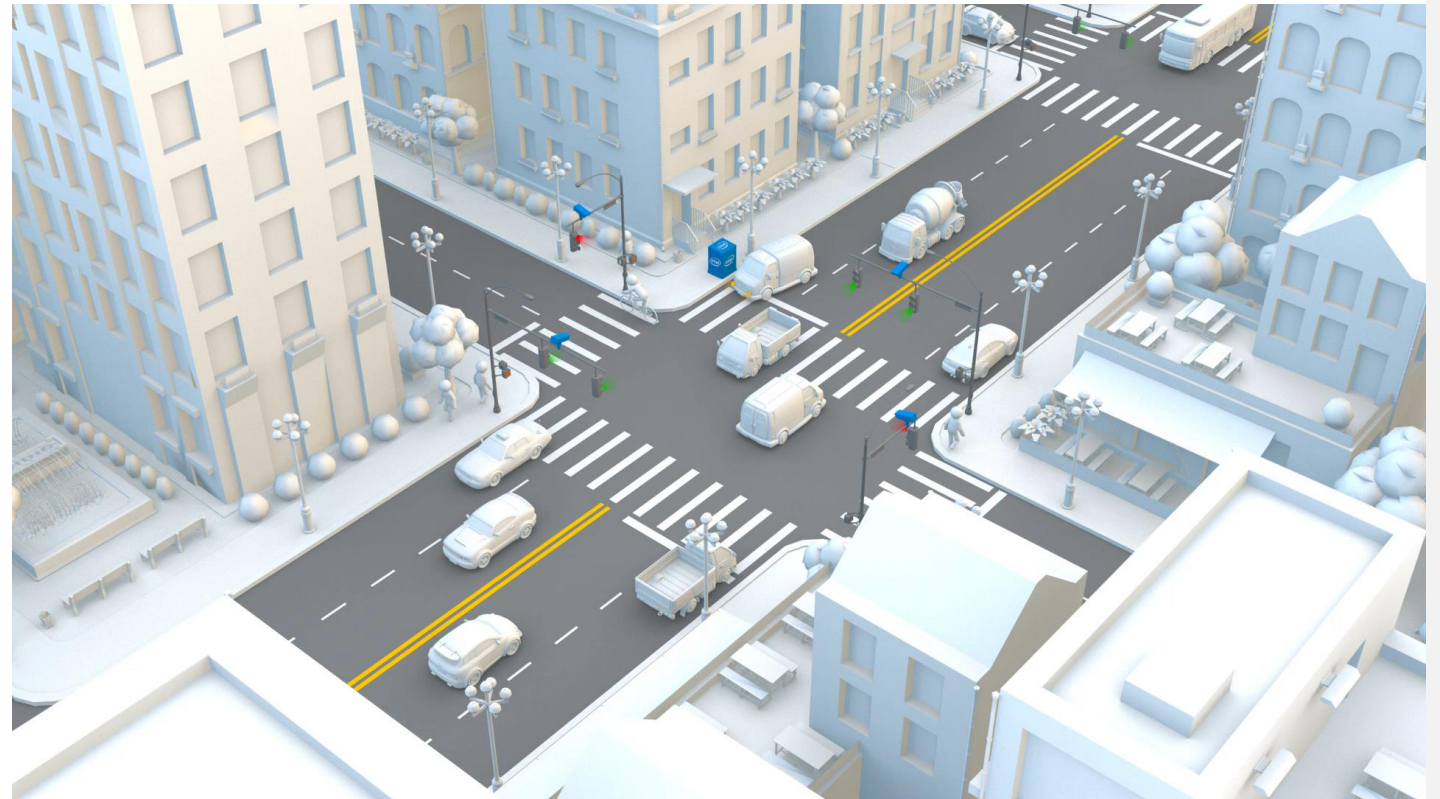


Largest Benefits from Infrastructure

VS. single automated car



- Additional independent source of perception
- Extended field of view
- Different perception vantage points
- Extended compute and energy envelope



<https://www.geospatialworld.net/news/mobileye-join-hands-enable-crowd-sourced-hd-mapping-automated-driving/>
<https://wtvox.com/fashion-innovation/the-future-of-driving-with-v2v-and-v2i-technology/>

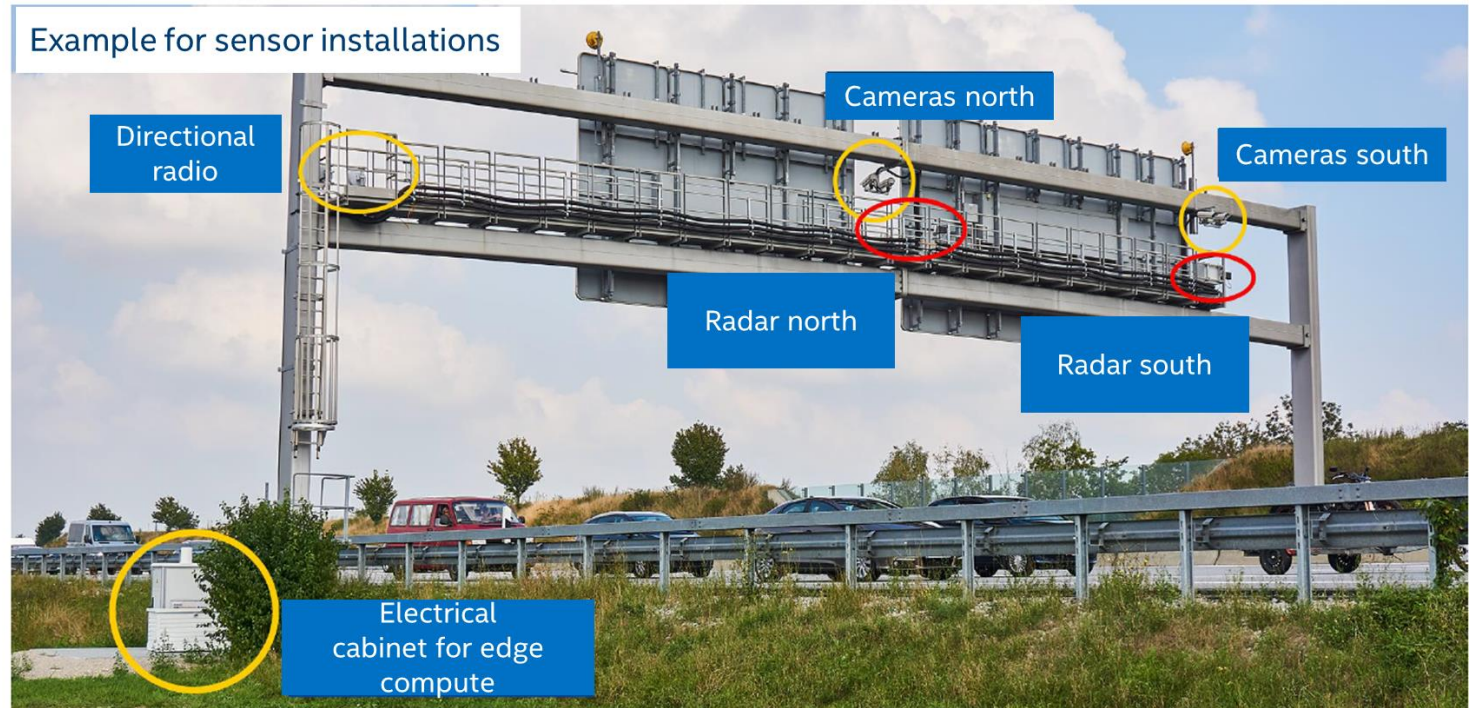
Providentia++



- Basis of the digitalized highway of the future: Real-Time Digital Twin for Smart Highway & Smart City

- Benefits:

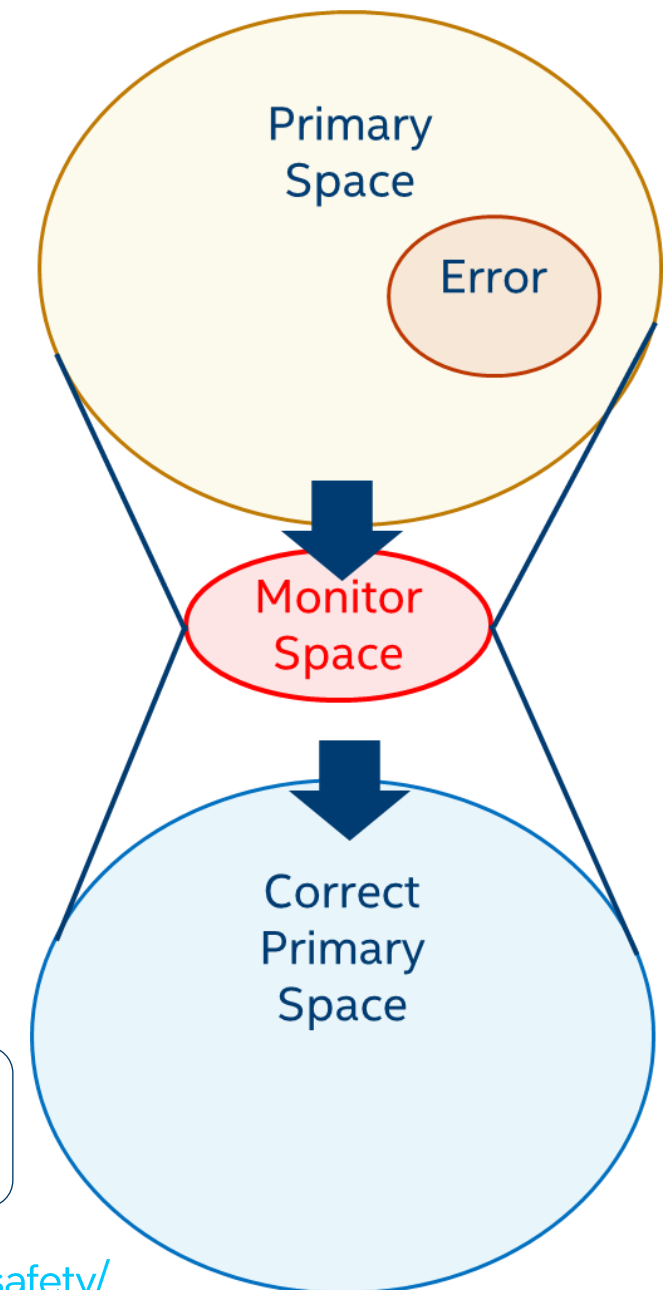
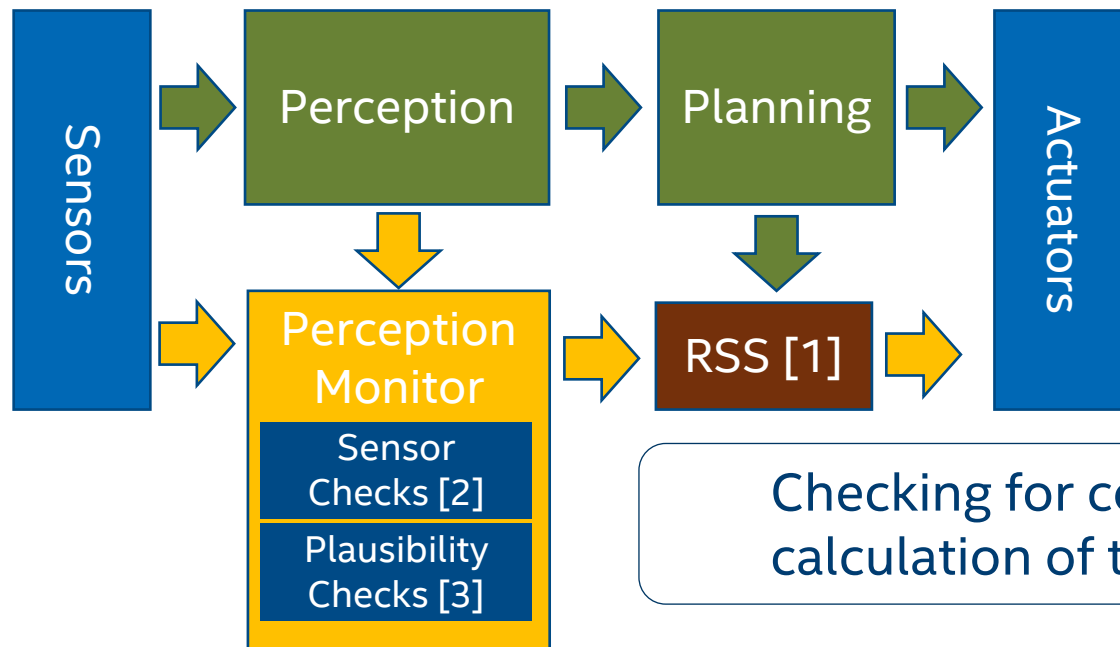
- Real life use case for dependability work
- Test bench for work of Intel Labs (ASMRL & DRL teams)
- Work with IOTG Autonomous Transportation and Infrastructure



Monitor Architecture at Board Level with Application and Application Monitor

How to implement a system that can monitor & recover function :

- requiring significant **less complexity**
- **not decreasing availability** of the primary channel

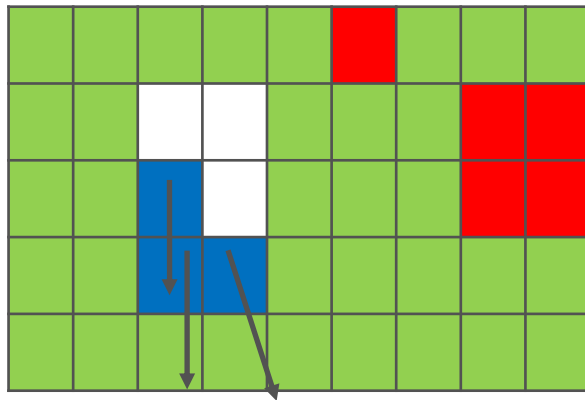


[1] <https://www.mobileye.com/responsibility-sensitive-safety/>

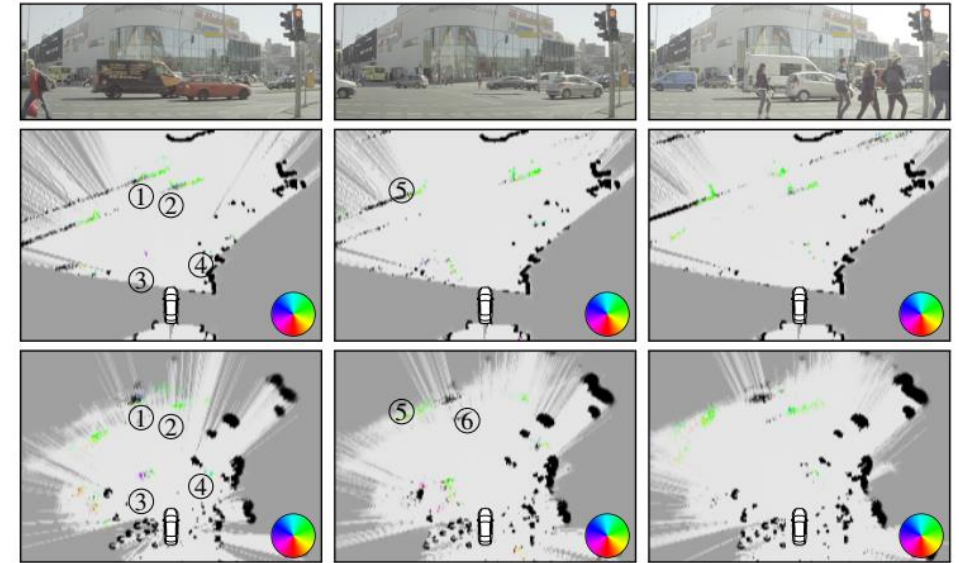
[2] <https://ieeexplore.ieee.org/iel7/9304518/9304528/09304571.pdf>

[3] <https://arxiv.org/pdf/2009.14756>

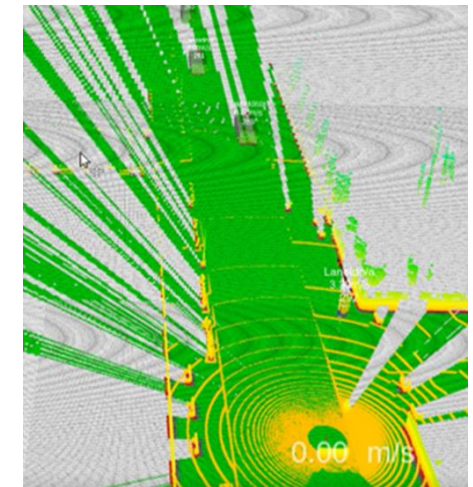
Dynamic Occupancy Grid



- Static Occupied Cells
- Dynamic Occupied Cells
- Free Cells
- Unknown



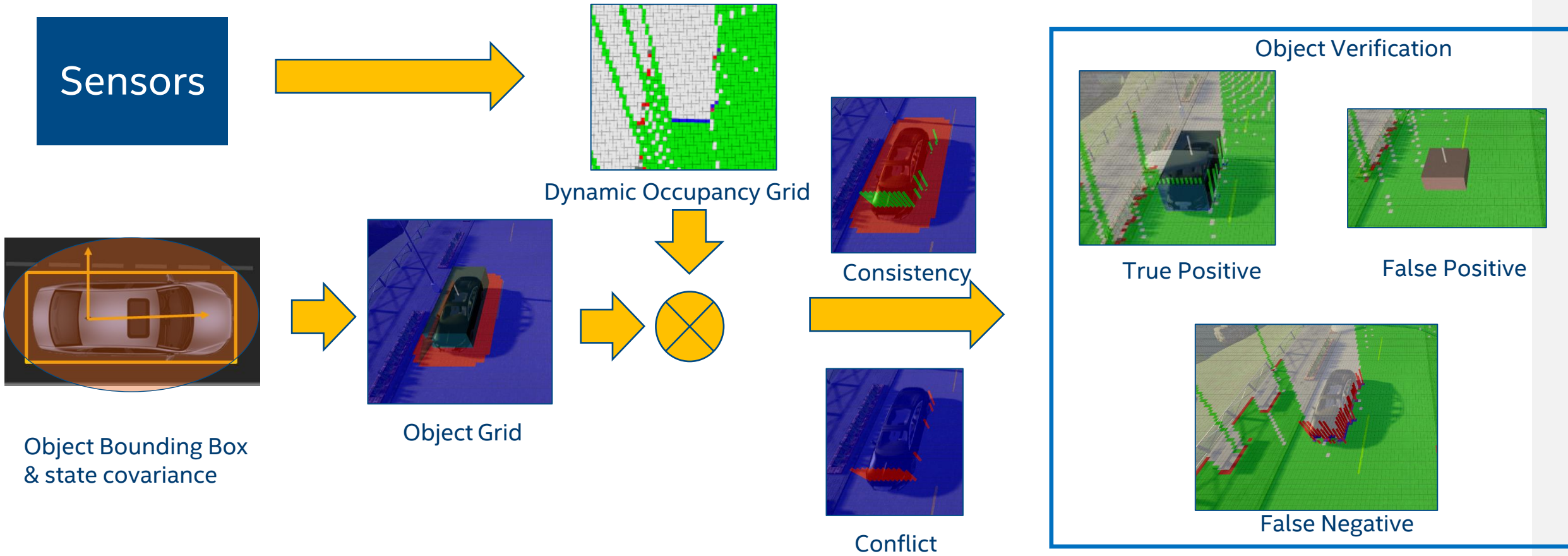
- Promising results in cluttered & dynamic environments
- Fusion of Lidar [1][2] and Radar sensor[3] information
- Classical algorithm* – redundancy towards ML algorithms



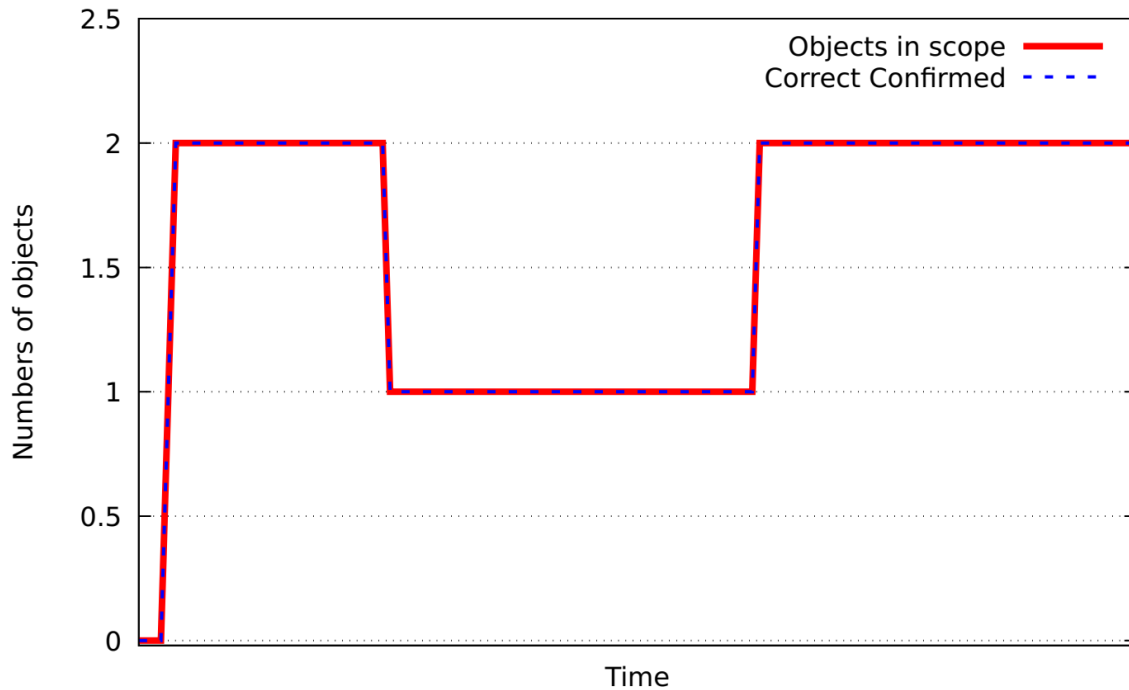
- [1] G. Tanzmeister and D. Wollherr, "Evidential Grid-Based Tracking and Mapping,"
- [2] D. Nuss, et al., "A random finite set approach for dynamic occupancy grid maps with real-time application"
- [3] Christopher Diehl, et al. "Radar-based Dynamic Occupancy Grid Mapping and Object Detection", ITSC 2020

Sensor Checks - Position

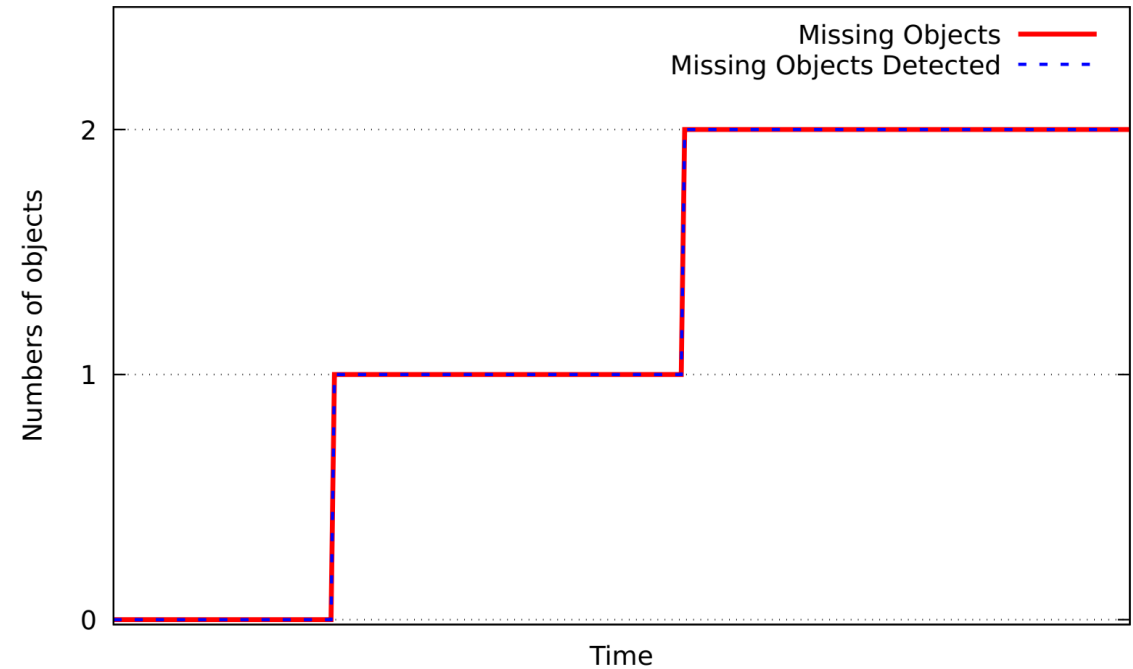
- Objects
 - Position
 - Velocity
 - Dimension



Results: True Positives & False Negatives



True Positives



False Negatives



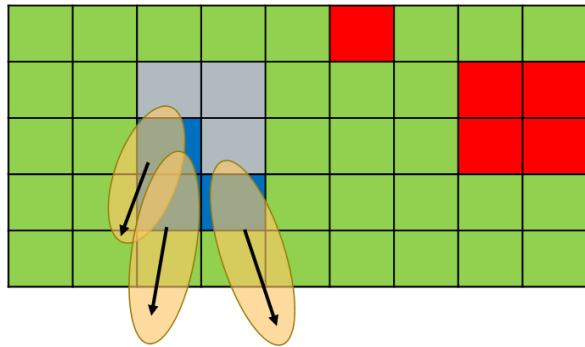
ITSC2020 – *Towards Online Environment Model Verification* | Cornelius Buerkle, Fabian Oboril & Kay-Ulrich Scholl

Velocity checks

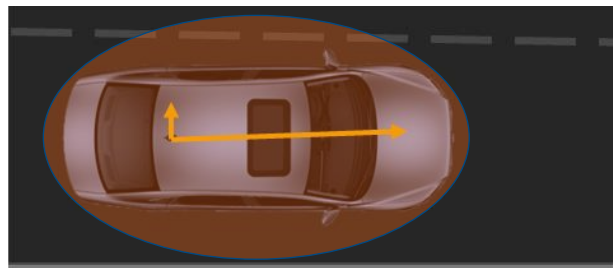
- Objects
 - Position
 - **Velocity**
 - Dimension

Sensor checks

Calculate similarity / distance of velocity distribution of each cell in object region with velocity distribution of object

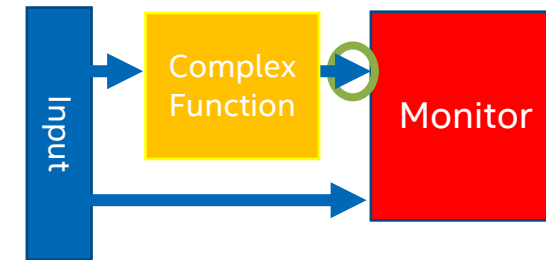


occupancy grid velocity information $V(\zeta)$

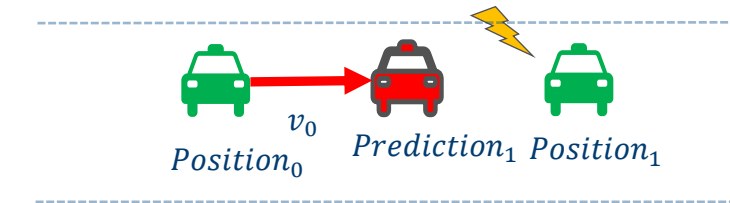


object velocity $\Sigma_v(o)$

Plausibility checks



Correct Velocity

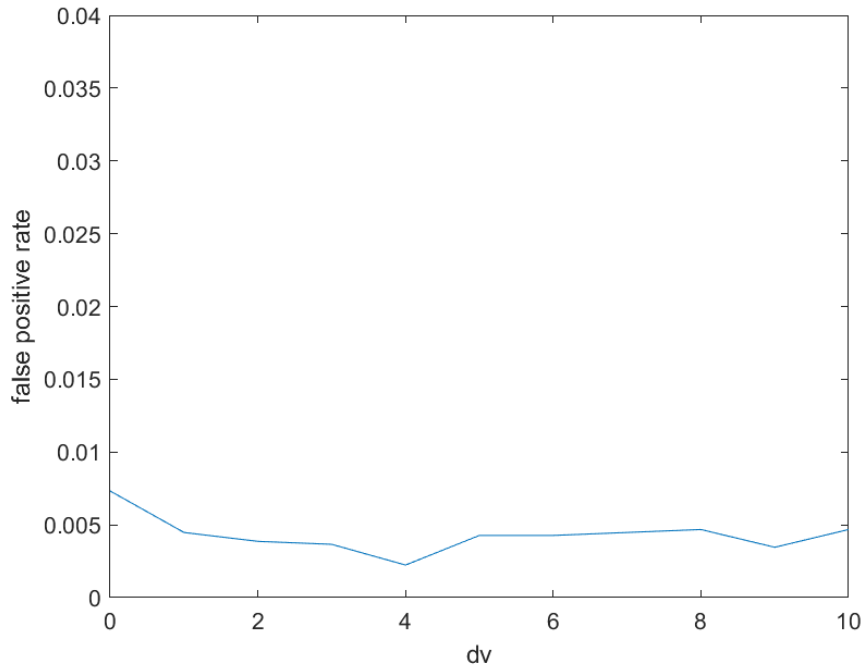


False Velocity

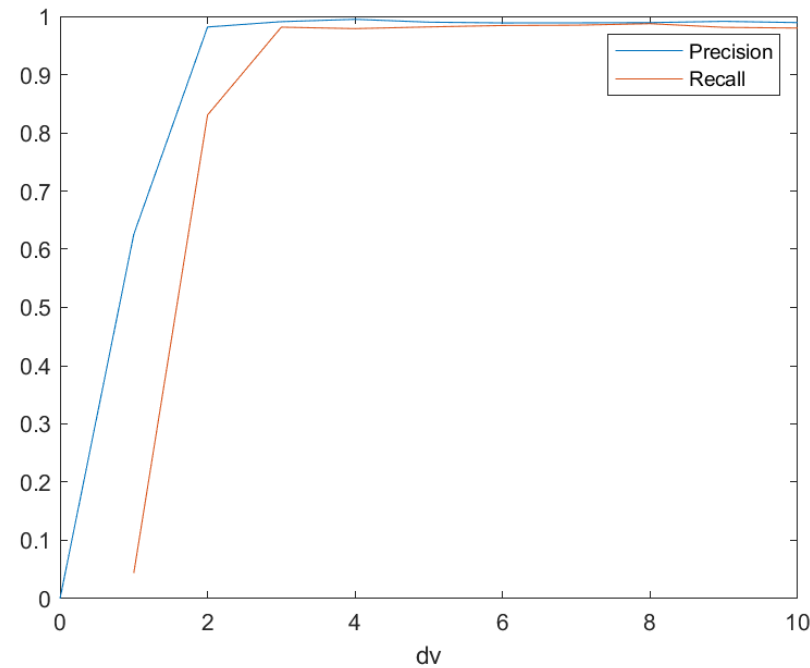
<https://arxiv.org/abs/2009.14756>

Results – Plausibility checks velocity

Chance of error injection per vehicle and time step: **0.25**
Scenario time: Nuscenes clip of **~125** time steps
Time intervals: **~0.2s**



- False positive rate (FPR) in absence of faults **~ 0.5%**
- Measured against **ground truth**, where avg velocity information estimates were added (pos diff over time from last time step).
- Rather insensitive to speed error.



- Increasing speed error faults dv injected. Here, $dv = 0$ means no faults injected
- **Recall ~ 0.98 and Precision ~ 0.98 with the given parameters for $dv \geq 3\text{m/s}$ (or $\sim 11\text{km/h}$)**

Addressing Safety at Different Levels ...

Main lane
Main function
Main value

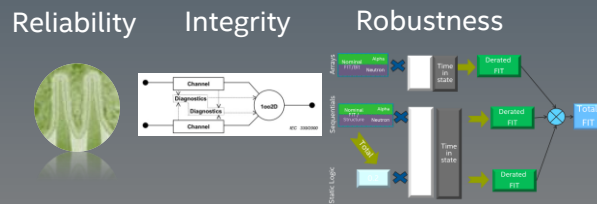
Monitor lane
Safety function
Safety value

CHIP LEVEL

Powerful computing e.g. to master perception challenge



High integrity and availability requires simple approaches

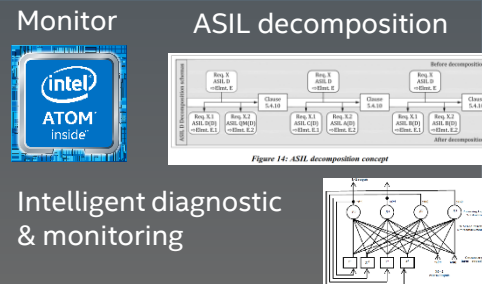


BOARD LEVEL

Powerful Computing and Accelerator

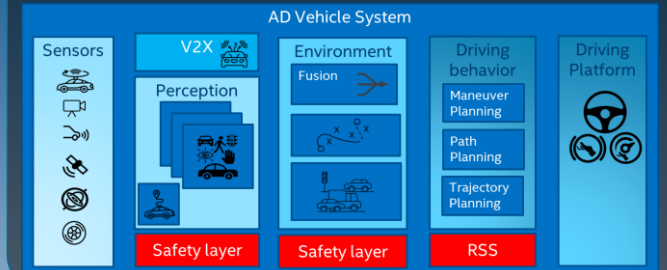


Safety Monitoring (Doer/Checker) for systematic and random hardware faults



AD APPLICATION LEVEL (SYSTEM LEVEL)

Complex AD Stack (with AI like DNN in perception, powerful multi-level)



Simple Checker App. – e.g. Responsibility Sensitive Safety (RSS)



... can sometimes be synergetic

Summary (Monitoring)

- Monitoring at different levels helps dependability (i.e. safety in this context)
- Application-Level Monitor approaches could help to solve the safety challenge for automated vehicles (AVs)
 - Showed initial realization for monitors of object information based on sensor and plausibility checks
 - Demonstrated feasibility with evaluation results in simulation
- Challenges:
 - Sensor checks rely on quality of dynamic occupancy grid and sensor preprocessing
 - Erroneous cells occupation decrease availability or error detection
 - Diversity arguments and proofs for effectiveness
 - Error detection effectiveness to be further investigated (incl. latent errors)

The Intel logo is centered on a solid blue background. It consists of the word "intel" in a white, lowercase, sans-serif font. A small blue square is positioned above the letter "i". To the right of the word "intel" is a white registered trademark symbol (®).