

Developing Safe Cyber-Physical Systems for Safety-Critical Applications

Presented by: Mansur M. Arief, Ph.D.

at the Mechanical and Aerospace Engineering Seminar, Nanyang Technological University, Singapore

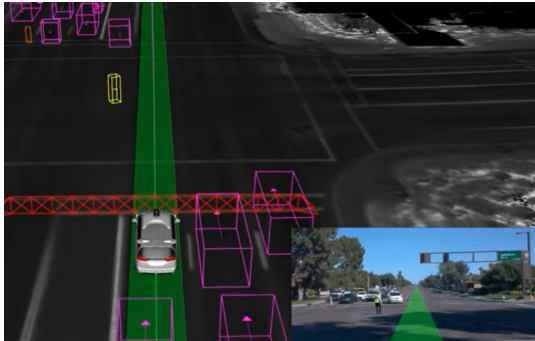
August 20, 2024

About Mansur

- **Postdoc, AeroAstro Department, Stanford, 2023-present**
 - Working with Mykel Kochenderfer at Stanford Intelligent Systems Lab (SISL) and Jef Caers at Mineral-X (Stanford Doerr School of Sustainability)
 - Working with PhD and master students, and RAs on AI for safety and sustainability
- **PhD in Mechanical Engineering, Carnegie Mellon, 2023**
 - Dissertation at Safe AI Lab: Certifiable Evaluation for Safe Intelligent Autonomy
 - Worked with Ding Zhao and Henry Lam (Columbia IEOR)
- **MSE, Industrial & Operations Engineering, University of Michigan, 2018**
- **BE, Industrial and Systems Engineering, Sepuluh Nopember Institute of Technology, Surabaya, 2014**



Cyber-Physical Systems (CPSs) are everywhere



Autonomous vehicles




Exploratory robots



**Aircraft collision
avoidance systems**

- Interacting with humans more intensively and collaboratively
- Making more important, even safety-critical, decisions


This is just the beginning...




Marco Pavone
Associate Professor at Stanford University and Director, Autonomous Vehicle Research at NVIDIA

Followers 8,611

Message




Marco Pavone reposted this



NVIDIA DRIVE
52,620 followers
3w · 🌐

+ Follow

Congrats to NVIDIA's [Marco Pavone](#) and Edward Schmerling and the team at Stanford for the RSS (Robotics: Science and Systems) 2024 Outstanding Paper Award on the topic of Real-Time Anomaly ...more

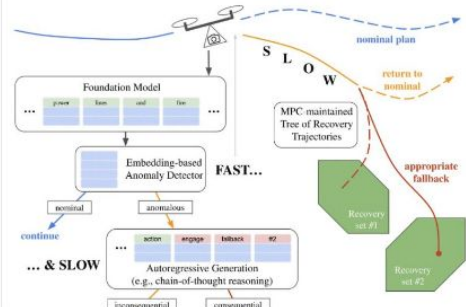


NVIDIA AI
1,005,059 followers
3w · 🌐


+ Follow

As AI is implemented in our daily lives, engaging with [#robots](#) and autonomous vehicles safely will become more important. RSS's best paper from [Stanford University](#) and [NVIDIA](#) presents a framework designed to improve the trustworthiness of dynamic robotic systems under resource and time constraints
<https://nvidia.ws/3YciRN7>

Congrats to the winners 🏆 [#RSS2024](#)




The diagram illustrates a system architecture for real-time anomaly detection and recovery. At the top, a drone icon represents the system. Below it, a 'Foundation Model' block contains 'prior', 'sens', and 'act' components. This feeds into an 'Embedding-based Anomaly Detector' which outputs 'nominal' or 'anomalous'. The 'anomalous' path leads to 'FAST...' (Autoregressive Generation, e.g., chain-of-thought reasoning), which then branches into 'inconsequential' and 'consequential' actions. The 'consequential' path leads to 'SLOW' (MPC-maintained Tree of Recovery Trajectories), which outputs 'nominal' or 'anomalous'. The 'anomalous' path leads to 'appropriate fallback' and 'Recovery set #1' and 'Recovery set #2'.



Conrad Tucker
Director of CMU-Africa| Professor of Mechanical Engineering & Machine Learning (Courtesy)| Board Member

Followers 4,290

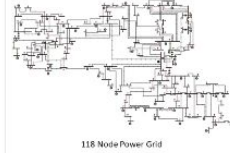
Message



Conrad Tucker · 1st
Director of CMU-Africa| Professor of Mechanical Engineering & ...
1yr · 🌐

Congrats to my Ph.D. student [James Cunningham](#) and our collaborators from the [Air Force Research Laboratory](#) Dr. [Alex A.](#), [David Ferris](#), and [Phillip Morrone](#) for our recent [IEEE](#) journal publication titled "A Deep Learning Game Theoretic Model for Defending Against Large Scale Smart Grid Attacks".

A short video-summary of the paper can be viewed below and the paper can be accessed here: <https://lnkd.in/gUA6-ZSJ>. The data and code for the model can be accessed on [GitHub](#) here: <https://lnkd.in/gu-BxVgH>
[#deeplearning](#), [#mechanicalengineering](#), [#gametheory](#), [#reinforcementlearning](#)

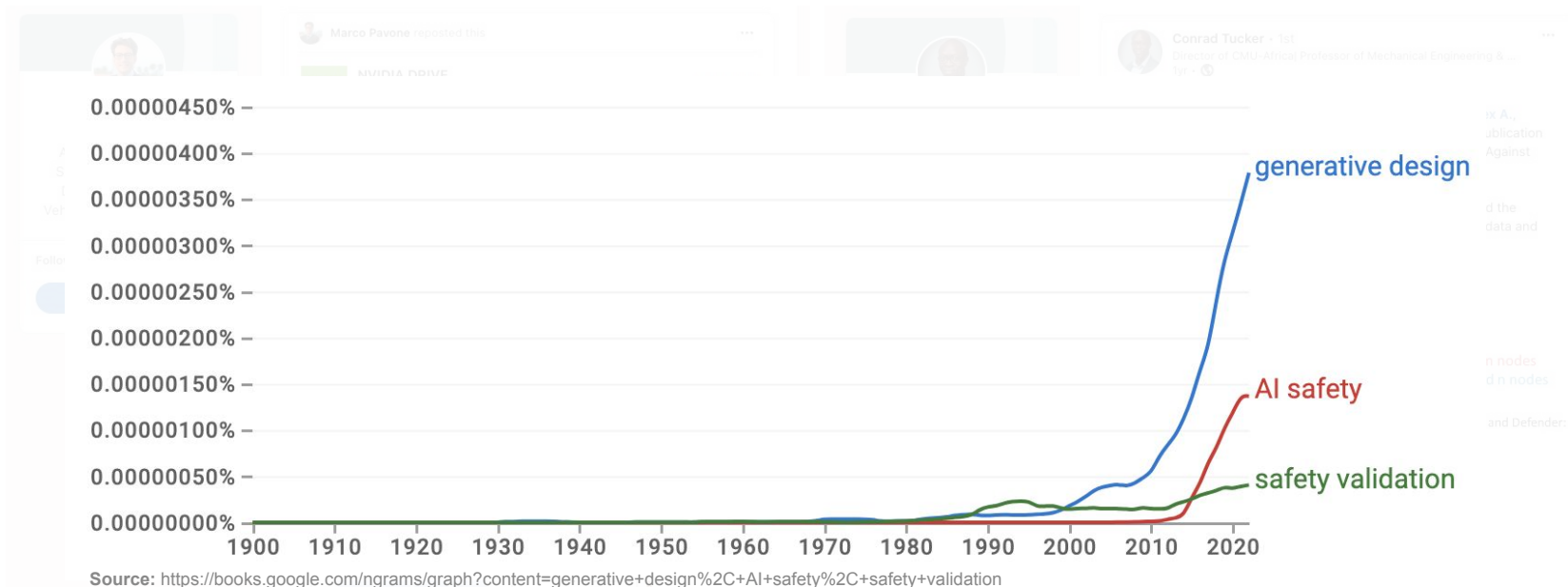


118 Node Power Grid

Attacker can attack n nodes
Defender can defend n nodes

Of actions for Attacker and Defender:
 $n=1$: 118
 $n=2$: 6093
 $n=3$: 260k
 $n=10$: 9.75×10^{13}

This is just the beginning...



but, safety engineering should catch up quickly!

The risk is real, the impact can be catastrophic!

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

Two Drug Possession Arrests

Person	Risk Level	Score
DYLAN FUGETT	LOW RISK	3
BERNARD PARKER	HIGH RISK	10

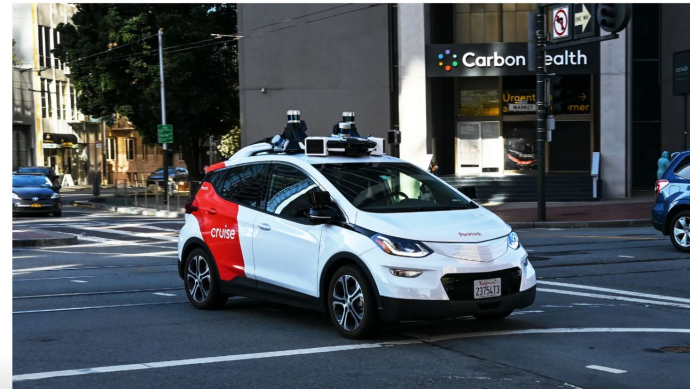
Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.

Source: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

BARZAN MARSHALL BUSINESS OCT 24, 2023 4:31 PM

GM's Cruise Loses Its Self-Driving License in San Francisco After a Robotaxi Dragged a Person

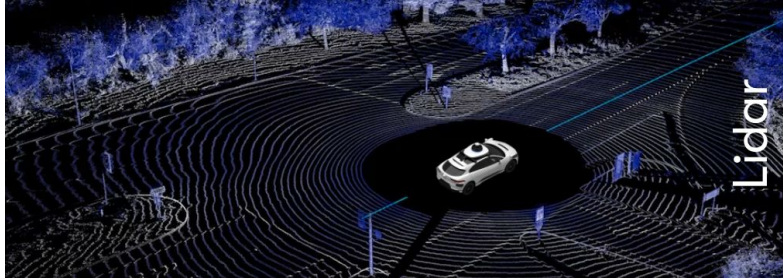
The California DMV says the company's autonomous taxis are "not safe" and that Cruise "misrepresented" safety information about its self-driving vehicle technology.



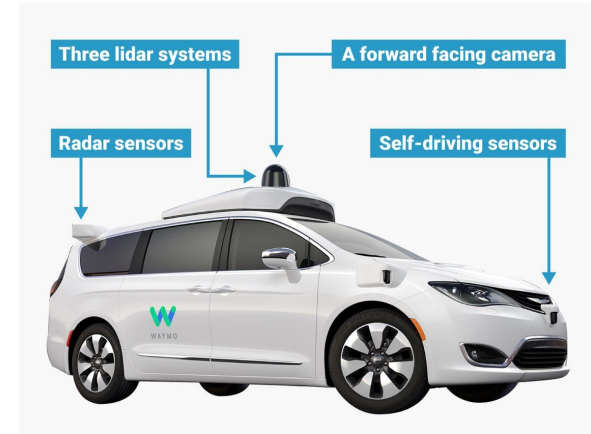
WIRED

Source: <https://www.wired.com/story/cruise-robotaxi-self-driving-permit-revoked-california/>

Modern CPS uses multimodal sensors,

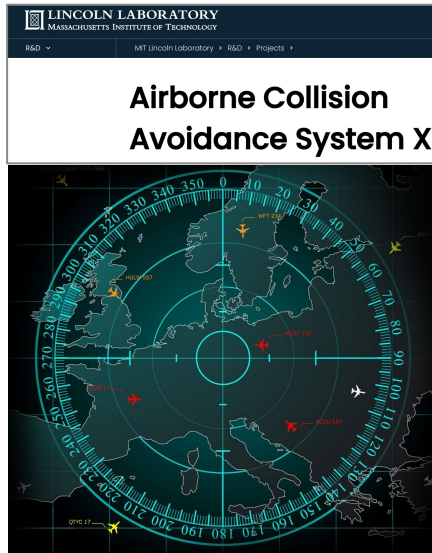


Source: <https://waymo.com/>



... and is robust to some degrees of uncertainty

ACAS-X



IATA safety statistics

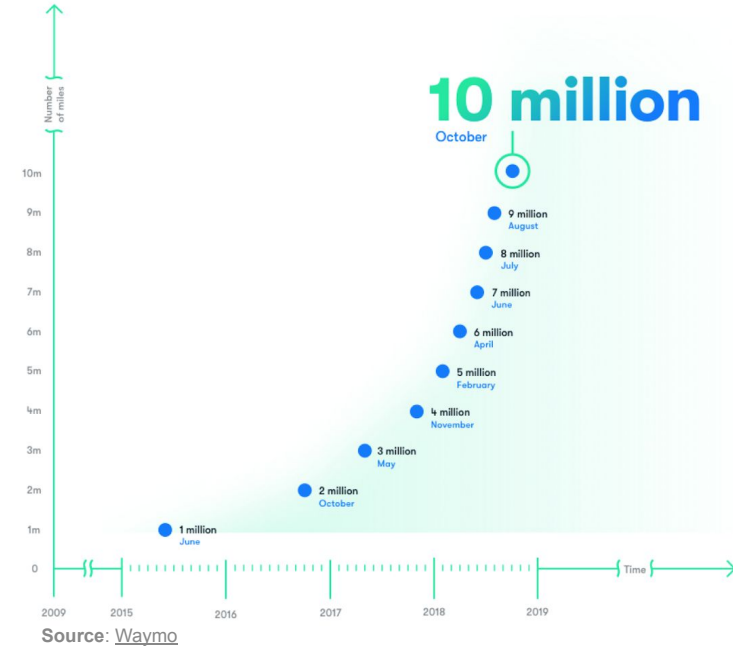
ACCIDENT TYPE	2023	2022	5-YEAR AVERAGE (2019-2023)
All accident rate (accidents per one million flights)	0.80 (1 accident every 1.26 million flights)	1.30 (1 accident every 0.77 million flights)	1.19 (1 accident every 0.88 million flights)
All accident rate for IATA member airlines	0.77 (1 accident every 1.30 million flights)	0.58 (1 accident every 1.72 million flights)	0.73 (1 accident every 1.40 million flights)
Total accidents	30	42	38
Fatal accidents	1 (0 jet and 1 turboprop)	5 (1 jet and 4 turboprop)	5
Fatalities	72	158	143
Fatality risk	0.03	0.11	0.11
IATA member airlines' fatality risk	0.00	0.02	0.04

Safe CPSs are challenging to evaluate



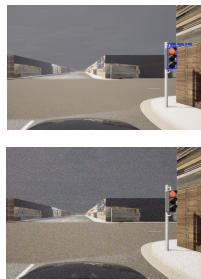
Main challenges, include:

- curse of dimensionality
- curse of rarity

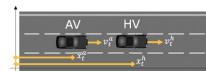
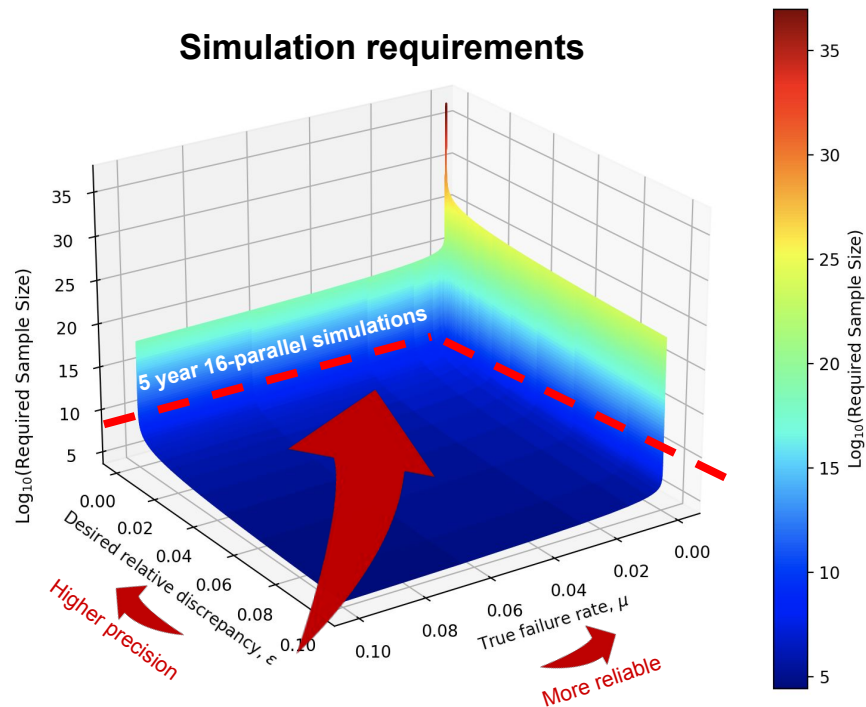


Airplane-level safety requires HUGE simulation runs,

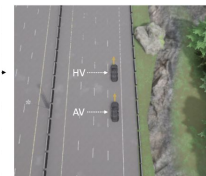
Simulation requirements



I ran simulations for about a month to compare 99.99% accuracy CV models.



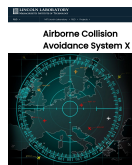
(a) Schematic diagram



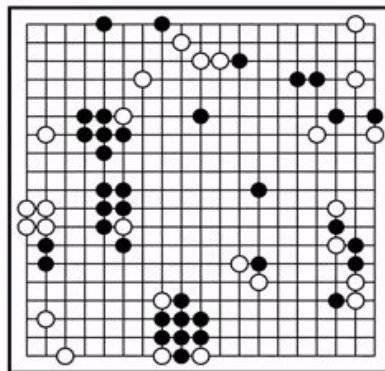
(b) CARLA topview camera

Even more for validating a 10^{-5} failure rate AV model.

Airplane-level safety requires HUGE simulation runs,



Mykel developed
ACAS-X in 2013



Accepted as standard after validation
in early 2020s (the method already
developed in several versions)

... and statistical and engineering rigor

FUNCTIONAL SAFETY SUPPORT THROUGHOUT THE DEVELOPMENT CYCLE

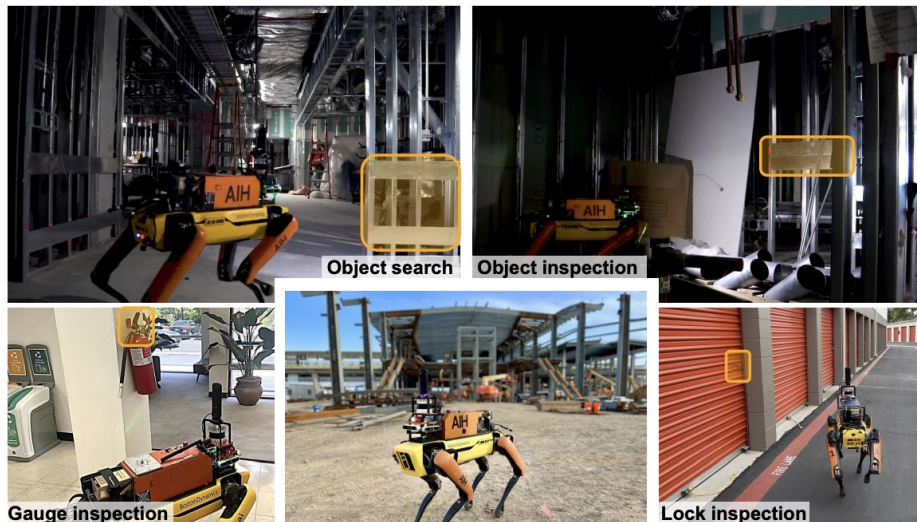


Mykel developed
ACAS-X in 2013



Accepted as standard after validation
in early 2020s (the method already
developed in several versions)

My Vision: Bring the airplane-level safety to CPS



Efficient inspection robots for semi-controlled indoor and outdoor area

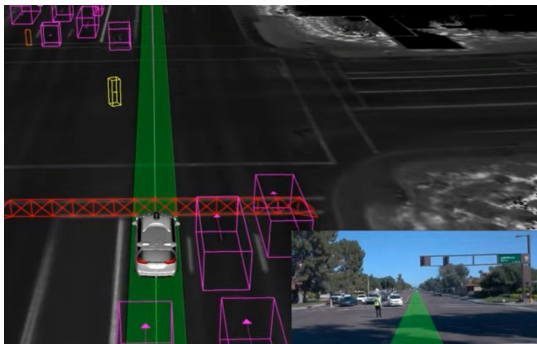


Multi-mobility collaboration for exploration

Ginting, M. F., Kim, S. K., Fan, D. D., Palieri, M., Kochenderfer, M. J., & Agha-Mohammadi, A. A. (2024). SEEK: Semantic Reasoning for Object Goal Navigation in Real World Inspection Tasks. *arXiv:2405.09822*.

Ginting, Muhammad Fadhil, Kyohei Otsu, Mykel J. Kochenderfer, and Ali-akbar Agha-mohammadi. "Capability-aware task allocation and team formation analysis for cooperative exploration of complex environments." IROS 2022.

Research Directions



**Rigorous and scalable
safety validation**

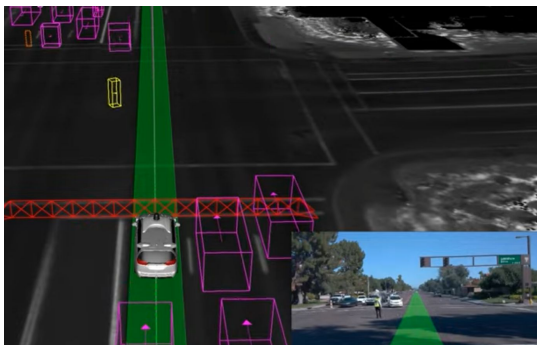


**Robust planning
& monitoring**



**Safety-centered
CPS development**

Research Directions



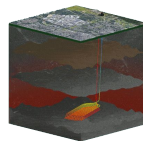
**Rigorous and scalable
safety validation**



Transportation



**Robust planning
& monitoring**



Sustainability



**Safety-centered
CPS development**

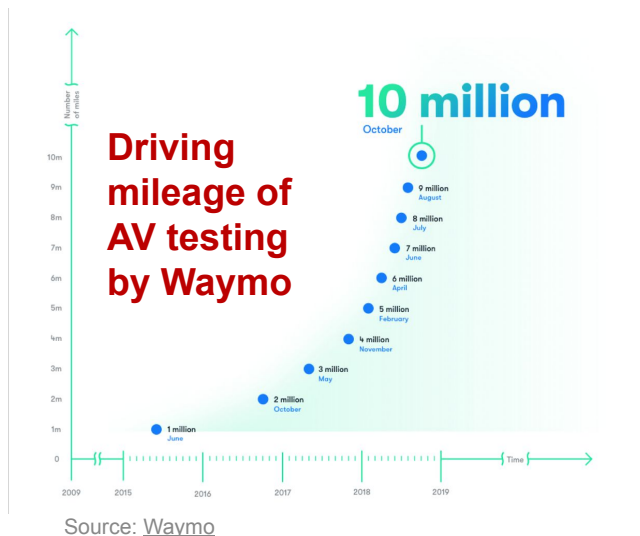


Manufacturing

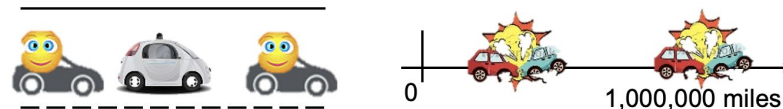
Application Areas

Rigorous and scalable safety validation

- If the failure rate is μ , smaller μ requires larger sample size.



Main reason:



Crashes happen extremely **rarely** (NHTSA, 2019)

How do we sample test scenarios more efficiently?

- **Objective:** Develop algorithms that can deal with
 - extreme rarity and high-dimensional inputs
- **Requirements:**
 - efficiency guarantee and efficient computation
- **Proposed algorithms:**
 - **Deep IS:** Deep Importance Sampling¹
 - **Deep-PrAE:** Deep Probabilistic Accelerated Evaluation²
 - **CERTIFY:** Computationally Efficient and Robust Evaluation of Safety³

¹[Arief, Mansur](#), Zhepeng Cen, Zhenyuan Liu, Zhiyuan Huang, Bo Li, Henry Lam, and Ding Zhao. "Certifiable Evaluation for Autonomous Vehicle Perception Systems Using Deep Importance Sampling (Deep IS)." In *Proceedings of the 2022 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2022. [[Link](#)]

²[Arief, Mansur](#), Zhiyuan Huang, Guru Koushik Senthil Kumar, Yuanlu Bai, Shengyi He, Wenhao Ding, Henry Lam, and Ding Zhao. "Deep Probabilistic Accelerated Evaluation: A Certifiable Rare-Event Simulation Methodology for Black-Box Autonomy." In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*. PMLR, 2021. [[Link](#)]

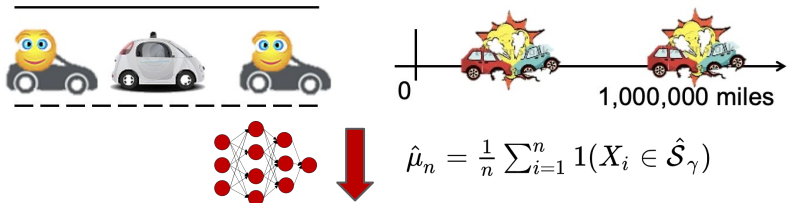
³[Arief, Mansur](#), Zhepeng Cen, Huan Zhang, Henry Lam, and Ding Zhao. "CERTIFY: Computationally Efficient Rare-failure Certification of Autonomous Vehicles." *Under review for IEEE T-IV*. [[Link](#)]

Importance Sampling (IS)

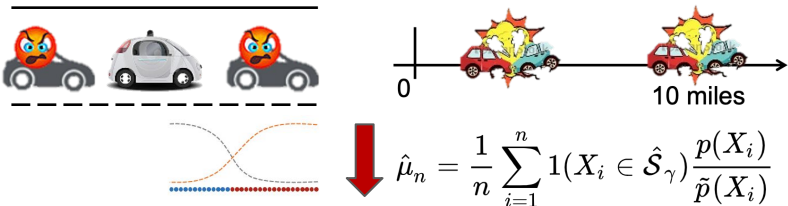
- **Importance Sampling (IS)** uses biased distribution to generate test cases and use importance weights to get unbiased results.

Importance Sampling (IS)

Naturalistic driving conditions:



Aggressive driving conditions:



Unbiased result

Key steps:

1. Start with normal driving
2. Learn the statistical model
3. Bias the statistics toward more aggressive driving
4. Use importance weights to obtain unbiased result
5. Return unbiased statistics

¹Arief, Mansur, Zhepeng Cen, Zhenyuan Liu, Zhiyuan Huang, Bo Li, Henry Lam, and Ding Zhao. "Certifiable Evaluation for Autonomous Vehicle Perception Systems Using Deep Importance Sampling (Deep IS)." In *Proceedings of the 2022 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2022. [\[Link\]](#)

Theoretical guarantees

- Crude technique sampling is **inadequate** to evaluate rare events (does not scale well in failure rarity)

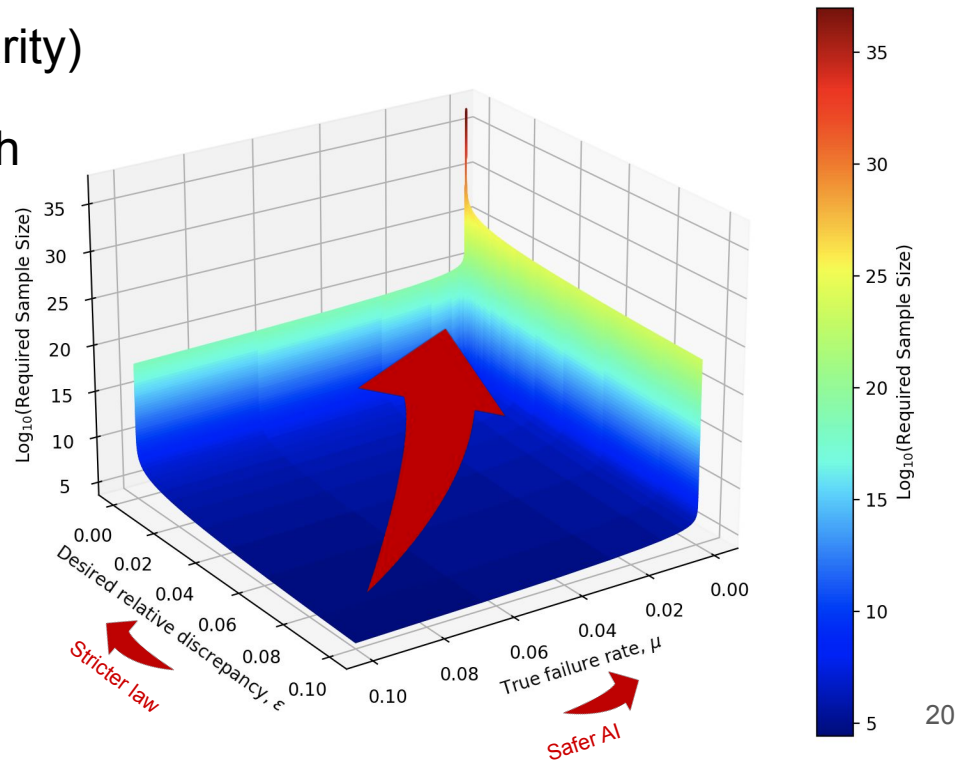
- Consider estimating a tiny μ with an estimator $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n Y_i$.

- A small ϵ & high confidence $1-\delta$

$\mathbb{P}(|\hat{\mu}_n - \mu| > \epsilon\mu) \leq \delta$
is achieved **only when**

$$n \geq \frac{\text{Var}(Y_i)}{\delta\epsilon^2\mu^2}.$$

- Thus, as $\mu \rightarrow 0, n \rightarrow \infty$.



Theoretical guarantees

- **Importance Sampling (IS)** uses proposal distribution \tilde{p} and computes

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \in \mathcal{S}_\gamma) L(X_i) = \frac{1}{n} \sum_{i=1}^n Y_i L(X_i),$$

$$L(X_i) = \frac{p(X_i)}{\tilde{p}(X_i)}. \Rightarrow \text{called the importance ratio}$$

Theoretical guarantees

- IS is provably unbiased

$$\begin{aligned}\mathbb{E}_{X \sim \tilde{p}}[\hat{\mu}_n] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \in \mathcal{S}_\gamma) L(X_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\mathbb{1}(X_i \in \mathcal{S}_\gamma) \frac{p(X_i)}{\tilde{p}(X_i)} \right] \\ &= \frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^d} \mathbb{1}(X_i \in \mathcal{S}_\gamma) \frac{p(X_i)}{\tilde{p}(X_i)} \tilde{p}(X_i) dX_i \\ &= \frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^d} \mathbb{1}(X_i \in \mathcal{S}_\gamma) p(X_i) dX_i \\ &= \mu.\end{aligned}$$

Theoretical guarantees

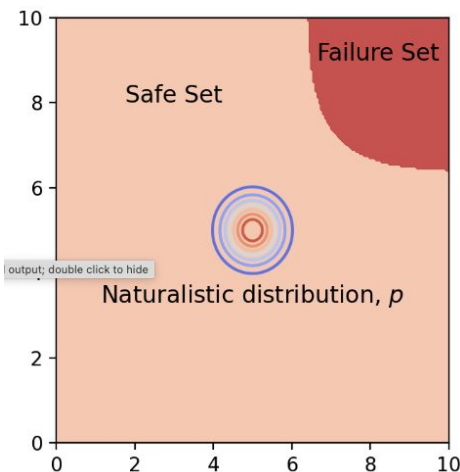
- **IS reduces variance** if the proposal distribution: $\tilde{p}(x) \propto \mathbb{1}(x \in \mathcal{S}_\gamma) p(x)$, i.e. the naturalistic distribution conditional on the failure set.
- **Cross Entropy (CE)** minimizes the KL-divergence between the proposal and this theoretically optimal distribution iteratively

$$\max_{\theta \in \Theta} \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbb{1}(X_i \in \mathcal{S}_\gamma) \frac{p(X_i)}{p_{\theta_j}(X_i)} \ln p_{\theta}(X_i)$$

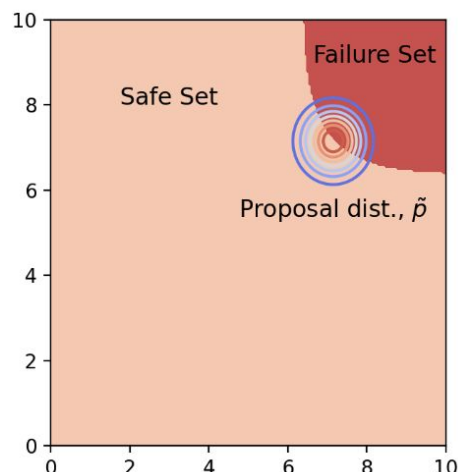
under some parametric class Θ .

What does it mean?

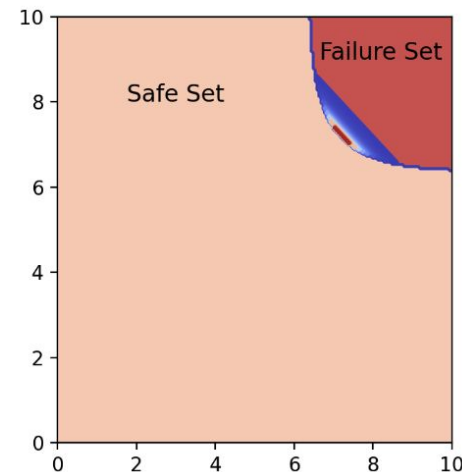
- Intuitively, IS **skews the distribution toward failures** and use likelihood ratio to compute an unbiased estimate.



Naturalistic conditions



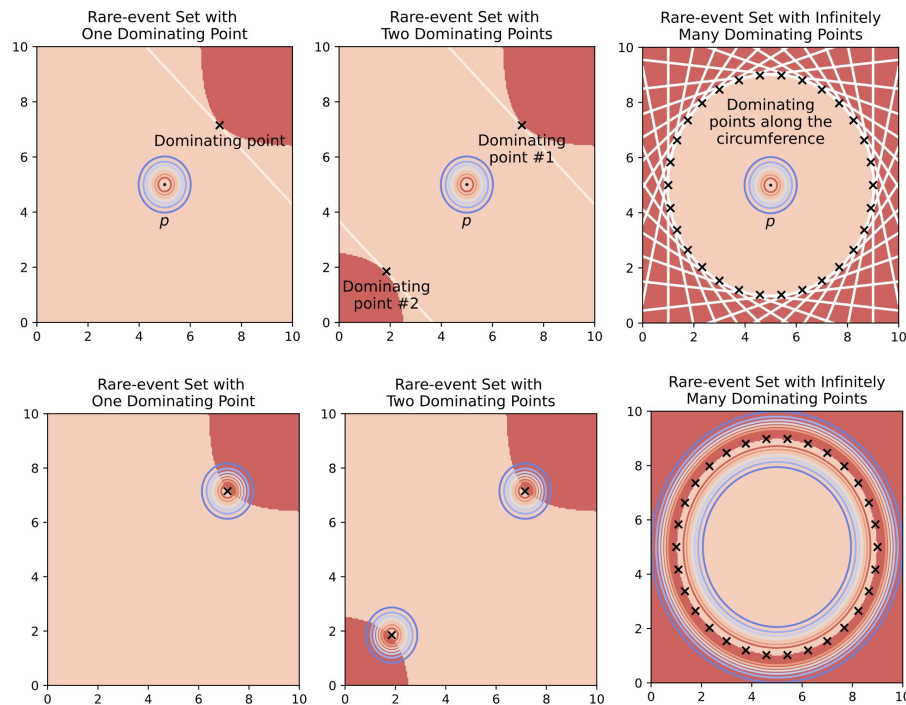
Skewed/aggressive conditions



Likelihood ratio conditional
on the failure set

What does it mean?

- Also applies for multiple failure modes.



Scaling to high dimensional problems

- Use neural net (NN) to approximate high-dimensional failure set

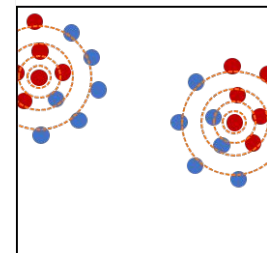
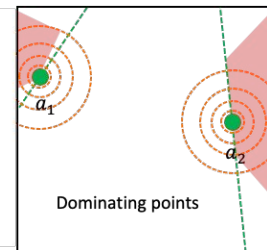
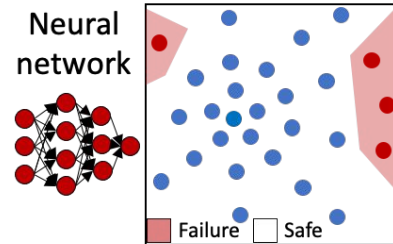
Benefits: Versatile, even to high-dimensions, given a sufficient training set

- Find dominating points using MIP reformulation

Benefits: Scalable in depth and complete, given ReLU-activated NNs

- Perform dominating-point-based IS and use NN predictions as labels

Benefits: Unbiased and faster (alleviating the need to run more simulations)



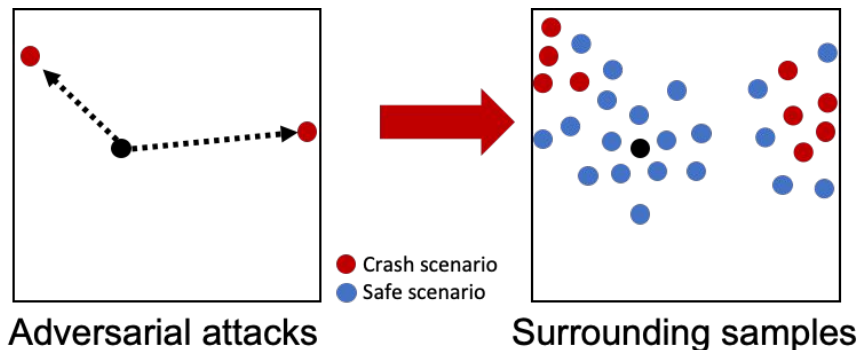
Deep IS: Unbiased, given an accurate approximation

- Suppose NN gives a set approximation $\hat{\mathcal{S}}_\gamma \approx \mathcal{S}_\gamma$ (the true failure rate) after training with n_1 samples. We have, with $n_2 = n - n_1$ samples,

$$\begin{aligned}
 \mathbb{E}_{X \sim \hat{p}}[\hat{\mu}_n] &= \mathbb{E} \left[\frac{1}{n_2} \sum_{i=1}^{n_2} \mathbb{1}(X_i \in \hat{\mathcal{S}}_\gamma) L(X_i) \right] \\
 &= \frac{1}{n_2} \sum_{i=1}^{n_2} \mathbb{E} \left[\mathbb{1}(X_i \in \hat{\mathcal{S}}_\gamma) \frac{\phi(\tilde{X}_i; \lambda, \Sigma)}{\sum_{a \in \hat{A}_\gamma} w_a \phi(\tilde{X}_i; a, \Sigma)} \right] \\
 &= \frac{1}{n_2} \sum_{i=1}^{n_2} \int_{\mathbb{R}^d} \mathbb{1}(X_i \in \hat{\mathcal{S}}_\gamma) \frac{\phi(\tilde{X}_i; \lambda, \Sigma)}{\sum_{a \in \hat{A}_\gamma} w_a \phi(\tilde{X}_i; a, \Sigma)} \sum_{a \in \hat{A}_\gamma} w_a \phi(\tilde{X}_i; a, \Sigma) dX_i \\
 &= \frac{1}{n_2} \sum_{i=1}^{n_2} \int_{\mathbb{R}^d} \mathbb{1}(X_i \in \hat{\mathcal{S}}_\gamma) \phi(\tilde{X}_i; \lambda, \Sigma) dX_i \\
 &= \frac{1}{n_2} \sum_{i=1}^{n_2} \mathbb{E}_{X \sim p} \mathbb{1}(X_i \in \hat{\mathcal{S}}_\gamma) \\
 &\approx \frac{1}{n_2} \sum_{i=1}^{n_2} \mathbb{E}_{X \sim p} \mathbb{1}(X_i \in \mathcal{S}_\gamma) \\
 &= \mu.
 \end{aligned}$$

Deep IS: Unlocks adversarial ML approaches

- Generate n_1 samples using adversarial attacks (FGSM, Boundary Attack) + surrounding samples.

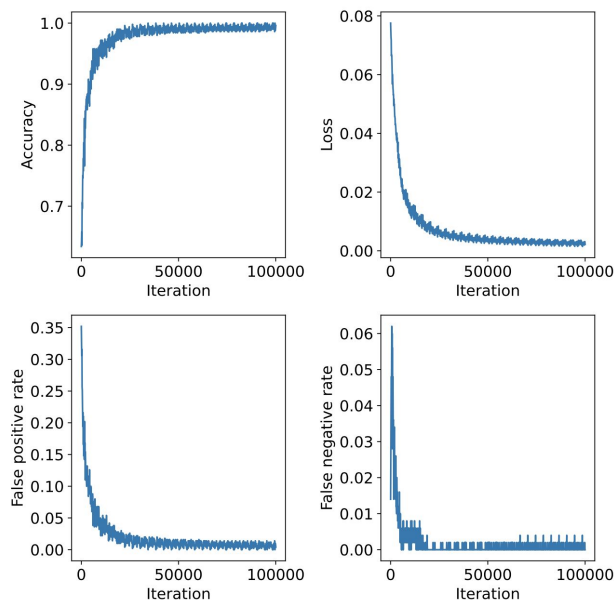


- Use log trick for the likelihood ratio during calculation

$$\log L(\tilde{X}_i) = \log \left(\frac{\phi(\tilde{X}_i; \lambda, \Sigma)}{\sum_{a \in \hat{A}_\gamma} w_a \phi(\tilde{X}_i; a, \Sigma)} \right) = \log \phi(\tilde{X}_i; \lambda, \Sigma) - \log \left(\sum_{a \in \hat{A}_\gamma} w_a \phi(\tilde{X}_i; a, \Sigma) \right) \quad 28$$

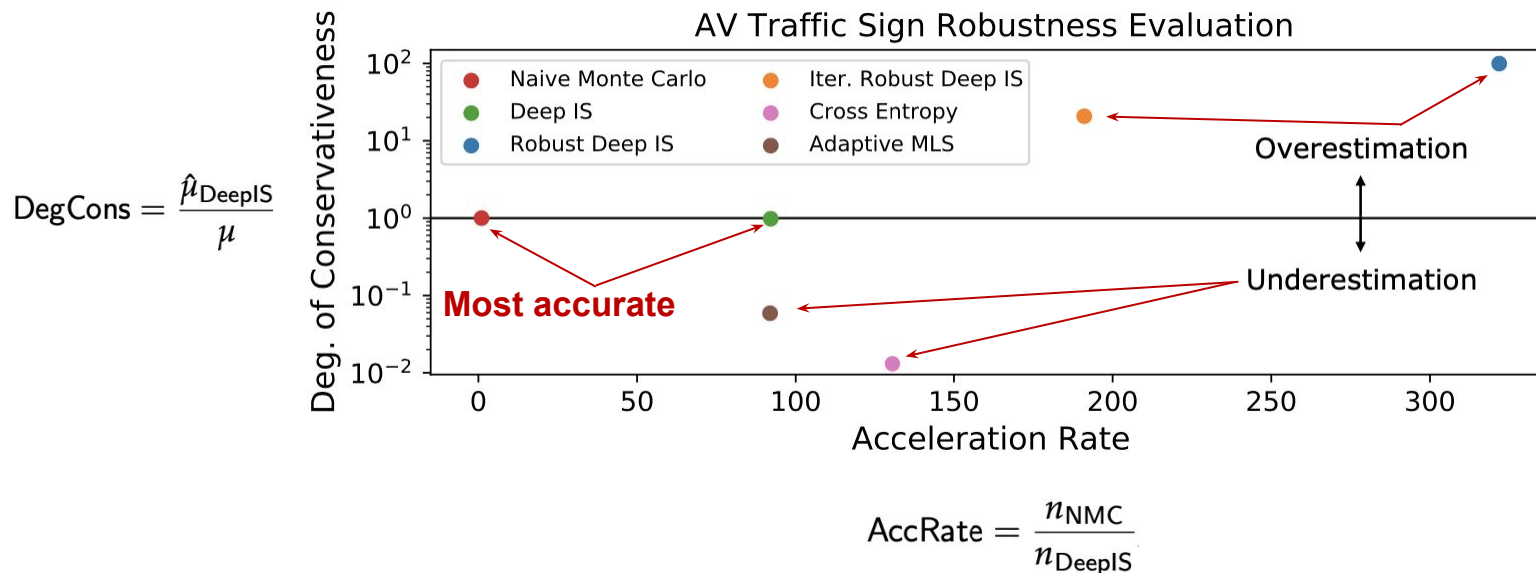
Deep IS: Numerical experiments

- **Deep IS classifier:** 4-layer feed-forward ReLU activated neural nets
- **Training:** 20,000 uniform Stage 1 samples, 512 batch size, Adam optimizer with L2 regularization (speeds up MIP by 20%).
- **ICP:** Terminates after 100 dominating points (some appear interpretable, most aren't)



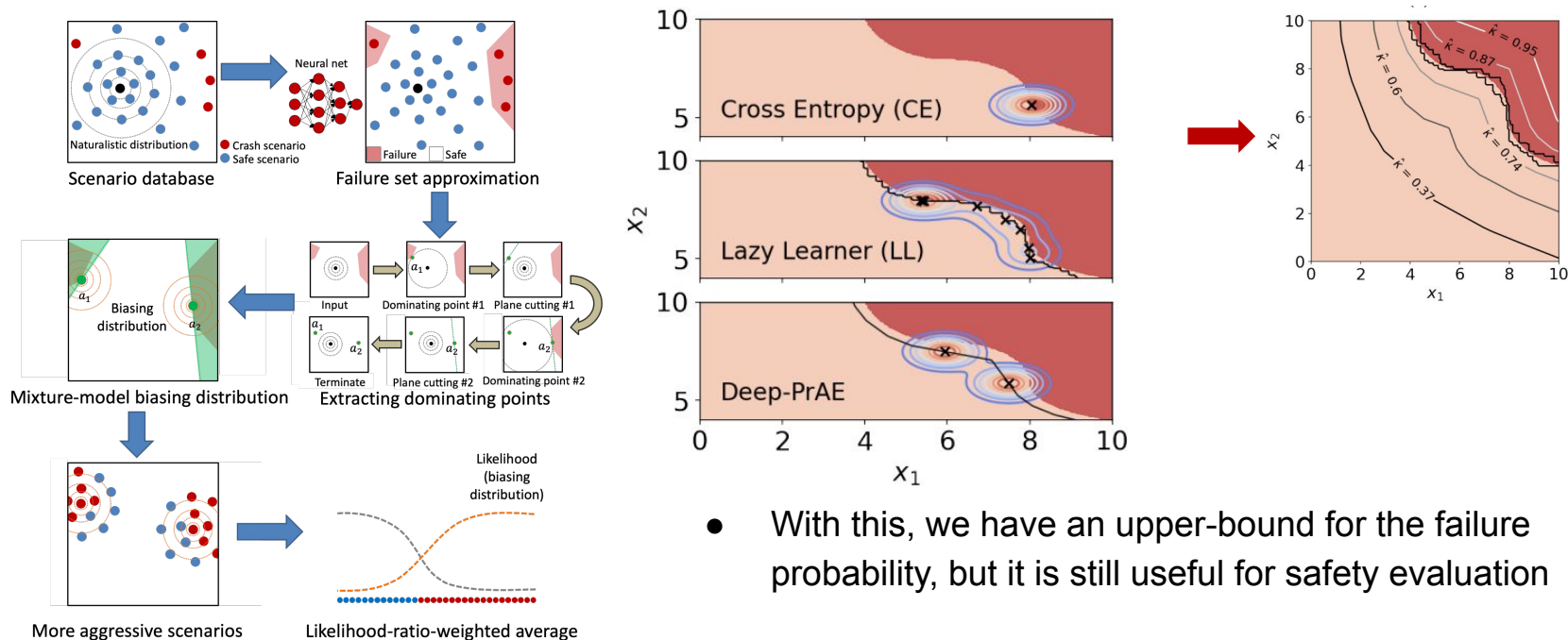
Deep IS: Numerical experiments

- Main result: Most accurate vs. other benchmarks (except huge NMC)



Further extensions: Deep-PrAE

- What if we have an error, can we prove efficiency? Yes, a conservative one!



- With this, we have an upper-bound for the failure probability, but it is still useful for safety evaluation

Robust planning & monitoring

- We cannot anticipate all corner cases during training.



In-context stop signs
during training



Rare, out-of-context signs
in the real world

Robust planning & monitoring

- We cannot anticipate all corner cases during training.

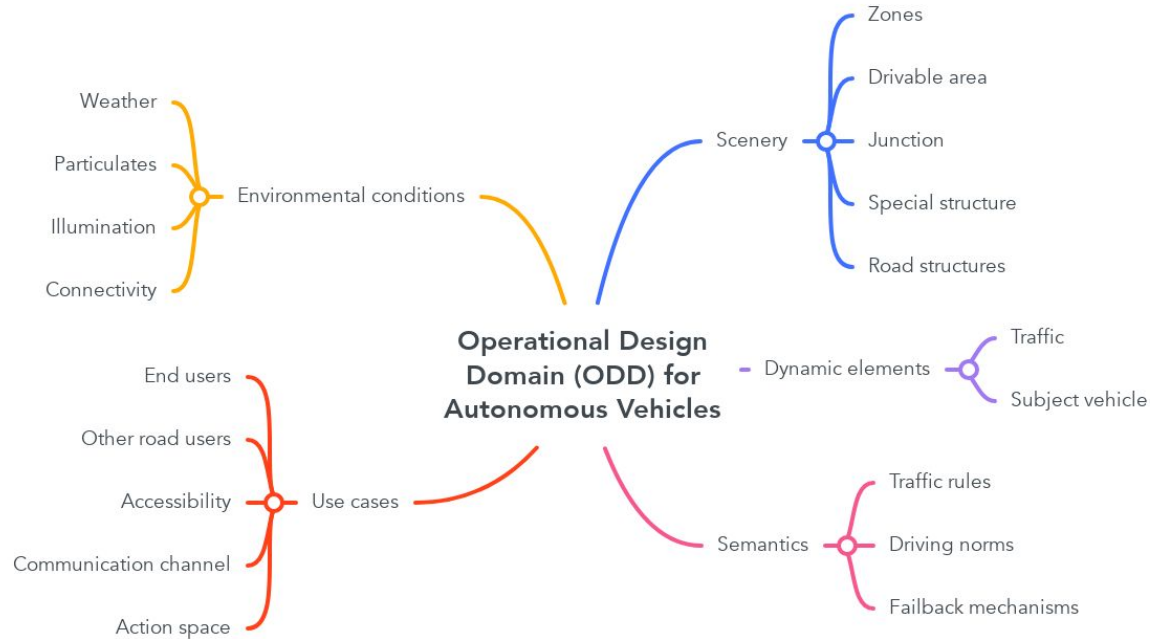


99.99% accuracy model
(success detection)



Same model
(misdetection when noisy)

Robust Operational Design Domain (ODD) Monitoring



ODD specifies the conditions for which the system **is designed** to function properly.

ODD-aware training and deployment improves safety



Importance Sampling-Guided Meta-Training for Intelligent Agents in Highly Interactive Environments

Mansur Arief, Mike Timmerman, Jiachen Li, David Isele, Mykel J. Kochenderfer. Under review.



Uncertainty Estimation & Out-Of-Model-Scope Detection Through Disentangled Concepts

Romeo Valentin, Sydney Katz, Dylan Asmar, Esen Yel, Mykel Kochenderfer

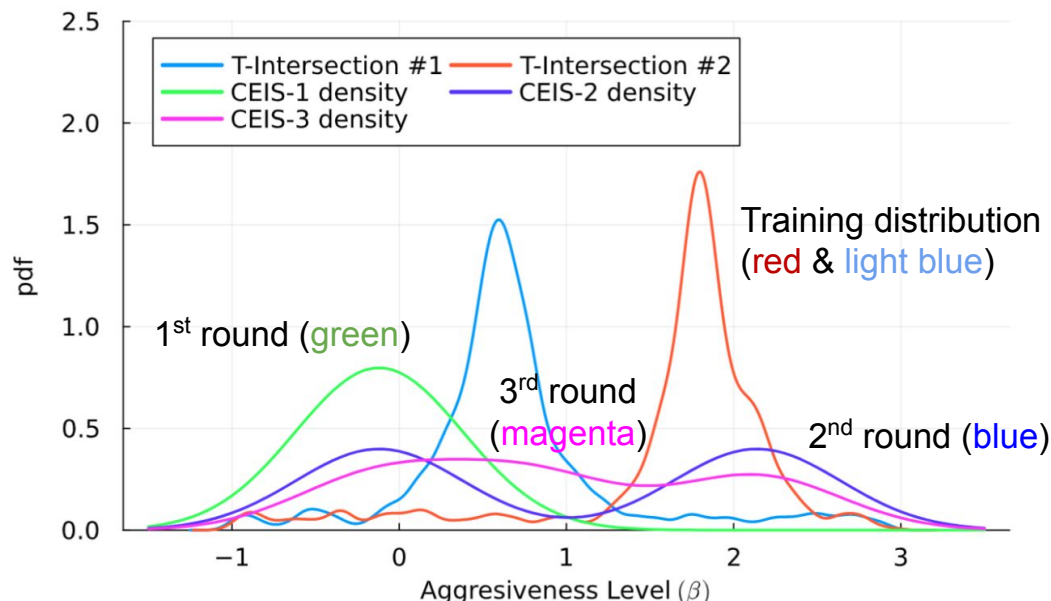


Efficient Safety Validation Using Meta-Learning

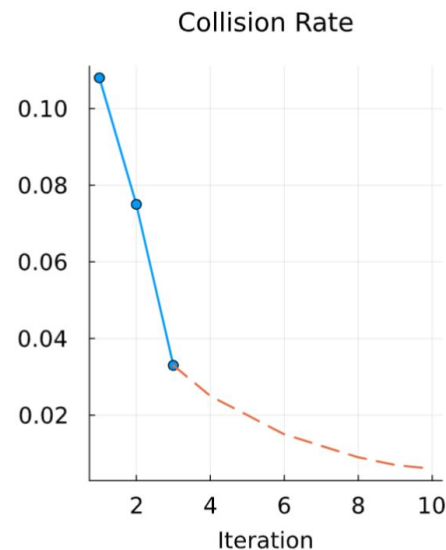
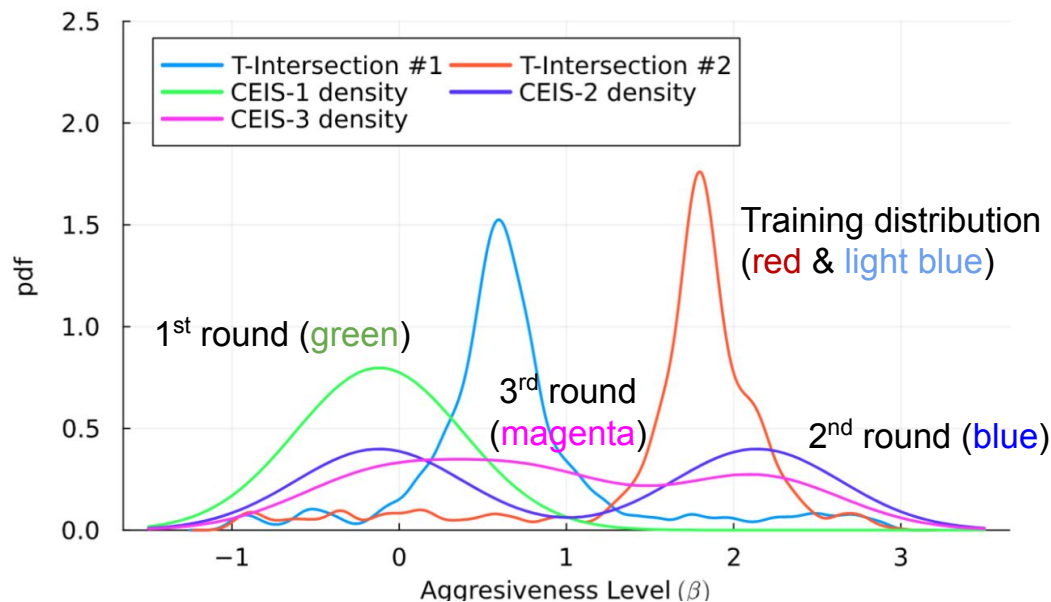
Marc R. Schlichting, Nina V. Boord, Anthony L. Corso, Mykel J. Kochenderfer. SAVME: Efficient Safety Validation for Autonomous Systems Using Meta-Learning. ITSC 2023.



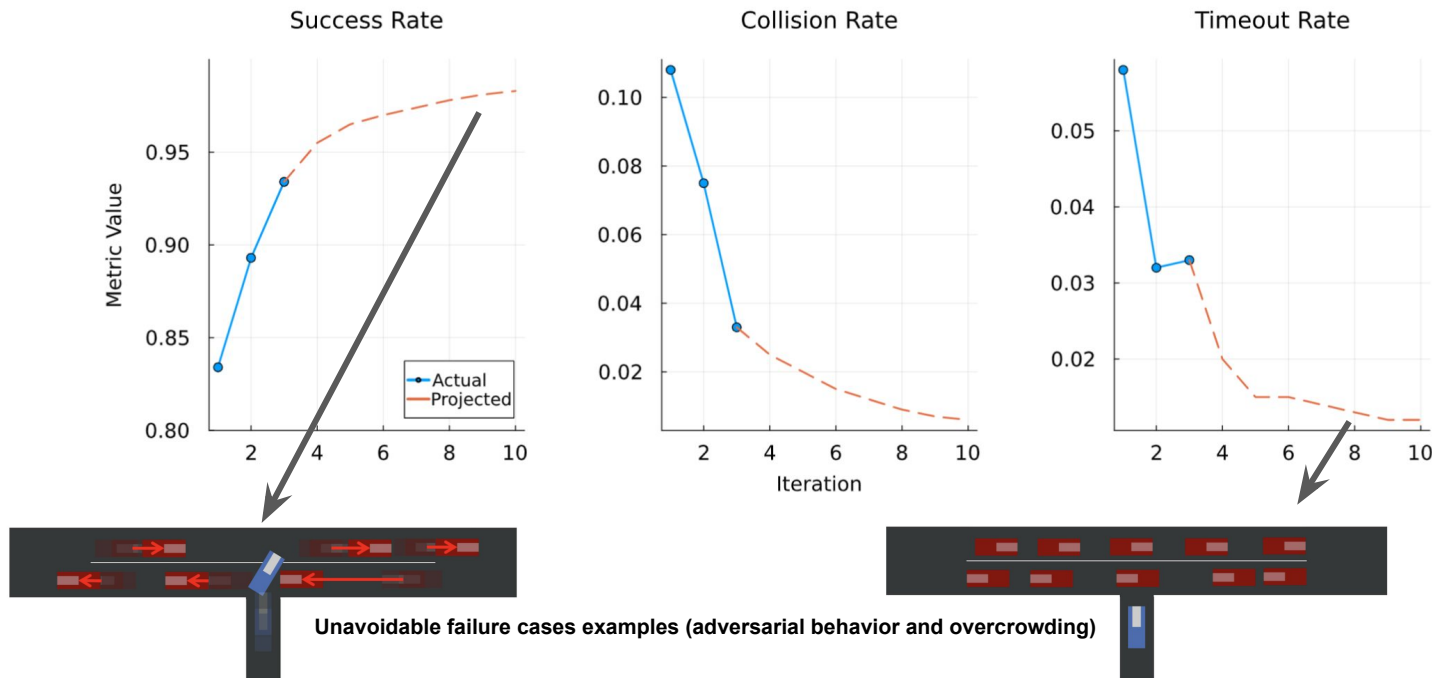
ODD-aware training and deployment improves safety



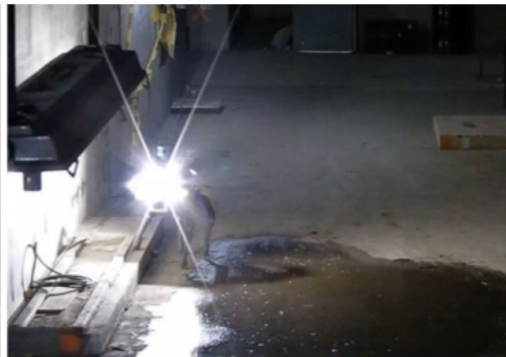
ODD-aware training and deployment improves safety



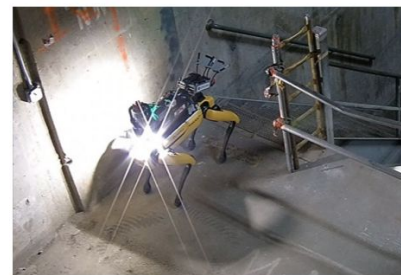
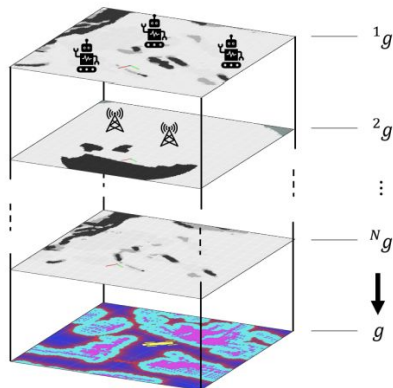
ODD-aware training and deployment improves safety



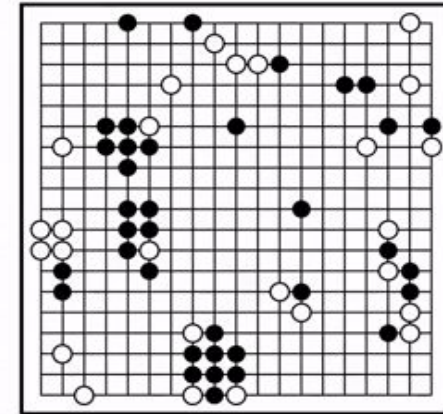
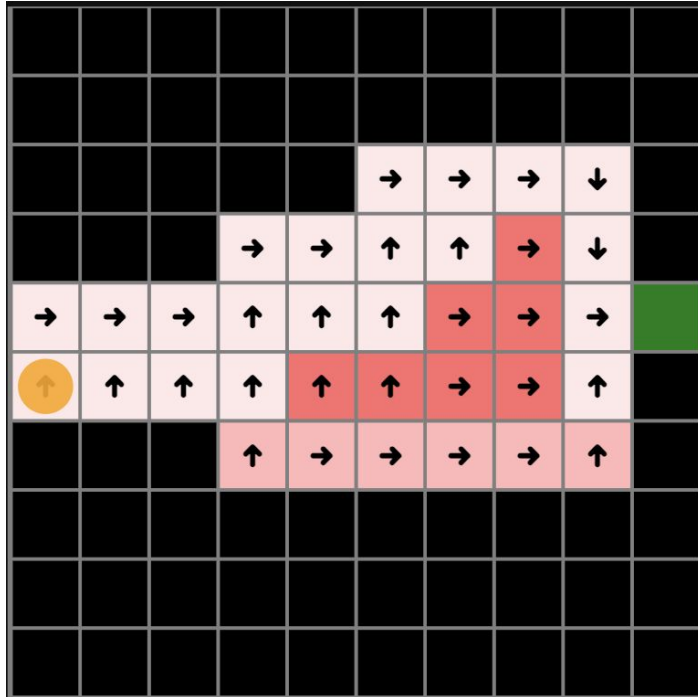
Safety-centered CPS Development



Main question:
Given uncertain and extreme conditions, how to explore safely and efficiently to fulfil mission objectives?



Applications: post mining, geosteering, blasting, etc.



Solved via AlphaGo-like simulations

Another challenge is vast outdoor exploration

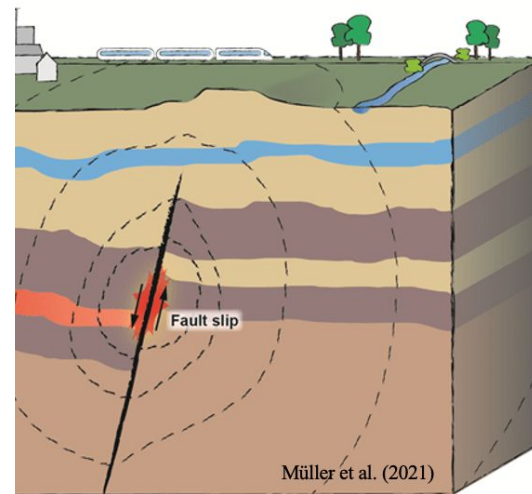


Sources of uncertainties: Noisy sensors, limited sensor range, vast area, moving obstacles

And, more importantly, safety!

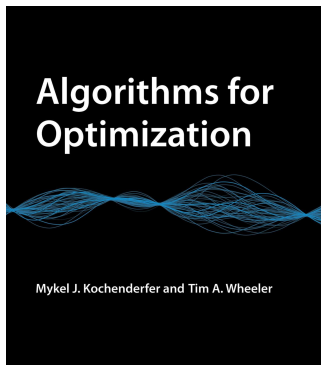


Safety risk for workers

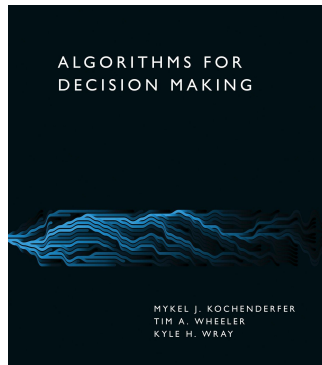


Risk of induced seismicity

Our approach toward safe intelligent autonomy



2019



2022



Coming soon!

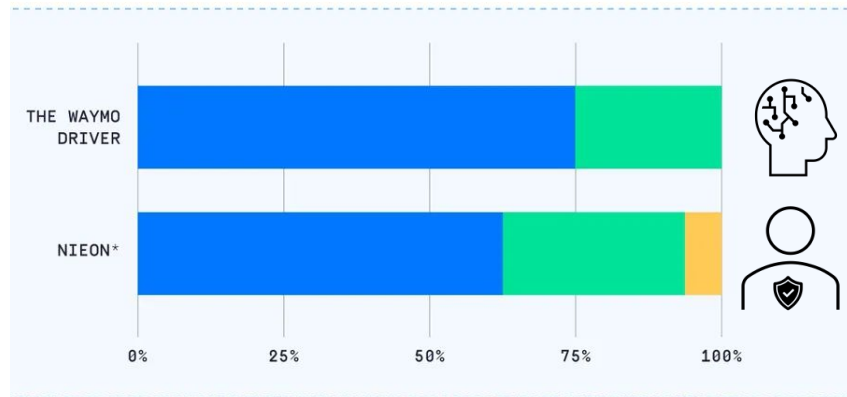
Soon: Algorithms for Validation book

Mykel Kochenderfer, Anthony Corso, Robert Moss, and Sydney Katz



AI systems have huge potential for improving safety

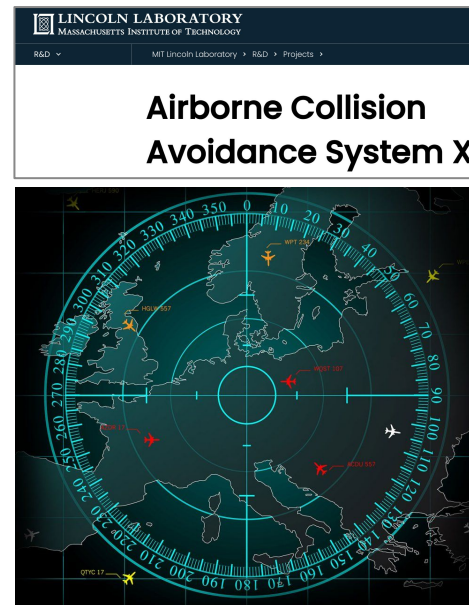
The Waymo Driver's collision avoidance performance in simulated tests



*NON-IMPAIRED, WITH EYES ALWAYS ON THE CONFLICT
HUMAN DRIVER THAT DOESN'T EXIST IN THE HUMAN POPULATION

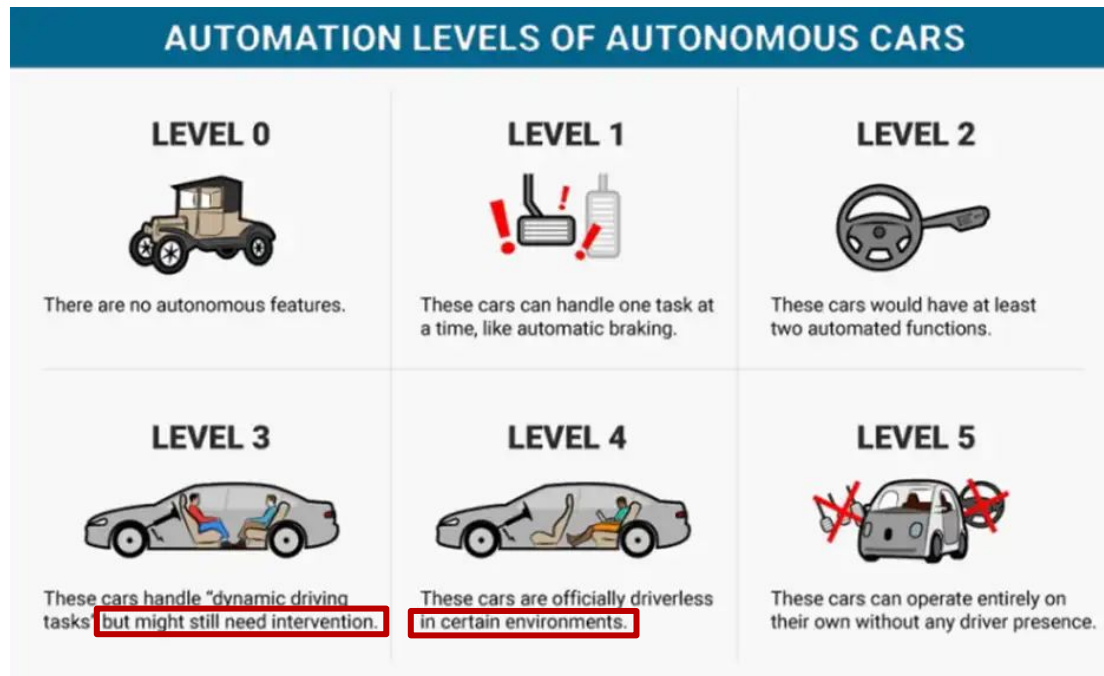
AVOIDED CRASH
MITIGATED CRASH
CRASH NOT MITIGATED

Source: <https://www.theverge.com/2022/9/29/23377219/waymo-av-safety-study-response-time-crash-avoidance>,
<https://waymo.com/waymo-one-san-francisco/>,



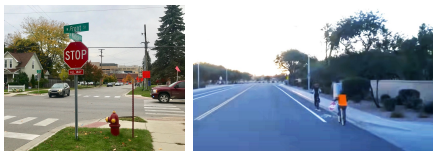
A next-generation collision avoidance system will help pilots and unmanned aircraft safely navigate the airspace.

But, we have to develop and deploy them cognizantly

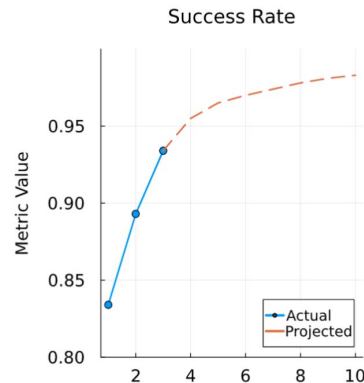


Collaboration opportunities

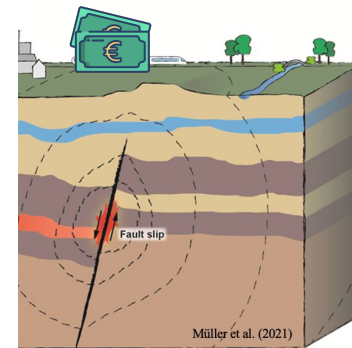
- How do we integrate airplane-level safety culture into the industry?



Runtime monitoring and
rigorous validation



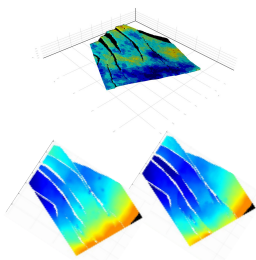
ODD-aware continuous
development



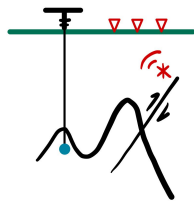
Risk-cognizant planning

Geothermal POMDP

Reservoir model

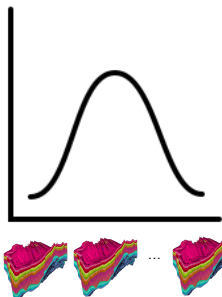


Well measurements



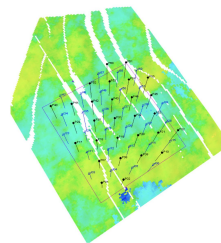
- Well temp
- Well pressure

Uncertainty



- Well temp
- Well pressure

Decisions to make



- Where to place wells?
- What rate of injection?

Goal: Max NPV, safety

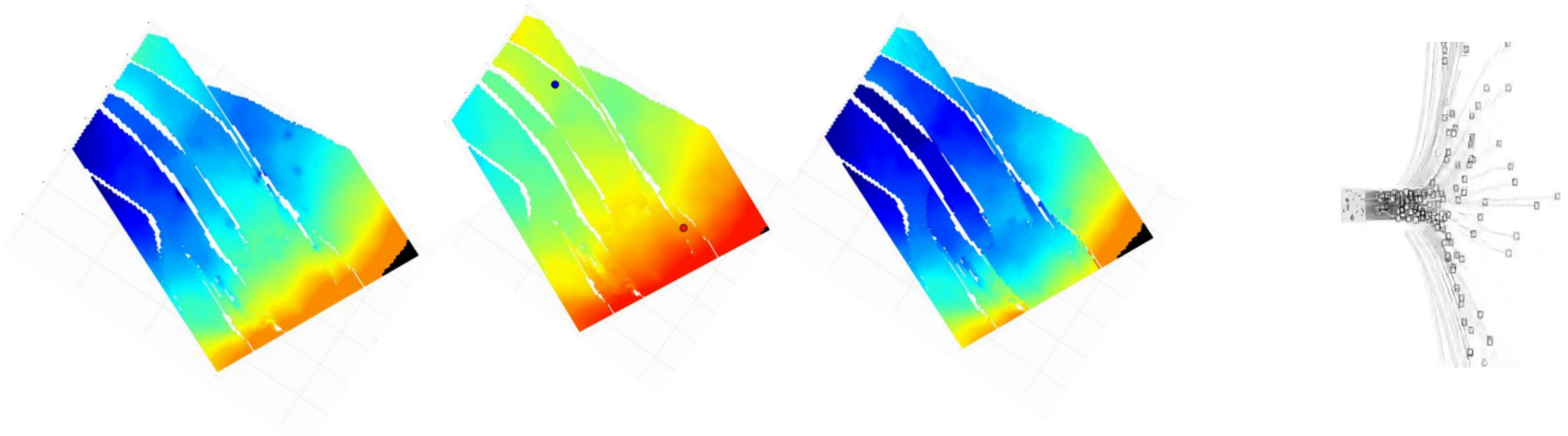


Reservoir simulator



- Our POMDP model ties together Earth and energy sciences, AI/data science, risk & safety, economics & business analysis

Our AlphaGo's approach for subsurface



```
save_config(sim_config, sim_config["FILEPATHS"]["config"])

println("Setting up action and running simulation...")
build_scenario(rw, sim_config)
run_intersect(sim_config)|

ro = ReservoirOutput(num_i=rw.num_i, num_j=rw.num_j, num_k=rw.num_k,
```

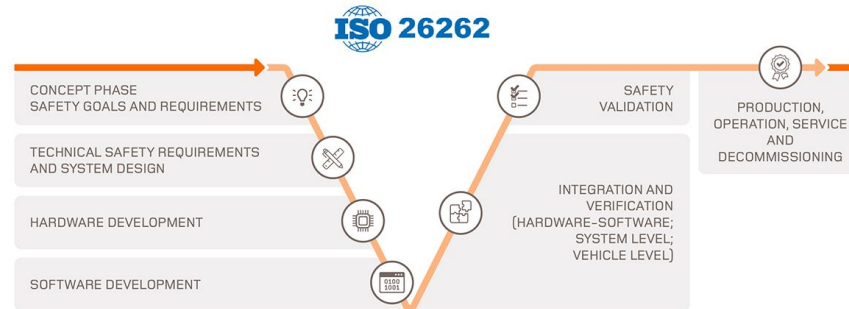
**We have run (automated)
~5k simulations to date...**

**= 1TB of data files (+1TB of
simulation files, compressed)**

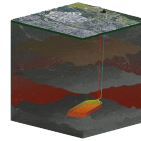
Unique Research Directions in MAE

Rigorous and scalable safety validation **Robust planning & monitoring** **Safety-centered CPS development**

FUNCTIONAL SAFETY SUPPORT THROUGHOUT THE DEVELOPMENT CYCLE



Transportation



Sustainability



Manufacturing

Thank you!

Mansur Arief

Email: mansur.arief@stanford.edu

Stanford Intelligent Systems Lab (SISL) and MineralX