



THE UNIVERSITY OF
MELBOURNE

Cyber-Physical System Security and Adversarial Machine Learning

Prof. Tansu Alpcan

Electrical & Electronic Engineering
The University of Melbourne

CSIRO/Data61 Seminar

20 July 2022



Outline

- **Securing Autonomous Vehicle Platoons**
 - V2V cyber-physical system security
 - Attacker model and defence method
 - Simulation results
- **Improving Adversarial Robustness Coding Theory**
 - Coding theory and adversarial robustness in DL
 - Effective Error-correcting output code (eECOC)
 - Neural Network Embedded Coding (NNEC)
 - Experimental results and analysis
- **Cyber(-Physical) Security Games**
- **Ongoing Research and Future Directions**





THE UNIVERSITY OF
MELBOURNE

Securing Autonomous Vehicle Platoons and V2V Networks

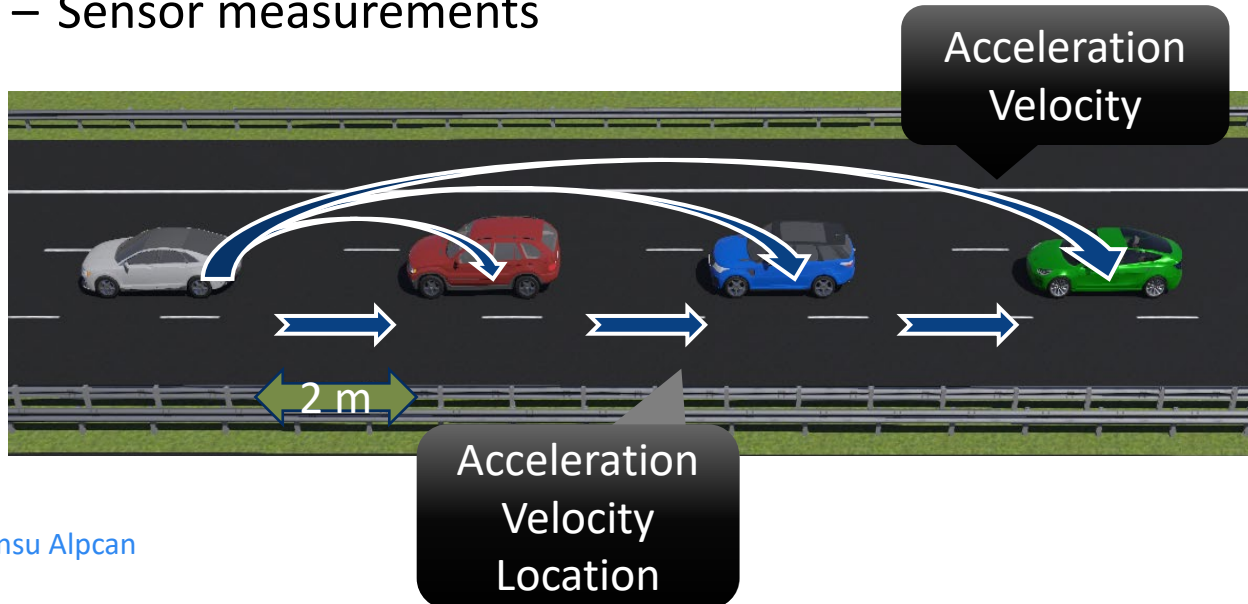
with Mr. Guoxin Sun
(PhD Student)

Simulation screenshot (Guoxin Sun)



Motivation

- **Autonomous vehicle platoons:**
 - A string of vehicles travelling as a single unit from an origin to a destination.
- **Platoon** maintains a **narrow inter-vehicle distance** and **relative velocity**, using
 - Wireless communication
 - Sensor measurements



Simulations with Webots and Sumo

Vehicle Platoon Model and Control

Information flow topology :

- Predecessor-leader following (PLF) – Each vehicle receives dynamics information from its immediate preceding vehicle and the leader vehicle.

Two control policies:

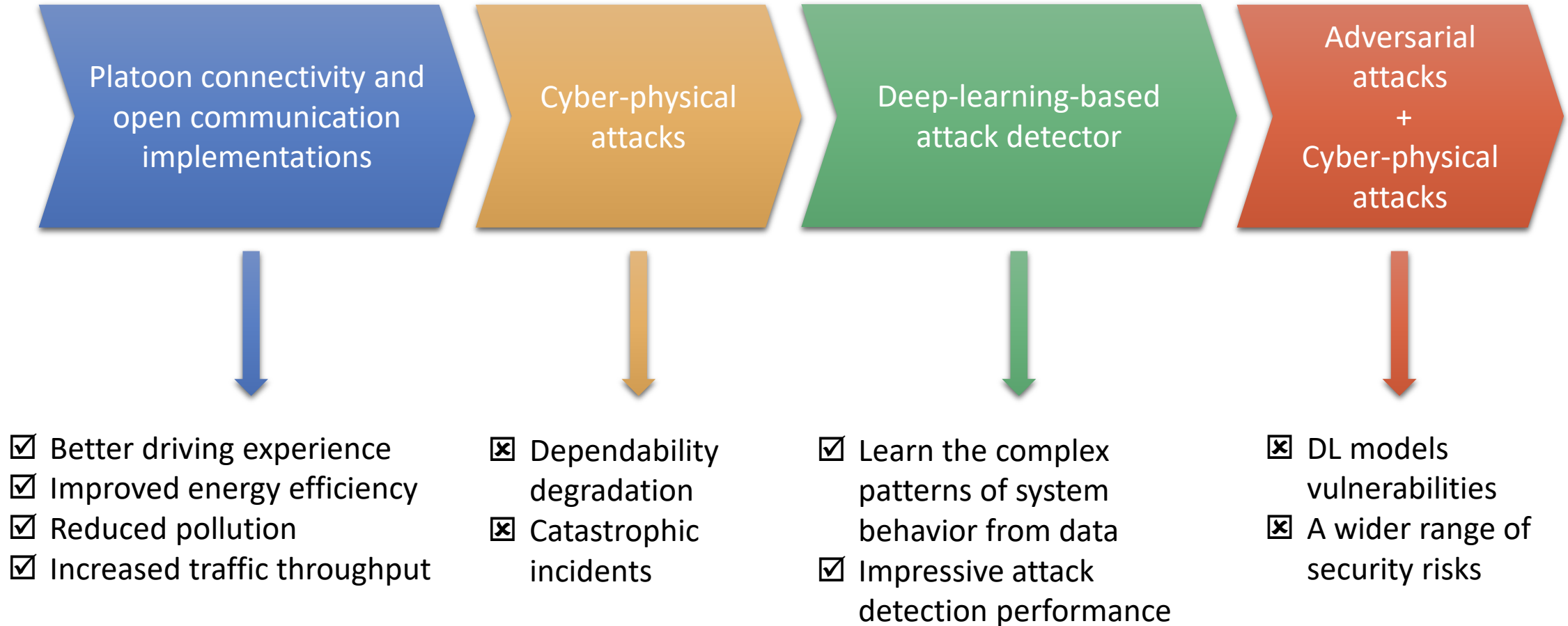
- Cooperative adaptive cruise control (CACC) – communication based
- Adaptive cruise control (ACC) – sensor based

Security Challenge:

- V2V communications in vehicle platoons are a target for cyber-physical attacks.



Problem Analysis



Attacker Model

False data injection (FDI): corrupts the content of wirelessly transmitted messages or sensor observations to cause performance degradation or catastrophic failure of safety-critical systems.

1. **Conventional Cyber-Physical Attacks** (to cause collisions in the platoon case)
 - I. **Vanilla False Data Injection Attack (v-FDI)** (same as FDI)
 - II. **Model-Aware False Data Injection Attack (m-FDI)** (the attack knows the underlying system model)
2. **Adversarially-Masked Cyber-Physical Attacks** (the attack deceives the anomaly detector)

Gradient based white-box adversarial attack: basic iterative method (BIM)

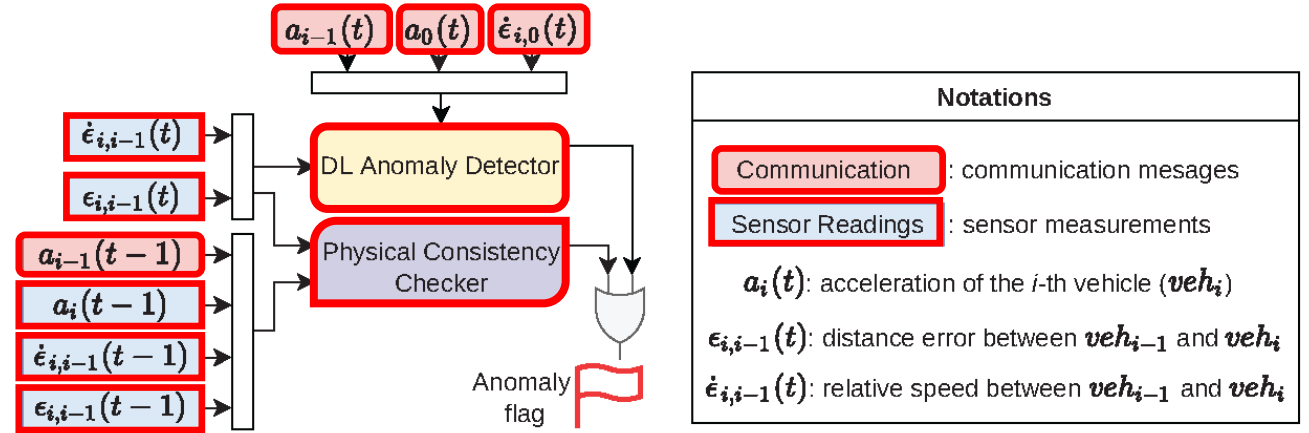
Table 1: Knowledge required by the attacker to conduct different attacks.

Attack Types \ Access to	Sensors	Communication	DL Model	System Model	Memory
v-FDI	✓	✓	✗	✗	✗
m-FDI	✓	✓	✗	✓	✗
v-FDI (adv. masked)	✓	✓	✓	✗	✓
m-FDI (adv. masked)	✓	✓	✓	✓	✓

Defense Method

DL-based Anomaly Detector: detects conventional cyber-physical attacks when existing modelling techniques fail to model the system accurately and reliably.

Physical Consistency Checker: assists in reporting adversarial perturbations to compensate for the deficiency of deep learning models



Algorithm 1 Double-Insured Anomaly Detection (DAD)

Input: Communication messages S and sensor readings R

Output: Anomaly flag

- 1: Initialization()
- 2: while Destination is not reached do
- 3: S and R are received and measured
- 4: $hist \leftarrow$ Load one-step history data
- 5: $flag1 \leftarrow$ AnomalyDetector($R, hist$)
- 6: $flag2 \leftarrow$ PhysicalConsistencyChecker($S, R, hist$)
- 7: if $flag1$ or $flag2$ is TRUE then
- 8: Anomaly flag \leftarrow Anomaly
- 9: else
- 10: Anomaly flag \leftarrow Normal
- 11: end if
- 12: end while

Double-Insured Anomaly Detection Algorithm

- The idea is to **detect anomalies using both deep learning pattern recognition and physical consistency check** using the underlying physical model (domain knowledge).
- It is important to consider (whenever possible) the **underlying physical system model** when addressing cyber-physical system security problems!

Algorithm 1 Double-Insured Anomaly Detection (DAD)

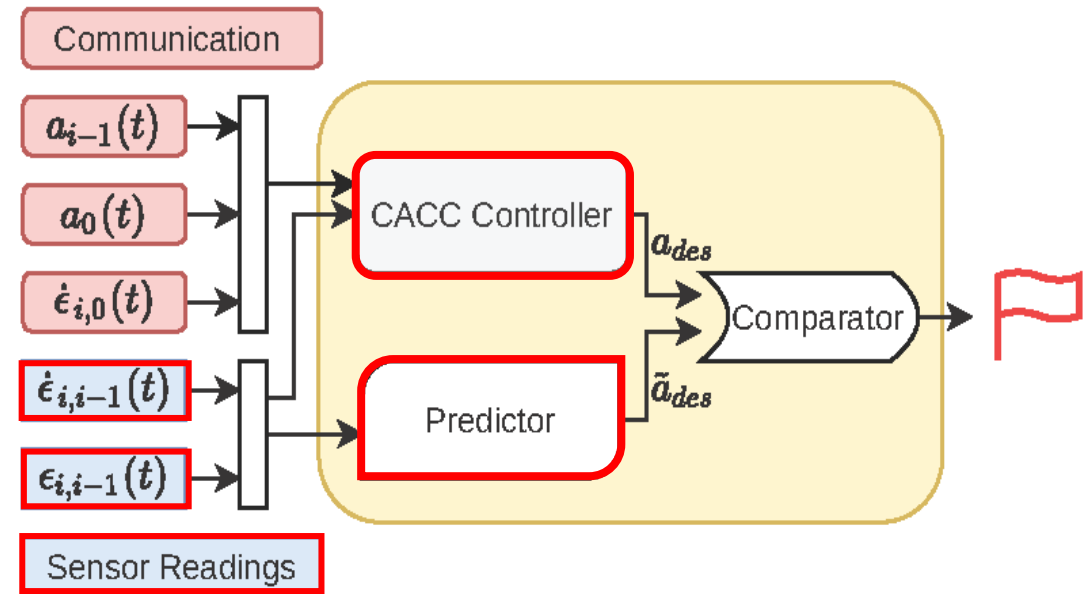
Input: Communication messages S and sensor readings R

Output: Anomaly flag

```
1: Initialization()
2: while Destination is not reached do
3:   Vehicle receives  $S$  and measures  $R$ 
4:    $hist \leftarrow$  Load one-step history data
5:    $flag1 \leftarrow AnomalyDetector(R, hist)$ 
6:    $flag2 \leftarrow PhysicalConsistencyChecker(S, R, hist)$ 
7:   if  $flag1$  or  $flag2$  is TRUE then
8:     Anomaly flag  $\leftarrow$  Anomaly
9:   else
10:    Anomaly flag  $\leftarrow$  Normal
11:  end if
12: end while
```

Data-Driven Anomaly Detector

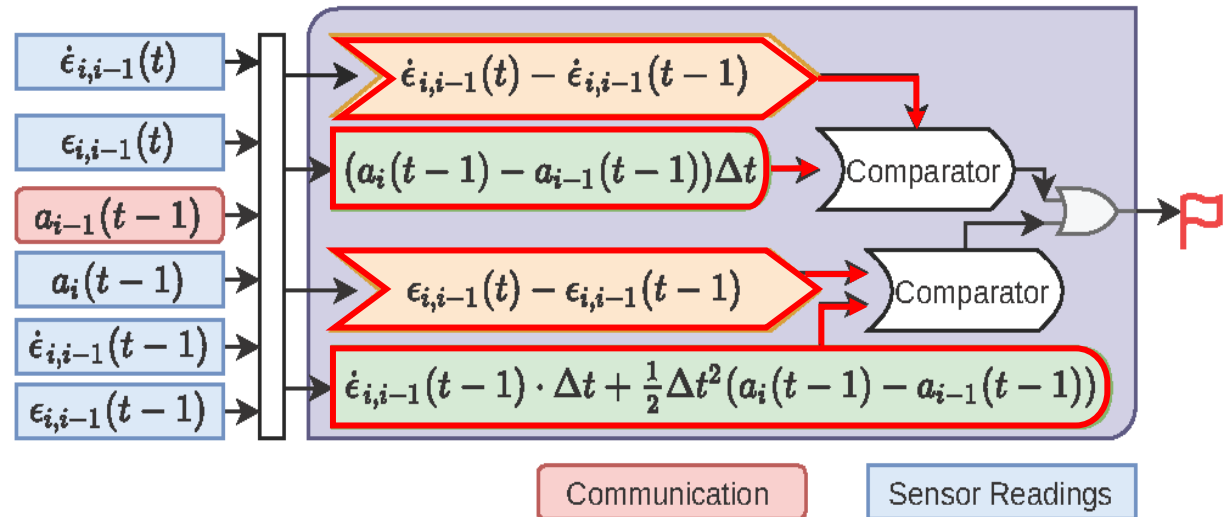
- Each vehicle obtains the relative **speed** and **distance** with respect to its predecessor
- **Only sensor measurements are used as detector input**
- **Predictor:**
 - Semi-supervised training
 - LSTM regression model to utilize the **temporal** information within the data
 - **Outputs** the expected desired acceleration value at the current time instance
- **Comparator**
 - **Computes the difference between the predictor output and controller output (which can be attacked!)**
 - **Raise anomaly flag when the deviation is large**



Physical Consistency Checker

- Corrupted controller inputs may not obey the underlying physical processes of the cyber-physical system.
- Simple kinematic model:

$$v_i(t) = v_i(0) + a_i t, \quad x_i(t) = x_i(0) + v_i(0)t + \frac{1}{2} a_i t^2$$
- It consists of **speed checker** and **distance checker**.
- Again, sensor- and communication-based (possibly attacked) results are compared.
- This model is domain specific. This defense approach can be generalized, for example, to power system domain or MLaaS.





Evaluation - Simulation Setup

- *Webots* simulation platform provides an efficient way of constructing different cyber-physical attacks and generate relevant training data.
- Platoon and traffic simulation
 - 4 **calibrated** BMW X5 **vehicles** (Webots models) on a highway segment
 - **Multiple sensors** such as Radar, transmitters and receivers
- **Traffic environment** includes
 - four types of vehicles (i.e., motorcycles, light-weight vehicles, trucks and trailers)
 - Various driving characteristics (i.e., cooperative or competitive)
- **Success criteria** are attack detection and inter-vehicular distance.

Conventional Cyber-physical Attacks

➤ Baseline attack detectors:

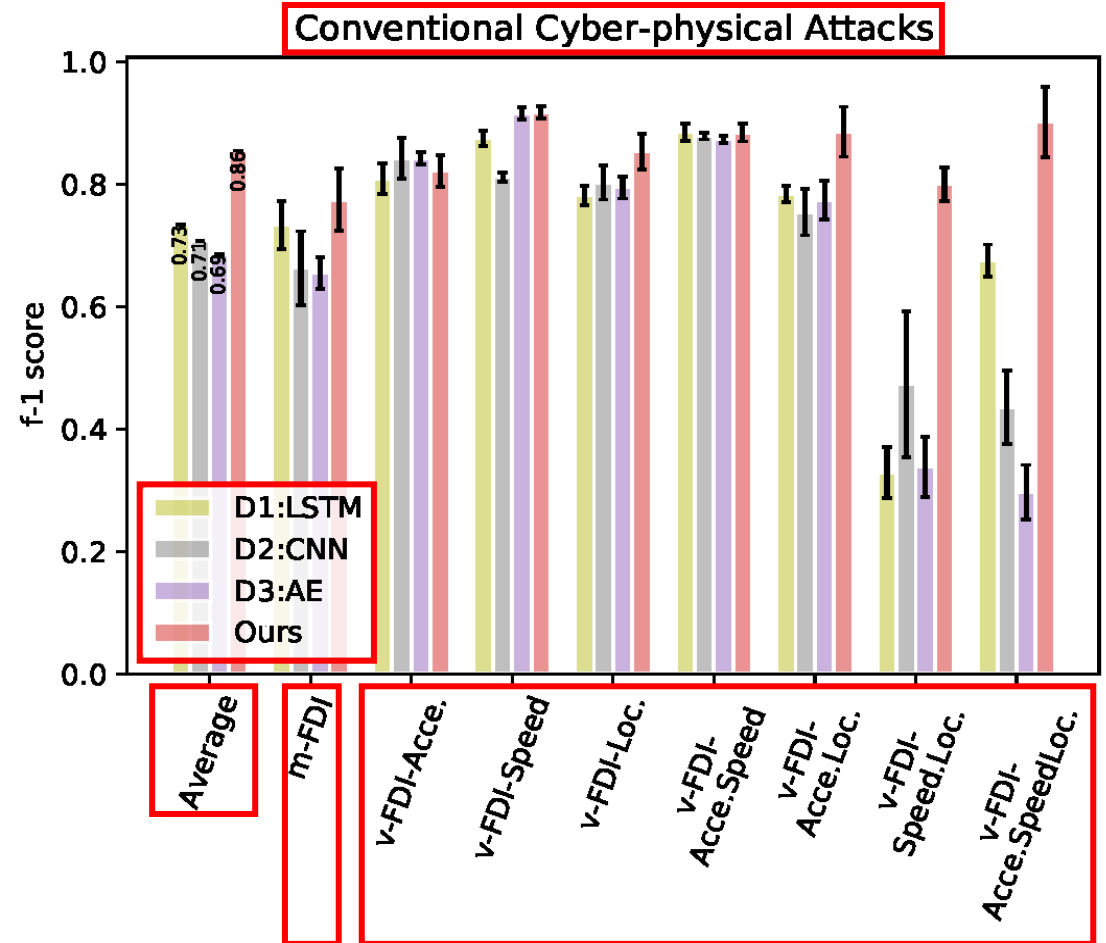
- D1: LSTM - Long Short-Term Memory
- D2: CNN - Convolutional Neural Networks
- D3: AE - Auto-Encoder

➤ Comparison metric:

- F1 score - the harmonic mean of the precision and recall

➤ Highlights:

- Data-driven detection methods in general perform well against such conventional attacks
- Our proposal slightly outperforms the baselines



Adversarially-masked Cyber-physical Attacks

➤ Baseline attack detectors:

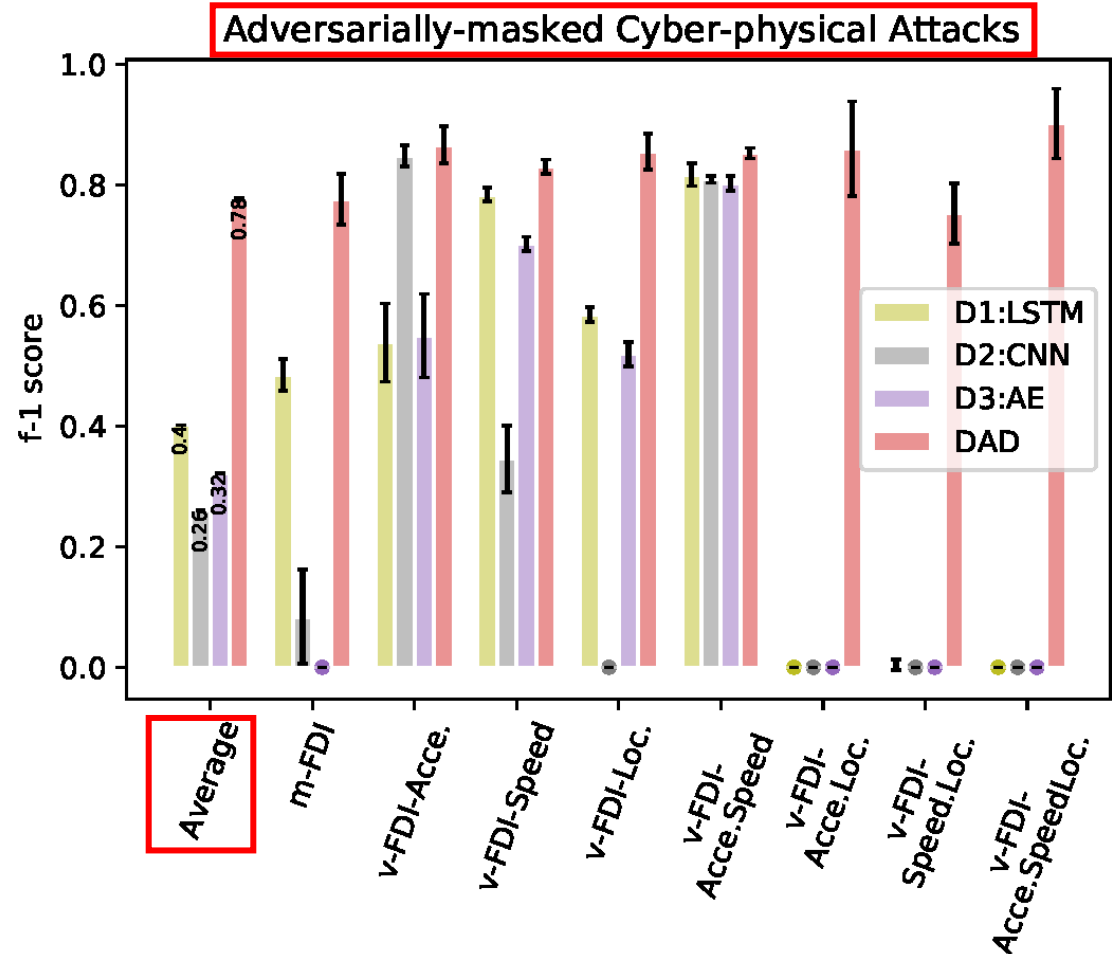
- D1: LSTM - Long Short-Term Memory
- D2: CNN - Convolutional Neural Networks
- D3: AE - Auto-Encoder

➤ Comparison metric:

- F1 score - the harmonic mean of the precision and recall

➤ Highlights:

- Baselines suffer from adversarial attacks
- Our proposal doubles the detection F1 score compared to the baselines



Model-aware False Data Injection (m-FDI)

- **Baseline attack detectors:**
 - D1: LSTM - Long Short-Term Memory
 - D2: CNN - Convolutional Neural Networks
 - D3: AE - Auto-Encoder
 - PCC - Physical Consistency Checker
 - D1: LSTM* - Adversarially retrained LSTM
- **Comparison metric:**
 - F1 score and Recall

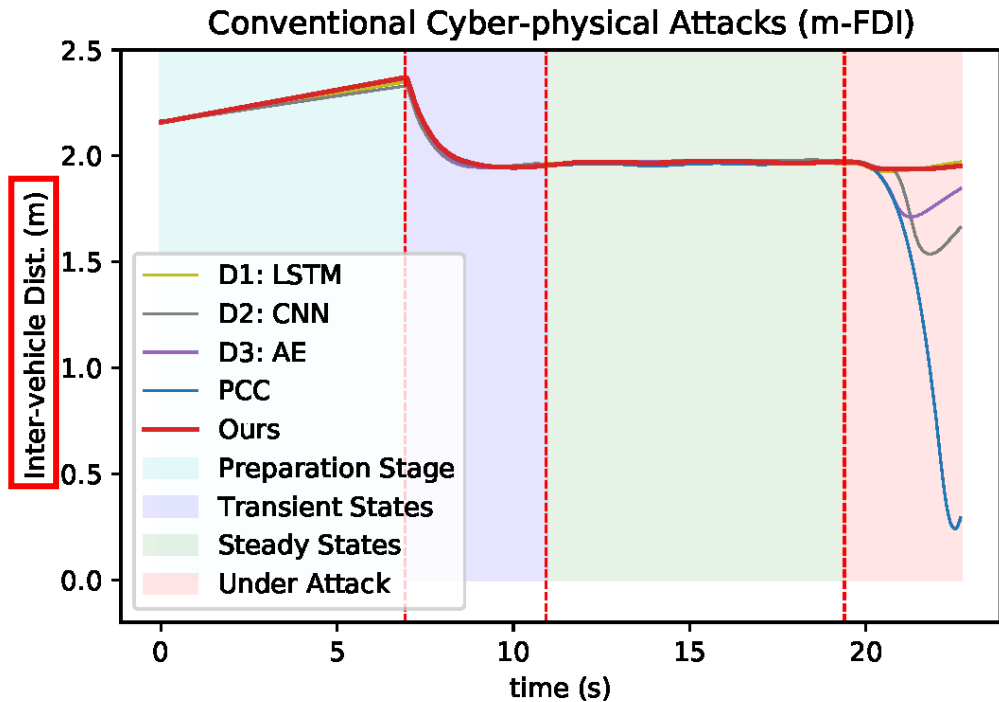
Table 2: Attack Detection Results against m-FDI with Different Detection Methods. * denotes adversarial training.

Attack	m-FDI		m-FDI (adv. masked)		
Defense	Rec	F1	Defense	Rec	F1
D1: LSTM	0.70	0.73	D1: LSTM	0.39	0.49
D2: CNN	0.57	0.66	D2: CNN	0.05	0.08
D3: AE	0.56	0.66	D3: AE	0.00	0.00
PCC	0.18	0.29	PCC	0.63	0.75
Ours: DAD	0.77	0.77	Ours: DAD	0.77	0.78
D1: LSTM*	0.70	0.73	D1: LSTM*	0.48	0.56
Ours: DAD*	0.75	0.76	Ours: DAD*	0.84	0.81

Results:

- Physics component (PCC) *when used alone* is defeated by powerful m-FDI attack.
- Our approach can be combined with existing adversarial defense approaches (e.g., adversarial training) to further enhance detection performance.

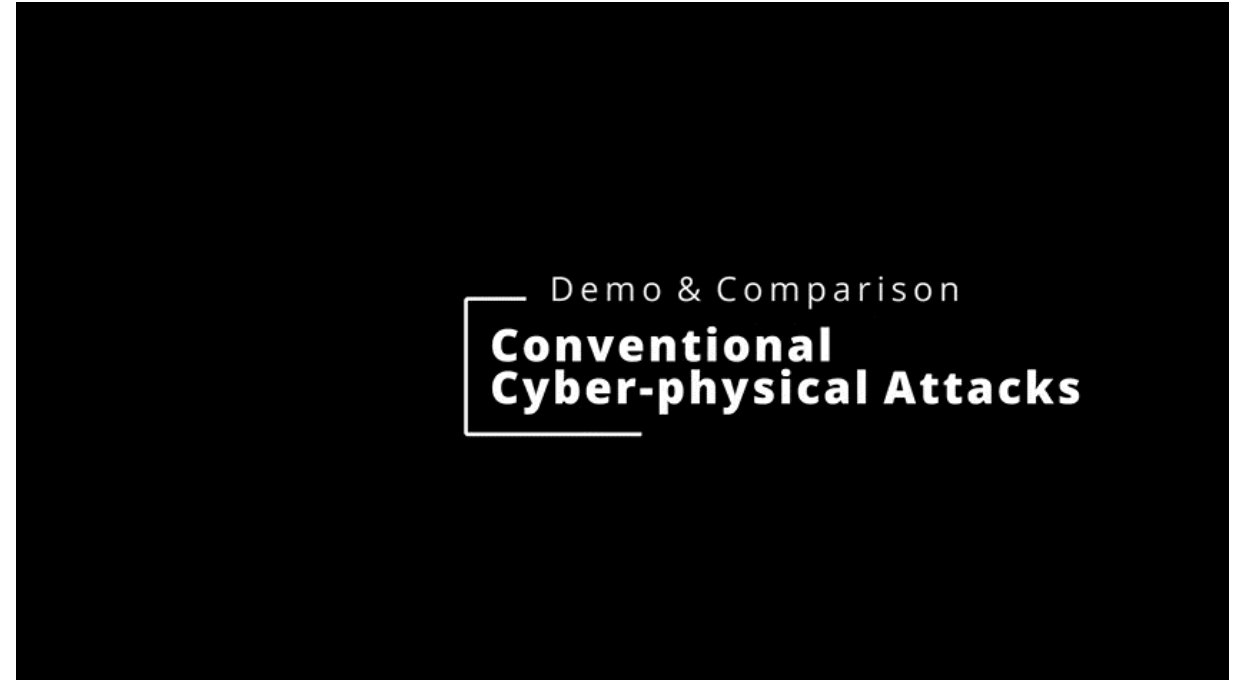
Simulation Demonstration



Comparison metric:

- Inter-vehicle distance – to measure the danger faced by the platoon

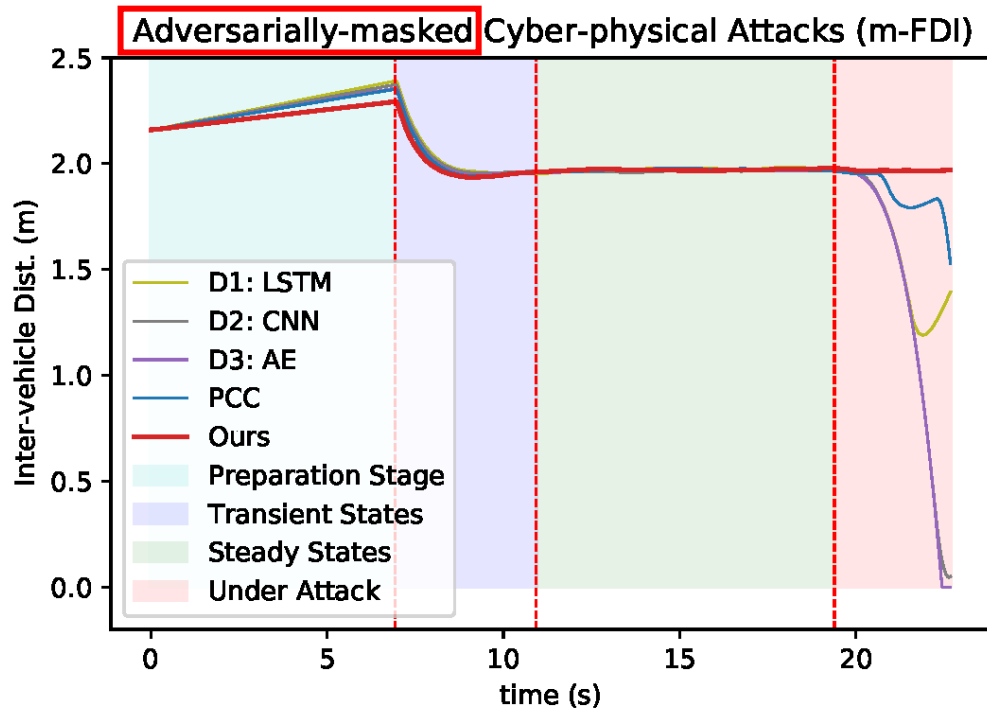
Tansu Alpcan



Results:

- Our approach results in nearly unnoticeable fluctuation throughout the entire attack period.

Simulation Demonstration



Comparison metric:

- Inter-vehicle distance – to measure the danger faced by the platoon

Highlights:

- Our approach results in the largest inter-vehicle distance throughout the entire attack period.



Contributions

- A novel **physics-enhanced data-driven attack detection system** for cyber-physical systems that leverages knowledge from both **data** and **physics**.
- Classical physics-modelling techniques can help to **mitigate the deficiency of deep learning**-based approaches, which extends the applicability of many state-of-the-art DL-based approaches for cyber-physical systems.
- As a demonstration, we successfully improve the security and dependability of **vehicle platoons**. Our defense system provides excellent detection performance against an informed white-box attacker.
- Our results are demonstrated using **sophisticated simulations**. It outperforms standard baseline attack detection methods and shows the **potential of application together with existing adversarial defense** techniques for better performance.



Ongoing and Future Work

- Applications of our approach to modern communication/computing systems, specifically Machine Learning as a Service (MLaaS).
- Domain-specific **physical models and systems-based approach** significantly better than pure machine learning techniques. Furthermore, our approach improves applicability of ML/DL to real-world scenarios.
- Further emphasis on **defence analysis and methods using game theory**, which has huge potential for cyber-physical system security.

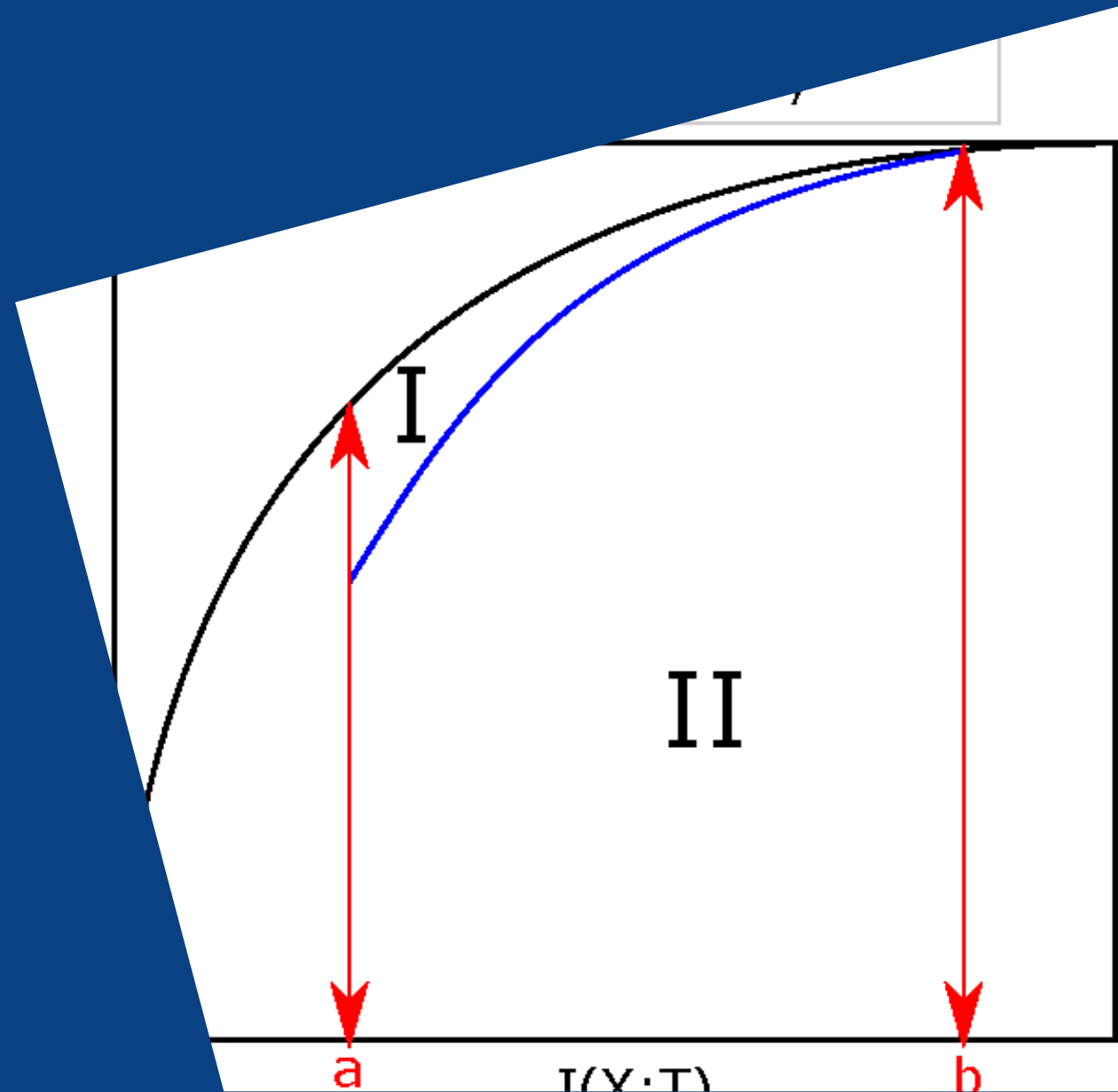


THE UNIVERSITY OF
MELBOURNE

Improving Adversarial Robustness Using Information and Coding Theories

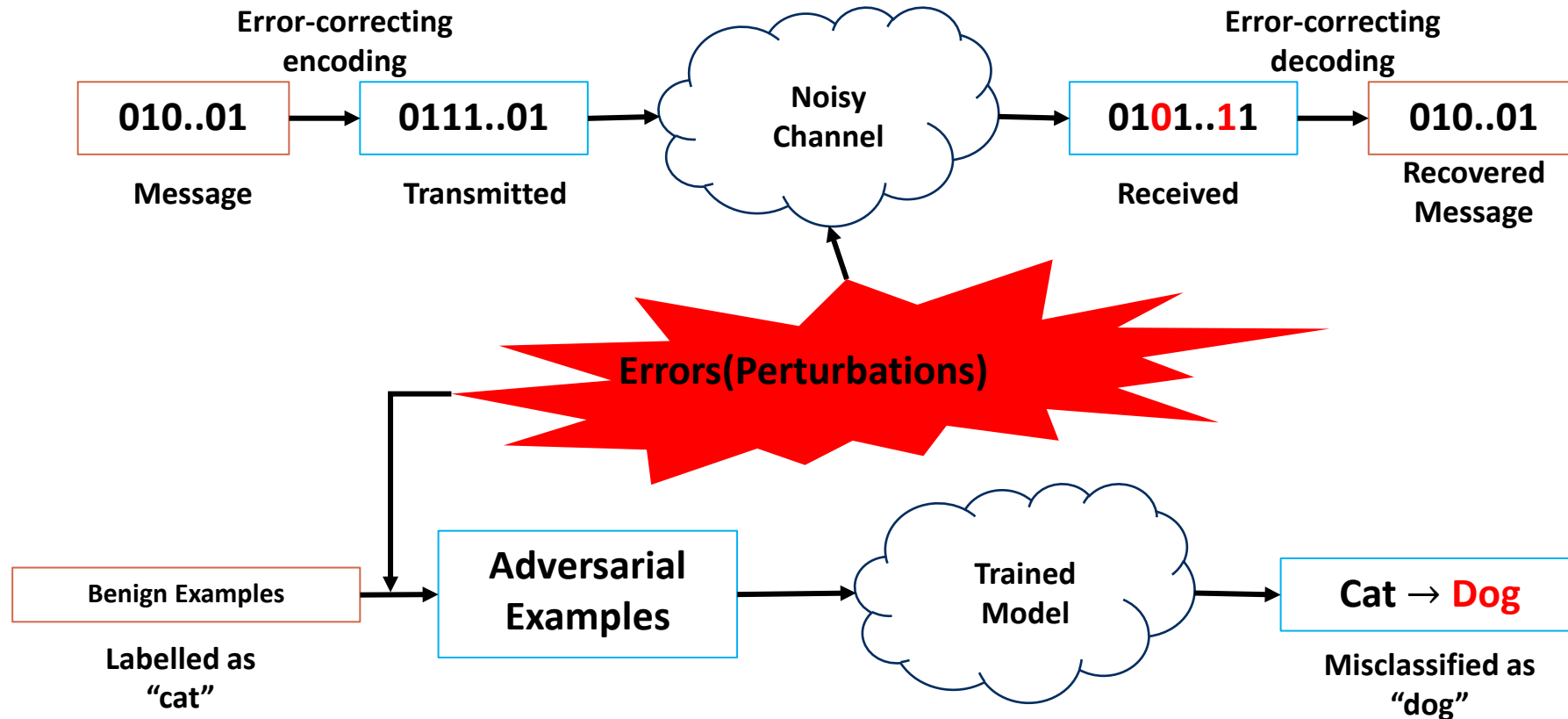
with Mr. Li Wan
(PhD Student)

Information Plane (Li Wan)



Motivation and Problem

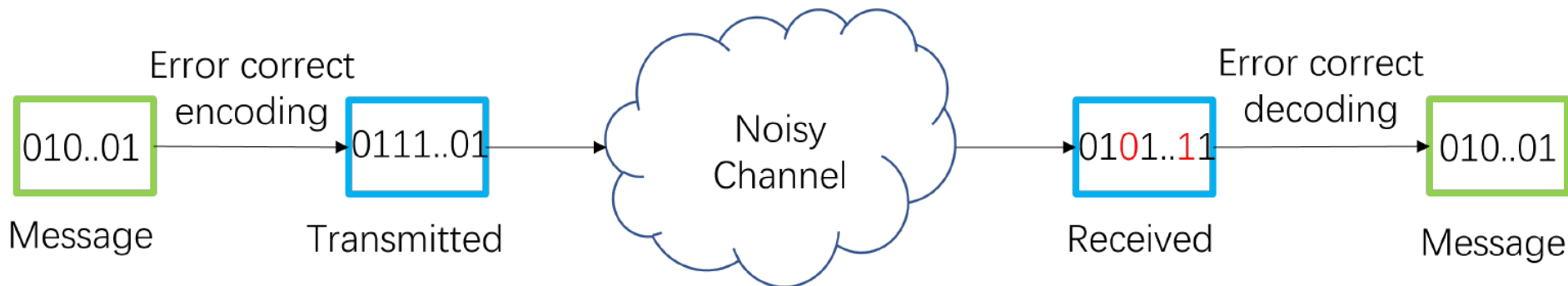
There are multiple interesting connections between **information theory and (adversarial) deep learning**, e.g. information bottleneck, error correcting output codes, etc.



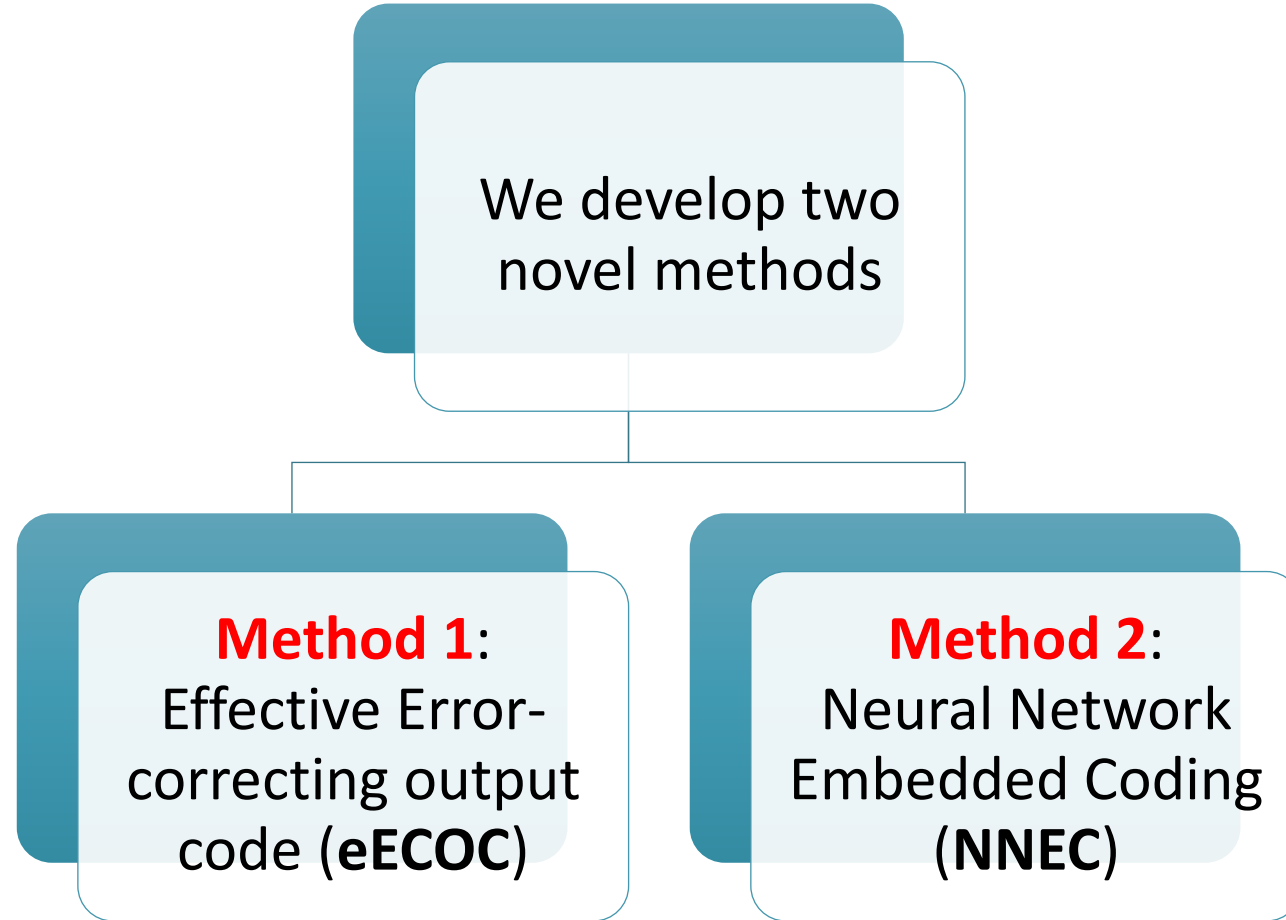
Research Questions

Research questions we try to address:

- How can we find **better Error-correcting output codes (ECOCs) to improve adversarial robustness** (building upon Verma, G. & Swami, A., 2019)?
- Is there an **efficient mapping between the selected codewords and classes**?
- How can we **encode the data within the layers of neural networks** to prevent accuracy drop while improving adversarial robustness?

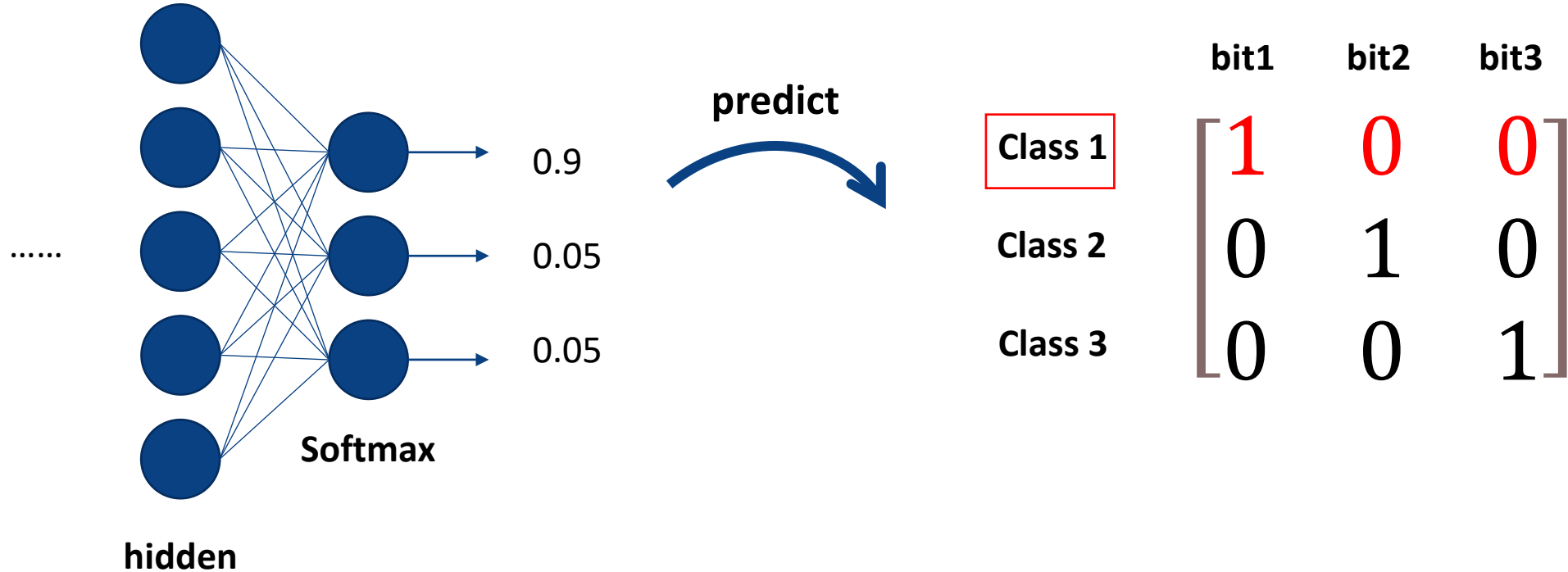


Contributions



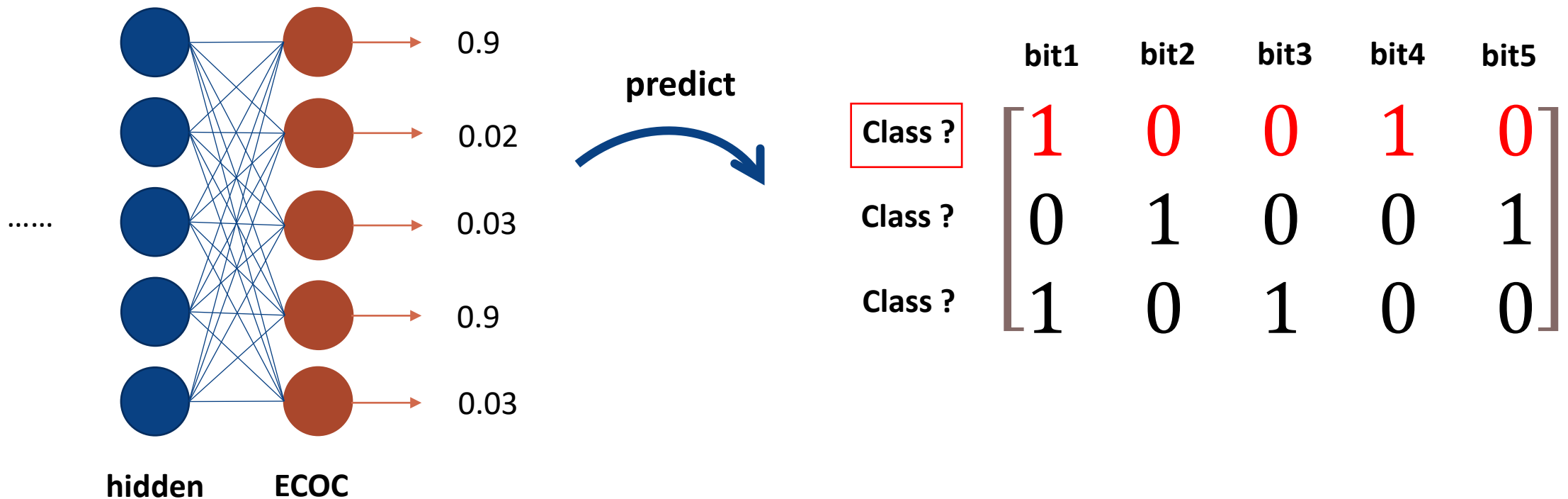
Effective Error-correcting output code (eECOC)

One-hot encoding and Softmax layer are widely used in classification tasks.

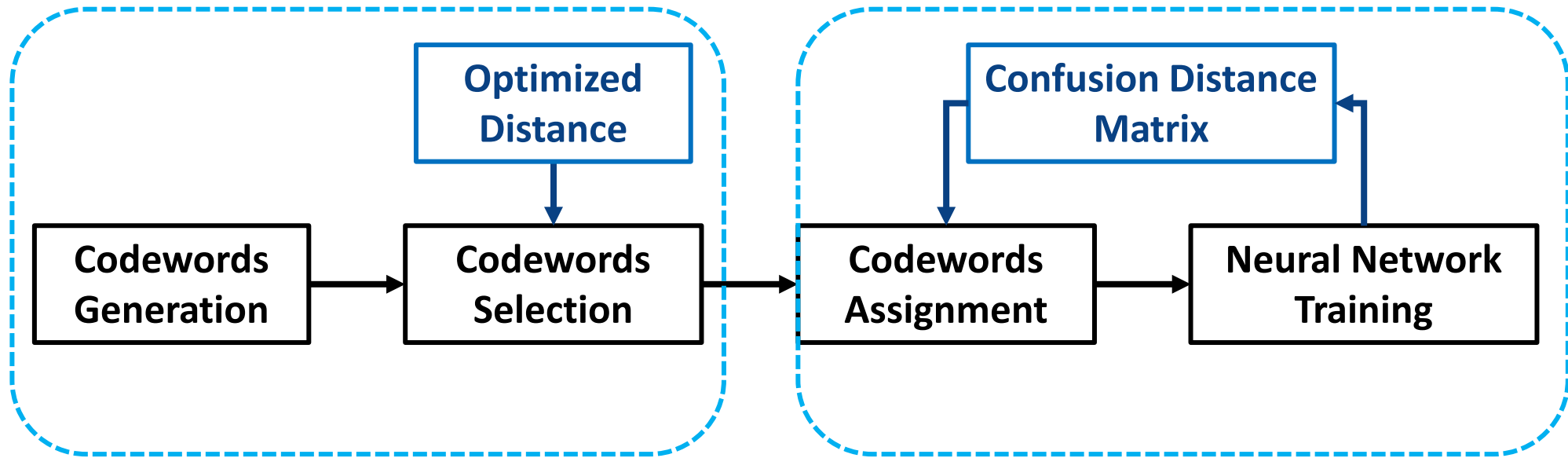


Effective Error-correcting output code (eECOC)

Error-correcting output codes replace the one-hot encoding.
 The ECOC layer replaces the softmax layer.



Effective Error-correcting output code (eECOC)



**Stage 1: A rule-based
codebook design**

**Stage 2: A codeword
assignment problem**

eECOC – Theory Basics

Hamming distance: In coding theory, the **Hamming distance** between any two binary codewords \mathbf{c} and $\hat{\mathbf{c}}$ denoted as $d(\mathbf{c}, \hat{\mathbf{c}})$ computes the number of different bits between two codewords. Therefore, the Hamming distance of a codebook \mathcal{C} is defined as:

$$d = \min\{d(\mathbf{c}, \hat{\mathbf{c}}) \mid \mathbf{c}, \hat{\mathbf{c}} \in \mathcal{C}, \mathbf{c} \neq \hat{\mathbf{c}}\}.$$

The Hamming distance of one-hot codebook is always 2.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

eECOC – Theory Basics

Theorem 1 (Error-correcting capability). *If the minimum distance of a codebook \mathcal{C} is d , a nearest neighbor decoder will always decode correctly when there are $\lfloor \frac{d-1}{2} \rfloor$ or fewer error.*

The error-correcting capability of one-hot codebook is 0.

$$[1 \ 0 \ 0] \rightarrow [1 \ 0 \ 1]$$

Class 1

Class 1 or 3?

We want to generate a code with large Hamming distance.

eECOC – Codebook design

Rule-based Heuristic Codebook Design: Given a dataset with M different classes, a codebook matrix $C \in \mathbb{R}^{M \times N}$, $N = 2^k$ should be generated based on the following rules:

- The elements in any one of the columns can not be the same.
(Discriminative power of each bit)
- The Hamming distance between any two codewords (rows) cannot be smaller than $2^{k-1} - 1$. (Guaranteed minimum distance)
- The codebook should have enough diversity to match the confusion distance matrix, while trying to maximise the Hamming distances between codewords. (For codeword assignment)



eECOC - Codeword Assignment

- The confusion distance matrix measures the separability between classes. Some classes are harder to distinguish while some classes are easier to classify.
- As for codebook, although we have a guaranteed minimum distance, some codewords have large distance to other codewords.
- **Intuition:** Assign the codeword with larger Hamming distance to the class that are harder to distinguish.
- **However, this problem is proved to be NP-hard.**
- A **greedy algorithm**, is used to find a sub-optimal solution.

eECOC - Codeword Assignment

An example of our proposed framework on $M = 10$ classes and $N = 16$ bits.

1. We start with a 16-bits Hadamard code.
2. Change bits starting with the first column.
3. Flip additional bits to increase the diversity of the distance with a guaranteed minimum distance.
4. Swap the rows (codes of classes) to match (Hamming distances) to the confusion distance matrix as much as possible.

codes

classes	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
	+	-	+	-	+	-	+	-	+	+	+	-	+	-	+
	+	+	-	-	+	+	-	-	+	+	-	-	+	+	-
	+	+	+	+	-	-	-	-	+	+	+	+	-	-	-
	+	-	+	-	-	+	-	+	+	-	+	-	-	+	-
	+	+	+	-	-	+	+	+	+	+	-	-	-	-	+
	+	-	-	+	-	+	+	-	+	-	-	+	-	+	+
	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-
	+	-	+	-	+	-	+	-	-	+	-	+	+	+	-



distance matrix

0	9	7	8	7	8	9	9	10	8
9	0	8	9	8	11	8	10	9	7
7	8	0	9	8	7	10	8	9	7
8	9	9	0	9	8	7	9	8	8
7	8	8	9	0	9	8	8	9	9
8	11	7	8	9	0	9	7	10	10
9	8	10	7	8	9	0	8	7	11
9	10	8	9	8	7	8	0	9	9
10	9	9	8	9	10	7	9	0	10
8	7	7	8	9	10	11	9	10	0

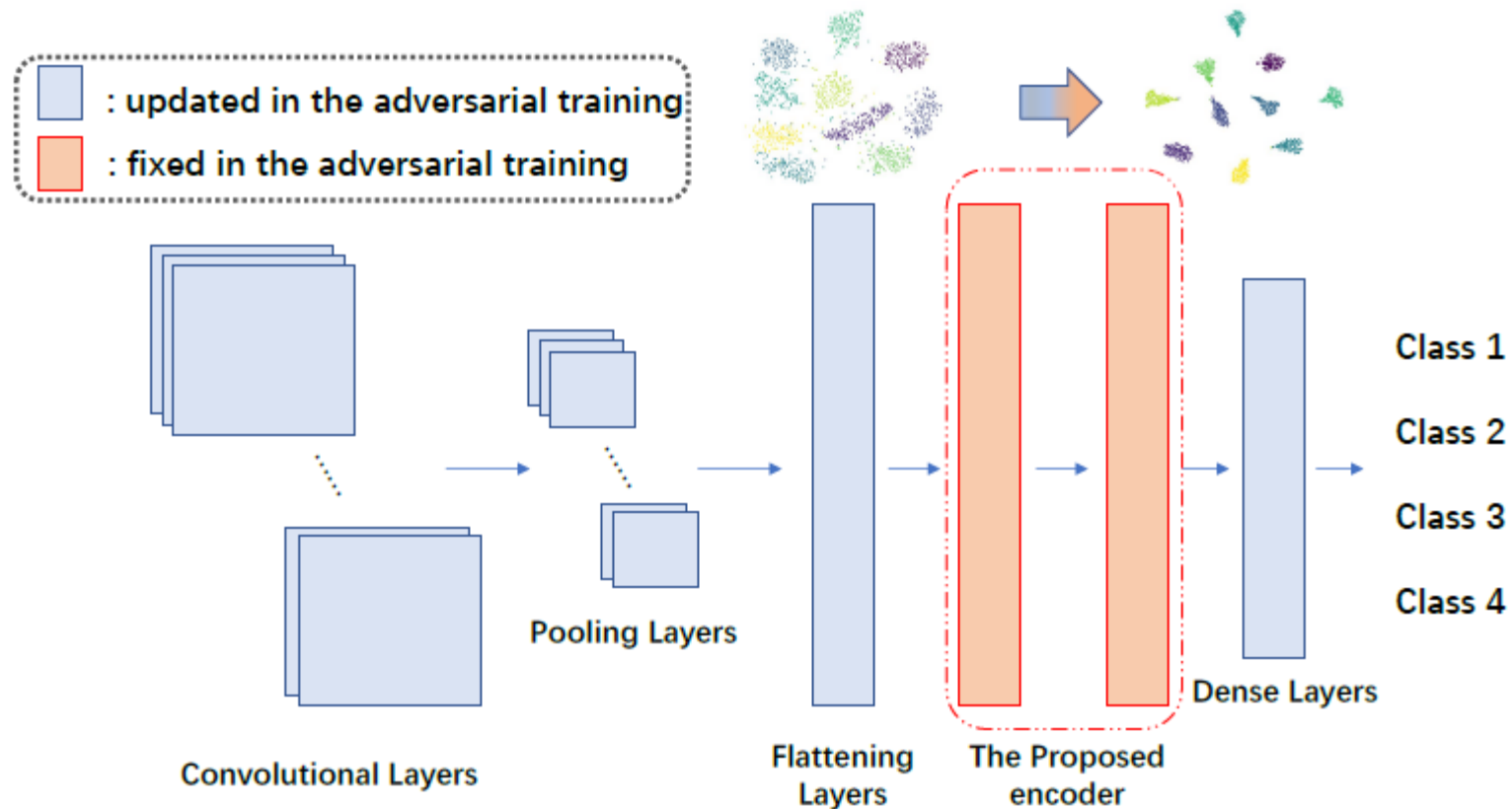


Neural Network Embedded Coding (NNEC)

- We aim to find an encoding scheme that can be embedded into neural network layered architecture (not necessarily to the output layer).
- The encoded features should be **well-separated inter-class** and **concentrated intra-class**.

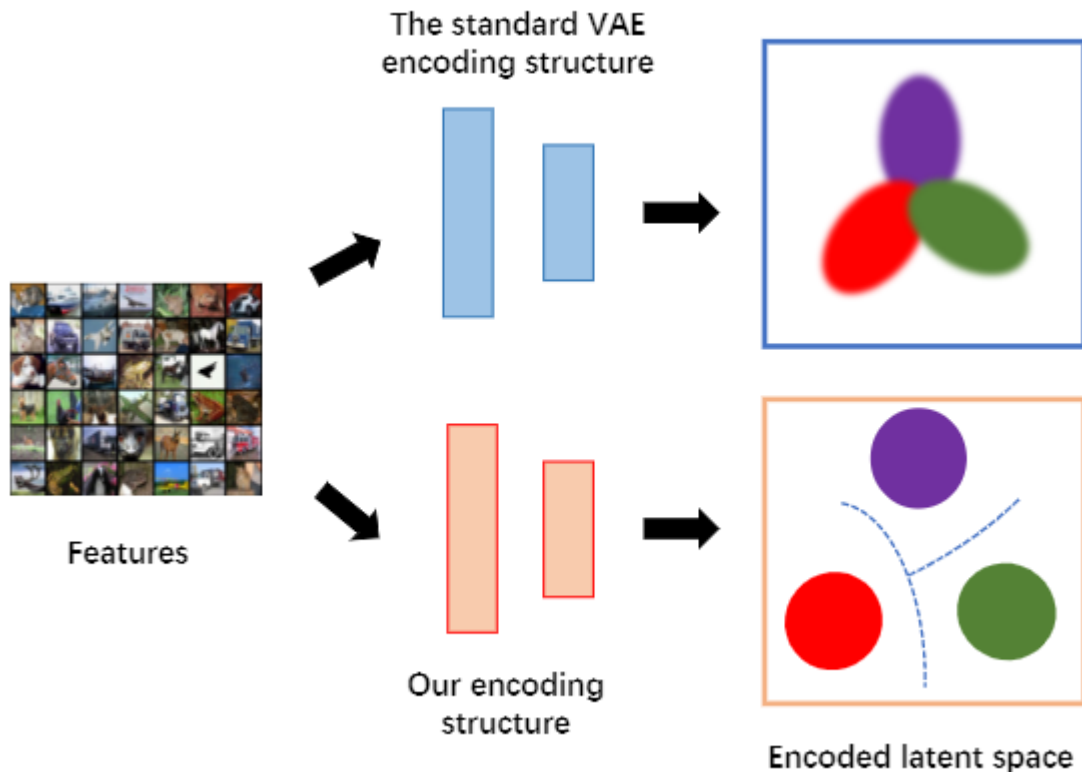
Question: How to encode features in continuous space with high dimension to achieve these goals?

Neural Network Embedded Coding (NNEC)



An overview of the proposed encoder. We use the transfer learning to embed the encoding part of the VAE into a neural network.

Neural Network Embedded Coding (NNEC)



Algorithm 1 Semi-supervised VAE training.

Input: Clean samples dataset (X, Y) , an initialized variational autoencoder (VAE) with latent space dimension d_z

Output: A trained encoding structure that can achieve adversarial robustness

- 1: Generate M number of clustering centers on a d_z dimensional sphere using Fibonacci lattice.
 - 2: Randomly pair up the clustering centers with each class
 - 3: **for** each training batch **do**
 - 4: Draw batch samples from dataset (X, Y)
 - 5: Pass the batch into the VAE
 - 6: Update the VAE based on Eq. (15)
 - 7: **end for**
 - 8: Output the encoding structure and trained parameters of the VAE
-

Coding theory helps us define the cluster centres in the latent space (here Fibonacci lattice)

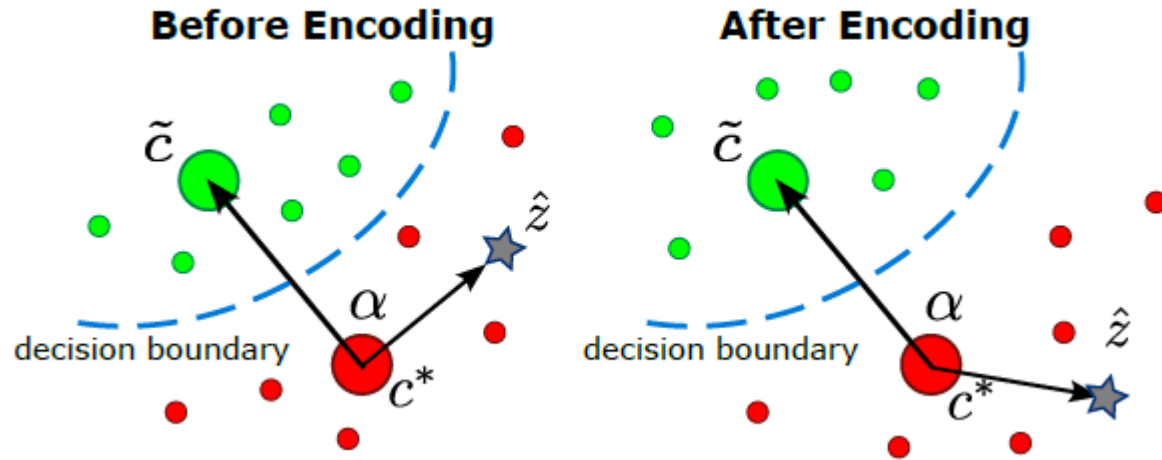
NNEC – Experimental Results

ADVERSARIAL ACCURACY (%) OF VARIOUS MODELS UNDER DIFFERENT ATTACKS.

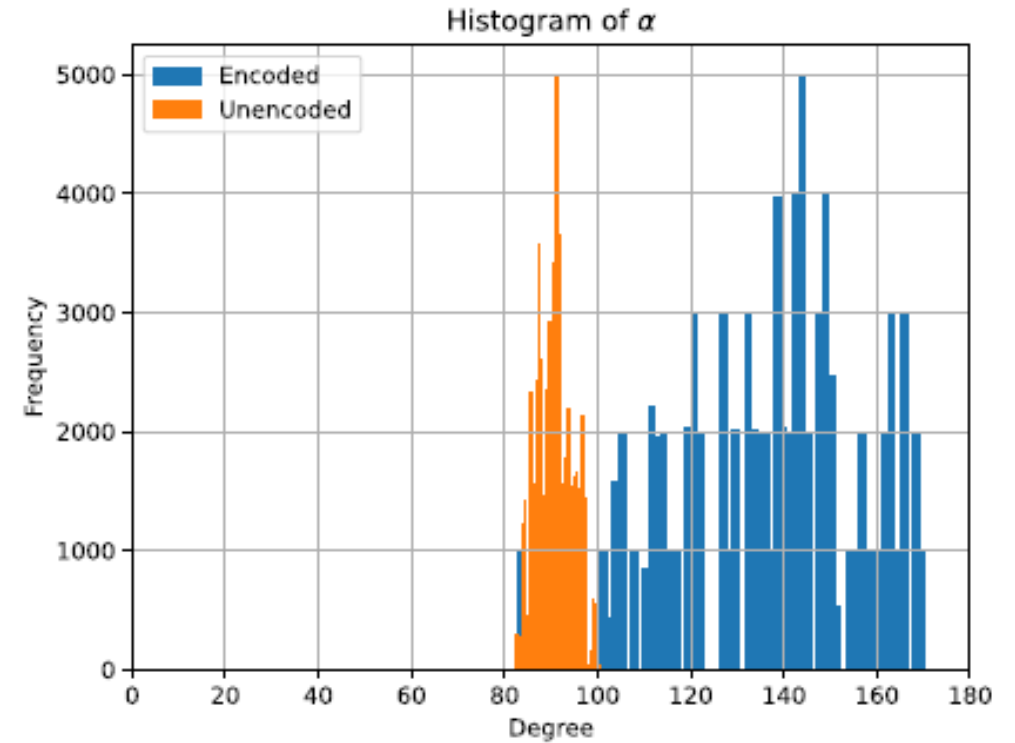
Models	MNIST($\epsilon = 0.3$)				FashionMNIST($\epsilon = 0.1$)				CIFAR-10($\epsilon = 0.03$)			
	Benign	FGSM	PGD	CW	Benign	FGSM	PGD	CW	Benign	FGSM	PGD	CW
Standard	99.24	49.44	3.40	4.29	92.15	13.67	0.00	3.09	92.62	8.43	3.12	6.94
ADV	99.00	92.95	93.20	97.96	83.25	75.62	70.79	66.61	81.24	45.13	41.41	63.17
NNEC	99.14	61.72	42.69	9.43	92.67	18.67	8.94	7.13	92.83	13.11	7.45	7.32
eECOC	98.60	92.12	90.21	98.57	90.42	72.43	71.69	67.23	89.21	50.13	44.72	60.37
ADV + NNEC	98.81	94.39	91.92	98.33	87.32	79.03	74.51	71.87	87.02	52.80	46.90	64.46
ADV + eECOC	96.34	92.81	92.24	96.29	81.13	78.92	72.31	73.13	82.69	49.96	46.83	61.29
NNEC + eECOC	99.16	92.54	91.03	98.45	90.87	73.82	70.96	68.92	90.01	50.67	45.19	62.84
ADR	99.35	91.38	94.52	97.57	84.21	77.94	72.57	70.64	82.43	51.26	42.75	62.84
Thermometer Encoding	99.17	91.94	92.74	97.69	86.15	80.13	73.52	72.36	86.95	53.74	47.59	63.91

We combine the coding approach with adversarial training to improve results.

NNEC – Analysis



An illustration of angle α between cluster centre and data points



The histogram of angle α before and after encoding

On average, our encoding scheme pushes the data points away from boundary!



THE UNIVERSITY OF
MELBOURNE

Cyber(-Physical) Security Games

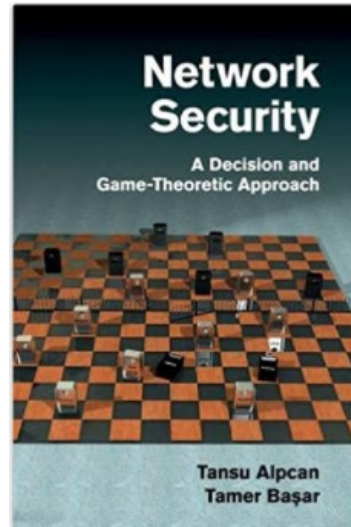
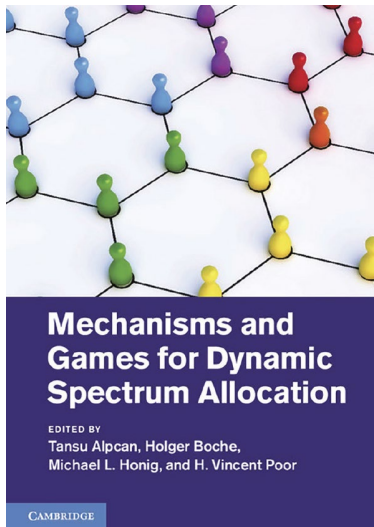
University of Melbourne, Old Arts Building. Parkville Campus

https://commons.wikimedia.org/wiki/File:Old_Arts_Building._Parkville_Campus_of_University_of_Melbourne.jpg



Game Theory in Engineering

- **Game Theory** provides a solid **quantitative and conceptual foundation** for analysing and developing **multi-agent decisions and systems**.
- **Successfully applied** to engineering resource allocation and security problems, specifically in communications, energy, and cybersecurity.



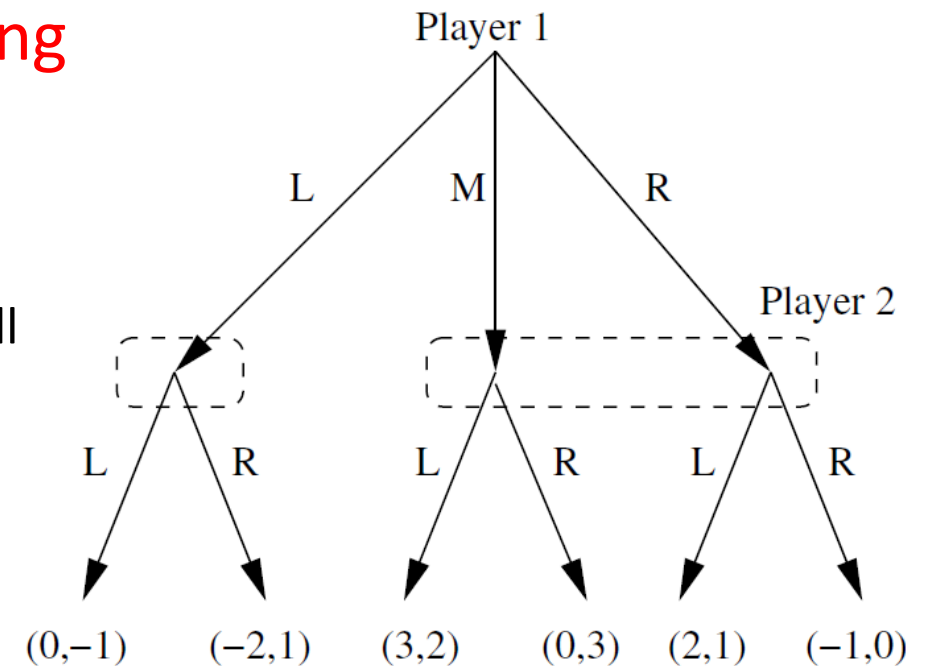
*13th edition of **GameSec**
Conference on Decision and Game Theory for Security*

Game Theory Basics

Game Theory studies multi-person decision making

A **strategic (non-cooperative) game** consists of:

- **players**, who are decision makers acting on their self-interest
- **actions** chosen from a strategy (action) space, which is the set of all actions available to player(s),
- **outcomes (pay-offs)**, which quantize gain or loss of players,
- **information structure (flow)**, characterizing how much each player knows about other's actions



Platoon Security Game-based Best Response

We model the interactions with a **non-cooperative cybersecurity game**:

Attacker:

$$\mathcal{A}^A := \{\mathbf{a}: \text{boiling frog attack}; \mathbf{na}: \text{not attacking}\}$$

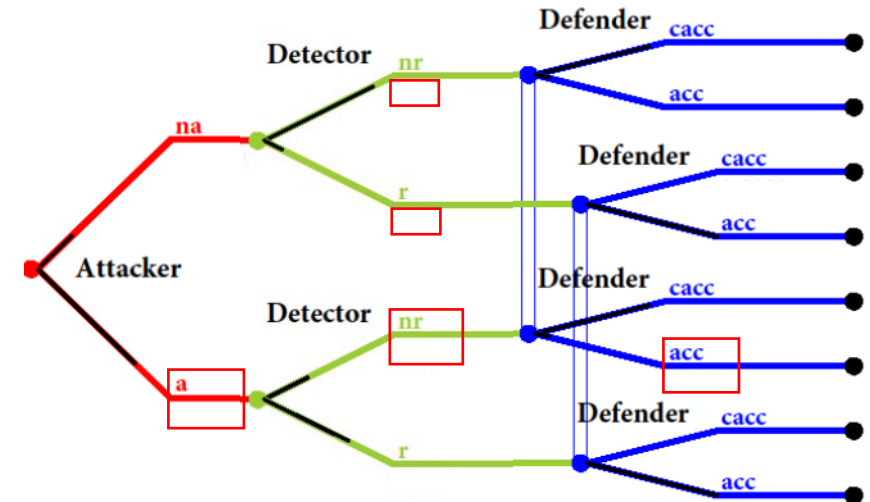
Anomaly detector:

$$\mathcal{C} := \{\mathbf{r}: \text{reporting an attack}; \mathbf{nr}: \text{not reporting an attack}\}$$

Defender:

$$\mathcal{A}^D := \{\mathbf{acc}: \text{ACC controller}; \mathbf{cacc}: \text{CACC controller}\}$$

Nash equilibrium solution is used to guide the decisions



ECML-PKDD 2021 article: Strategic mitigation against wireless attacks on autonomous platoons, Guoxin Sun, Tansu Alpcan, Benjamin Rubinstein and Seyit Camtepe

Simulation Results

Attack: communication link between vehicle 2 and 3 is compromised

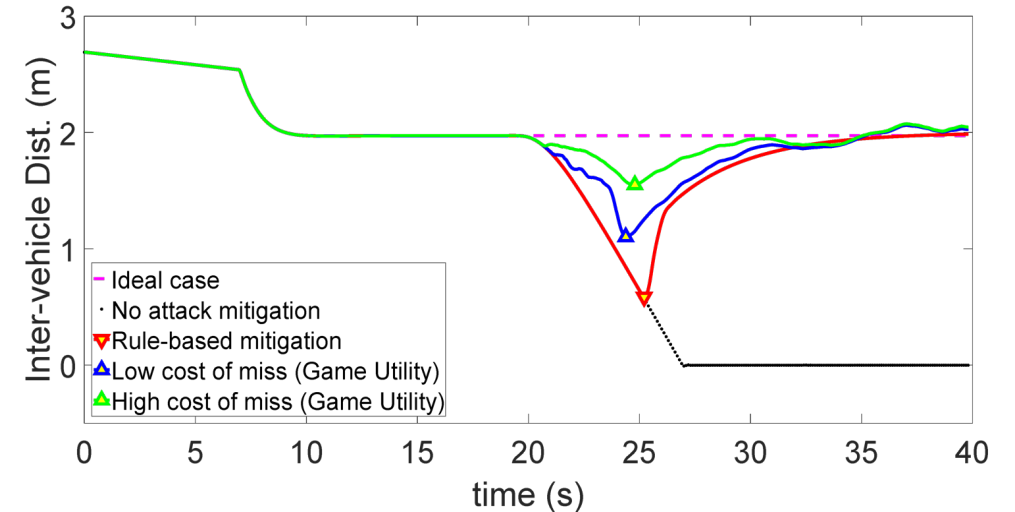
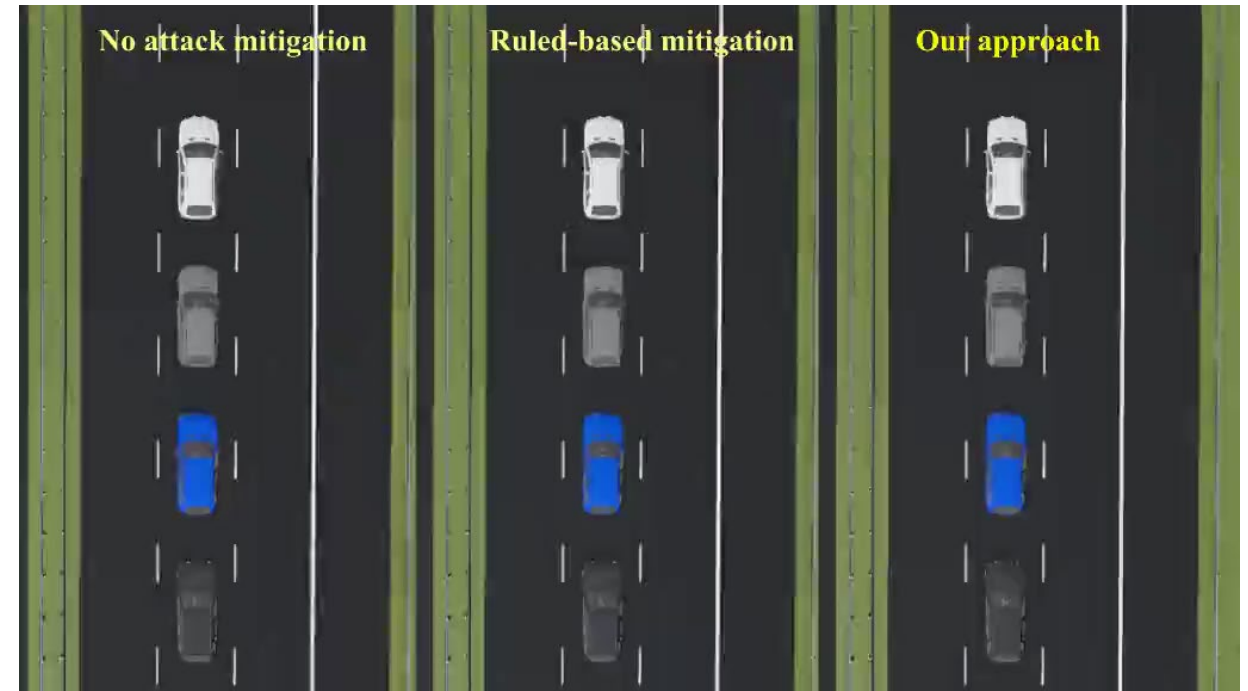
Comparison metric: intervehicle distance between vehicle 2 and vehicle 3

Observation:

- Our proposed defense framework not only avoids a collision but also results in a much safer situation.

Other evaluation scenarios:

- Defense against greedy and rational attackers
- Realistic driving scenario
- Comparison of players' utilities





Optimisation, Game Theory, and Learning

- **Game theory helps better understand player incentives, actions, and strategies.**
- Distributed systems are connected to each with a variety of wired/wireless communication technologies resulting in complex networked systems
→ interaction between decision makers.
- Agents share various resources
→ competition for available resources (resource allocation).
- Optimisation methods help Agents and entire systems to decide on their optimal actions
→ global solutions as ideal benchmark and individual agent best responses.
- Machine learning methods help closing modelling gaps and provide flexibility
→ adaptive and practical solutions using data-oriented approaches.
- Attackers and Defenders continuously improve their attacks and defences
→ adversarial behaviour is modelled using security games.



THE UNIVERSITY OF
MELBOURNE

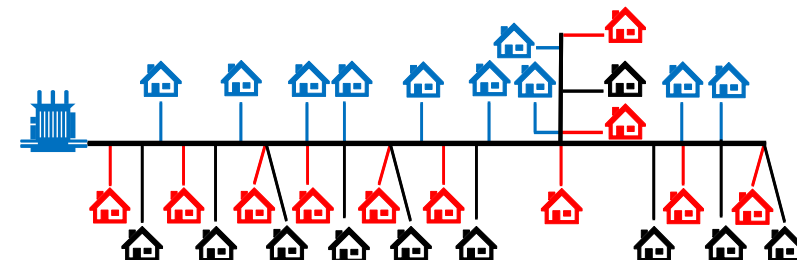
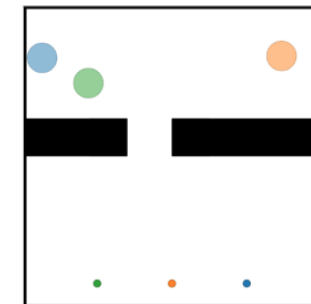
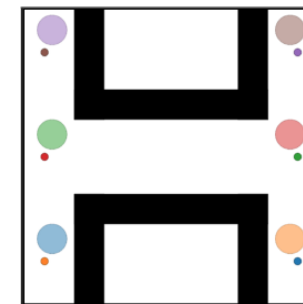
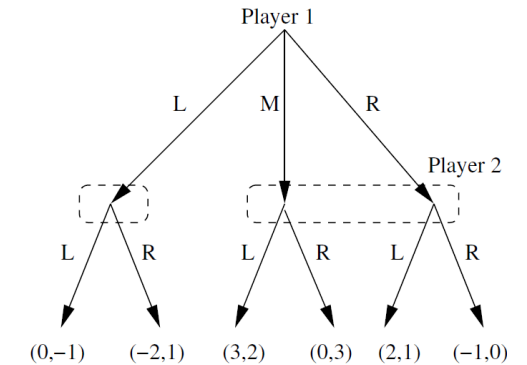
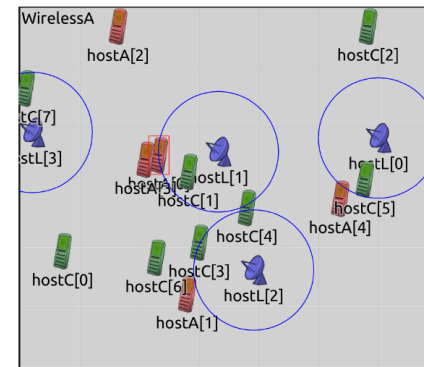
Ongoing Research and Future Directions

Simple game (Tansu Alpcan)



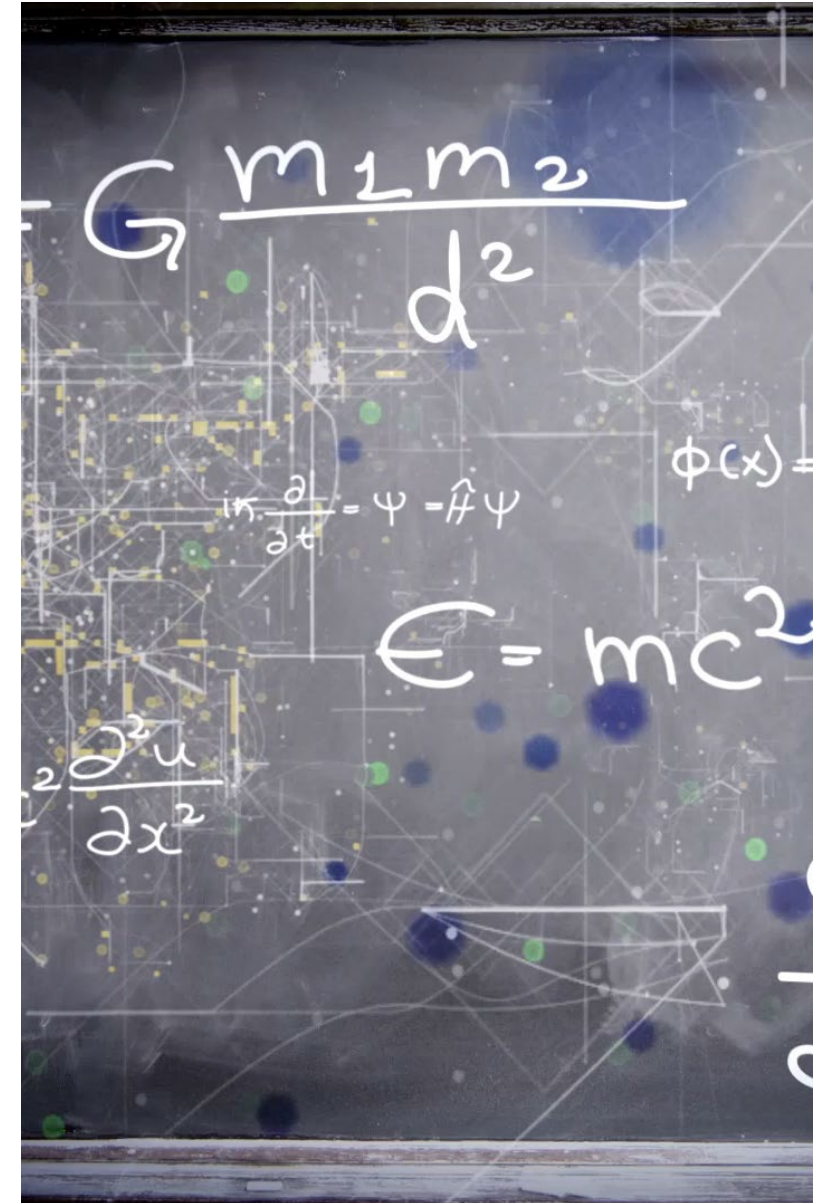
Ongoing Research

- Distributed Anomaly Detection for Cognitive Radio Networks
- (Adversarial) Machine Learning and Game Theory
- Model-based hybrid Reinforcement Learning
- Identification of power distribution networks using deep learning



Future Research

- Solving **large-scale games** for cybersecurity and cyber-physical systems
- Distributed optimisation and machine learning for E2E QoE in emerging network/IoT applications
- (Adversarial) Machine Learning for robustness and security





Thanks to my PhD students, collaborators and funding agencies!

Thanks for funding support

- Project with **CSIRO/Data61** (PhD top-up for Guoxin)
- AdvML Project with DSTG and CSIRO/Data61
- Australian Research Council (ARC) Linkage project (with NGC, USA)

My valuable students and collaborators

PhD students: [Mr. Guoxin Sun](#)
[and Mr. Li Wan](#)

***Dr. Seyit Camtepe** (CSIRO/Data61, Australia), Prof. Ben Rubinstein, Prof. Margreta Kuijper, (University of Melbourne, Australia) Prof. Emanuele Viterbo (Monash Univ., Australia)*



Publications

- G. Sun, T. Alpcan, B. I. P. Rubinstein and S. Camtepe, "Securing Cyber-Physical Systems: Physics-Enhanced Adversarial Learning for Autonomous Platoons," to appear in Proc. of ECML PKDD 2022.
- G. Sun, T. Alpcan, B. I. P. Rubinstein and S. Camtepe, "A Communication Security Game on Switched Systems for Autonomous Vehicle Platoons," *2021 60th IEEE Conference on Decision and Control (CDC)*, 2021, pp. 2690-2695, doi: 10.1109/CDC45484.2021.9683543.
- Sun, G., Alpcan, T., Rubinstein, B.I.P., Camtepe, S. (2021). Strategic Mitigation Against Wireless Attacks on Autonomous Platoons. In: Dong, Y., Kourtellis, N., Hammer, B., Lozano, J.A. (eds) Machine Learning and Knowledge Discovery in Databases. Applied Data Science Track. ECML PKDD 2021. Lecture Notes in Computer Science(), vol 12978. Springer, Cham. https://doi.org/10.1007/978-3-030-86514-6_5
- Li Wan; Tansu Alpcan; Emanuele Viterbo; Margreta Kuijper, "Efficient Error-correcting Output Codes for Adversarial Learning Robustness," *IEEE International Conference on Communications (ICC) 2022, Best Paper Award*.
- L. Wan, T. Alpcan and M. Kuijper, "Interpretable Dictionary Learning Using Information Theory," *GLOBECOM 2020 - 2020 IEEE Global Communications Conference*, 2020, pp. 1-6, doi: 10.1109/GLOBECOM42002.2020.9322557.
- Multiple journal/conference papers under review or being submitted these days.



THE UNIVERSITY OF
MELBOURNE

Thank you
Questions?

