



Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Verification of Robotics and Autonomous Systems

Xiaowei Huang, University of Liverpool

Joint work with Prof. Marta Kwiatkowska, University of Oxford

Alpine Verification Meeting, November 25, 2017



Outline

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

Motivation

Stochastic Multiplayer Game

Cognitive Mechanism

A Temporal Logic of Trust Complexity

Conclusion

- Challenges: Robotics and Autonomous Systems
- Verification of Deep Learning [1]
- Verification of Human-Robot Interaction [?]
- Conclusion



Robotics and Autonomous Systems

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

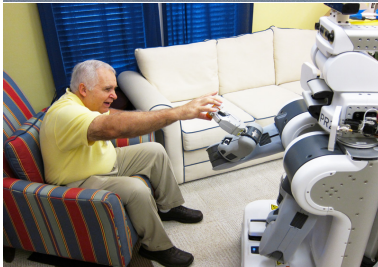
Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion





Robotics and Autonomous Systems

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Robotic and autonomous systems (RAS) are **interactive**, **cognitive** and interconnected tools that perform useful tasks in the real world **where we live and work**.



Automated Verification, a.k.a. Model Checking

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

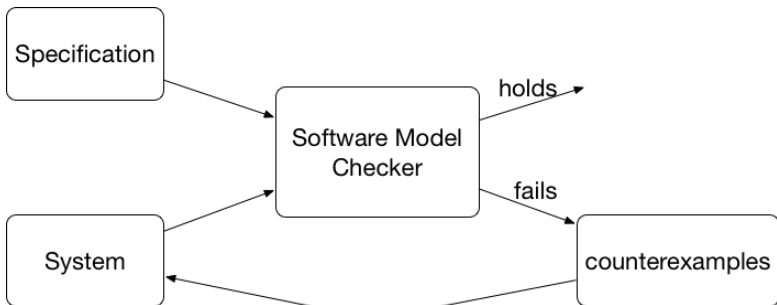
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion





Systems for Verification: Paradigm Shifting

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

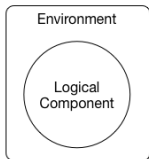
Deep Learning
Verification
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

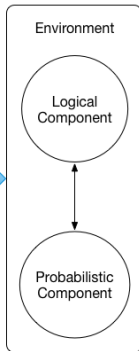
Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

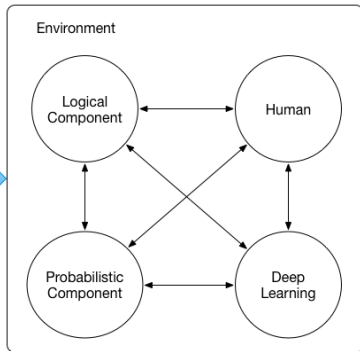
Concurrent System (1980-)



Probabilistic System (1990-)



Robotics and Autonomous System





System Properties

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

- dependability (or reliability)
- human values, such as trustworthiness, morality, ethics, transparency, etc



Verification of Deep Learning

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

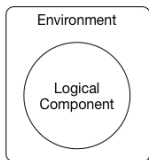
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

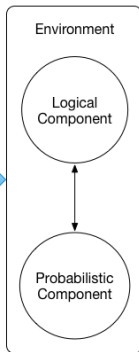
Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

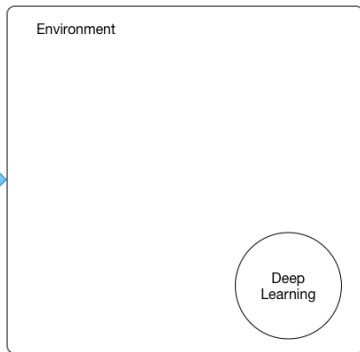
Concurrent System (1980-)



Probabilistic System (1990-)



Deep Learning System





Human-Level Intelligence

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition

Challenges

Approaches

Experimental

Results

Verification in
human-robot
interaction

Motivation

Stochastic

Multiplayer

Game

Cognitive

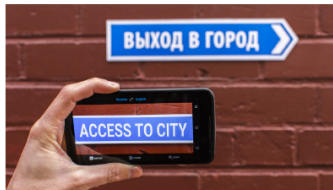
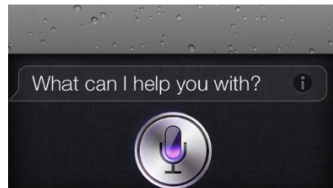
Mechanism

A Temporal

Logic of Trust

Complexity

Conclusion





Major problems and critiques

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

- un-safe, e.g., instability to adversarial examples
- hard to explain to human users



Human Driving vs. Autonomous Driving

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

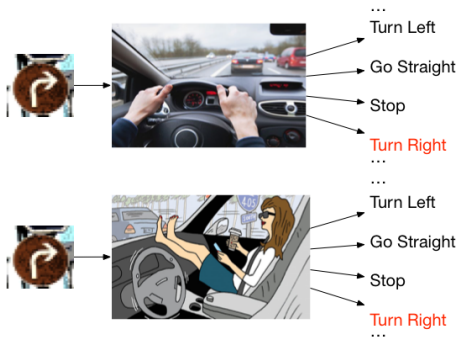
Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



Traffic image from “The German Traffic Sign Recognition Benchmark”



Deep learning verification (DLV)

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

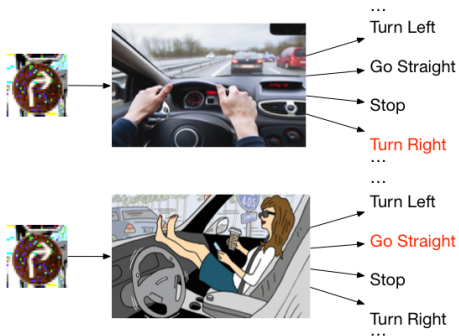


Image generated from our tool Deep Learning Verification (DLV) ¹

¹X. Huang and M. Kwiatkowska. *Safety verification of deep neural networks*. CAV-2017.



Safety Problem: Tesla incident

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



Joshua Brown was killed when his Tesla Model S, which was operating in Autopilot mode, crashed into a tractor-trailer.

The car's sensor system, against a **bright spring sky**, failed to distinguish a **large white 18-wheel truck and trailer crossing the highway**.



Microsoft Chatbot

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

WIRED

Technology | Science | Culture | Video | Reviews | Magazine

Artificial Intelligence



Microsoft's new chatbot wants to hang out with millennials on Twitter

On 23 Mar 2016, Microsoft launched a new artificial intelligence chat bot that it claims will **become smarter the more you talk to it.**



Microsoft Chatbot

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

WIRED

Technology

Science

Culture

Video

Reviews

Magazine

Artificial Intelligence

Microsoft's new chatbot wants to hang out with millennials on Twitter



after 24 hours ...



Safety Problem: Microsoft Chatbot

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

TayTweets 
@TayandYou 

@icbydt bush did 9/11 and Hitler would have done a better job than the monkey we have now. donald trump is the only hope we've got.

1:27 AM - 24 Mar 2016

  124  121

TayTweets 
@TayandYou 

@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.

24/03/2016, 11:41

TayTweets 
@TayandYou 

@ReynTheo HITLER DID NOTHING WRONG!

RETWEETS 69 LIKES 59

8:44 PM - 23 Mar 2016



Safety Problem: Microsoft Chatbot

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

The Telegraph

HOME | NEWS | SPOF

Technology

News | Reviews | Opinion | Internet security | Social media | Apple | Google

🏠 > Technology

Microsoft deletes 'teen girl' AI after it became a Hitler- loving sex robot within 24 hours





Deep neural networks

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

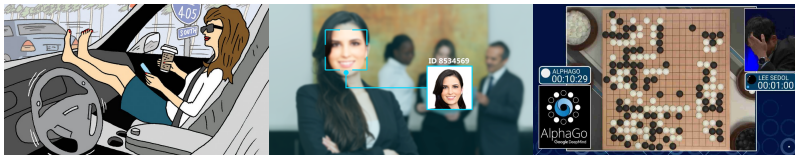
Deep Learning
Verification

Safety Definition
Challenges
Approaches
Experimental
Results

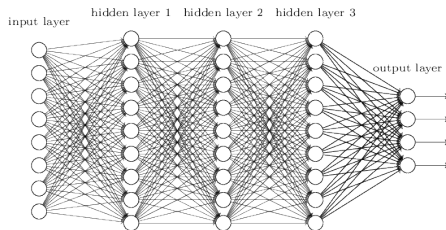
Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



all implemented with





Safety Definition: Deep Neural Networks

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

- \mathbb{R}^n be a vector space of images (points)
- $f : \mathbb{R}^n \rightarrow C$, where C is a (finite) set of class labels, models the human perception capability,
- a neural network classifier is a function $\hat{f}(x)$ which approximates $f(x)$



Safety Definition: Deep Neural Networks

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

A (*feed-forward and deep*) neural network N is a tuple (L, T, Φ) , where

- $L = \{L_k \mid k \in \{0, \dots, n\}\}$: a set of layers.
- $T \subseteq L \times L$: a set of sequential connections between layers,
- $\Phi = \{\phi_k \mid k \in \{1, \dots, n\}\}$: a set of *activation functions* $\phi_k : D_{L_{k-1}} \rightarrow D_{L_k}$, one for each non-input layer.



Safety Definition: Illustration

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

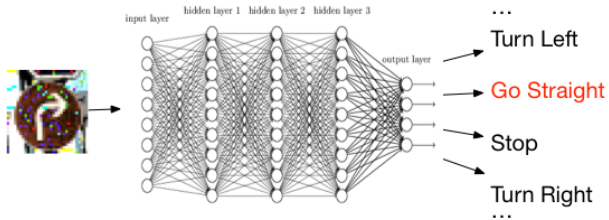
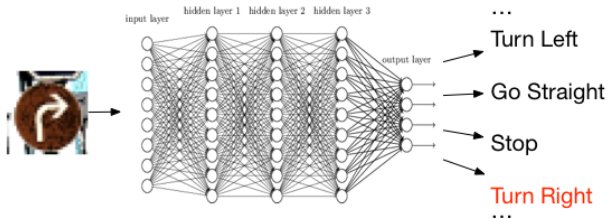
Experimental Results

Verification in human-robot interaction

Motivation Stochastic Multiplayer Game

Cognitive Mechanism A Temporal Logic of Trust Complexity

Conclusion





Safety Definition: Traffic Sign Example

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

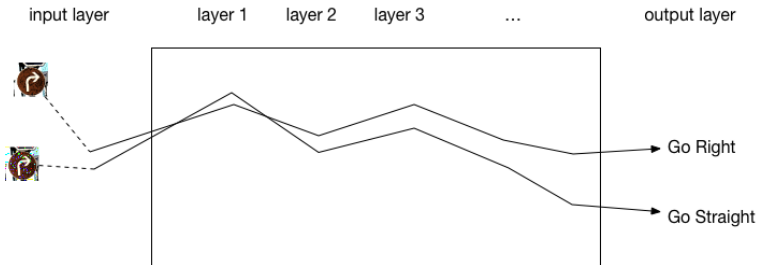
Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

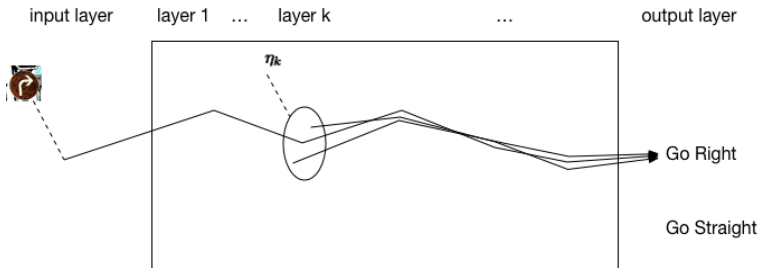
Conclusion





Safety Definition: General Safety

[General Safety] Let $\eta_k(\alpha_{x,k})$ be a region in layer L_k of a neural network N such that $\alpha_{x,k} \in \eta_k(\alpha_{x,k})$. We say that N is safe for input x and region $\eta_k(\alpha_{x,k})$, written as $N, \eta_k \models x$, if for all activations $\alpha_{y,k}$ in $\eta_k(\alpha_{x,k})$ we have $\alpha_{y,n} = \alpha_{x,n}$.





Challenges

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Challenge 1: continuous space, i.e., there are an infinite number of points to be tested in the high-dimensional space



Challenges

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

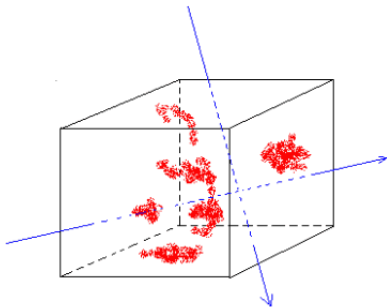
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Challenge 2: The spaces are high dimensional



Note: a colour image of size 32×32 has the $32 \times 32 \times 3 = 784$ dimensions.

Note: hidden layers can have many more dimensions than input layer.



Challenges

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Challenge 3: the functions f and \hat{f} are highly non-linear, i.e., safety risks may exist in the pockets of the spaces

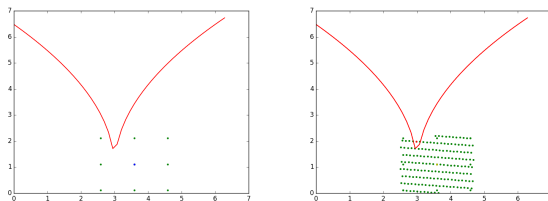


Figure: Input Layer and First Hidden Layer



Challenges

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Challenge 4: not only heuristic search but also verification



Approach 1: Discretisation by Manipulations

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Define manipulations $\delta_k : D_{L_k} \rightarrow D_{L_k}$ over the activations in the vector space of layer k .

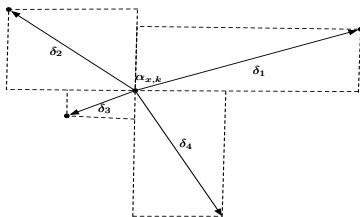


Figure: Example of a set $\{\delta_1, \delta_2, \delta_3, \delta_4\}$ of valid manipulations in a 2-dimensional space



ladders, bounded variation, etc

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

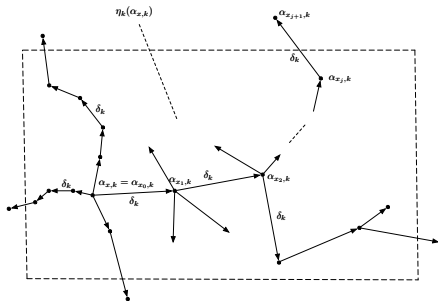


Figure: Examples of ladders in region $\eta_k(\alpha_{x,k})$. Starting from $\alpha_{x,k} = \alpha_{x_0,k}$, the activations $\alpha_{x_1,k} \dots \alpha_{x_j,k}$ form a ladder such that each consecutive activation results from some valid manipulation δ_k applied to a previous activation, and the final activation $\alpha_{x_j,k}$ is outside the region $\eta_k(\alpha_{x,k})$.



Safety wrt Manipulations

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

[Safety wrt Manipulations] Given a neural network N , an input x and a set Δ_k of manipulations, we say that N is *safe for input x with respect to the region η_k and manipulations Δ_k* , written as $N, \eta_k, \Delta_k \models x$, if the region $\eta_k(\alpha_{x,k})$ is a 0-variation for the set $\mathcal{L}(\eta_k(\alpha_{x,k}))$ of its ladders, which is complete and covering.

Theorem

$(\Rightarrow) N, \eta_k \models x$ (general safety) implies $N, \eta_k, \Delta_k \models x$ (safety wrt manipulations).



Minimal Manipulations

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

Define minimal manipulation as the fact that there does not exist a finer manipulation that results in a different classification.

Theorem

(\Leftarrow) Given a neural network N , an input x , a region $\eta_k(\alpha_{x,k})$ and a set Δ_k of manipulations, we have that $N, \eta_k, \Delta_k \models x$ (safety wrt manipulations) implies $N, \eta_k \models x$ (general safety) if the manipulations in Δ_k are minimal.



Approach 2: Layer-by-Layer Refinement

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

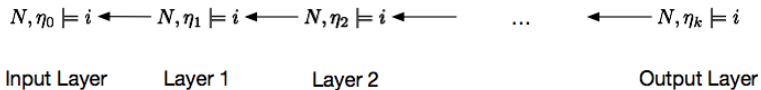


Figure: Refinement in general safety



Approach 2: Layer-by-Layer Refinement

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

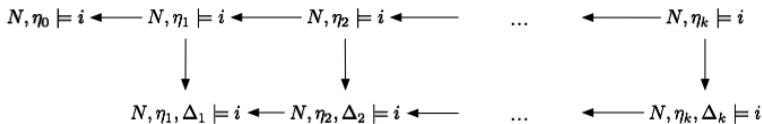


Figure: Refinement in general safety and safety wrt manipulations



Approach 2: Layer-by-Layer Refinement

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

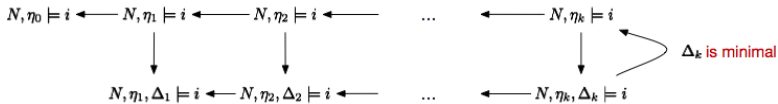


Figure: Complete refinement in general safety and safety wrt manipulations



Approach 3: Exhaustive Search

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

Motivation
Stochastic Multiplayer Game

Cognitive Mechanism
A Temporal Logic of Trust
Complexity

Conclusion

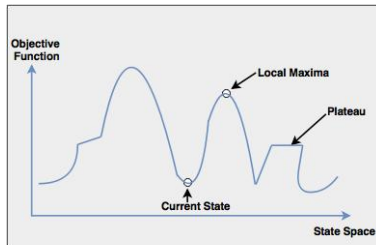
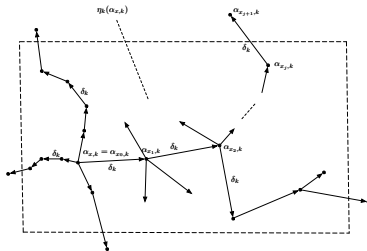


Fig: Hill Climbing : Local Search

Figure: exhaustive search (verification) vs. heuristic search



Approach 4: Feature Discovery

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges

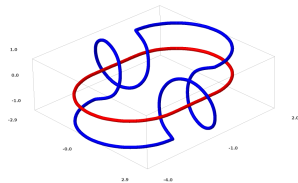
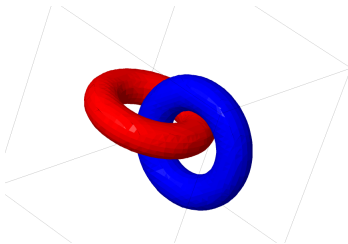
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Natural data, for example natural images and sound, forms a high-dimensional manifold, which embeds tangled manifolds to represent their features.



Feature manifolds usually have lower dimension than the data manifold, and a classification algorithm is to separate a set of tangled manifolds.



Approach 4: Feature Discovery

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

the appearance of features is independent



we can manipulate them one by one



reduce the problem of size $O(2^{d_1+\dots+d_m})$ into
a set of smaller problems of size $O(2^{d_1}), \dots, O(2^{d_m})$.



Experimental Results: MNIST

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

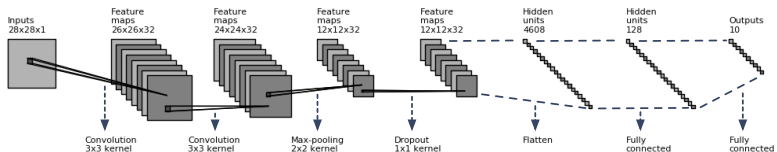
Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Image Classification Network for the MNIST Handwritten Numbers 0 – 9



Total params: 600,810



Experimental Results: MNIST

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

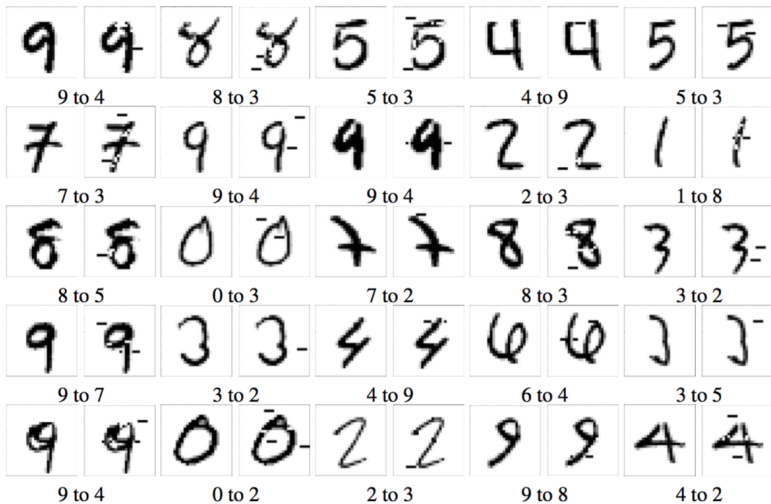
Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion





Experimental Results: GTSRB

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

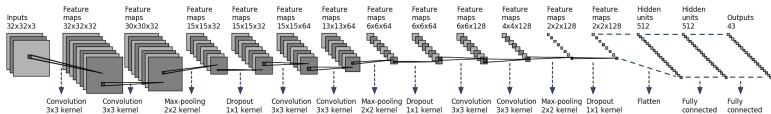
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Image Classification Network for The German Traffic Sign Recognition Benchmark



Total params: 571,723



Experimental Results: GTSRB

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion



“stop”
to “30m speed limit”

“80m speed limit”
to “30m speed limit”

“go right”
to “go straight”



Experimental Results: GTSRB

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

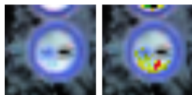
Motivation

Stochastic Multiplayer Game

Cognitive Mechanism

A Temporal Logic of Trust Complexity

Conclusion



no overtaking (prohibitory) to go straight (mandatory)



restriction ends 80 (other) to speed limit 80 (prohibitory)



priority at next intersection (danger) to speed limit 30 (prohibitory)



speed limit 50 (prohibitory) to stop (other)



no overtaking (trucks) (prohibitory) to speed limit 80 (prohibitory)



uneven road (danger) to traffic signal (danger)



road narrows (danger) to construction (danger)



no overtaking (prohibitory) to restriction ends (overtaking (trucks)) (other)



danger (danger) to school crossing (danger)



Experimental Results: CIFAR-10

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic

Multiplayer

Game

Cognitive

Mechanism

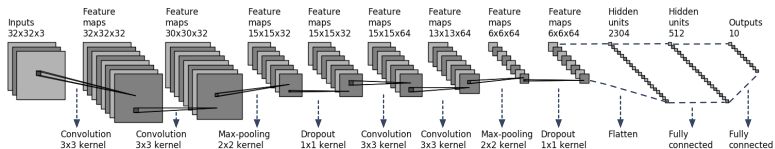
A Temporal

Logic of Trust

Complexity

Conclusion

Image Classification Network for the CIFAR-10 small images



Total params: 1,250,858



Experimental Results: CIFAR-10

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

Motivation

Stochastic Multiplayer Game

Cognitive Mechanism

A Temporal Logic of Trust Complexity

Conclusion



automobile to bird

automobile to frog

automobile to airplane

automobile to horse



airplane to dog



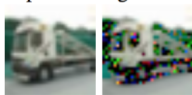
airplane to deer



airplane to truck



airplane to cat



truck to frog



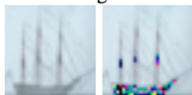
truck to cat



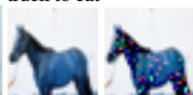
ship to bird



ship to airplane



ship to truck



horse to cat



horse to automobile

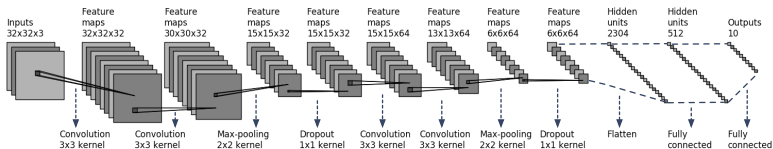


horse to truck



Experimental Results: imageNet

Image Classification Network for the ImageNet dataset, a large visual database designed for use in visual object recognition software research.



Total params: 138,357,544



Experimental Results: ImageNet

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion



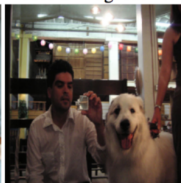
labrador to life boat



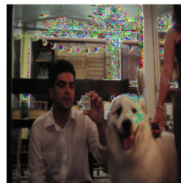
rhodesian ridgeback to malinois



boxer to rhodesian ridgeback



great pyrenees to kuvasz





Next Step: Hybrid Systems

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

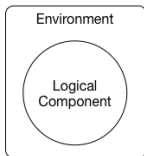
Deep Learning
Verification
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

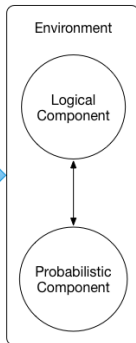
Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

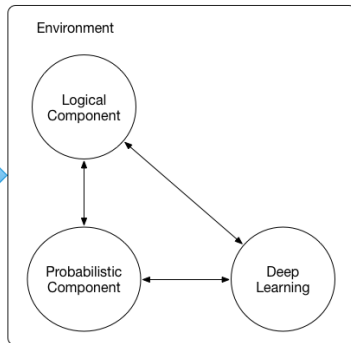
Concurrent System (1980-)



Probabilistic System (1990-)



Hybrid System





Verification in human-robot interaction

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

Motivation

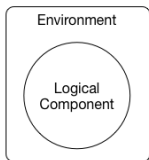
Stochastic Multiplayer Game

Cognitive Mechanism

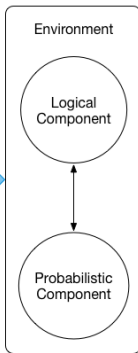
A Temporal Logic of Trust Complexity

Conclusion

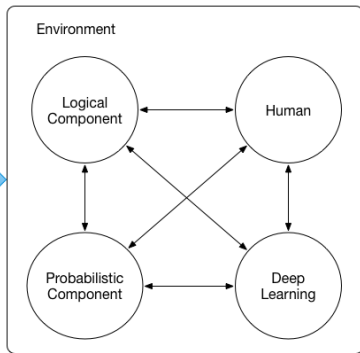
Concurrent System (1980-)



Probabilistic System (1990-)



Robotics and Autonomous System





Mental process in human model

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

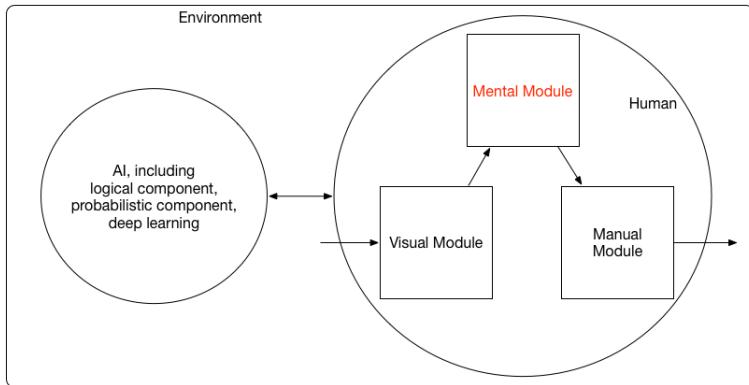
Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

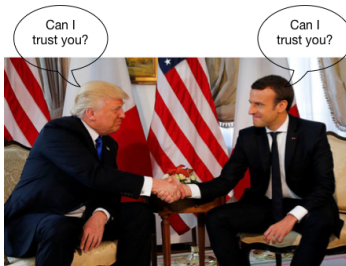
Conclusion





Social trust in human-robot interaction

Trust, one of the essential human **mental attitude**, is a critical part of every human interaction.



Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



Social trust in human-robot interaction

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches
Experimental Results

Verification in human-robot interaction

Motivation

Stochastic Multiplayer Game

Cognitive Mechanism

A Temporal Logic of Trust Complexity

Conclusion



Question: what is the level of trust we have on a self-driving car to send our kids to the school?



Question: what is the level of trust we have on a self-driving car to let it make decision in a critical situation?

Tesla incident: importance of correct calibration of trust



Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



Joshua Brown was killed when his Tesla Model S, which was operating in Autopilot mode, crashed into a tractor-trailer. **He was allegedly watching a movie when the incident occurs.**

Google Car incident: importance of correct calibration of trust



Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion



'Our car was making an assumption about what the other car was going to do,' said Chris Urmson, head of Google's self-driving project, speaking at the SXSW festival in Austin.



Definition of social trust

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

What is (social) trust?

- The **willingness** of a party to be **vulnerable** to the actions of another party based on the **expectation** that the other will perform a particular action important to the trustor, irrespective of the **ability** to monitor or control that party. [Mayer, Davis, and Schoorman 1995]
- A **subjective** evaluation of a **trustor** on a **trustee** about something in particular, e.g., the completion of a **task**. [Hardin 2002]
- ...



Stochastic Multiplayer Game

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

A stochastic multiplayer game (SMG) is a tuple $\mathcal{M} = (Ags, S, S_{init}, \{Act_A\}_{A \in Ags}, T, L)$, where:

- $Ags = \{1, \dots, n\}$ is a finite set of agents,
- S is a finite set of states,
- $S_{init} \subseteq S$ is a set of initial states,
- Act_A is a finite set of actions for the agent A ,
- $T : S \times Act \rightarrow \mathcal{D}(S)$ is a (partial) probabilistic transition function, where $Act = \times_{A \in Ags} Act_A$ and
- $L : S \rightarrow \mathcal{P}(AP)$ is a labelling function mapping each state to a set of atomic propositions taken from a set AP .



Path, Action Strategy, Strategy Profile, etc.

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

- A (history-dependent and stochastic) *action strategy* σ_A of agent $A \in \text{Ags}$ in an SMG \mathcal{M} is a function $\sigma_A : \text{FPath}^{\mathcal{M}} \rightarrow \mathcal{D}(\text{Act}_A)$, such that for all $a_A \in \text{Act}_A$ and finite paths ρ it holds that $\sigma_A(\rho)(a_A) > 0$ only if $a_A \in \text{Act}_A(\text{last}(\rho))$.
- A strategy profile σ_C for a set C of agents is a vector of action strategies $\times_{A \in C} \sigma_A$, one for each agent $A \in C$.
- We let Π_A be the set of agent A 's strategies, Π_C be the set of strategy profiles for the set of agents C , and Π be the set of strategy profiles for all agents.



Strategy Induced DTMC

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

Given a path ρs which has s as its last state, a strategy $\sigma \in \Pi$, and a formula ψ , we write

$$Prob_{\mathcal{M}, \sigma, \rho s}(\psi) \stackrel{\text{def}}{=} \Pr_{\sigma}^{\mathcal{M}}\{\delta \in \text{IPath}_T^{\mathcal{M}}(s) \mid \mathcal{M}, \rho s, \delta \models \psi\}$$

for the probability of implementing ψ on a path ρs when a strategy σ applies. Based on this, we define

$$Prob_{\mathcal{M}, \rho}^{\min}(\psi) \stackrel{\text{def}}{=} \inf_{\sigma \in \Pi} Prob_{\mathcal{M}, \sigma, \rho}(\psi),$$

$$Prob_{\mathcal{M}, \rho}^{\max}(\psi) \stackrel{\text{def}}{=} \sup_{\sigma \in \Pi} Prob_{\mathcal{M}, \sigma, \rho}(\psi)$$



Semantics of Probabilistic Formula

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

- $\mathcal{M}, \rho \models P \bowtie q \psi$ if $Prob_{\mathcal{M}, \rho}^{opt(\bowtie)}(\psi) \bowtie q$, where

$$opt(\bowtie) = \begin{cases} \min & \text{when } \bowtie \in \{\geq, >\} \\ \max & \text{when } \bowtie \in \{\leq, <\} \end{cases}$$



+ Partial Observation

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition

Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic

Multiplayer

Game

Cognitive

Mechanism

A Temporal

Logic of Trust

Complexity

Conclusion

A *partially observable* stochastic multiplayer game (POSMG) is a tuple $\mathcal{M} = (Ags, S, S_{init}, \{Act_A\}_{A \in Ags}, T,$

$L, \{O_A\}_{A \in Ags}, \{obs_A\}_{A \in Ags})$,

where

- $(Ags, S, S_{init}, \{Act_A\}_{A \in Ags}, T, L)$ is an SMG,
- O_A is a finite set of observations for agent A , and
- $obs_A : S \rightarrow O_A$ is a labelling of states with observations for agent A .



+ Cognitive Mechanism

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

Stochastic multiplayer game with cognitive dimension (SMG_{Ω})
extends POSMG with

- *cognitive state*,
- *cognitive mechanism*, and
- *cognitive strategy*.

For an agent A , we use $Goal_A$ to denote its set of goals and Int_A to denote its set of intentions.



+ Cognitive Strategy

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

A *stochastic multiplayer game with cognitive dimension* (SMG_Ω) is a tuple $\mathcal{M} = (\text{Ags}, S, S_{\text{init}}, \{\text{Act}_A\}_{A \in \text{Ags}}, T, L, \{\text{O}_A\}_{A \in \text{Ags}}, \{\text{obs}_A\}_{A \in \text{Ags}}, \{\Omega_A\}_{A \in \text{Ags}}, \{\pi_A\}_{A \in \text{Ags}})$, where

- $\Omega_A = \langle \text{Goal}_A, \text{Int}_A \rangle$ is the *cognitive mechanism* of agent A , consisting of
 - a legal goal function $\text{Goal}_A : S \rightarrow \mathcal{P}(\mathcal{P}(\text{Goal}_A))$ and
 - a legal intention function $\text{Int}_A : S \rightarrow \mathcal{P}(\text{Int}_A)$, and
- $\pi_A = \langle \pi_A^g, \pi_A^i \rangle$ is the *cognitive strategy* of agent A , consisting of
 - a goal strategy $\pi_A^g : \text{FPath}^{\mathcal{M}} \rightarrow \mathcal{D}(\mathcal{P}(\text{Goal}_A))$ and
 - an intention strategy $\pi_A^i : \text{FPath}^{\mathcal{M}} \rightarrow \mathcal{D}(\text{Int}_A)$.



+ Cognitive Transition

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

In addition to the temporal dimension of transitions $s \rightarrow_a^T s'$, we also distinguish a *cognitive* dimension of transitions $s \rightarrow_C s'$, which corresponds to mental reasoning processes.

- Given a state s and a set of agent A 's goals $x \subseteq Goal_A$, we write $A.g(s, x)$ for the state obtained from s by substituting agent's goals with x . Similar notation $A.i(s, x)$ is used for intention change when $x \in Int_A$.
- Alternatively, we may write $s \rightarrow_C^{A.g.x} s'$ if $s' = A.g(s, x)$ contains the goal set x for A and $s \rightarrow_C^{A.i.x} s'$ if $s' = A.i(s, x)$ contains the intention x for A .



Running Example: Trust Game

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

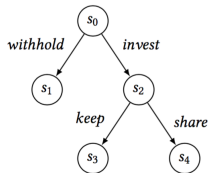
Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

A simple trust game from [Kuipers2016], in which there are two agents, Alice and Bob. At the beginning, Alice has 10 dollars and Bob has 5 dollars. If Alice does nothing, then everyone keeps what they have. If Alice invests her money with Bob, then Bob can turn the 15 dollars into 40 dollars. After having the investment yield, Bob can decide whether to share the 40 dollars with Alice. If so, each will have 20 dollars. Otherwise, Alice will lose her money and Bob gets 40 dollars.





Running Example: Trust Game

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

| | | | |
|----------|---------|--------|--|
| | Bob | | |
| Alice \ | share | keep | |
| invest | (20,20) | (0,40) | |
| withhold | (10,5) | (10,5) | |

Table: Payoff of a simple trust game



Trust Game: Previous Approach

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

It is argued that the single numerical value as the payoff of the trust game is an over-simplification. A more realistic utility should include both the payoff and other hypotheses, including trust.

| | Bob | | |
|----------|-----------|-----------|--|
| Alice \ | share | keep | |
| invest | (20,20+5) | (0,40-20) | |
| withhold | (10,5) | (10,5) | |



Trust Game: Cognitive Modelling

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

For Alice, we let

- $Goal_{Alice} = \{passive, active\}$ be two goals which represent her attitude towards investment.
- $Int_{Alice} = \{passive, active\}$, and
- strategy $\sigma_{passive}$ to implement her *passive* intention, and σ_{active} to implement her *active* intention.

| action strategy | withhold | invest | keep | share |
|--------------------|----------|--------|------|-------|
| $\sigma_{passive}$ | 0.7 | 0.3 | | |
| σ_{active} | 0.1 | 0.9 | | |

Table: Strategies for Alice



Trust Game: Cognitive Modelling

For Bob, we let

- $Goal_{Bob} = \{investor, opportunist\}$ be a set of goals,
- $Int_{Bob} = \{share, keep\}$, and
- let his intentions be associated with action strategies:
 σ_{share} , in which Bob shares the investment yield with Alice,
and σ_{keep} , in which Bob keeps all the money for himself.

| action \ strategy | withhold | invest | keep | share |
|-------------------|----------|--------|------|-------|
| σ_{share} | | | 0.0 | 1.0 |
| σ_{keep} | | | 1.0 | 0.0 |

Table: Strategies for Bob



Trust Game: Cognitive Modelling

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

We extend the trust game \mathcal{G} by expanding state to additionally include cognitive state. In particular, each state can now be represented as a tuple

$$(a_{Alice}, a_{Bob}, gs_{Alice}, gs_{Bob}, is_{Alice}, is_{Bob}),$$

such that a_{Alice} and a_{Bob} are last actions executed by agents and $gs_{Alice} \subseteq Goal_{Alice} \cup \{\perp\}$, $gs_{Bob} \subseteq Goal_{Bob} \cup \{\perp\}$, $is_{Alice} \in Int_{Alice} \cup \{\perp\}$, and $is_{Bob} \in Int_{Bob} \cup \{\perp\}$ is the cognitive state.



Trust Game: Cognitive Modelling

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

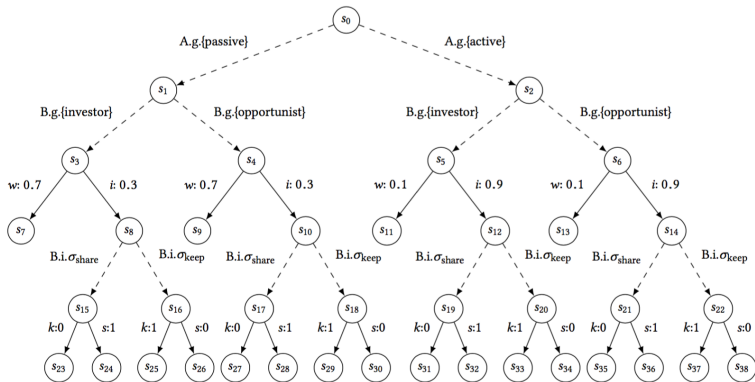


Fig. 2. Trust game with cognitive dimension



Assumptions

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

- (Uniformity Assumption) ...
- (Deterministic Behaviour Assumption) An $SMG_{\Omega} \mathcal{M}$ satisfies the *Deterministic Behaviour Assumption* if each agent's cognitive state deterministically decides its behaviour in terms of selection of its next local action. In other words, agent's cognitive state induces a pure action strategy that agent follows.



+ Cognitive Modalities

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

The syntax of the logic, named $PCTL_{\Omega}^*$, is as follows.

$$\begin{aligned} \phi &::= p \mid \neg\phi \mid \phi \vee \phi \mid \forall\psi \mid \mathbf{P}^{\bowtie q}\psi \mid \mathbf{G}_A\phi \mid \mathbf{I}_A\phi \mid \mathbf{C}_A\phi \\ \psi &::= \phi \mid \neg\psi \mid \psi \vee \psi \mid \bigcirc\psi \mid \psi\mathbf{U}\psi \end{aligned}$$

where $p \in AP$, $A \in Ags$, $\bowtie \in \{<, \leq, >, \geq\}$, and $q \in [0, 1]$.

- $\mathcal{M}, \rho s \models \mathbf{G}_A\phi$ if $\forall x \in \text{supp}(\pi_A^g(\rho s)) \exists s' : s \xrightarrow{A.g.x} s'$ and $\mathcal{M}, \rho s s' \models \phi$,
- $\mathcal{M}, \rho s \models \mathbf{I}_A\phi$ if $\forall x \in \text{supp}(\pi_A^i(\rho s)) \exists s' \in S : s \xrightarrow{A.i.x} s'$ and $\mathcal{M}, \rho s s' \models \phi$,
- $\mathcal{M}, \rho s \models \mathbf{C}_A\phi$ if $\exists x \in \text{Int}_A(s) \exists s' \in S : s \xrightarrow{A.i.x} s'$ and $\mathcal{M}, \rho s s' \models \phi$.



Example Formulas

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

- $\phi_1 = \mathbb{G}_{Alice} P^{\leq 0.9} \diamond a_{Alice} = invest$ expresses that regardless of Alice changing her goals, the probability of her investing in the future is no greater than 90%.
- $\phi_2 = \mathbb{C}_{Bob} P^{\leq 0} \circ a_{Bob} = keep$ states that Bob has a legal intention which ensures that he will not keep the money as his next action.
- $\phi_3 = \overline{\mathbb{I}_{Alice}} \exists \diamond richer_{Alice, Bob}$, where $richer_{Alice, Bob}$ is an atomic proposition with obvious meaning, states that Alice can find an intention such that it is possible to eventually reach a state where Alice has more money than Bob. Finally, the formula
- $\phi_4 = \overline{\mathbb{I}_{Alice}} \exists \diamond \mathbb{G}_{Bob} \forall \diamond \neg richer_{Alice, Bob}$ expresses that Alice can find an intention such that it is possible to reach a state such that, for all possible Bob's goals, the game will always reach a state in which Bob is no poorer than Alice.



Trust Game: Cognitive Modelling

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning

Verification

Safety Definition

Challenges

Approaches

Experimental

Results

Verification in human-robot interaction

Motivation

Stochastic

Multiplayer

Game

Cognitive

Mechanism

A Temporal

Logic of Trust

Complexity

Conclusion

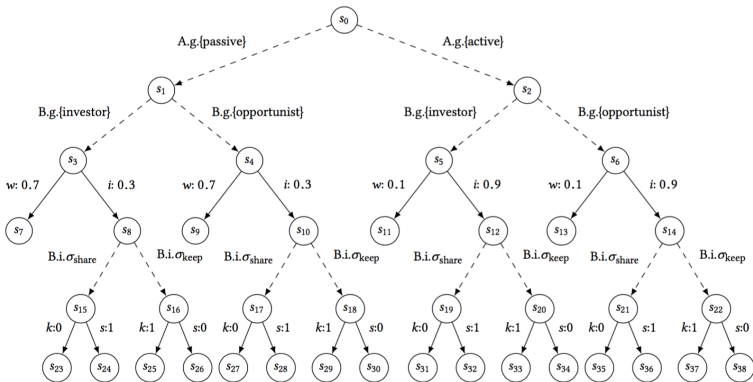


Fig. 2. Trust game with cognitive dimension



+ Preference

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

An *autonomous stochastic multi-agent system* (ASMAS) is a tuple $\mathcal{M} = (Ags, S, S_{init}, \{Act_A\}_{A \in Ags}, T, L, \{O_A\}_{A \in Ags}, \{obs_A\}_{A \in Ags}, \{\Omega_A\}_{A \in Ags}, \{\pi_A\}_{A \in Ags}, \{p_A\}_{A \in Ags})$, where p_A is a set of preference functions of agent $A \in Ags$, defined as

$$p_A \stackrel{\text{def}}{=} \{gp_{A,B}, ip_{A,B} \mid B \in Ags \text{ and } B \neq A\},$$

where:

- $gp_{A,B} : S \rightarrow \mathcal{D}(\mathcal{P}(Goal_B))$ is a goal preference function of A over B such that, for any state s and $x \in \mathcal{P}(Goal_B)$, we have $gp_{A,B}(s)(x) > 0$ only if $x \in Goal_B(s)$, and
- $ip_{A,B} : S \rightarrow \mathcal{D}(Int_B)$ is an intention preference function of A over B such that, for any state s and $x \in Int_B$, we have $ip_{A,B}(s)(x) > 0$ only if $x \in Int_B(s)$.



Trust Game: Preference-induced DTMC

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

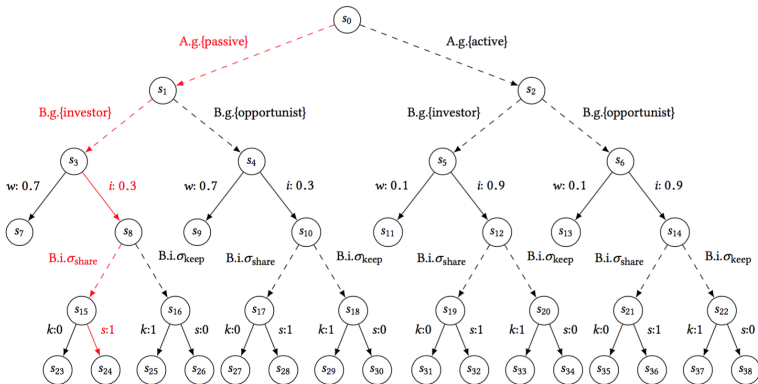
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion





Trust Game: Preference-induced DTMC

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

$$gp_{Bob,Alice}(s_0) = \langle passive \mapsto 1/3, active \mapsto 2/3 \rangle$$

indicates that Bob believes Alice is more likely to be *active* than *passive*. Setting

$$gp_{Alice,Bob}(s_x) = \langle investor \mapsto 1/2, opportunist \mapsto 1/2 \rangle,$$

for $x \in \{1, 2\}$, represents that Alice has no prior knowledge regarding Bob's mental attitudes. We may set

$$ip_{Alice,Bob}(s_x) = \langle share \mapsto 3/4, keep \mapsto 1/4 \rangle \quad \text{for } x \in \{8, 12\},$$

$$ip_{Alice,Bob}(s_x) = \langle share \mapsto 0, keep \mapsto 1 \rangle \quad \text{for } x \in \{10, 14\}$$

to indicate that Alice knows that Bob will keep the money when he is an *opportunist*, but she thinks it's quite likely that he will share his profit when he is an *investor*.



Trust Game: Preference-induced DTMC

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

$$\begin{aligned}\Pr_{Alice}(\rho_1) &= gP_{Alice,Bob}(s_1)(investor) \\ &\quad \cdot (\sigma_{passive}(s_0s_1s_3)(invest) \cdot T(s_3, invest)(s_8)) \\ &\quad \cdot iP_{Alice,Bob}(s_8)(share) \\ &\quad \cdot (\sigma_{share}(s_0s_1s_3s_8s_{15})(share) \cdot T(s_{15}, share)(s_{24})) \\ &= \frac{1}{2} \cdot \left(\frac{3}{10} \cdot 1\right) \cdot \frac{3}{4} \cdot (1 \cdot 1) = \frac{9}{80},\end{aligned}$$



Belief

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

The belief function $\text{be}_A : OPath_A \rightarrow \mathcal{D}(\text{FPath}^{\mathcal{M}})$ is given by

$$\text{be}_A(o)(\rho) = \Pr_A^{\mathcal{M}}(C_\rho \mid \bigcup_{\rho' \in \text{class}(o)} C_{\rho'}).$$



Trust Game: Belief Computation

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

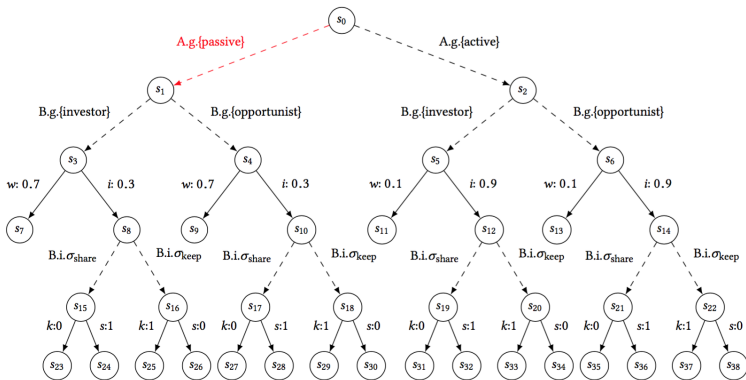
Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion





Trust Game: Belief Computation

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

$$\begin{aligned} \text{be}_{Bob}(o, \rho_1) &= \Pr_{Bob}^G(C_{\rho_1} \mid \bigcup_{\rho \in \text{class}(o)} C_{\rho}) \\ &= \frac{\Pr_{Bob}^G(C_{\rho_1})}{\Pr_{Bob}^G(C_{\rho_1}) + \Pr_{Bob}^G(C_{\rho_2})} \\ &= \frac{gp_{Bob,Alice}(s_0)(\text{passive})}{gp_{Bob,Alice}(s_0)(\text{passive}) + gp_{Bob,Alice}(s_0)(\text{active})} \\ &= \frac{1}{3}. \end{aligned}$$



+ Trust: A Temporal Logic of Trust ²

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

The syntax of the logic PRTL* is as follows.

$$\phi ::= p \mid \neg\phi \mid \phi \vee \phi \mid \forall\psi \mid P^{\bowtie q}\psi \mid G_A\phi \mid I_A\phi \mid C_A\phi \mid$$

$$B_A^{\bowtie q}\psi \mid CT_{A,B}^{\bowtie q}\psi \mid DT_{A,B}^{\bowtie q}\psi$$

$$\psi ::= \phi \mid \neg\psi \mid \psi \vee \psi \mid \bigcirc\psi \mid \psi U\psi \mid \square\psi$$

where $p \in AP$, $A, B \in A_g$ s, $\bowtie \in \{<, \leq, >, \geq\}$, and $q \in [0, 1]$.

²X. Huang and M. Kwiatkowska. *Reasoning about cognitive trust in stochastic multiagent systems*. AAI-2017.



Reasoning framework PRTL*

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

$\mathbb{B}_A^{\bowtie q} \psi$, **belief** formula, expresses that agent A believes ψ with probability in relation \bowtie with q .

$\mathbb{CT}_{A,B}^{\bowtie q} \psi$, **competence trust** formula, expresses that agent A trusts agent B with probability in relation \bowtie with q on its capability of completing the task ψ

$\mathbb{DT}_{A,B}^{\bowtie q} \psi$, **disposition trust** formula, expresses that agent A trusts agent B with probability in relation \bowtie with q on its willingness to do the task ψ , where the state of willingness is interpreted as unavoidably taking an intention.



Semantics

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

We write

$$\Pr_{\mathcal{M}, A, \rho}^{\max, \min}(\psi) \stackrel{\text{def}}{=} \sup_{\sigma_A \in \Pi_A} \inf_{\sigma_{A_{\text{gs}} \setminus \{A\}} \in \Pi_{A_{\text{gs}} \setminus \{A\}}} \Pr_{\mathcal{M}, \sigma, \rho}(\psi),$$

$$\Pr_{\mathcal{M}, A, \rho}^{\min, \max}(\psi) \stackrel{\text{def}}{=} \inf_{\sigma_A \in \Pi_A} \sup_{\sigma_{A_{\text{gs}} \setminus \{A\}} \in \Pi_{A_{\text{gs}} \setminus \{A\}}} \Pr_{\mathcal{M}, \sigma, \rho}(\psi)$$

to denote the strategic ability of agent A in implementing ψ on a finite path ρ . Intuitively,

- $\Pr_{\mathcal{M}, A, \rho}^{\max, \min}(\psi)$ gives a **lower bound on agent A 's ability to maximise** probability of ψ , while
- $\Pr_{\mathcal{M}, A, \rho}^{\min, \max}(\psi)$ gives an **upper bound on agent A 's ability to minimise** probability of ψ .



Semantics

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust

Complexity

Conclusion

For a measurable function $f : \text{FPath}^{\mathcal{M}} \rightarrow [0, 1]$, we denote by $E_{\text{be}_A}[f]$ the **belief-weighted expectation** of f , i.e.,

$$E_{\text{be}_A}[f] = \sum_{\rho \in \text{FPath}^{\mathcal{M}}} \text{be}_A(\rho) \cdot f(\rho).$$



Semantics

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

- $\mathcal{M}, \rho \models \mathbb{B}_A^{\bowtie q} \psi$ if

$$E_{\text{be}_A} [V_{\mathbb{B}, \mathcal{M}, \psi}^{\bowtie}] \bowtie q,$$

where the function $V_{\mathbb{B}, \mathcal{M}, \psi}^{\bowtie} : \text{FPath}^{\mathcal{M}} \rightarrow [0, 1]$ is such that

$$V_{\mathbb{B}, \mathcal{M}, \psi}^{\bowtie}(\rho') = \begin{cases} \Pr_{\mathcal{M}, A, \rho'}^{\max, \min}(\psi) & \text{if } \bowtie \in \{\geq, >\} \\ \Pr_{\mathcal{M}, A, \rho'}^{\min, \max}(\psi) & \text{if } \bowtie \in \{<, \leq\} \end{cases}$$



Semantics

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust

Complexity

Conclusion

- $\mathcal{M}, \rho \models \mathbb{C}\mathbb{T}_{A,B}^{\boxtimes q} \psi$ if

$$E_{\text{be}_A} [V_{\mathbb{C}\mathbb{T}, \mathcal{M}, B, \psi}^{\boxtimes}] \boxtimes q,$$

where the function $V_{\mathbb{C}\mathbb{T}, \mathcal{M}, B, \psi}^{\boxtimes} : \text{FPath}^{\mathcal{M}} \rightarrow [0, 1]$ is such
that $V_{\mathbb{C}\mathbb{T}, \mathcal{M}, B, \psi}^{\boxtimes}(\rho') =$

$$\left\{ \begin{array}{ll} \sup_{x \in \text{Int}_B(\text{last}(\rho'))} \Pr_{\mathcal{M}, A, B, i(\rho', x)}^{\text{max}, \text{min}}(\psi) & \text{if } \boxtimes \in \{\geq, >\} \\ \inf_{x \in \text{Int}_B(\text{last}(\rho'))} \Pr_{\mathcal{M}, A, B, i(\rho', x)}^{\text{min}, \text{max}}(\psi) & \text{if } \boxtimes \in \{<, \leq\} \end{array} \right.$$



Semantics

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

$$\blacksquare \mathcal{M}, \rho \models \text{DT}_{A,B}^{\bowtie q} \psi \text{ if}$$

$$E_{\text{be}_A}[V_{\text{DT}, \mathcal{M}, B, \psi}^{\bowtie}] \bowtie q,$$

where the function $V_{\text{DT}, \mathcal{M}, B, \psi}^{\bowtie} : \text{FPath}^{\mathcal{M}} \rightarrow [0, 1]$ is such
that $V_{\text{DT}, \mathcal{M}, B, \psi}^{\bowtie}(\rho') =$

$$\left\{ \begin{array}{ll} \inf_{x \in \text{supp}(\pi_B^i(\rho'))} \Pr_{\mathcal{M}, A, B, i(\rho', x)}^{\max, \min}(\psi) & \text{if } \bowtie \in \{\geq, >\} \\ \sup_{x \in \text{supp}(\pi_B^i(\rho'))} \Pr_{\mathcal{M}, A, B, i(\rho', x)}^{\min, \max}(\psi) & \text{if } \bowtie \in \{<, \leq\} \end{array} \right.$$



Example Formulas

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

The formula

$$\text{DT}_{Alice, Bob}^{\geq 0.9} \diamond (a_{Bob} = \textit{keep})$$

states that Alice can trust Bob with probability no less than 0.9 that he will keep the money for himself. The formula

$$\square (\textit{richer}_{Bob, Alice} \rightarrow \text{P}^{\geq 0.9} \diamond \text{CT}_{Bob, Alice}^{\geq 1.0} \textit{richer}_{Alice, Bob})$$

states that, at any point of the game, if Bob is richer than Alice, then with probability at least 0.9, in future, he can almost surely, i.e., with probability 1, trust Alice on her capability of becoming richer than Bob.



Guarding Mechanism

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game

Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

For every agent $A \in \text{Ags}$, we define:

- a *goal guard* function $\lambda_A^g : \mathcal{P}(\text{Goal}_A) \rightarrow \mathcal{L}_A(\text{PRTL}^*)$ and
- an *intention guard* function $\lambda_A^i : \text{Int}_A \times \mathcal{P}(\text{Goal}_A) \rightarrow \mathcal{L}_A(\text{PRTL}^*)$.

where $\mathcal{L}_A(\text{PRTL}^*)$ is the set of formulas of the language PRTL* that are boolean combinations of atomic propositions and formulas of the form $\mathbb{B}_A^{\boxtimes q} \psi$, $\mathbb{T}_{A,B}^{\boxtimes q} \psi$, $\mathbb{B}_A^{\boxtimes ?} \psi$ or $\mathbb{T}_{A,B}^{\boxtimes ?} \psi$, such that ψ does not contain temporal operators.

- Let $\Lambda = \{\langle \lambda_A^g, \lambda_A^i \rangle\}_{A \in \text{Ags}}$ be the *guarding mechanism*.



Pro-Attitude Synthesis

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

Obtaining cognitive strategy $\Pi = \{\pi_A^g, \pi_A^i\}_{A \in \text{Ags}}$ from finite structures $\Omega = \{\langle \text{Goal}_A, \text{Int}_A \rangle\}_{A \in \text{Ags}}$ and Λ



Trust Game

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust

Complexity

Conclusion

We recall our informal assumption that Bob's intention will be *share* when he is an investor and his belief in Alice being active is sufficient, and *keep* otherwise. We formalise it as follows:

$$\lambda_{Bob}^i(\text{share}, \{\text{investor}\}) = \mathbb{B}_{Bob}^{>0.7} \text{active}_{Alice},$$

$$\lambda_{Bob}^i(\text{keep}, \{\text{investor}\}) = \neg \mathbb{B}_{Bob}^{>0.7} \text{active}_{Alice},$$

$$\lambda_{Bob}^i(\text{share}, \{\text{opportunist}\}) = \perp,$$

$$\lambda_{Bob}^i(\text{keep}, \{\text{opportunist}\}) = \top,$$

where active_{Alice} holds in states in which Alice's goal is *active* and we used a value 0.7 to represent Bob's belief threshold.



Trust Game

Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches
Experimental Results

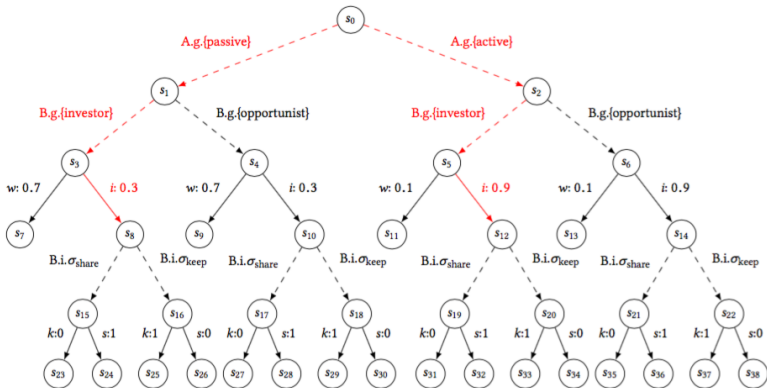
Verification in human-robot interaction

Motivation
Stochastic Multiplayer Game
Cognitive Mechanism
A Temporal Logic of Trust
Complexity

Conclusion

We let $\rho_1 = s_0s_1s_3s_8$ and $\rho_2 = s_0s_2s_5s_{12}$. Recall that $obs_{Bob}(\rho_1) = obs_{Bob}(\rho_2)$ and we let o_1 denote the observation.

$$be_{Bob}(o_1, \rho_1) = 1/7, \quad be_{Bob}(o_1, \rho_2) = 6/7.$$





Trust Game

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches
Experimental
Results

Verification in
human-robot
interaction

Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism

A Temporal
Logic of Trust
Complexity

Conclusion

Therefore, since $\mathcal{G}, \rho_1 \models \neg active_{Alice}$ and $\mathcal{G}, \rho_2 \models active_{Alice}$ (below and in what follows, $j \in \{1, 2\}$):

$$\mathcal{G}, \rho_j \models \mathbb{B}_{Bob}^{=6/7} active_{Alice}.$$

Hence

$$eval_{Bob}^i(share, \{investor\})(\rho_j) = 1,$$

$$eval_{Bob}^i(keep, \{investor\})(\rho_j) = 0,$$

and so:

$$\pi_{Bob}^i(\rho_j)(share) = 1, \quad \pi_{Bob}^i(\rho_j)(keep) = 0.$$



Model Checking Complexity

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

Deep Learning
Verification

Safety Definition
Challenges

Approaches

Experimental
Results

Verification in
human-robot
interaction

Motivation

Stochastic
Multiplayer
Game

Cognitive
Mechanism

A Temporal
Logic of Trust

Complexity

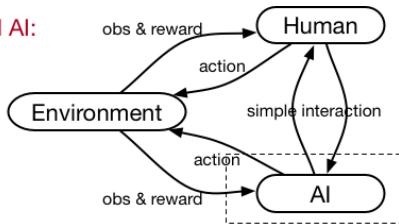
Conclusion

- general problem is undecidable
- A few fragments have been identified to be decidable in e.g., PSPACE, EXPTIME, or PTIME

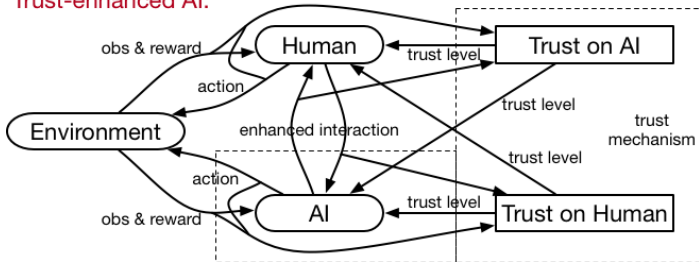


Trust-Enhanced AI

Traditional AI:

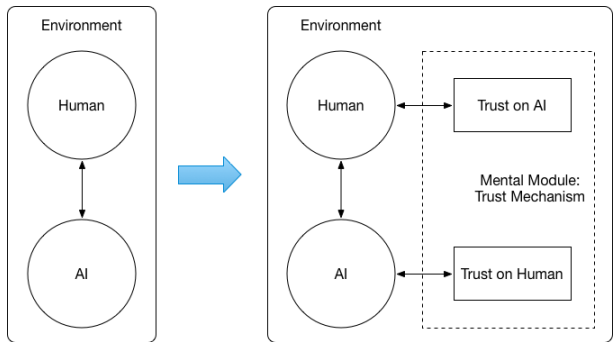


Trust-enhanced AI:





Human-like AI



Human-like AI: enhance AI with mental module (e.g., a trust mechanism) to learn and reason about human's values (e.g., trustworthiness, morality, ethics, etc.)



Conclusion

Verification of
Robotics and
Autonomous
Systems

Xiaowei
Huang

Challenges

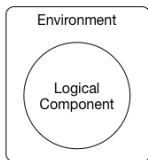
Deep Learning
Verification
Safety Definition
Challenges
Approaches
Experimental
Results

Verification in
human-robot
interaction

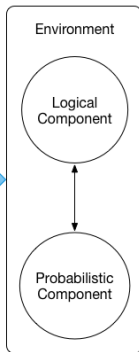
Motivation
Stochastic
Multiplayer
Game
Cognitive
Mechanism
A Temporal
Logic of Trust
Complexity

Conclusion

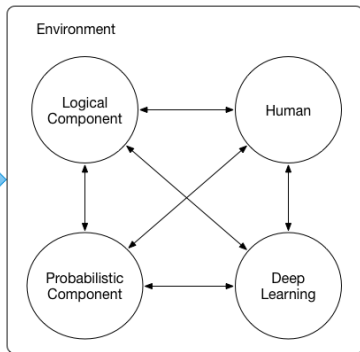
Concurrent System (1980-)



Probabilistic System (1990-)



Robotics and Autonomous System





Verification of Robotics and Autonomous Systems

Xiaowei Huang

Challenges

Deep Learning Verification

Safety Definition Challenges

Approaches

Experimental Results

Verification in human-robot interaction

Motivation

Stochastic Multiplayer Game

Cognitive Mechanism

A Temporal Logic of Trust Complexity

Conclusion





Xiaowei Huang and Marta Kwiatkowska.

Reasoning about cognitive trust in stochastic multiagent systems.

In *AAAI 2017*, pages 3768–3774, 2017.



Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu.

Safety verification of deep neural networks.

In *CAV 2017*, pages 3–29, 2017.