

# Trustworthy Artificial Intelligence

Sargur N. Srihari  
srihari@buffalo.edu

# Topics in Trustworthy AI

- Scope of AI
- Five Subdisciplines of AI
- Four Aspects of Trustworthy AI

# Definition of AI

- AI enables performing tasks requiring human cognition/decision-making
  - Thus understanding thought, intelligent behavior and their implementation is useful
- AI endeavor is multidisciplinary, encompassing
  - Perceive, learn, reason, communicate, act in the world
  - Exhibit flexibility, resourcefulness, creativity, real-time responsiveness, and long-term reflection
  - Use a variety of representation or reasoning approaches
  - Demonstrate competence in complex environments and social contexts

# Subdisciplines of AI

- Following is a definition of AI and its principle disciplines
  - Core AI
  - Biologically-inspired AI
  - Perception and Communication
  - Embodied AI
  - Trustworthy AI
  - AI/ML Theory
  - AI Hardware Architectures
- Following slides describe each subdiscipline

# 1. Core AI

- Theory and methods for these target abilities
  - Learning, abstraction, and inference
  - Architectures for intelligence, multiagent systems
- ML has made great advances through
  - algorithms, computing power, and growing data
- Other technologies
  - Knowledge representation
  - Logical and probabilistic reasoning,
  - Planning,
  - Search, constraint satisfaction, and optimization

## 2. Biologically-inspired AI

- Models may be inspired by living systems
  - Connectionism, behavior, emergence
- Computational neuroscience
  - Theory of computation in the nervous system
- Behavioral and cognitive science
  - Typical of human perceptual, motor, and cognitive processes and their interactions

# 3. Perception and Communication

- Computer vision
  - Methods to sense and reason about visual world
- Human language technologies (also called NLP, NLU)
  - To analyze, produce, translate, and respond to human text and speech

## 4. Embodied AI

- Intelligent systems may be able to act upon the world through embodiment
- Robotics is closely aligned with but not identical to embodied AI
- An embodied AI may be a robot
- This solicitation does not include *teleoperated robots* or *industrial robots* that simply repeat programmed patterns of motion

# 5. Trustworthy AI

- AI amplifies human capabilities to accomplish individual and collective goals,
  - Need to assess benefits, effects, and risks
  - How human, technical, and contextual aspects of systems interact to shape those effects
- Relevant research areas
  - Explainable and trustworthy AI
  - Validation of AI-enabled systems
  - AI safety, security, and privacy (including, for example, AI for security and security for AI),
  - Role of emotion and affect in the design and perception of AI

# Aspects of Trustworthy AI

- Trust is key in adoption of AI for economic growth and innovations to benefit society
- Today, ability to understand AI decisions and measure their trustworthiness is limited
- Aspects of trustworthy AI include:
  - a) Trusted to function *Reliably*
  - b) Trusted to be able to *Explain* conclusions
  - c) Trusted not violate *Privacy*
  - d) Trusted not to exhibit socially *Harmful Bias*

# (a) Trusted to Function Reliability

- Must go beyond performance on a test set
  - Mission-critical applications may require guarantees of worst-case solution quality and time to solution
  - With room for error, we can provide probabilistic guarantees – but for safety-critical applications, guarantee needs to be precise and explicit
    - E.g., an algorithm may be reliable in the face of random noise, but not in the face of an adversarial attack
  - Research challenges include
    - detecting malicious inputs, understanding of model characteristics that might be taken advantage of by malicious attackers, richer measures of reliability, and verifying robustness formally and experimentally.

# (b) Trusted to be able to Explain

- Explanations are vital in decision making
  - Establishing human trust in the outcome requires the exchange of reasons for that outcome
  - Explanations must be in terms as appropriate to the task and as needed by users.
  - Explanations help pinpoint errors or data
- Research challenges include finding ways to make "black box" AI systems explainable
  - Models and frameworks for learning and reasoning that are both inherently explainable and powerful
  - Integrating psychology, cognitive science, to better understand and acceptability of an explanation

# (c) Trusted not to violate privacy

- Privacy of people whose data are used to train
  - Not leak personal information during/after training
- Approaches to enhancing privacy in ML
  - Federated learning, homomorphic encryption, differential privacy, and synthetic datasets
- Solutions needed to allow people to establish boundaries on the use of their personal data

## (d) Trusted not to have harmful bias

- Harmful bias arises from incorrect assumptions
  - when training data doesn't reflect full domain
  - when system fails to respect legal and social norms
- Research challenges include
  - Creating notions of fairness in decisions
  - Contexts in which they are achievable or impossible
  - Detecting and measuring data bias
  - Preventing "amplified" bias by learning algorithms
  - Teaching human rules against unlawful bias

# Requirements of Trustworthy AI

