

ARTIFICIAL INTELLIGENCE

The Very Idea

Vasant G. Honavar

Dorothy Foehr Huck and J. Lloyd Huck Chair in Biomedical Data Sciences and Artificial Intelligence Professor of Data Sciences, Informatics, Computer Science, Bioinformatics & Genomics and Neuroscience Director, Artificial Intelligence Research Laboratory

Director, Center for Artificial Intelligence Foundations and Scientific Applications Associate Director, Institute for Computational and Data Sciences Pennsylvania State University

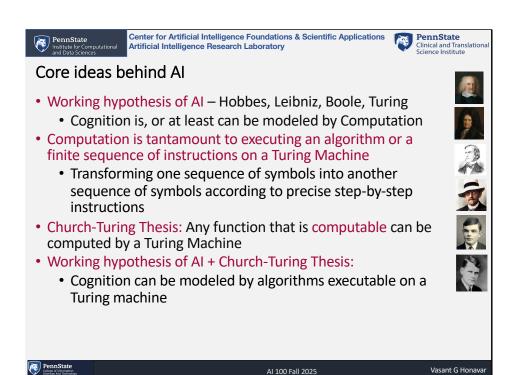
vhonavar@psu.edu http://faculty.ist.psu.edu/vhonavar http://ailab.ist.psu.edu



AI 100 Fall 2025

Vasant G Honavar

PennState
Clinical and Translationa



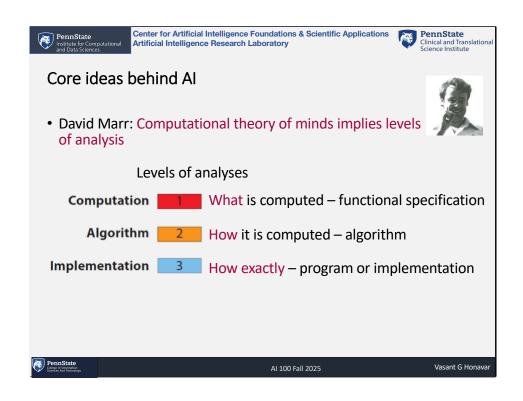


Core ideas behind Al

- The working hypothesis of AI implies a computational theory of minds
- · We will have a theory of
 - Problem-solving if we can devise algorithms that can solve problems
 - Game-playing if we can devise algorithms that play games
 - Reasoning if we can devise algorithms that reason from assumptions to conclusions
 - Learning if we can devise algorithms for learning from experience
 - Language if we can devise algorithms that effectively communicate using language
 - Cooperation if we have a computational model of multi-agent collaboration
 - Creativity if we can devise algorithms that exhibit creativity



AI 100 Fall 2025



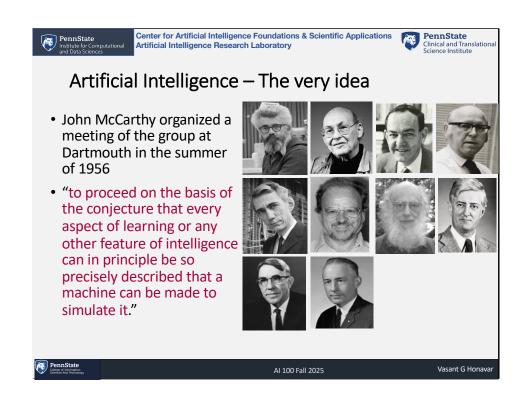


The focus of this course

- In this course, our focus is on
 - what is being computed by machines that behave as if they have minds
 - how it is computed by an algorithm
 - but not necessarily how exactly it is being computed that is, realized by a program
- Church-Turing thesis says that a given algorithm can have many implementations
- Details of implementation of the algorithms of the mind on specific physical substrates is left to other courses
 - Al implementation using programs written in a specific programming language executed on computers
 - Neuroscience implementation in brains



AI 100 Fall 2025



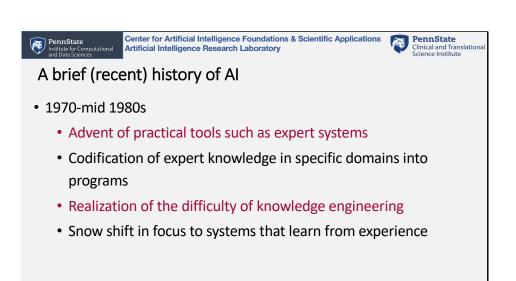


A brief (recent) history of AI

- Birth of artificial intelligence (1956)
- Early demonstrations of artificial intelligence and the publication of Computers and Thought (1959)
- 1960-1970
 - Optimism fueled by early success on some problems thought to be hard (e.g., theorem proving)
 - Slow progress on many problems thought to be easy (e.g., vision, language);
 - Fragmentation of AI into sub-areas focused on problemsolving, knowledge representation and inference, vision, planning, language processing, learning, etc.

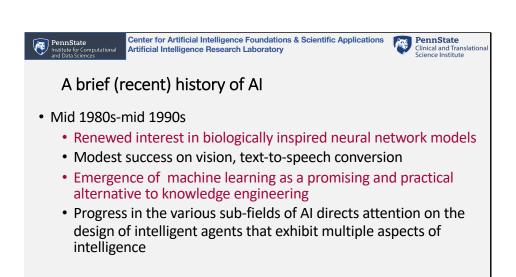


AI 100 Fall 2025



PennState
College of Information

AI 100 Fall 2025



PennState
College of Information
Sciences And Technolog

AI 100 Fall 2025



A brief (recent) history of A

- Mid 1990s-2010
 - The advent of the World-Wide-Web, big data, and computing advances lead to AI systems trained on massive amounts of data
 - Major breakthroughs in learning theory offer insights that lead to practical advances, e.g., kernel machines, in machine learning
 - Large data coupled with machine learning yield breakthroughs in many applications: information retrieval, computer vision, information extraction, robotics, among others
 - Successes in design of intelligent agents shifts attention to multiagent systems – inter-agent communication, coordination, and multi-agent organizations

PennState
College of Information
Sciences And Technology

AI 100 Fall 2025



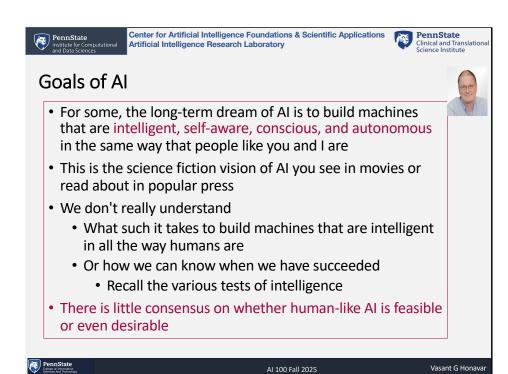
- 2020-present
 - Large language models capture public interest and imagination

progress on computer vision, natural language processing

- Advances in powerful AI systems that could automate aspects of human intellectual work highlight ethical concerns
- Goals of AI shift from automating intelligent behavior to augmenting and extending human intellect and abilities

PennState
College of Information
Sciences And Extendiby

AI 100 Fall 2025





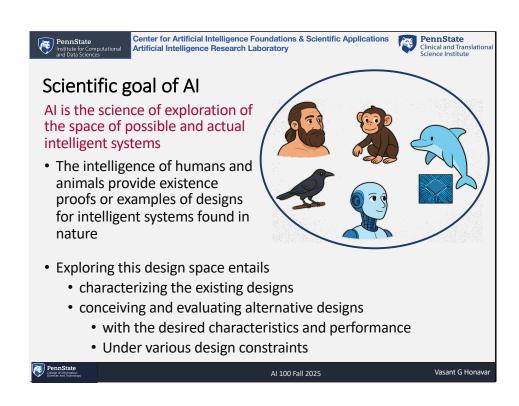
Scientific goal of Al

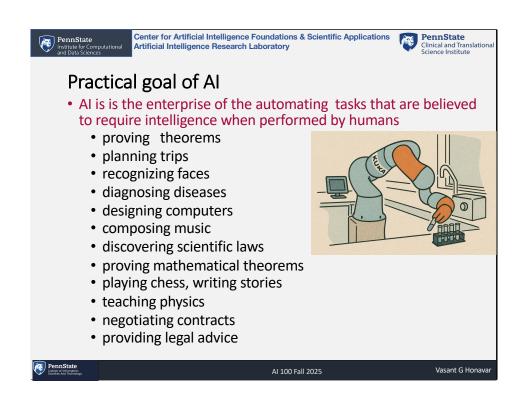
The central scientific goal of AI is to understand intelligence by building and studying computational models of intelligence

- · We will have a theory of
 - Problem-solving if we can devise algorithms that can solve problems
 - Game-playing if we can devise algorithms that play games
 - Reasoning if we can devise algorithms that reason from assumptions to conclusions
 - Learning if we can devise algorithms for learning from experience
 - Language if we can devise algorithms that effectively communicate using language
 - Cooperation if we have a computational model of multi-agent collaboration
 - Creativity if we can devise algorithms that exhibit creativity



AI 100 Fall 2025





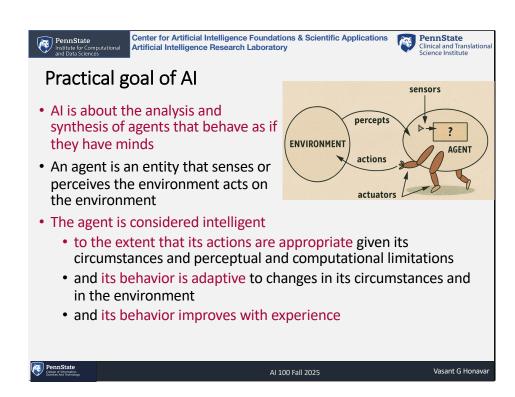


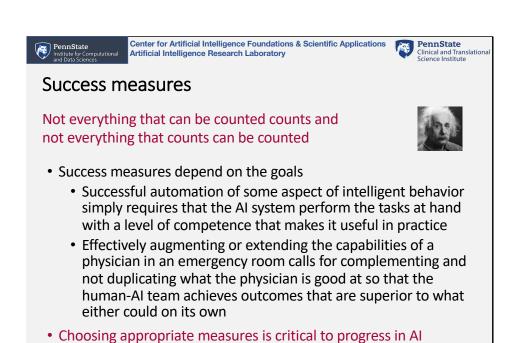
Practical goal of Al

- Al is about augmenting and extending human intelligence, problem-solving and creativity
 - An Al physician's assistant helps medical practitioners make better decisions
 - A search engine augments human memory
 - Natural language translation systems help people communicate across linguistic barriers
 - An Al writer's assistant can act as a muse
 - Al-powered scientist's assistants can help identify promising hypotheses to pursue, optimal experiments to run, and help analyze and interpret experimental results



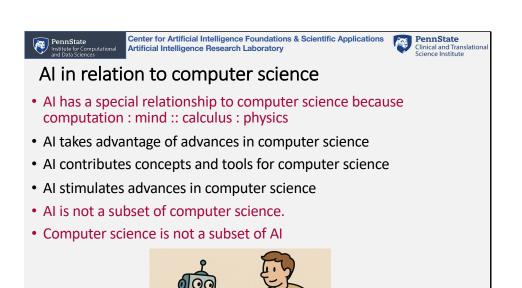
Al 100 Fall 2025





AI 100 Fall 2025





COMPUTER SCIENCE

AI 100 Fall 2025

Vasant G Honavar

20



Al in relation to psychology

- Al is often seen as a sibling of psychology
- Psychology is concerned with studies of human and animal behavior
- Al is concerned with computational models or artifacts that exhibit intelligent behavior
- Al is not committed to human-like mechanisms or human-like implementation of such mechanisms
- Computational models from AI influence research in psychology
- Findings and insights from psychology inform the design of AI models



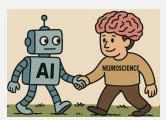
PennState
College of Information
Sciences And Technology

AI 100 Fall 2025



Al in relation to neuroscience

- AI has deep connections to neuroscience
- Neuroscience is about brains
- Al is about computational theories of brain function
- Al is informed by or at least inspired by findings in neuroscience
- Advances in machine learning offer new perspectives on old questions in neuroscience



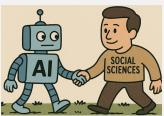
PennState
College of Information
Sciences And Technology

AI 100 Fall 2025



Al in relation to social sciences

- Insofar as AI is concerned with the design of intelligent agents that
 - act rationally
 - interact with other agents
- Al draws on ideas from economics, game theory, organizational theory, decision theory, and other areas of social sciences
- Al systems can inform studies that drive advances in the social sciences



PennState
College of Information
Sciences And Technology

AI 100 Fall 2025



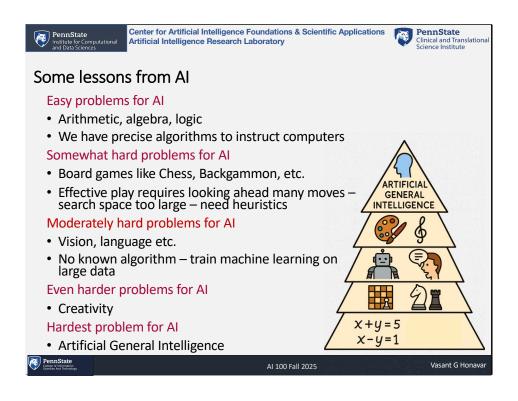
Al in relation to engineering

- Insofar as AI is concerned with the design of intelligent artifacts, it both contributes to, and draws on advances in engineering
- Al advances have resulted in practical tools for
 - configuring computer systems
 - Diagnosing faults in machinery
 - software agents that scour the Internet for information on demand
 - Intelligent systems for planning and scheduling
 - Computer-aided design tools in many engineering disciplines
 - Self-driving automobiles, Smart buildings, smart robots, etc.



PennState
College of Information
Sciences And Technology

AI 100 Fall 2025





Example: Automated driving

- Each year, 1 million people die in auto accidents worldwide
- Over 90% of auto accidents are caused by human error
- · Automating driving could save countless lives, reduce injuries
- Automating driving could result in significant job loss
- Should we automate driving?

Similar examples abound in other areas

- AI and robotics have unleashed the 4th industrial revolution
- Majority of jobs we have today will be lost or dramatically changed due to advances in AI, robotics, and automation



AI 100 Fall 2025



Al and economy

- AI will disrupt all areas of life
 - Transportation, accounting, healthcare, journalism, banking
 - What will happen to the workers who once occupied those jobs?
 - Will new jobs be created?
 - How can workers get trained for the new jobs?
 - If AI automates almost all work, what will humans do?
 - Do we need new economic models?



AI 100 Fall 2025

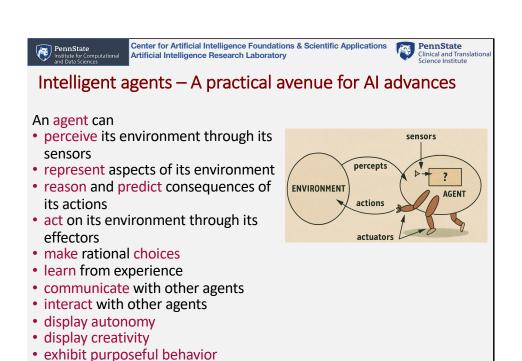


Moral Imperative of Al

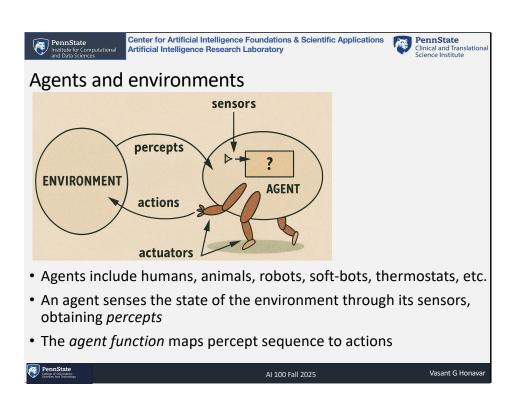
- How can we maximize the benefits of AI while minimizing the harms of AI?
- How do we build AI systems that can augment and extend human abilities?
- How do we build AI systems that can explain their decisions?
- How do we ensure that AI systems are safe?
- How do we ensure that AI systems are accountable?
- How do we ensure that AI systems are fair?
- How do we ensure that AI systems adhere to human ethical and social norms?

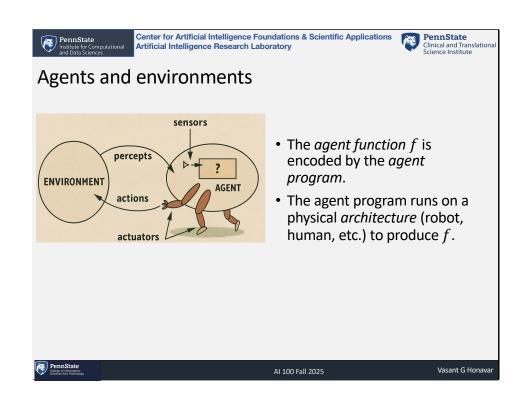


AI 100 Fall 2025



AI 100 Fall 2025

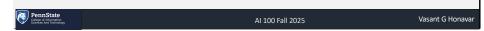


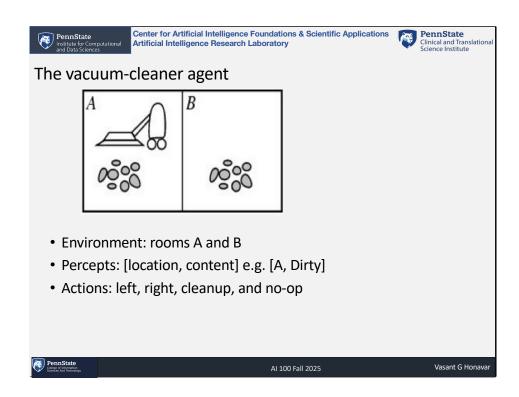


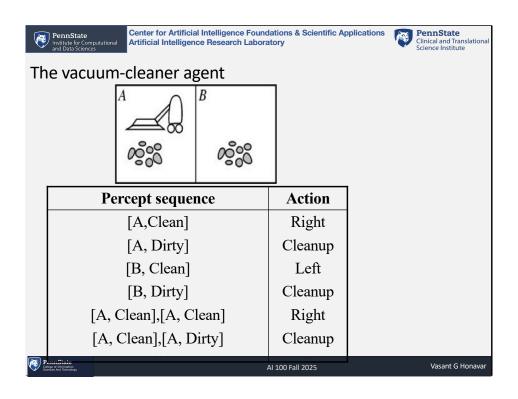


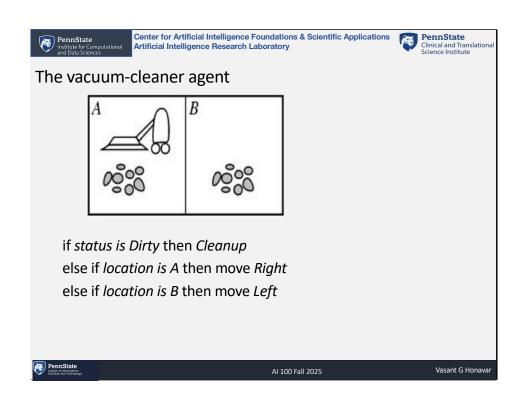
Al Thermostat

- One of the simplest agents we can imagine is a simple thermostat
- The thermostat senses the ambient temperature using its temperature sensor
- If the sensed temperature is greater than the preset temperature, it turns on the air conditioner
- If the sensed temperature is less than equal to the preset temperature, it shuts off the air conditioner if it is on
- If the sensed temperature is less than the preset temperature, it turns on the heater
- If the sensed temperature is greater than equal to the preset temperature, it shuts off the heater if it is on











- How agents ought to behave is a topic of debate in moral philosophy
- Al adopts a form of consequentialism
 - whether or not an action is the right one depends on its consequences relative to what we value
- Example:
 - Most people would agree that lying is wrong
 - But if telling a lie would help save a person's life, consequentialism says it's the right thing to do

PennState
College of Information
Sciences And Technology

AI 100 Fall 2025

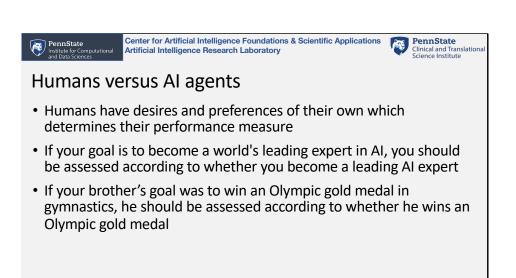


How can we relate actions to consequences?

- An agent it generates actions in response to the percepts it receives, thereby causing the environment to go through a sequence of states
- Example: the thermostat agent would ensure that the room temperature is maintained to its preset value
- If the consequence of the agent's actions is desirable, then the agent has performed well
- How does an agent know the desirability of its actions?
 - Performance measure
- An agent's performance in its environment is evaluated by the chosen performance measure



AI 100 Fall 2025



PennState
College of Information

AI 100 Fall 2025



Center for Artificial Intelligence Foundations & Scientific Applications Artificial Intelligence Research Laboratory



Humans versus AI agents

- Machines, unlike humans, do not have desires and preferences of their own
- The performance measure for an AI agent must be, initially at least, specified by
 - the AI agent's designer, or
 - · its users, or more broadly,
 - the society at large

"If we use, to achieve our purposes, a mechanical agency with whose operation we cannot interfere effectively, we had better be quite sure that the purpose put into the machine is the purpose which we really desire."

 How do we ensure that the AI system adheres to human and societal values?



AI 100 Fall 2025



Specifying good performance measures is hard

- Specifying the performance measure to reflect precisely how the agent out to behave from an individual or societal point of view is highly non-trivial
- Consequently, the more autonomous and more powerful an AI agent is, the greater the concern about ensuring that its

performance measure is aligned with human and societal values Vasant G Honavar AI 100 Fall 2025



Specifying good performance measures is hard

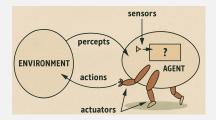
- When AI agents are designed to serve the needs of multiple users,
 - we end up with a piece of software, copies of which will serve different individuals
- We cannot possibly anticipate in advance the goals and preferences of each individual
- Even if we could, custom-designing of agents for each individual is likely to be highly impractical
- · Al agents should
 - accommodate uncertainty about the actual performance measure against which they will be assessed and
 - refine the measure over time
 - through their interactions with their respective users



AI 100 Fall 2025



How an agent should behave



What an agent ought to do at any given time depends on:

- The performance measure that defines the criterion of success against which the agent is evaluated
- The agent's prior knowledge of its environment
- The actions that the agent has at its disposal
- The agent's percept sequence up to that time



AI 100 Fall 2025

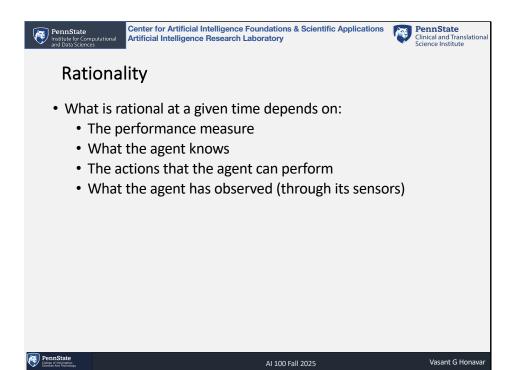


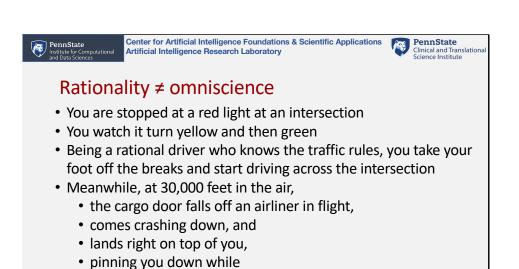
Rational Agents

- For each possible percept sequence, a rational agent should select an action that is expected to maximize its performance measure, given
 - the evidence provided by the percept sequence
 - · whatever built-in knowledge the agent has and
 - the actions at its disposal
- Examples of performance measures
 - the amount of dirt cleaned within a certain time
 - how clean the floor is
 - · the amount of dirt cleaned per unit of electricity used
- What is your performance measure?
- Who decides what the performance measure should be?
 - · Internal drives
 - External rewards



AI 100 Fall 2025





other vehicles crash into you causing a pileup

Was your behavior rational?

• If so, why?

• If not, why not?

PennState
College of Information
Sciences And Technology

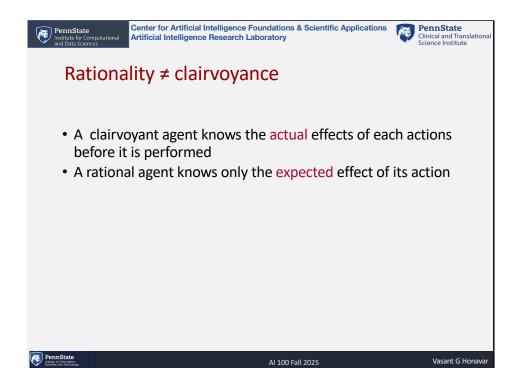
AI 100 Fall 2025

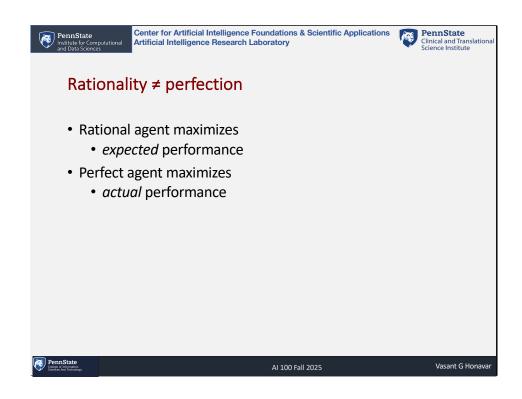


- the expected state of the world resulting from its action
- Your behavior at the intersection was rational based on what you knew
 - You are not omniscient
 - You couldn't possibly be expected to know that the cargo door of the airplane flying overhead would come crashing down on you

PennState
College of Information
Sciences And Technology

AI 100 Fall 2025







Is the thermostat agent rational?

- The performance measure awards one point for each hour the temperature is maintained within ± 2 degrees Celsius from the preset temperature
- What the agent knows:
 - The preset temperature in degrees Celsius.
 - If the air conditioner is turned on, it cools the room, causing a decrease in the temperature;
 - if the heater is turned on, it heats the room, causing an increase in the temperature;
 - The heater, once turned on, will remain on until it is turned off;
 - The air conditioner and the heater cannot both be on at the same time.
- The available actions are to turn the air conditioner on or off, and the heater on or off, or do nothing (leave things the way they are).





Is the thermostat agent rational?

- The thermometer or room temperature sensor is operational and provides the correct temperature readings
- An air conditioner status sensor tells the thermostat whether the air conditioner is on
- The heater status sensor tells the thermostat whether the heater is on
- If the sensed temperature is greater than the preset temperature, the thermostat turns on the air conditioner
- If the sensed temperature is less than equal to the preset temperature, it shuts off the air conditioner if it is on
- If the sensed temperature is less than the preset temperature, it turns on the heater
- If the sensed temperature is greater than equal to the preset temperature, it shuts off the heater if it is on



AI 100 Fall 2025



Is the new thermostat agent rational?

- Suppose the thermostat is additionally equipped with room occupancy sensor that tells it whether the room is occupied by people
- The performance measure awards the agent
 - \bullet one point for each hour the temperature is maintained within $\pm~2$ degrees Celsius of the preset temperature whenever the room is occupied; and
 - -10 points if either the heater or air conditioner are found to be on when the room is unoccupied by people
- Under these circumstances, does the agent function described previously ensure that the agent is rational? Why or why not?
- How can you restore rationality?

