**Symbols**
**Logic**
**Replace**
**Disappointment**

**Neurons**
**Probability**
**Augment**
**Doomsday**

# Where will the AI Pendulum Swing Next?

**Subbarao Kambhampati**
Arizona State University

Video of the talk available at
http://rakaposhi.eas.asu.edu/ai-pendulum.html

Pop Quiz: Magellan, the explorer, went around the world three times. On one of his trips, he died. Which trip did he die?
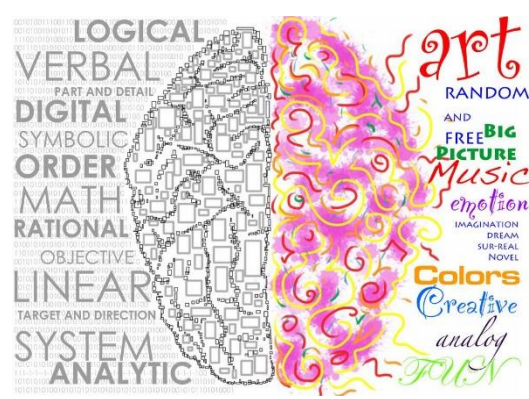
[A 3min montage video of accomplishments of AI]

# What is "Intelligence" anyway?

Clearly that indefinable quality that *you* have and your bozo friends don't…

Magellan, the explorer, went around the world three times. On one of his trips, he died.

Question: Which trip did he die in?

# Many Intelligences..





- Perceptual tasks that seem to come naturally to us
  - Form the basis for the Captchas..
    - But rarely form the basis for our own judgements about each other's intelligence

- Cognitive/reasoning tasks
  - That seem to be what we get tested in in SAT etc

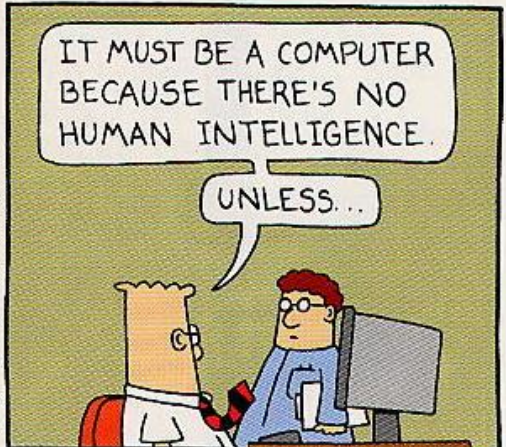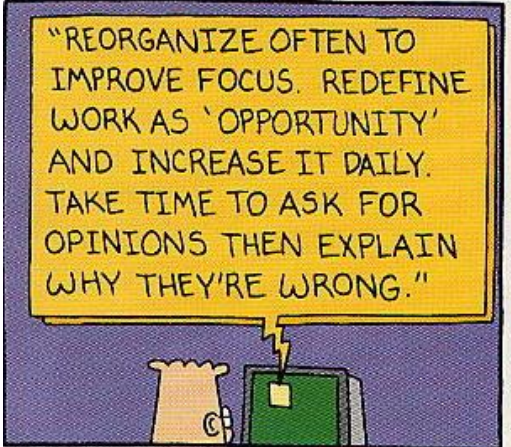- Emotional Intelligence
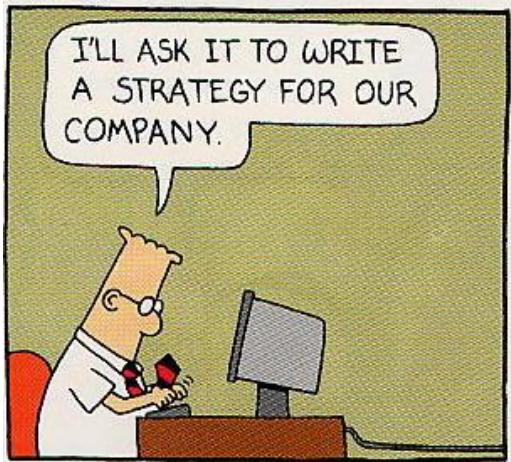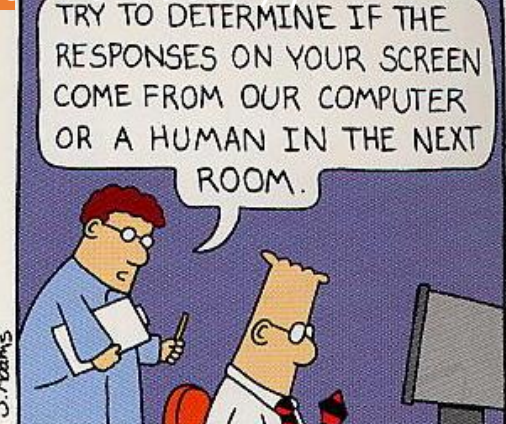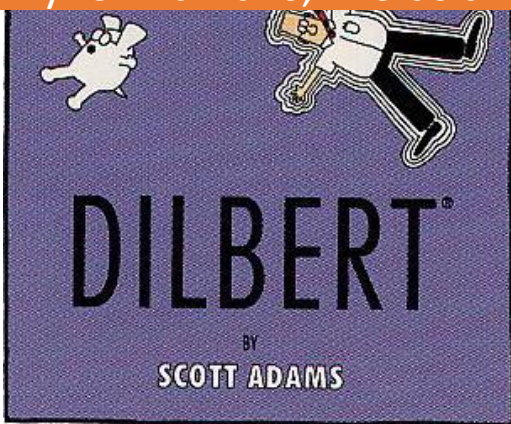
- Social Intelligence..

# When is a computer Intelligent?

- When it does tasks that, when done by a human, would be seen as requiring intelligence..
  - Nice circular definition ☺

# AI's progress towards intelligence

- 80's --- Expert systems
  - Rule-based systems for many businesses
- 90's --  Reasoning systems
  - Dethroned Kasparov
- 00's: Perceptual tasks
  - Speech recognition common place!
  - Image recognition has improved significantly
- Current:  Connecting reasoning and perception

Notice the contrast.. Human babies master perception before they get good at reasoning tasks!

7

# If you want to know limits of AI, look at the Captcha's!

- AI could imitate experts earlier than it could imitate 4 year olds..

I'm not a robot
reCAPTCHA

Type the text

250

Verify

I'm not a robot
reCAPTCHA

Select all images below that match this one:

Verify

**Qualifying question**
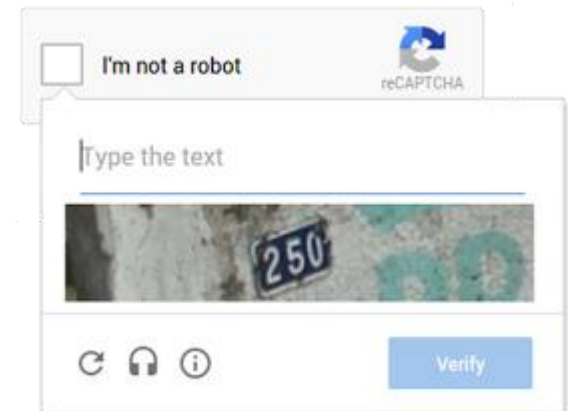
Just to prove you are a human, please answer the following math challenge.

Q: Calculate:

$$\frac{\partial}{\partial x}\left[6\cdot\sin\left(x-\frac{\pi}{2}\right)+3\cdot\cos\left(2\cdot x-\frac{\pi}{2}\right)\right]\Big|_{x=\pi}$$

A: 

*mandatory*

Note: If you do not know the answer to this question, reload the page and you'll (probably) get another, easier, question.

Password (required)
********
Birthday (required)
March | 31 | 1981
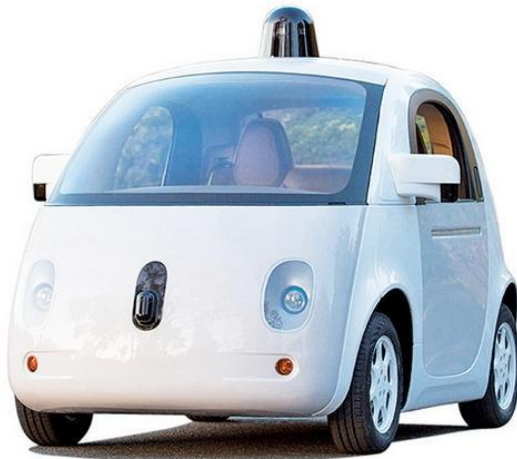Human test (required)
Type in the text you see in the box below.

Sorry, your text and the image didn't match. Please try again.

Read (really!)
☑ I have read and agree to the Terms of Use and Privacy Policy.

# Still Elusive Commonsense

- When did Magellan Die?

**Symbols**
**Logic**
**Replace**
**Disappointment**

**Neurons**
**Probability**
**Augment**
**Doomsday**

# Symbols ←→Neurons

# Symbols or Neurons?

**Neurons**

- Clearly, brain works by neurons.

**Symbols**

- But, from Greeks on, human knowledge has been codified in symbolic fashion





Qn: Should AI researchers look at symbols or neurons as the substrate?
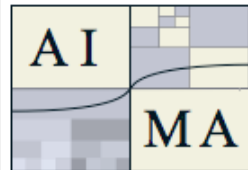
13

# Symbols or Neurons?

- "A physical symbol system has the necessary and sufficient means for general intelligent action.

    *--Allen Newell &*

    *Herbert Simon*

- Symbols are Luminiferous Aether of AI

    *—Geoff Hinton*

# Artificial Intelligence: A Modern Approach

## (Third edition) by Stuart Russell and Peter Norvig

The leading textbook in Artificial Intelligence.
Used in over **1300** universities in over **110** countries.
The 22nd most cited computer science publication on Citeseer (and 4th most cited publication of this century).

## What's New

- **Free Online AI course**, Berkeley's CS 188, offered through edX.

## Comments and Discussion

- Comments from readers
- Errata list (errors in the book)
- **AIMA-talk** discussion list, open to all

## AI Resources on the Web

- AI On the Web, a list of over 900 links
- AI Books in many categories
- AI courses that are using AIMA (1200 schools)

## Online Code Repository

- Pseudo-code algorithms from the book in pdf.
- Online code in Lisp, Python, Java etc.
- Data for the online code
- Online demos (Java applets and Javascript)
- The OpenNERO 3D multiagent simulator

## For the Instructor

- AI Instructor's Resource Page
- Lecture slides coming soon

## Table of Contents

# Deep networks were in deep hibernation for most of recent past.. But clearly the pendulum swung their way now



World Cloud of IJCAI-16 Submission Titles

# Prediction?



- AI systems were good at reasoning tasks (the SAT stuff..) before it became good at the perception tasks (vision, language understanding etc.)

- The successes on the perception front did have a lot to do with neural architectures
  - This doesn't necessarily imply that pendulum would stay at the neural end
    - Deep learning has been good until now  on non-cognitive tasks; extending their success seems to require reasoning
      - Example: Success of AlphaGo

- tldr; We might want our aether…. ☺



Google DeepMind calls the self-guided method reinforced (sic) learning, but it's really just another word for "deep learning," the current AI buzzword.
          -IEEE Spectrum (!!) 1/27/16

# Logic←→Probability

# Does Tweety Fly?

**Logic**
- Bird(x) => Fly(x)
- Bird(Tweety)
- ?
- But if I tell you Tweety is an ostrich? A magical ostrich?
- Non-monotonic logic

**Probability**
- P(TF|TB) =0.99
- P(TF|TB&TO) =0.4
- P(TF|TB&TO&TMO)= 0.8
- Posterior probabilities
  - Bayes Rule
  - P(A|B)=P(B|A)*P(A)/P(B)

# IN DEFENCE OF LOGIC

## P. J. Hayes
Essex University
Colchester, U.K.

### Introduction

Modern formal logic is the most successful precise language ever developed to express human thought and inference. Measured across any reasonably broad spectrum, including philosophy, linguistics, computer science, mathematics and artificial intelligence, no other formalism has been anything like so successful. And yet recent writers in the AI field have been almost unanimous in their condemnation of logic as a representational language, and other formalisms are in a state of rapid development.

I will argue that most of this criticism misses the point, and that the real contribution of logic is not its usual rather sparse syntax, but the semantic theory which it provides. AI is as much in need now of good semantic theories with which to compare formalisms as it always has been. I will also re-examine the procedural/declarative controversy and show how regarding representational languages as programming languages has, ironically, made procedural ideas as vulnerable to the old theorem-proving paradigm was. I will argue that the contrast between assertional and procedural languages is false: we have rather two kinds of subject-matter than two kinds of language.

This paper is deliberately polemical in tone. Much has been written from the proceduralist point of view. It's time the other arguments were put.

### Logic is not a programming system

It will, and has been, said that to defend logic is to adopt a reactionary position. Logic has been tried (in the late sixties) and found wanting; now it has been superceded by better systems, in particular, procedural languages such as uPLANNER [17] , CONNIVER f18] and more recently KRL [2].

But logic is false: we have rather two kinds of subject-matter than two kinds of language.

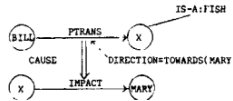performs inferences: some of its processes are the making of inferences.

But two different systems may be based on the same notion of inference and the same representational language. The inference structure of the language used by a system does not depend on the process structure. In particular, a system may have a logical inference structure - may be making deductively valid inferences - without being a classical uniform theorem-prover which just "grinds lists of clauses together".

### What logic is: the extensional analysis of meaning

One of the first tasks which faces a theory of representation is to give some account of what a representation or representational language means. Without such an account, comparisons betwer representations or languages can only be very superficial. Logical model theory provides such an analysis.

Suppose it is claimed that:

```
                    IS-A:FISH
  (BILL)  PTRANS   (  x  )
         CAUSE      DIRECTION=TOWARDS(MARY)
  ( x )   IMPACT   (MARY)
```

means that Bill hit Mary with a fish (to take a representative example), or that:

((DO(^AGENT)*BADTHING)CAUSE(^AGENT)DISPLAY
(tfNEGATIVEEMOTION)))

means that people often seem upset when bad things happen (to take another). How could one judge whether they really do mean those things? What would count as a specification of their meanings? Several answers can be suggested.

IJCAI 1977

---

## In Defense of Probability

Peter Cheeseman
SRI International
333 Ravenswood Ave., Menlo Park, California 94025

In this paper, it is argued that probability theory, when used correctly, is sufficient for the task of reasoning under uncertainty. Since numerous authors have rejected probability as inadequate for various reasons, the bulk of the paper is aimed at refuting these claims and indicating the sources of error. In particular, the definition of probability as a measure of belief rather than a frequency ratio is advocated, since a frequency interpretation of probability drastically restricts the domain of applicability. Other sources of error include the confusion between relative and absolute probability, the distinction between probability and the uncertainty of that probability. Also, the interaction of logic and probability is discusses and it is argued that many extensions of logic, such as "default logic" are better understood in a probabilistic framework. The main claim of this paper is that the numerous schemes for representing and reasoning about uncertainty that have appeared in the AI literature are unnecessary—probability is all that is needed.

### 1 Introduction

A glance through any major AI publication shows that an overwhelming proportion of papers are concerned with what might be described as the logical approach to inference and knowledge representation. It now widely accepted that many knowledge representations can be mapped into (first order) predicate calculus, and the corresponding inference procedures can be reduced to a type of controlled logical deduction. However, examples of human reasoning (judgements) are full of such terms as "probably", "most", "usually" etc., showing that many patterns of human reasoning are *not* logical in form, but intrinsically probabilistic.

The claim that many patterns of human reasoning are probabilistic does not mean that the underlying "logic" of such patterns cannot be axiomatized. On the contrary, a basis for such an axiomatization is given in section 3. The claim is that when such an exercise is performed, the resulting patterns of inference are different in form from those found in analogous logical deductions. A characteristic dif-

inference paths ("proofs") connecting the evidence to the hypothesis of interest must be examined and "combined", while in logic it is sufficient to establish a single path between the axioms and the theorem of interest. Also, the output is different, the former includes at least one numerical measure, the latter simply true or false.

Unfortunately, the logical style of reasoning is so prevalent in AI that many have attempted to force intrinsically probabilistic situations into a logical straight-jacket with predictable limited success. Two conspicuous examples of this are "Default Logic" [19] and "Non-Monotonic Logic" [15] discussed in more detail below. These methods are appropriate for dealing with some situations where limited knowledge is available. The same cannot be said for those who invent new theories for reasoning under uncertainty, such as "Certainty Factors", "Schafer/Dempster Theory", "Confirmation Theory", "Fuzzy Logic", "Endorsements" etc.

These theories will be shown below to be at best unnecessary and at worst misleading (not to mention confusing to the poor novice faced with so many possibilities). Each one is an attempt to circumvent some perceived difficulty of probability theory, but as shown below these difficulties exist only in the minds of their inventors. However, these supposed difficulties are common misconceptions of probability, generally springing from the inadequate frequency interpretation. A major aim of this paper is to put forward the older view (Bayes, Laplace etc.), that probability is a measure of belief in a proposition given particular evidence. This definition avoids the difficulties associated with the frequency definition and answers the objections of those who felt compelled to invent new theories.

An analogy can be draw between the situation in AI in the late 1970s, where Pat Hayes, in a paper entitled "In Defence of Logic" [10], found it necessary to take a broadside at the proliferation of new representation languages (with associated inference procedures) that purported to solve difficulties with the logical approach. He showed that far from being "nonlogical" it is possible to cast such languages into an equivalent logical form, and by doing so provide a clear semantics. In addition, he pointed out the obvious but unpopular fact that logic has been around for

IJCAI 1985

---

if $E$ splits two nodes $x$ and $y$, then the probabilities to be assigned to these nodes are conditionally independent given $E$.

Given these definitions, the conditional independence assumptions implicit in influence diagrams seem to be reasonable ones—if the area in a diagram like that in Figure 12.1 or Figure 12.2 correspond to the causal connections in our domain, so that the color of the light can only affect our mood via the possibility of our getting a ticket, the conclusion that sanctions the derivation of (12.5) from (12.4) seems a valid one. The reason that influence diagrams are of such interest to the probabilistic community is that they provide a compact, effective, and useful way to represent the wealth of independence assumptions needed by practical probabilistic systems.

There is another way to look at this as well. In order to evaluate all of the probabilities in the traffic example of Figure 12.1, we need to know the probability

$$pr(c \land t \land m \land l) \qquad [12.6]$$

for each choice of color $c$, ticket possibility $t$ (yes or no; did I get a ticket or not?), mood $m$ (good or bad), and loss of license $l$ (yes or no). As we have already remarked, there are twenty-four of these probabilities, and they are potentially constrained only by the requirement that they sum to 1.

Of course, we know that we can always rewrite (12.6) as

$$pr(c) \cdot pr(t|c) \cdot pr(m|c \land t) \cdot pr(l|c \land t \land m)$$

The conditional independence assumptions associated with the influence diagram allow us to rewrite this in the simpler form

$$pr(c) \cdot pr(t|c) \cdot pr(m|t) \cdot pr(l|t) \qquad [12.7]$$

Once again, only ten values are needed to evaluate the various instances of (12.7)—one probability for each color, three for the probability of getting a ticket as a function of color, and two each to give my mood and chance of losing my license depending on whether I've gotten a ticket or not. Either way we think of it, the lesson is the same:

*Influence diagrams allow us to conveniently represent the conditional independence assumptions used to reduce the amount of information needed by a probabilistic reasoner.*

### 12.4 ARGUMENTS FOR AND AGAINST PROBABILITY IN AI

My final aim in this chapter is to discuss the philosophical questions underlying the application of probability to AI. After all, the success of PROSPECTOR is not necessarily evidence that probabilities have a funda-

Published 1993

---

### 12.4 ARGUMENTS FOR AND AGAINST PROBABILITY IN AI

My final aim in this chapter is to discuss the philosophical questions underlying the application of probability to AI. After all, the success of PROSPECTOR is not necessarily evidence that probabilities have a funda-
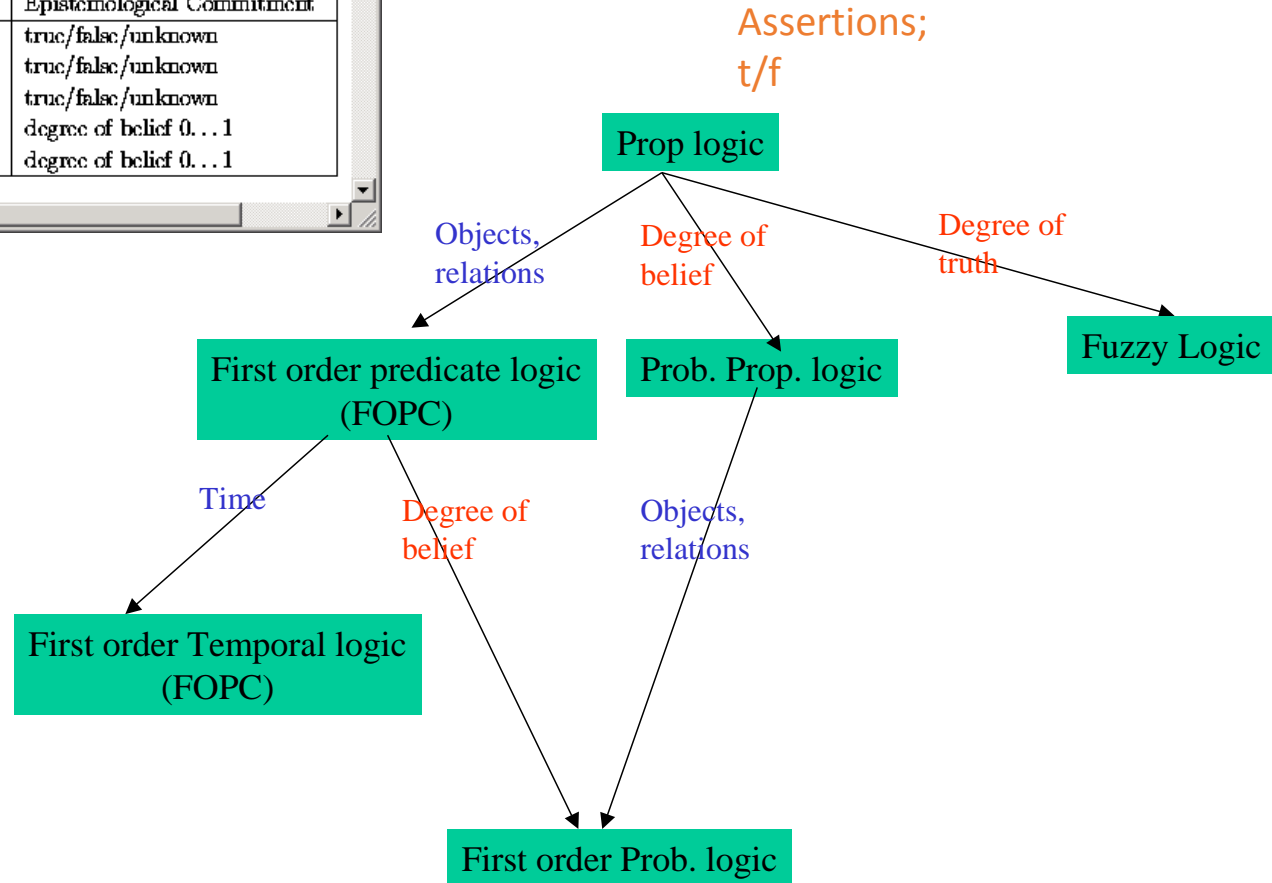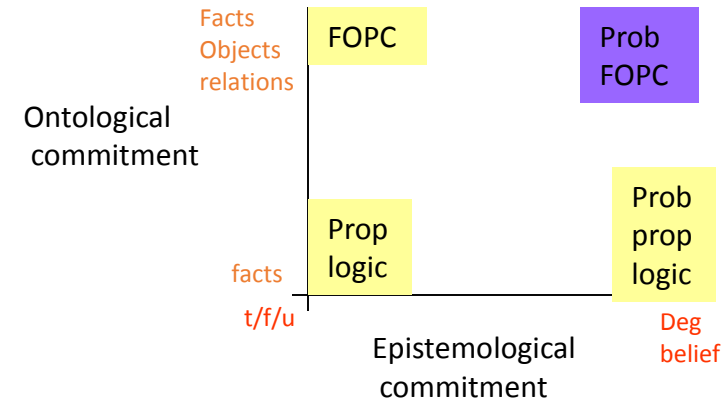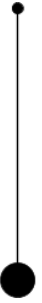
## Types of logic

Logics are characterized by what they commit to as "primitives"

Ontological commitment: what exists—facts? objects? time? beliefs?

Epistemological commitment: what states of knowledge?

| Language | Ontological Commitment | Epistemological Commitment |
|---|---|---|
| Propositional logic | facts | true/false/unknown |
| First-order logic | facts, objects, relations | true/false/unknown |
| Temporal logic | facts, objects, relations, times | true/false/unknown |
| Probability theory | facts | degree of belief 0...1 |
| Fuzzy logic | degree of truth | degree of belief 0...1 |

Facts
Objects
relations

FOPC

Prob
FOPC

Ontological
commitment

Prob
prop
logic

facts

Prop
logic

t/f/u

Epistemological
commitment

Deg
belief

Assertions;
t/f

Prop logic

Objects, relations

Degree of belief

Degree of truth

First order predicate logic (FOPC)

Prob. Prop. logic

Fuzzy Logic

Time

Degree of belief

Objects, relations

First order Temporal logic (FOPC)

First order Prob. logic

# Replace←→Augment

# AI's Curious Ambivalence to humans..

- Our systems seem happiest
    - either far away from humans
    - or in an adversarial stance with humans







*You want to help humanity, it is the people that you just can't stand…*

# What happened to Co-existence?

- Whither McCarthy's advice taker?

- ..or Janet Kolodner's house wife?

- …or even Dave's HAL?
    - (with hopefully a less sinister voice)

**HAAI**
**Human-aware AI**

Special Theme: Human Aware AI

# Planning: The Canonical View

Full Problem Specification



**PLANNER**

Fully Specified Action Model

Fully Specified Goals

Completely Known (Initial) World State

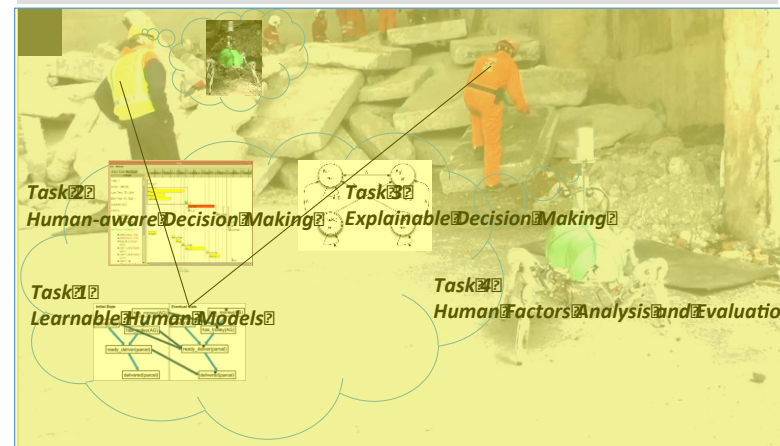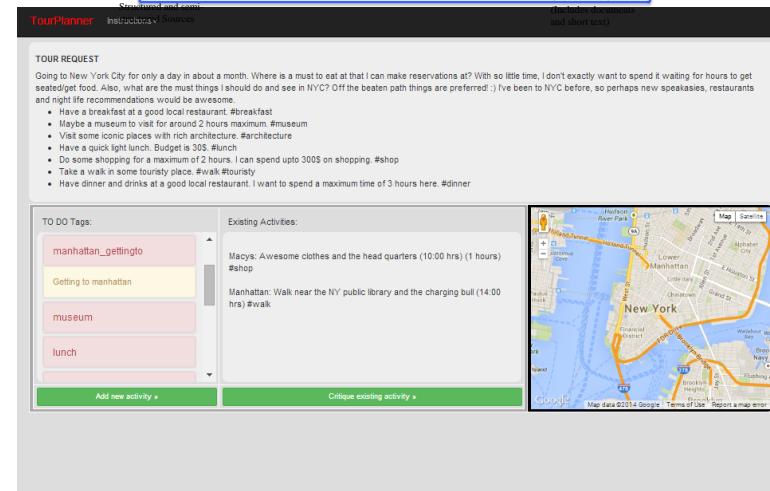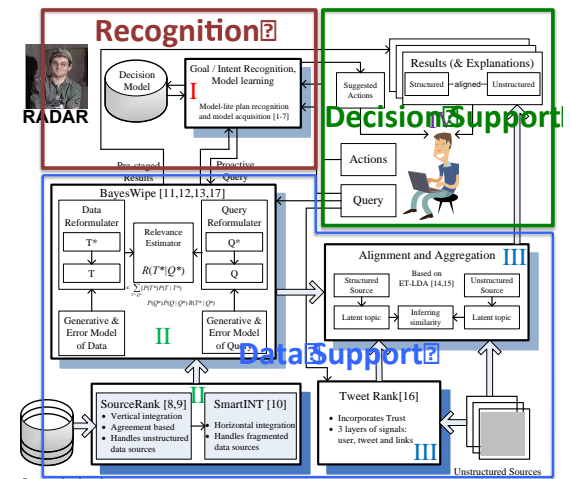Plan (Handed off for Execution)

Assumption:
- → Complete Action Descriptions
- → Fully Specified Preferences
- → All objects in the world known up front
- → One-shot planning

Allows planning to be a pure inference problem

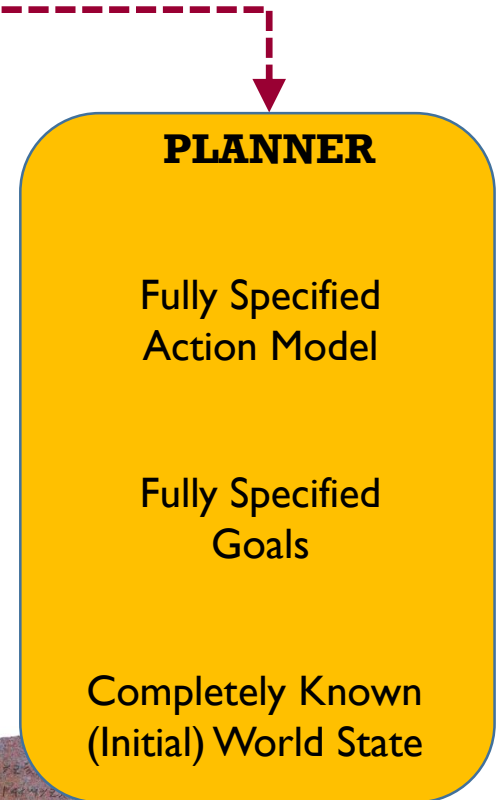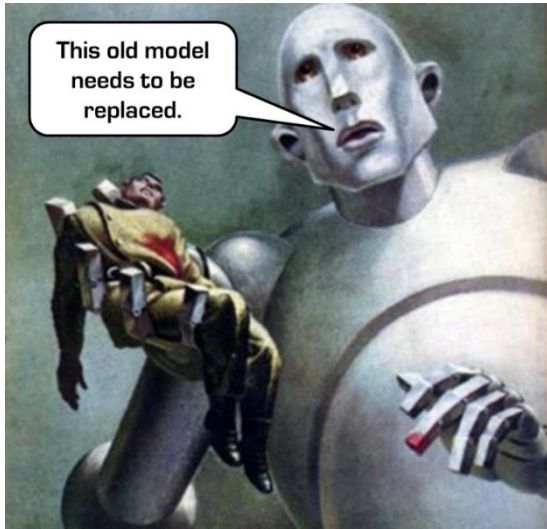☹ But humans in the loop can ruin a really a perfect day ☹

26

# Human-in-the-Loop Planning

- In many scenarios, humans are part of the planning loop, because the planner:

  - Needs to plan to avoid them

    - Human-Aware Planning

  - Needs to provide decision support to humans

    - Because "planning" in some scenarios is too important to be left to automated planners

    - "Mixed-initiative Planning"; "Human-Centered Planning"; "Crowd-Sourced Planning"

  - (May need) help _from_ humans

    - Mixed-initiative planning; "Symbiotic autonomy"

  - Needs to team with them

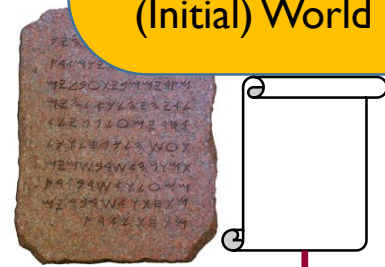    - Human-robot teaming; Collaborative planning

# Planning: The Classical View



This old model needs to be replaced.

Full Problem Specification

## PLANNER

Fully Specified Action Model

Fully Specified Goals

Completely Known (Initial) World State

Plan (Handed off for Execution)

Violated Assumptions:
→ Complete Action Descriptions (Split knowledge)
→ Fully Specified Preferences (uncertain users)
→ Packaged planning problem (Plan Recognition)
→ One-shot planning (continual revision)
**Planning is no longer a pure inference problem** ☹

☹ But humans in the loop can ruin a really a perfect day ☹

28

# Human-in-the-Loop Planning & Decision Support

## AAAI 2015 Tutorial

rakaposhi.eas.asu.edu/hilp-tutorial

**Subbarao Kambhampati**
Arizona State University
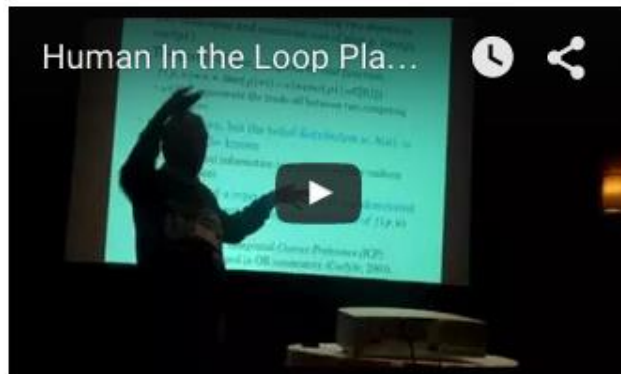
**Kartik Talamadupula**
IBM T.J. Watson Research Center

AAAI-15 Austin, Texas USA
The First *Winter* AI Conference!

**Materials**

[Tutorial Slides (Final version, as given) [PDF]](#)

# Challenges in Human-in-the-loop Planning

- Interpret what humans are doing based on incomplete human and domain models (Modeling)
  - Plan/goal/intent recognition
- Plan with incomplete domain models (Decision Making)
  - Robust planning/execution support with "lite" models
  - Proactive teaming support
- Explanations/Excuses  (Interaction/Communication)
  - How should the human and robot coordinate
- Understand effective interactions between humans and machines (Evaluation)
  - Human factor study

# Planning for Human-Robot Teaming

## Open World Goals

- › When to start sensing?
  - › Indicator to start sensing
- › What to look for?
  - › Object type
  - › Object properties
- › When to stop sensing?
  - › When does the planner know the world is closed?
- › Why should the robot sense?
  - › Does the object fulfill a goal?
  - › What is the reward? Is it a bonus?

[Talamadupula, Benton et al., ACM TIST 2010]

## Problem Updates
[TIST10]

**Assim** **Info**

## *Planning for*

### Replanning for Changing Worlds

- › New Information
  - › Sensors
  - › Human teammate
- › New Goals
  - › Orders: Humans
  - › Requests
- › Requirement
  - › New plan that works in new world (state)
  - › Achieves the changed goals

[Talamadupula et al. AAAI10]

## Problem Specification

When is a plan "Explainable" to the human in the loop?

- The robot generates its plan of action using its model $M_R$
- The human "interprets" this plan in light of her understanding of the Robot's model $M^*_R$
- $M_R$ and $M^*_R$ can be quite different..
- Differences can be a result of:
  - ◇ Different capabilities (e.g., possible actions)
  - ◇ Different knowledge (e.g., level of modeling)
  - ◇ Different interpretation of behaviors (e.g., plans) interacting with the world -- more than just trajectory planning!

$$\underset{\pi_{M_R}}{\arg\min}\; cost(\pi_{M_R}) + \alpha \cdot dist(\pi_{M_R}, \pi_{\mathcal{M}^*_R})$$

But, alas, $M^*_R$ is not known!

M

Re

## Model Updates
(via natural language)

- › "To go into a room when you are at a closed door, push it one meter."
  - › Precondition: "you are at a closed door"
  - › Action definition: "push it one meter"
  - › Effect: "go into a room"
- › NLP Module
  - i. Reference resolution
  - ii. Parsing
  - iii. Background knowledge
  - iv. Action submission (to planner)

[Cantrell, Talamadupula et al., HRI 2012]

[In collaboration with hrilab, Tufts University]

s

[IROS14]

# Interpretable AI… (Symbols/Neurons Redux)

- We humans may be made of neurons, but we seem to care a "lot" about comprehensibility and "explanations"

- If we want AI systems to work with us, they better handle this
  - This is an important challenge for the neural architectures
    - What do those middle layers represent?
      - Hinton says that (eventually?) we can just connect them to language generator networks and in effect "ask them"..

# Spock or Kirk?: Should AI have emotions?



- By dubbing "acting rational" as the definition of AI, we carefully separated the AI enterprise from "psychology", "cognitive science" etc.

- But pursuit of HAAI pushes us right back into these disciplines (and more)
  - Making an interface that improves interaction with humans requires understanding of human psychology..
    - E.g. studies showing how programs that have even a rudimentary understanding of human emotions fare much better in interactions with humans

# Do we really know what (sort of assistance) humans want?

# Proactive Help Can be Disconcerting!



The Sentence Finisher

*We dance round in a ring and suppose,*
*But the Secret sits in the middle and knows.*

# Human-human Teaming Analysis in Urban Search and Rescue

Simulated search task (Minecraft) with human playing role of USAR robot

- 20 internal/external dyads tested
- Conditions of autonomous/intelligent or remotely controlled robot
- Differences in SA, performance, and communications

# Disappointment ← → Doomsday



Thursday January 31, 2013

37

# Musk, Wozniak and Hawking urge ban on warfare AI and autonomous weapons

More than 1,000 experts and leading robotics researchers sign open letter warning of military artificial intelligence arms race

VB

NEWS EVEN

## Netflix's Hastings: Battle for machines and genetically m

CHRIS O'BRIEN    JANUARY 18, 2016 3:41 AM
TAGS: AI, GENETICS, NETFLIX, REED HASTINGS

Image Credit: Flickr/epSos .de

Before Reed Hastings cofounded a little company called Netflix, which is now changing the way we watch TV, he was an artificial intelligence engineer.

AI has come a long way since Hastings got his masters from Stanford University in 1988. But he still follows developments in the field closely. And during a conversation on stage today at the DLD Conference in Munich, Germany, Hastings said he was far less worried about looming threats of an AI-triggered apocalypse than are many other observers, such as Tesla's Elon Musk.

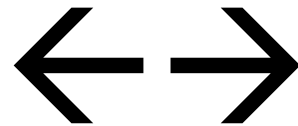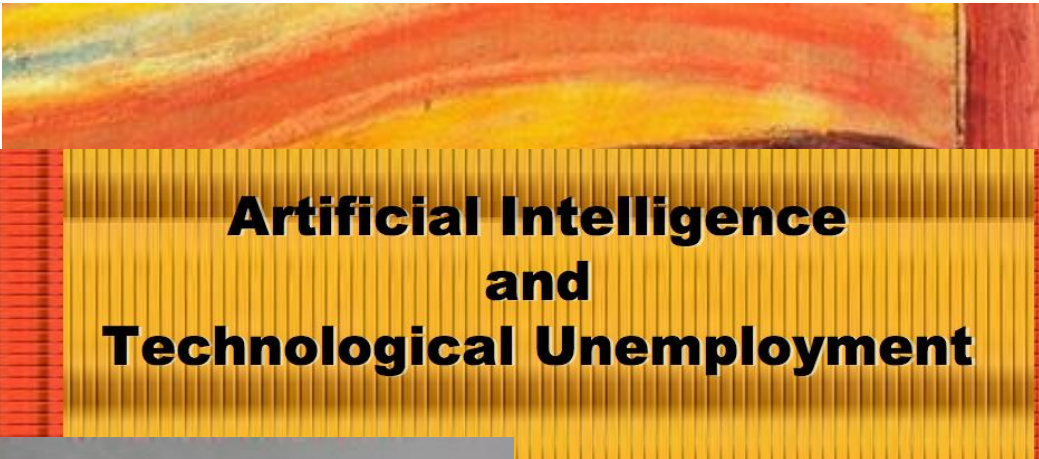"Some people worry about what happens when machine intelligence is too strong," Hastings said. "That's like worrying about our Mars colony and people

## Artificial Intelligence and Technological Unemployment

America is the only country that went from barbarism to decadence without civilization in between

~ Oscar Wilde ~

*AI is the only technology that is going from disappointment to deadly without touching beneficial.. (?)*

www.StatusMind.com

### Press Releases

As Tax Season is Set to Take off, Taxhub™ a NYC Startup, Offers Disruptive New Take on the Personal Income Tax Filing Industry

Zero Gravity Solutions, Inc. Signs Space Act Agreement with NASA Ames Research Center

16

**the key to** **nizing other** **planets.** But the renowned physicist, whose recent lecture will be broadcast next week, does not think that will happen soon.

BBC News ↗

# Why we don't need to over-worry…

# Captain America to the Rescue?

# Why we don't need to over-worry…

- We already have autonomous systems; making them intelligent can't be bad!
- We get to design AI—we don't need to imbue them with the same survival instincts
- The way to handle possible problems with AI are to allow multiple AIs
- Technological unemployment is a big concern..
  - ..but even here, the opinion is divided
    - "biased advantages" vs. "rising tide lifts all boats"

# Suppose Evil AI is right around the corner.. How do we stop it?

- What wont work
  - Renunciation
  - Tight regulation
  - Fierce internal programming
- What works?
  - More AI!



Forbes / Tech

The One Thing We Need To Stop Robots From Achieving World Domination

What constraints to AI and machine learning algorithms are needed to prevent AI from becoming a dystopian threat to humanity? originally appeared on Quora: The best answer to any question.

Answer by David Brin, author of The Postman, Earth and The Transparent Society, on Quora:

It is, of course, wise and beneficial to peer ahead for potential dangers and problems — one of the central tasks of high-end science fiction. Alas, detecting that a danger lurks is easier than prescribing solutions that can prevent it. Take the plausibility of malignant AI, remarked-upon recently by luminaries ranging from Stephen Hawking to Elon Musk. Indeed, my own novels contain some chilling warnings about failure modes with our new, cybernetic children.

It is called Competition.

If you fear a super smart, Skynet level AI getting too clever for us and running out of control, then give it rivals who are just as smart but who have a vested interest in preventing any one AI entity from becoming a would-be God.

It is how the American Founders used constitutional checks and balances to prevent runaway power grabs by our own leaders, for the first time in the history of varied human civilizations. It is how companies prevent market warping monopoly, that is when markets are truly kept flat-open-fair.

Alas, this is a possibility almost never portrayed in Hollywood sci fi – except on the brilliant show Person of Interest – wherein equally brilliant computers stymie each other and this competition winds up saving humanity.

# AI & Unemployment:

If machines can do
Everything that people
can, then what will
people do?



MOTIVATION

If a Pretty Poster and a Cute Saying are All it Takes to Motivate You,
You Probably have a Very Easy Job. The Kind Robots Will be Doing Soon.

www.despair.com

# AI & Unemployment

- Taxi Drivers
- Factory workers
- Journalists
- Doctors (??)
- Cocktail Waiter

Technology

**Intelligent Machines: The jobs robots will steal first**

By Jane Wakefield
Technology reporter

14 September 2015 | Technology

THE SECOND MACHINE AGE

WORK, PROGRESS, AND PROSPERITY
IN A TIME OF
BRILLIANT TECHNOLOGIES

ERIK BRYNJOLFSSON
ANDREW McAFEE

# The many good things AI can bring to the society

- Assistive technologies
  - Elder care; care for the disabled;
  - cognitive orthotics
    - Personal Digital Assistants
      - ("Not Eric Schmidt")
- Accident free driving..
- Increased support for diversity
  - Language translation technologies (real life Babel Fish!)
- … <many many others>



BABEL FISH

DIGESTIVE NERVE CHORD   ENERGY ABSORPTION FILTER
TELEPATHIC EXCRETOR   BRAIN   GAS BLADDER
OLFACTORY BULB

EXTENDABLE
NERVE SIGNAL
SENSOR   LIVER   DIGESTION   CONSCIOUS   UNCONSCIOUS FREQUENCY SENSORS
GILL RAKERS   HEART   FREQUENCY SENSORS

THE BABEL FISH IS SMALL, YELLOW, LEECHLIKE,
AND PROBABLY THE ODDEST THING IN THE UNIVERSE.
IT FEEDS ON BRAIN WAVE ENERGY, ABSORBING ALL

# Summary

- What is Intelligence

- Progress of AI

- The pendulum swings in AI
  - Symbols – Neurons
  - Logic  -- Probability
  - Replace – Augment
    - Spock –Kirk
  - Disappointment – Doomsday

# The Fundamental Questions Facing Our Age

- Origin of the Universe
- Origin of Life
- Nature of Intelligence



"To know your future you must know your past"

— George Santayana

Predictions are hard,
especially about the future

--Niels Bohr

rakaposhi.eas.asu.edu/cse471/

Apps  Tasks  AAAI-16 Tutorial For...  British Airways - Bo...  Dec  session  Nishant Jain | Linked...  Vivid dreaming can ...  SD Card for surface ...  slow  IJCAI-2016-PC-Mem...  e  avi  midtown  »

# CSE 471/598 Lecture Notes (Spring 2012)

## Current Offering: Fall 2015; Friday 1:30PM--4PM. LSE 106

## ****Here is the Class Schedule with Videos for Fall 2015 Offering****

## Piazza site for Fall 2015 Offering

**Google "rao ai"**

**Additional pointers:**

- Check out what students say about the last offering of this course in **CEAS student evaluations**.
- Check out what students "learned" from the last offering of this course **(in their own words)**.

- Intro; Intelligent agent design [R&N Ch 1, Ch 2] (sound).

  - L1 [Jan 5, 2012] Video (~4gig) and Audio of the lecture. Administrivia; the space odyssey *son et lumeire*; Using it as a vehicle to do an interactive overview of AI developments; Definitions of AI--and why thinking humanly, thinking rationally and acting humanly do not quite provide the right fit as general definitions for AI enterprise.

  - L2 [Jan 10, 2012] Video (~4gig) and Audio of the lecture. Rational agency; performance metrics; percept/action/goal/environment types; agent designs and how they motivate the course topics.
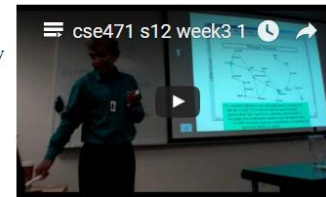
- Atomic (Problem Solving) Agents [R&N Ch 3 3.1--3.5]

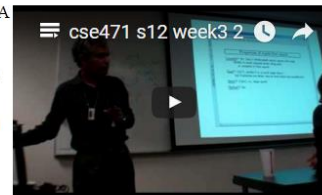  - L3 [Jan 12, 2012] Video (~4gig) and Audio of the lecture. Atomic,

    propositional and relational representations--their tradeoffs (and how functions can allow representation of infinite state spaces compactly.). Atomic agent design. How search is at the heart--and child-generator, goal-test functions are needed. Blind vs. Informed search. Before continuing to single state search algorithms, thinking of the effect of enviornment accessibility. Multiple-initial state search or belief-state search. Without sensing (conformant planning) and with sensing (contingent planning; the medicate problem).

  - L4 [Jan 17, 2012] Video (~4gig) and Audio of the lecture. Exploration/exploitation tradeoffs. Setting state-spaces for atomic agents. Search for single-state atomic agents. World state vs. Search node difference. Viewing search algorithms in terms of their queuing functions. Dept-First/Breadth-first searches and how they differ in optimiality vs. space consumption.

  - L5 [Jan 19, 2012] Video (~4gig) and Audio of the lecture. A discussion of the uniform search tree model. A slow and comprehensive discussion of blind search strategies BFS, DFS, Depth limited DFS and IDDFS and their tradeoffs. A discussion on graph vs. tree search. A discussion on handling duplicate expansions with closed list vs. ancestor checking.

- Informed Search

  - L6 [Jan 20, 2012; for the class of 1/24] Video (~2.5 gig),

**Sidebar navigation:**

Course Overview

Homeworks

Lecture Notes

Metaphors

Projects

Mail Archive

Blog

AI@ASU

Lisp-in-a-box (Free Lisp for PC**)
Free book on lisp
Lisp vs. Scheme

S09 Notes

F07 Notes

F06 Notes
F06 Blog
F06 Mail Archive

F03 Notes

F03 Acquired Wisdom

F'00 Schedule

Related Courses

Rao Kambhampati

063635