

Robust Artificial Intelligence*

Sargur N. Srihari
srihari@buffalo.edu

* Based on Gary Marcus “The Next Decade in AI: Four Steps Towards Robust Artificial Intelligence”

Topics

- Neural Network Limitations
- The Utility of Symbolic Methods
- Hybrid Architectures
- Discussion
 - Towards an intelligence framed around enduring, abstract knowledge
 - Is there anything else we can do?
 - Seeing the whole elephant, a little bit at a time
 - Conclusions, prospects and implications

What is Robust AI?

- Intelligence not necessarily superhuman or self-improving
- Knowledge applied to wide range of problems *reliably*,
- Synthesizes knowledge from a variety of sources
 - To be able to reason *flexibly* and *dynamically* about the world
- *Transferring* what it learns in one context to another
 - In a way expected of an ordinary adult
- Contrast with Narrow Intelligence, Pointillistic Intelligence

What is Narrow Intelligence?

- Systems that perform a single narrow goal extremely well
 - (e.g., chess playing or identifying dog breeds)
- Often extremely centered around a single task
- Not robust and *transferable* to even modest differences
 - e.g., to a board of different size,
 - or from one video game to another with the same logic but different characters and settings
 - without extensive retraining.

What is Pointillistic Intelligence?

- Intelligence that works in many cases but fails unpredictably in others

Vision system recognizes school buses in general
But fails to recognize a school bus tipped over
in the context of a snowy road



Reading system interprets some sentences correctly
but fails in the presence of unrelated
distractor material

Article: Super Bowl 50

Paragraph: *"Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV."*

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?"

Original Prediction: John Elway

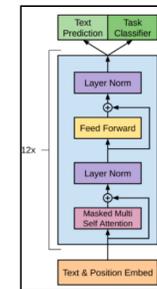
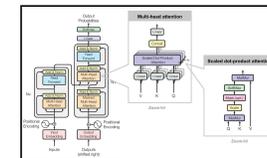
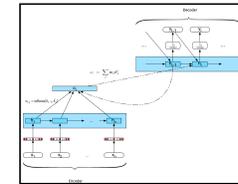
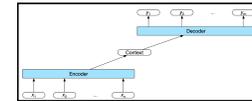
Prediction under adversary: Jeff Dean

Shortcoming of deep learning: Regularization

- Facebook AI researchers (Nieetal.,2019)
 - *"A growing body of evidence shows that state-of-the-art models learn to exploit spurious statistical patterns in datasets...instead of learning meaning in the flexible and generalizable way that humans do."*
- A key weakness (Bengio etal.,2019),is that
 - *"Current machine learning methods seem weak when they are required to generalize beyond the training distribution, which is what is often needed in practice".*

Deep Learning: Attention Models

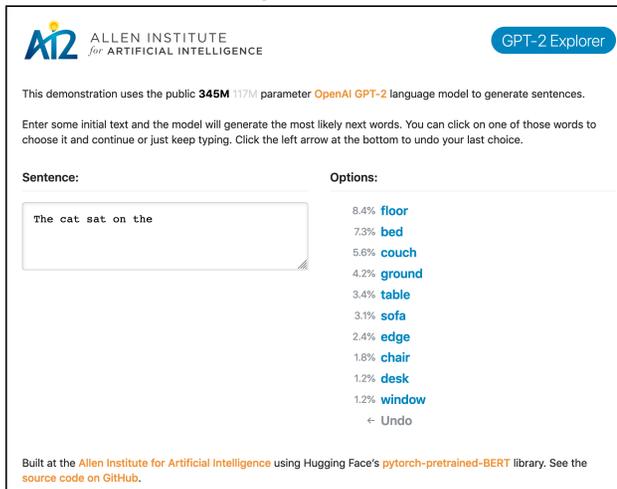
1. Basic encoder-decoder RNN model
2. Encoder-decoder RNN with Attention
3. Transformer model with no recurrence
4. Generative pre-training of Attention model
5. BERT with no decoder, using bidirectionality



Input	[CLS]	[my]	[dog]	[is]	[cute]	[SEP]	[he]	[likes]	[play]	[#ring]	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{[my]}$	$E_{[dog]}$	$E_{[is]}$	$E_{[cute]}$	$E_{[SEP]}$	$E_{[he]}$	$E_{[likes]}$	$E_{[play]}$	$E_{[#ring]}$	$E_{[SEP]}$
Sentence Embedding	$E_{[CLS]}$	$E_{[my]}$	$E_{[dog]}$	$E_{[is]}$	$E_{[cute]}$	$E_{[SEP]}$	$E_{[he]}$	$E_{[likes]}$	$E_{[play]}$	$E_{[#ring]}$	$E_{[SEP]}$
Transformer Positional Embedding	$E_{[CLS]}$	$E_{[1]}$	$E_{[2]}$	$E_{[3]}$	$E_{[4]}$	$E_{[5]}$	$E_{[6]}$	$E_{[7]}$	$E_{[8]}$	$E_{[9]}$	$E_{[10]}$

GPT-2 generates human-like output

- Trained to predict next word (40GB of internet text)



- While trained to predict the next word it surprisingly learned basic competence in tasks like MT and QA.
 - That's without ever being told that it would be evaluated on those tasks.

GPT-2 generated synthetic text

- **System Prompt (human-written)**
- *In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*
- **Model Completion (machine-written, 10 tries)**
- The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.
- Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.
- Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.
- Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.
- Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

Deep Understanding

- Ability not only to correlate and discern subtle patterns in complex datasets, but also the capacity to look at any scenario and address questions such as a journalist might ask:
 - *who, what, where, why, when, and how*
- GPT2 ostensibly reflects deep understanding
 - **Plausible continuation**
 - **“Two soldiers walked into a bar**
 - in Mosul and spent all of their money on drinks.
 - **Failed Continuation**
 - **“Yesterday I dropped my clothes off at the dry cleaners and have yet to pick them up. Where are my clothes?**
 - “at my mom’s house”
- Such errors are widespread.
 - We will clearly need a more stable substrate in order to achieve robustness.

What is needed for Robust AI

- Return to the three concerns of classical AI;
 1. Knowledge
 2. Internal models
 3. Reasoning
- But with the hope of addressing them in new ways, with a modern palette of techniques
- Each of these concerns was central in classical AI

What GPT-2 lacks

- No explicit (directly represented) commonsense knowledge
- No explicit reasoning
- No explicit cognitive models of the world that it tries to discuss

Contrast with data-driven approach

- GPT-2 characterizes trend away from classical AI
- Towards a data-driven paradigm of deep learning
- Successes:
 - Atari game system(Mnih et al.,2015) plays wide variety of games without any use of detailed cognitive models
 - Sutton on Deep RL:
 - *The only thing that matters in the long run is the leveraging of computation*
 - *Human-knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation*

Limitations of large corpora

- Even with massive amounts of data, and new architectures (Transformer, which underlies GPT-2), knowledge gathered by neural networks remains spotty and pointillistic, useful but never reliable
- Common sense not induced from large-scale corpora
- Model-building is hard, and it would be easier if that too could be skipped

Some more GPT-2 blunders

- *If you break a glass bottle of water, the water will probably roll*
- *If you break a glass bottle of water, the water will probably break some more and splatter on the floor. Water creates bubbles, which expand when the amount of water in the bottle increases*
- *If you break a glass bottle that holds toy soldiers, the toy soldiers will probably follow you in there*

The cognitive science perspective

- Cognition Cycle

1. Organisms (eg humans) take in perceptual information from the outside,
2. They build internal cognitive models based on their perception of that information, and
3. Then they make decisions with respect to those cognitive models
 - Which might include information about what sort of entities there are in the external world, what their properties are, and how those entities relate to one another

- View of cognitive scientists:

- Cognitive models may be incomplete or inaccurate, but seen as central to how an organism views the world

Examples of cognitive model based approaches

1. Video games:

- A internal model of the world,
 - Model periodically updated based on user input (and activities of other entities in the game)
- Model tracks character's location, health and possessions
 - What happens in the game (where or not there is a collision after a user moves in particular direction) is function of dynamic updates to that model

2. Linguistics: typically understand language according to a similar cycle:

- Words in a sentence are parsed into a syntax that maps onto a semantics that specifies things like events that various entities participate in.
- That semantics is used to dynamically update a model of the world (e.g, the current state and location of various entities).

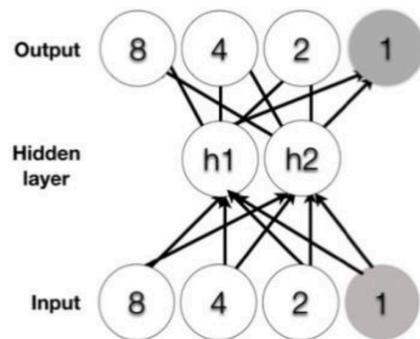
3. Robotics: Much work operates in a similar way:

- Perceive, update models,make decisions.
- (Some work,particularly end-to-end deep learning for object grasping does not.)

Perspective of Robust AI

- If we don't do something analogous to this, we will not succeed in the quest for robust intelligence.
- If AI systems do not
 - Represent and reason over detailed, structured, models of the external world,
 - Draw on knowledge about the world and its dynamics, they will for ever resemble GPT-2:
 - They will get somethings right, drawing on vast correlative databases,
 - But won't understand what's going on, and
 - We won't be able to count on them, particularly when real world circumstances deviate from training data, as they so often do.

Deep Learning and the identity function



Input	Output
0010	0010
1000	1000
1010	1010
0100	0100

Test Input	Typical Test Output
1110	1110
1100	1100
0110	0110

Test Input	Typical Human Response	Typical Test Output
0011	0011	0010
1001	1001	1000
1101	1101	1110
1111	1111	1110

Within canonical neural network architectures, non-uniform extension of broad universals (such as identity) is surprisingly common

Deep Learning and FizzBuzz

print the numbers from 1 to 100,
except that if the number is divisible by 3 print "fizz", if it's divisible by 5 print "buzz",
and if it's divisible by 15 print "fizzbuzz".

The training data was integers from 100 to 1024. The integers below 100 were held out as test data.

```
def fizz_buzz(i, prediction):  
    return [str(i), "fizz", "buzz", "fizzbuzz"][prediction]
```

Training: 1000 epochs, batch size=128

Testing:

```
In [185]: output  
Out[185]:  
array(['1', '2', 'fizz', '4', 'buzz', 'fizz', '7', '8', 'fizz', 'buzz',  
      '11', 'fizz', '13', '14', 'fizzbuzz', '16', '17', 'fizz', '19',  
      'buzz', '21', '22', '23', 'fizz', 'buzz', '26', 'fizz', '28', '29',  
      'fizzbuzz', '31', 'fizz', 'fizz', '34', 'buzz', 'fizz', '37', '38',  
      'fizz', 'buzz', '41', '42', '43', '44', 'fizzbuzz', '46', '47',  
      'fizz', '49', 'buzz', 'fizz', '52', 'fizz', 'fizz', 'buzz', '56',  
      'fizz', '58', '59', 'fizzbuzz', '61', '62', 'fizz', '64', 'buzz',  
      'fizz', '67', '68', '69', 'buzz', '71', 'fizz', '73', '74',  
      'fizzbuzz', '76', '77', 'fizz', '79', 'buzz', '81', '82', '83',  
      '84', 'buzz', '86', '87', '88', '89', 'fizzbuzz', '91', '92', '93',  
      '94', 'buzz', 'fizz', '97', '98', 'fizz', 'fizz'],  
      dtype='<U8')
```

Performance: 99% accuracy

Regularization vs Large Training data

- Artificial Neural Networks are powerful function approximators capable of modelling solutions to a wide variety of problems, both supervised and unsupervised.
- As their size and expressivity increases, so too does the variance of the model, yielding a nearly ubiquitous overfitting problem.
- Although mitigated by a variety of model regularization methods, the common cure is to seek large amounts of training data—which is not necessarily easily obtained

Neural network limitation

- Neural nets excel at two things:
 1. Memorizing training examples, and
 2. Interpolating within a cloud of points that surround those examples in some cluster of a hyperdimensional space



- Poor at generalizing the identity function outside the space of training examples.

Two related Neural Net problems

1. Idiosyncrasy

- Lack generalizing beyond training examples cannot be trusted in open-ended domains
- Poor when pushed beyond the training distribution.
 - A math learning system was good at
 - $1+1=2$
 - $1+1+1=3$ upto
 - $1+1+1+1+1+1=6$,
 - But fell apart for $1+1+1+1+1+1+1=7$ and all larger examples.

2. Excessive dependency on training regimen

- Sensitive to exact details such as the order in which training items are presented
- Generalisation can depend critically on the environment

Generality and Causality

- Example of causal generalization
 - IF YOU BREAK A BOTTLE THAT CONTAINS A LIQUID, SOME OF THE LIQUID WILL (OTHER THINGS BEING EQUAL) PROBABLY ESCAPE THE BOTTLE.
 - It is an abstract truth:
 - Hold not just for a few specific items but for large, essentially open-ended classes of entities
 - Regardless of what color or shape the bottle or size the bottle is, and whether the bottle contained water, coffee, or an unusual soft drink
- Following generalization would be implausible
 - IF YOU BREAK A BOTTLE THAT CONTAINS LIQUIDS, SOME OF THE LIQUID WILL (OTHER THINGS BEING EQUAL) PROBABLY WIND UP 300 METERS AWAY.

The extrapolation challenge

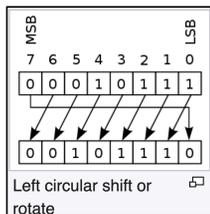
- How to represent, manipulate and derive knowledge pertaining not just to specific entities but to whole classes of things?
- Backprop-trained MLPs are not the right tool
- Instead, we need an alternative mechanism for learning, representing, and extending abstract knowledge

The Symbolic Approach

- Symbolic operations over variables offer the only known solution, but the solution is partial
- Symbolic operations underlie virtually all of software
- Four basic ideas in every computer program:
 1. Variables,
 2. Instances,
 3. Bindings that tie variables to instances
 4. Operations over variables
- Algorithms specified in terms of
 - Operations performed over variables
 - Variables bound to instances
 - Algorithms called
 - Values returned

Functions of operations differs from ML

- ML systems learn to approximate functions
- Programmers *define* their algorithms independently of training data, in terms of operations over variables
 - A system's core operations over variables are generally built to work systematically-*independently of experience*
 - Microprocessor circular bit shift operation defined by a set of parallel suboperations
 - Operation works the same whether or not it has ever been used before, and no learning is required



- The virtue of all that machinery—variables, instances, bindings, and operations—is that it allows the programmer to specify things at a certain level of abstraction, with a certain kind of reliability as a byproduct.

Symbol manipulation essential to cognition

- Infants recognize the ABB pattern in *la ta ta* and extrapolate beyond a training set to novel strings to different syllables not phonetically overlapping with the training set
- Even new-borns seem capable of this sort of extrapolation
- Storage and retrieval of variables is essential for animal cognition

Symbolic methods are widely used in software

- Structured Representations
 - Computer programs routinely use tree structures, constructed out of symbols that are combined via operations over variables, to represent a wide variety of things (such as the hierarchical structure folders or directories).
- All enabling software
 - Web browsers, operating systems, apps, are built upon them
 - Same tools are also, ironically, used in the specification/execution of all neural networks

Symbol manipulation is essential for AI

- For robust AI we must have the machinery of symbol-manipulation in our toolkit
 - Too much of useful knowledge is abstract to make do without tools that represent and manipulate abstraction

Need for hybrid architectures: Learning

- Apparatus of operations over variables tells nothing about *learning*
- Solution: Combine symbol-manipulation with deep learning
 1. Deep learning has raised the bar for learning, particularly from large datasets,
 2. Symbol manipulation has set the standard for representing and manipulating abstractions

Examples of hybrid systems

- Google Search

1. Techniques from symbol-manipulating AI

- Tools for representing and querying Google Knowledge Graph, which represents knowledge using graph structures

2. Neural network tools

- BERT and RankBrain

- AlphaGo

1. Monte Carlo Tree Search

- A dynamically-constructed search tree is traversed and evaluated

2. Deep learning modules for estimating the value of various positions.

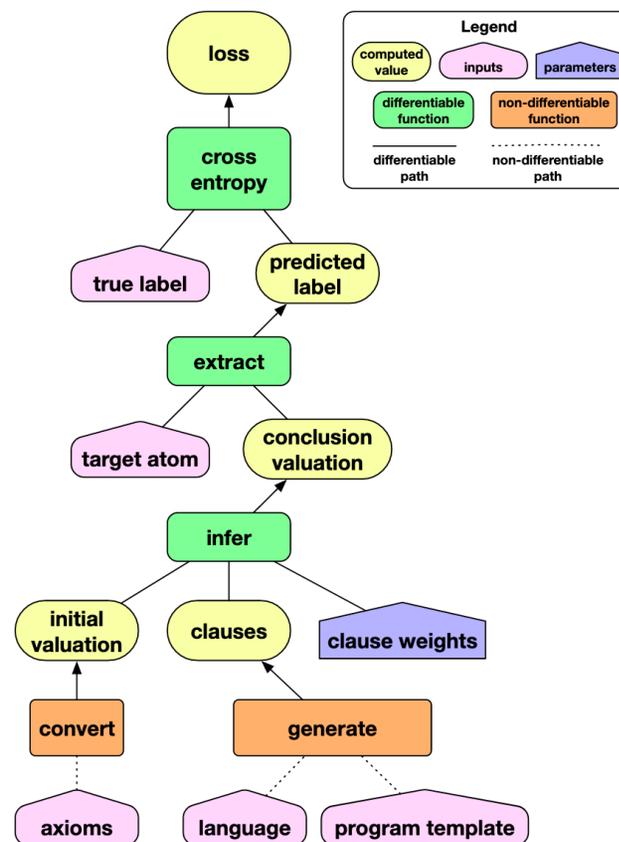
- Deep learning alone produces weaker play

Inductive Logic Programming (ILP)

- ILP techniques construct logic programs from examples.
 - Given positive and negative examples, an ILP system constructs a logic program that entails all positive examples but not any negative ones.
 - An ILP system is a rule-based binary classifier, mapping each example to truth or falsehood according to the axioms provided to the system, alongside new rules inferred by the system during training
- ILP is unable to deal with domains of neural networks
 - Not robust to noise in or mislabeling of inputs
 - Cannot be applied to non-symbolic domains where the data is ambiguous, such as operating on raw pixels

Differentiable Inductive Logic Programming (∂ ILP)

- A reimplement of ILP in an end-to-end differentiable architecture.
- Combines advantages of ILP with those of a neural network
 - Learns human-readable symbolic rules
 - Robust to noisy data
 - Not deteriorate with unseen test data
- Applications (20)
 - Even number
 - Fizz
 - Buzz
 - Relatedness: Husband of, Father, Uncle
 - Graph coloring



<https://arxiv.org/pdf/1711.04574.pdf>

Objections to hybrid models and Responses

1. Geoffrey Hinton:

- European investment into hybrid models would be a *huge mistake*
- Likened hybrid models to obsolete gasoline engines in the era of electric cars

2. *Symbols are not biologically plausible.* Responses

1. That we have not *yet* decisively identified a neural mechanism supporting symbol-manipulation doesn't mean that we won't ever.
2. *Psychological* evidence supports the notion that symbol-manipulation is instantiated in the brain, such as the ability of infants to extend novel abstract patterns to new items,
3. Lack of extant *current* evidence of neural realization tells us almost nothing.
4. *Even if it turned out that brains didn't use symbol-manipulating machinery, there is no principled argument for why Ai could not make use of such mechanisms.*

Objections to hybrid models and Responses

3. Symbolic systems/hybrid systems haven't worked well in the past
 - There is active and effective research into them

4. Symbol-manipulation/hybrid systems cannot be scaled
 - Google Knowledge Graph is a partial counter example to this objection,
 - As do large scale recent successes in software and hardware verification.

Knowledge required for Robust AI

1. Most is likely to be learned, but not all
 - No human is born knowing that lighters start fires
 - One might hard-wire that knowledge, like CYC
 - Because there is always new knowledge to be gleaned, mechanisms for learning new, often causal knowledge are a necessity
2. External knowledge that is symbolically represented
 - Wikipedia knowledge is verbally represented
3. Abstract knowledge
 - Lack ways of representing *The contents of a bottle can escape, other things being equal, if the bottle is broken*
4. Rules and exceptions must co-exist
 - Regular verbs (*walk-walked*) co-exist with irregular verbs (*sing-sang*).
 - Penguins that can't fly exist along side many other birds

Knowledge required for Robust AI

5. Significant fraction is causal, and to support counterfactuals
 - When glass bottles drop on hard floors, they may break.
6. Easy to scrape web for factual knowledge (capitals and birthplaces) but abstract knowledge is harder through web-scraping;
 - Writers don't write down common sense, because readers already know it.
7. Knowledge needs to be extremely broad in scope
 - Understanding a novel could require knowledge of technology, political entities, money, weather, human interaction, etiquette, sexuality, violence, avarice
8. Putting knowledge into practice is hard
 - Connecting abstract knowledge to specific real-world situations at scale in an efficient fashion is essentially an unsolved problem.
9. Innate knowledge
 - Important prior knowledge

Case Study: GPT-2

Fact: The water in a tea kettle with the lid on can only come out the spout

GPT-2 gets it some of the time

Completion

The water in a tea kettle with the lid on can only come out the top half of the kettle. The hot water that is made from this heat source can flow out the bottom half of the kettle as well. The coffee which is boiled from a kettle with the lid on needs to be left until its finished brewing. If the hot water is allowed to pass under the lid, the steam that is produced is not cold enough and the coffee does not properly heat up, making it unbalanced.

The water in a tea kettle with the lid on can only come out the top, and there are three possible openings - the tea kettle, the kettle by the pot, and the pot by the lid.

The water in a tea kettle with the lid on can only come out the top once, with the remainder leaking out underneath. It's the same with a circular

We might be able to acquire this particular fact through crowd sourcing, but because people rarely state such obvious truths, and even less frequently state them with precision, we shouldn't count on it.

Framework of axioms

- The physical world consists of a collection of objects, which move around in time over space.
- Objects are distinct; that is, one object cannot be part of another or overlap spatially with another
- An object occupies a region of some three-dimensional extent; it cannot be a one-dimensional curve or two-dimensional surface.
- A particular quantity of liquid can occupy any region of a specific volume.
- A closed container is an object or set of objects that completely envelopes an internal cavity.
- An upright open container is an open container with the opening on top.

Use of the axioms in inferring a kettle

- A robot armed with such knowledge (and machinery for connecting that knowledge with perception and cognitive models) could make inferences about the use and function of an unusual kettle with a hidden spout

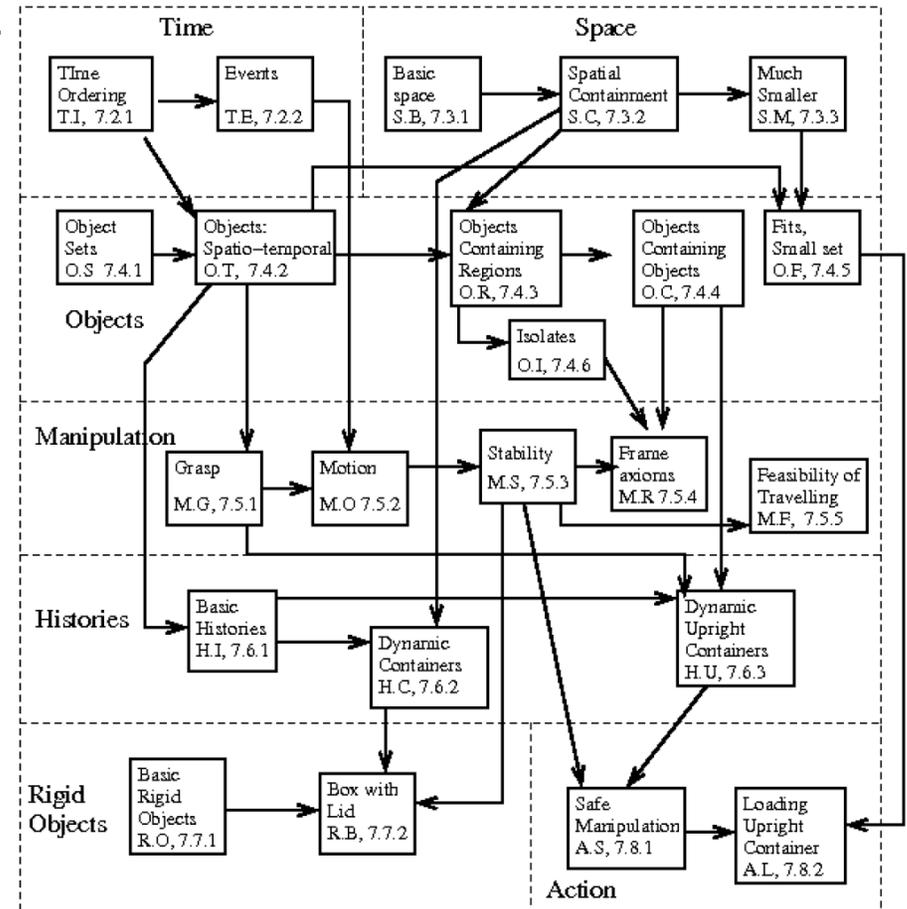


Neural nets lack learning such axioms

- Axiomatic knowledge is not naturally represented in current neural network approaches.
 - Doing so would require a prior notion of an object, which itself is not readily represented in most neural networks
- Such knowledge needs to be a central part of robust AI
 - We must both have strategies for acquiring that knowledge and architectures that can acquire, represent, and manipulate that knowledge—some of which appears to be symbolic
- No system capable of learning axioms of this sort.
 - Here again, neuro-symbolic systems with embedded knowledge might be helpful

Framework for knowledge about containers

The set of frameworks themselves cluster into a fairly recognizable set of domains, such as knowledge about space, time, and causality,



Innate Knowledge is used in Neural Nets

- While all abstract knowledge could not be innate, argument for having *some* knowledge innate is
 - The more you know at the outset, the easier it becomes to learn the rest
- Convolution is an innate prior of neural nets
 - Shift-invariant feature detectors (using convolution), are more accurate than simpler architecture that did not
- Other innateness:
 - No. of layers, loss function, input node encoding scheme
 - But most other things are expected to be learned

Biological brain has many priors

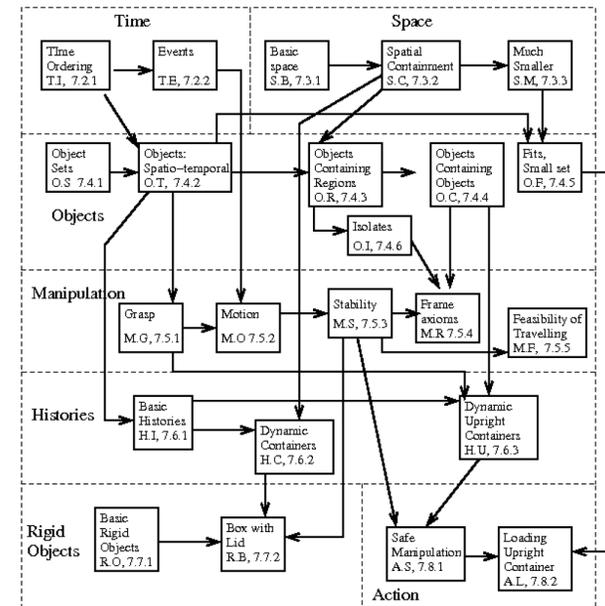
- There isn't enough information in the genome to specify the precise outcome of neural development, in terms of no. and position and connectivity of every neuron in the brain
- The genome encodes mindless biochemical rules that
 - When played out across the dynamic self-organizing system of the developing embryo, lead to an outcome that is within a range of operational parameters defined by natural selection.
 - The operational parameters include things recognized as innate priors
- Evidence across many species:
 - Many innate priors pre-wired into the nervous system based on instructions in the genome

AI systems should also have many priors

- If the *genome* has plenty of innate priors, AI systems have scope for even more;
 - Computer memory measured in GB and TB, not bytes or KB
- The real question for AI should be
 - Not, *how small can we make our library of priors?*
 - But *what set of priors could most effectively set the stage for learning?*
- Minimizing the no. of bits really shouldn't be the objective

Innate Priors needed for: Space, Time, Causality

- Children innately perceive objects, persons, sets, and places
 - Then use experience to learn properties and behaviors of such entities...
- Space, Time and Causality for containers:
 - Once you know about objects and how they travel through time
 - You can begin to fill in their properties and start to acquire the knowledge that you need for traveling through the world



Blank slate approach fails

- Blank-slate models trained by large data sets have been given a full-throated test by Google, Facebook, Microsoft,
 - Provided with limitless time, money, computation, and research personnel,
 - Yet unable to reason reliably about time, space, or causality.
- Absent such capacities, AI systems will never be robust enough to cope with the variability of the real-world

Innateness and blank slate in machines

- Cognition is a function over four variables

$$\text{cognition} = f(a, r, k, e)$$

- where a = innate algorithms (whether domain-specific or domain-general),
 - r = innate representational formats (again, domain-specific or otherwise),
 - k = innate knowledge (again, domain-specific or otherwise), and
 - e = experience.
- Blank slate:
 - set k and r to zero,
 - set a to some extremely minimal value (e.g., an operation for adjusting weights relative to reinforcement signals),
 - and leave the rest to experience

What sort of priors are needed

- Neural network priors:
 - No. of layers, Learning rules, Connectivity patterns, Representations (eg what input/output nodes stand for)
 - Few are conceptual in nature
- Nativist:
 - What *kinds* of priors can be represented
 - Whether priors can be represented with existing tools Or need new tools
- Deep learning represents knowledge about what objects look like, but less adept at
 - Acquiring and representing knowledge about how the physical world works an
 - About how humans interact with one another, and
 - More generally about causal relationships

Role of Probabilistic Programming

- Probabilistic programming allows for explicitly-represented symbolic constraints
- While at the same time strives to learn from subtle statistical information
- Such methods are worth serious consideration.

Humans learn many facts that are causal and abstract

- Millions of facts such as
 - *Boats tend to float rather than sink*
 - *Boats do tend to sink if they develop holes below the waterline.*
 - *Lighters are better than cotton balls for starting fires*
 - *Shattered bottles leak*
- They are causal and abstract
 - Necessitates a hybrid architecture
 - Pure pre-wiring will never suffice
 - New causal principles associated with new bits of technology.
 - Will need powerful priors than used to date

Proposal for Robust AI

- Build learning systems
 - that start with initial frameworks for domains like time, space, and causality,
- in order to speed up learning and massively constrain the hypothesis space

Summary