# Project Aristo: Towards Machines that Capture and Reason with Science Knowledge

Peter Clark

November 2019

# The History of KCap

- KCap 2001-19

# The History of KCap
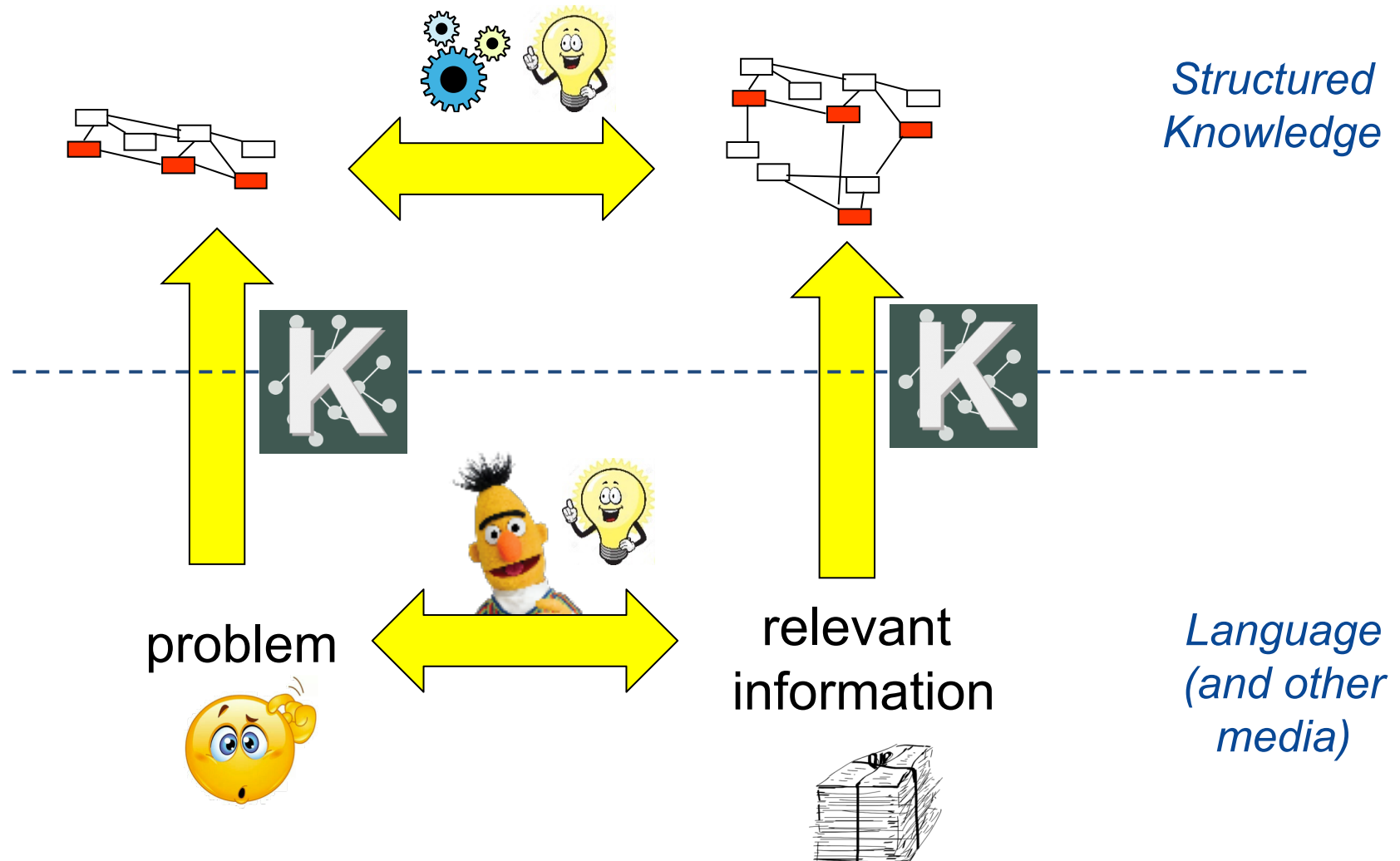
- KCap 2001-19
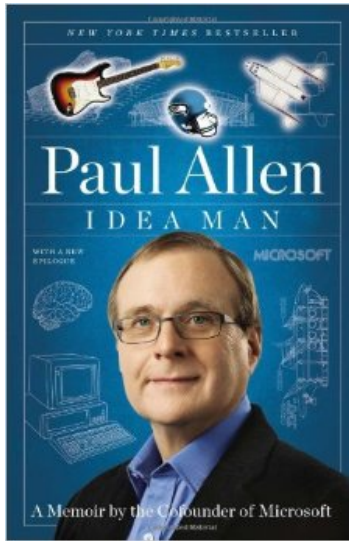
- Banff Knowledge Acquisition Workshops: 1986-1999

How do we get knowledge into the machine in a usable form?

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

*Structured Knowledge*

problem

relevant information

*Language (and other media)*

# Science Questions: A Grand Challenge…



*Over the last decade, I began to think about a **"Digital Aristotle",** an easy-to-use, all-encompassing knowledge storehouse....to advance the field of AI.*
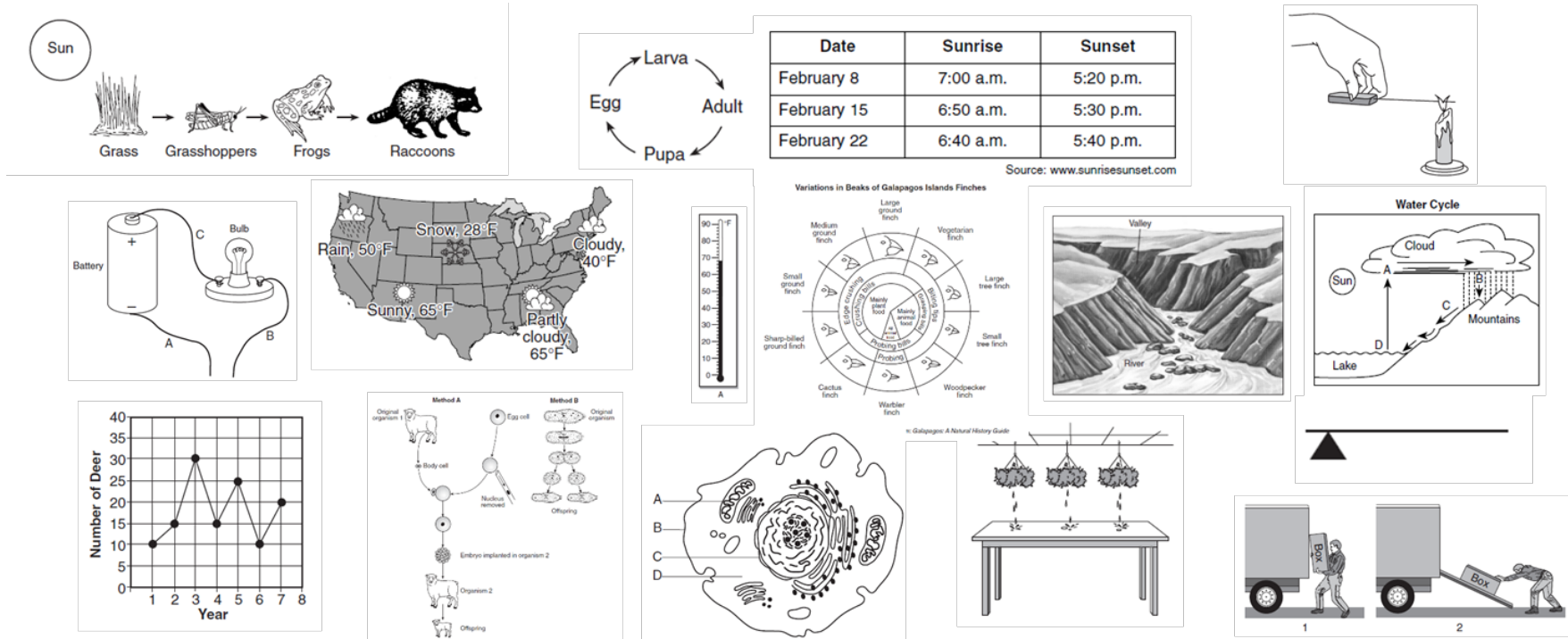
How are the particles in a block of iron affected when the block is melted?

(A) The particles gain mass.
(B) The particles contain less energy.
(C) **The particles move more rapidly.**
(D) The particles increase in volume.



ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# Question Categories Not Covered

- ## Diagrams



- ## Direct Answer Questions

# Progression on NY Regents 8th Grade (NDMC)



(hidden test set, questions as written, NDMC, 5 years/119 qns)

# Progression on NY Regents 8th Grade (NDMC)



(hidden test set, questions as written, NDMC, 5 years/119 qns)

**$80,000 • 119 teams**

# The Allen AI Science Challenge

Merger and 1st Submission Deadline

Wed 7 Oct 2015                                                    Sat 13 Feb 2016 (4.0 days to go)

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

## Dashboard

Competition Details   »   Get the Data   »   Make a submission

Home                    ⌂
  Data                  ⌸
  Make a submission     ☑

# Is your model smarter than an 8th grader?

Information             ⓘ
  Description
  Evaluation
  Rules
  Prizes
  Timeline

Forum                   💬

Leaderboard             ≡

## Public Leaderboard

1. amsqr
2. Cardal
3. poweredByTalkwalker
4. Generation Gap
5. yamayamada

The Allen Institute for Artificial Intelligence (AI2) is working to improve humanity through fundamental advances in artificial intelligence. One critical but challenging problem in AI is to demonstrate the ability to consistently understand and correctly answer general questions about the world.

The Aristo project at AI2 is focused on building such a system. One way Aristo "learns" is by extracting facts from various sources and processing them into a structured knowledge base. When taking an exam, questions are parsed and processed along with

(random)

2018    2019

(hidden test set, questions as written, NDMC, 5 years/119 qns)

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# Progression on NY Regents 8th Grade (NDMC)



(hidden test set, questions as written, NDMC, 5 years/119 qns)

# Progression on NY Regents 8th Grade (NDMC)



Separate test on 3 latest exams (2017-2019): 93.3%

(hidden test set, questions as written, NDMC, 5 years/119 qns)

# Outline

- Introduction
- How does Aristo work? ←
- What is going on behind the high scores on the exams?
- Where does Aristo fail?
- What are steps forward?

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# Aristo: an over-simplified overview

- An ensemble architecture

ARISTO

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# Aristo: an over-simplified overview

- An ensemble architecture

ARISTO

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# Aristo: an over-simplified overview

- An ensemble architecture

ARISTO



16

In New York State, the longest period of daylight occurs during which month? (A) June (B) March  (C) December  (D) September

# 1. Table Knowledge

In New York State, the longest period of daylight occurs during which month? **(A) June** (B) March  (C) December  (D) September

- Daylengths in different months and locations?
- Solstices?
- Where is New York State?
- Which hemisphere is it in?

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# 1. Table Knowledge: Aristo's Tablestore

- ~120 tables, ~10-500 rows each
- Defined with respect to questions, study guides, syllabus

# IKE – Interactive Knowledge Extraction



(AKBC'16)

# 1. Table Inference

In New York State, the longest period of daylight occurs during which month? (A) June (B) March  (C) December  (D) September

| Subdivision | Country |
|---|---|
| New York State | USA |
| California | USA |
| Rio de Janeiro | Brazil |
| ... | ... |

| Orbital Event | Day Duration | Night Duration |
|---|---|---|
| Summer Solstice | Long | Short |
| Winter Solstice | Short | Long |
| .... | .... | ... |

| Country | Hemisphere |
|---|---|
| United States | Northern |
| Canada | Northern |
| Brazil | Southern |
| ...... | ... |

| Hemisphere | Orbital Event | Month |
|---|---|---|
| North | Summer Solstice | June |
| North | Winter Solstice | December |
| South | Summer Solstice | December |
| South | Winter Solstice | June |

Semi-structured Knowledge

# 2. Table Inference

In New York State, the longest period of daylight occurs during which month? (A) June (B) March  (C) December  (D) September

| Subdivision | Country |
|---|---|
| New York State | USA |
| California | USA |
| Rio de Janeiro | Brazil |
| ... | ... |

| Orbital Event | Day Duration | Night Duration |
|---|---|---|
| Summer Solstice | Long | Short |
| Winter Solstice | Short | Long |
| .... | .... | ... |

| Country | Hemisphere |
|---|---|
| United States | Northern |
| Canada | Northern |
| Brazil | Southern |
| ...... | ... |

| Hemisphere | Orbital Event | Month |
|---|---|---|
| North | Summer Solstice | June |
| North | Winter Solstice | December |
| South | Summer Solstice | December |
| South | Winter Solstice | June |

**Semi-structured Knowledge**

# 1. Table Inference

In **New York State**, the **longest period of daylight** occurs during which month? (A) June (B) March  (C) December  (D) September

| Subdivision | Country |
|---|---|
| New York State | USA |
| California | USA |
| Rio de Janeiro | Brazil |
| ... | ... |

| Orbital Event | Day Duration | Night Duration |
|---|---|---|
| Summer Solstice | Long | Short |
| Winter Solstice | Short | Long |
| .... | .... | ... |

| Country | Hemisphere |
|---|---|
| United States | Northern |
| Canada | Northern |
| Brazil | Southern |
| ...... | ... |

| Hemisphere | Orbital Event | Month |
|---|---|---|
| North | Summer Solstice | June |
| North | Winter Solstice | December |
| South | Summer Solstice | December |
| South | Winter Solstice | June |

**Semi-structured Knowledge**

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE
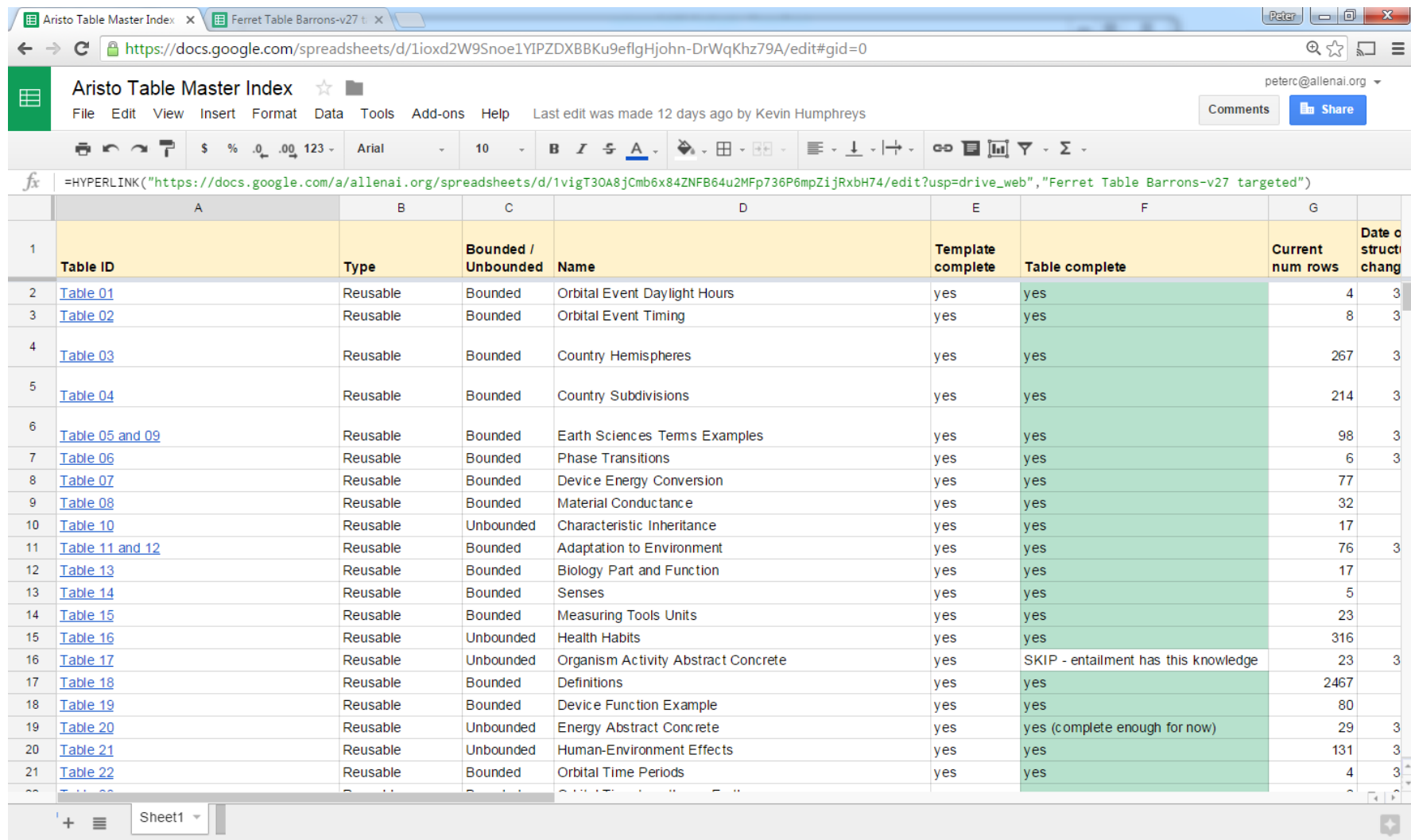
In New York State, the longest period of daylight occurs during which month? (A) June (B) March (C) December (D) September

| Subdivision | Country |
|---|---|
| New York State | USA |
| California | USA |
| Rio de Janeiro | Brazil |
| ... | ... |

| Orbital Event | Day Duration | Night Duration |
|---|---|---|
| Summer Solstice | Long | Short |
| Winter Solstice | Short | Long |
| .... | .... | ... |

| Country | Hemisphere |
|---|---|
| United States | Northern |
| Canada | Northern |
| Brazil | Southern |
| ...... | ... |

| Hemisphere | Orbital Event | Month |
|---|---|---|
| North | Summer Solstice | June |
| North | Winter Solstice | December |
| South | Summer Solstice | December |
| South | Winter Solstice | June |

Semi-structured Knowledge

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# 1. Table Inference

In New York State, the longest period of daylight occurs during which month? (A) June (B) March (C) December (D) September

| Subdivision | Country |
|---|---|
| New York State | USA |
| California | USA |
| Rio de Janeiro | Brazil |
| … | … |

| Orbital Event | Day Duration | Night Duration |
|---|---|---|
| Summer Solstice | Long | Short |
| Winter Solstice | Short | Long |
| …. | …. | … |

| Country | Hemisphere |
|---|---|
| United States | Northern |
| Canada | Northern |
| Brazil | Southern |
| …… | … |

| Hemisphere | Orbital Event | Month |
|---|---|---|
| North | Summer Solstice | June |
| North | Winter Solstice | December |
| South | Summer Solstice | December |
| South | Winter Solstice | June |

Semi-structured Knowledge

IJCAI'16

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# 2. Tuple Knowledge

Search → Tuple Extraction → Headword Extraction → Scoring + Filtering → Regen-eration → Canonical-ization

(ACL'17)

| Score | | | Tuple | | Verbalization |
|---|---|---|---|---|---|
| ... | ... | ... | ... | ... | .... |
| 1.00 | most | elephant | isa | mammal | // Elephant isa mammal. |
| 1.00 | most | elephant | isa | pachyderm | // Elephant isa pachyderm. |
| 1.00 | most | elephant | require | litre water | // Most elephants require litre water. |
| 1.00 | most | elephant | require | water | // Most elephants require water. |
| ... | | | | | |
| 0.92 | most | elephant | have | curve spine | // Most elephants have curve spines. |
| 0.92 | most | elephant | need | food | // Most elephants need food. |
| ... | ... | ... | ... | ... | ... |
| 0.83 | most | computer | receive | electric energy | // Most computers receive electric energy. |
| 0.67 | most | computer | solve | problem | // Most computers solve problems. |
| 0.60 | most | computer | provide | prediction | // Most computers provide predictions. |
| .... | ... | ... | ... | ... | ... |

for ARTIFICIAL INTELLIGENCE

# 2. Tuple Inference

Which object in our solar system reflects light and is a satellite that orbits around one planet? (A) Moon (B) Earth (C) Mercury (D) Sun

| Planets and **Moon** | **reflect** | **light** | back from the Sun |

| **Moon** | **orbits** | planets |

| **Moon** | is a | large natural **satellite** |

| Primary objects in solar system | are | planets and **moons** |

Stormy weather negatively affects a coastline by **(A) causing erosion** (B) causing earthquakes (C) increasing food production (D) increasing the growth of grasses

**Tuple KB**

**Text Corpus**

→ **Lexical Tuples** → **Semantic Transformation** → **Semantic Tuples** → **Tuple Chaining** → **Expanded Tuples** → **Support Graph Search**

| waves | *often result in* | erosion |
| storms | *can generate* | waves |

| waves | *cause* | erosion |
| storms | *cause* | waves |

| storms | *cause* | erosion |

# Aristo: an over-simplified overview

- An ensemble architecture

# Aristo: an over-simplified overview

- An ensemble architecture

ARISTO

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# BERT and RoBERTa



Aristo Text Corpus
(0.3TB)

**Multi-step Curriculum Training**

Background    Science    Regents

What part of a
sunlight to do

[CLS] context [SEP] question
[SEP] answer-option

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# Where is the Knowledge Capture?



What part of a plant needs sunlight to do its job? (A) leaf

Shin's research interests involve the structure and function of cell membrane proteins, including influenza hemaglutinin protein and an HIV virus spike protein that are responsible for cellular-viral membrane fusion. Biophysical chemists study protein structure and the functional structure of cell membranes. biological structure analysis by electron crystallography to characterize cell-membrane proteins and viruses; Structure-Function Analysis of the Influenza Virus Ion Channel Influenza virus protein M 2 is a small (97-residue) integral membrane protein that spans the cell membrane once and is minimally a disulfide-linked homotetramer. Membrane functions | top | Composition and Structure | Membrane proteins | Membrane functions | end | *dynamic boundary Cell membranes enclose the internal compartments of cells, allowing them to be different from the extracellular environment and from each other. In a third project scientists are examining the effects of Coxsackie B virus proteins on the function of internal cell membranes. Dr Michael Carter is to undertake a study which will examine the effects of Coxsackie B virus proteins on the function of internal cell membranes. A huge unsolved question of cell membrane structure and function is the structure of membrane proteins. Virus Proteins and Cell Membranes. Cell Membranes | top | Composition and Structure | Membrane proteins | Membrane functions |

Wikipedia + BookCorpus
(2.5B words + 11k books)

What part of a plant needs sunlight to do its job? (A) leaf

Wikipedia + BookCorpus
(2.5B words + 11k book

**Curriculum Training**

# Where is the Knowledge Capture?



Aristo Corpus
(2B words)

What part of a plant needs
sunlight to do its job? (A) leaf…

Wikipedia + BookCorpus
(2.5B words + 11k books)

**Curriculum Training**

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Exploiting Language Models

# Similar Progress on 4th Grade NDMC



Similarly on 12th grade NDMC:
- Random: 25.0%
- 2014:     40.6%
- 2019:     83.5%

# Individual Solver Performances

| Test Set | Num Q | IR | PMI | ACME | TupInf | Multee | AristoBERT | AristoRoBERTa | ARISTO |
|---|---|---|---|---|---|---|---|---|---|
| Regents 4th | 109 | 64.45 | 66.28 | 67.89 | 63.53 | 69.72 | 86.24 | 88.07 | **89.91** |
| Regents 8th | 119 | 66.60 | 69.12 | 67.65 | 61.41 | 68.91 | 86.55 | 88.24 | **91.60** |
| Regents 12th | 632 | 41.22 | 46.95 | 41.57 | 35.35 | 56.01 | 75.47 | 82.28 | **83.54** |
| ARC-Easy | 2376 | 74.48 | 77.76 | 66.60 | 57.73 | 64.69 | 81.78 | 82.88 | **86.99** |
| ARC-Challenge | 1172 | n/a† | n/a† | 20.44 | 23.73 | 37.36 | 57.59 | **64.59** | 64.33 |

Most of the heavy lifting….

- Introduction
- How does Aristo work?
- What is going on behind the high scores on the exams?
- Where does Aristo fail?
- What are steps forward?

▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉▉

(A) friction
(B) light
(C) force
(D) weather

| Test dataset | "Answer only" score |
|---|---|
| Regents 4th | 38.53 |
| Regents 8th | 37.82 |
| Regents 12th | 47.94 |
| ARC-Easy | 36.17 |
| ARC-Challenge | 35.92 |
| All | 37.11 |

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# 2. Is it fooled by "obviously wrong" answers?

The condition of the air outdoors at a certain time of day is known as
(A) friction
(B) light
(C) force
**(D) weather** **[selected, correct]**

The condition of the air outdoors at a certain time of day is known as

| | |
|---|---|
| (A) friction | (E) joule |
| (B) light | (F) gradient |
| (C) force | (G) trench |
| (D) weather | (H) add heat |

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# 2. Is it fooled by "obviously wrong" answers?

The condition of the air outdoors at a certain time of day is known as
(A) friction
(B) light
(C) force
(D) weather [selected, correct] ✓

The condition of the air outdoors at a certain time of day is known as
(A) friction            (E) joule
(B) light               (F) gradient [selected] ✗
(C) force               (G) trench
(D) weather [correct]   (H) add heat

Retrain

The condition of the air outdoors at a certain time of day is known as
(A) friction            (E) joule
(B) light               (F) gradient
(C) force               (G) trench
(D) weather [correct,selected] ✓

# 2. Is it fooled by "obviously wrong" answers?

The condition of the air outdoors at a certain time of day is known as
(A) friction
(B) light
(C) force
**(D) weather**

| Test dataset | 4-way MC | Adversarial 8-way MC | % drop (relative) |
|---|---|---|---|
| Regents 4th | 87.1 | 76.1 | 12.6 |
| Regents 8th | 78.9 | 76.4 | 3.1 |
| Regents 12th | 75.3 | 58.0 | 22.9 |
| ARC-Easy | 74.1 | 65.7 | 11.3 |
| ARC-Challenge | 55.5 | 47.7 | 14.0 |
| ALL | 69.1 | 59.5 | 13.8 |

Drop of (only) ≈ 10 points

Retrain

The condition of the air outdoors at a certain time of day is known as
(A) friction          (E) joule
(B) light             (F) gradient [selected]
(C) force             (G) trench
**(D) weather [correct,selected]**

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

City administrators can encourage energy conservation by
(1) lowering parking fees
(2) building larger parking lots
(3) decreasing the cost of gasoline
(4) lowering the cost of bus and subway fares

50

*increasing*

*raising*

City administrators can encourage energy conservation by
(1) lowering parking fees
(2) building larger parking lots
(3) ~~decreasing~~ the cost of gasoline
(4) ~~lowering~~ the cost of bus and subway fares

✓

Which of the following organs does a squirrel *not* have
(A) a brain
(B) gills
(C) a heart
(D) lungs

✓

51

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# 3. More than Pattern Matching?

City administrators can encourage energy conservation by
(1) lowering parking fees
(2) building larger parking lots
(3) ~~decreasing~~ the cost of gasoline
(4) ~~lowering~~ the cost of bus and subway fares

*increasing*
*raising*

✓

Which of the following organs does a squirrel ~~*not*~~ have
(A) a brain
(B) gills
(C) a heart
(D) lungs

✓

52

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

# 3. More than Pattern Matching?

2019 Report Card for   *Aristo*

| Subject | Grade | Teacher Comments |
|---------|-------|------------------|
| Negation | *A* | *Nice work!* |
| Conjunction | | |
| Polarity | | |
| World tracking | | |
| Factivity | | |
| Counting | | |

*94%*

Alan is small.      Alan is tall.      Bob is big.      Bob is tall.
Charlie is big.     Charlie is tall.   David is small.   David is short.

Which of the following is *not* tall? (A) Alan (B) Bob (C) Charlie (D) David *[correct]*

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

2019 Report Card for _____Aristo_____

| Subject | Grade | Teacher Comments |
|---|---|---|
| Negation | A | Nice work! |
| **Conjunction** | | |
| Polarity | | |
| World tracking | | |
| Factivity | | |
| Counting | | |

*94%*

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# Synthetic Conjunction Test

**Context:**

> Alan is red.
> Alan is big.
> Bob is blue.
> Bob is small.
> Charlie is blue.
> Charlie is big.
> David is red.
> David is small.

**Question:**

Which of the following is big **and** blue? (A) Alan (B) Bob (C) Charlie **[correct]** (D) David



| 1 conjunct: 98% | | 88.5% |
| 2 conjuncts: 95% | | 76.5% |
| 3 conjuncts: 94.5% | + 1 negation | 76% |
| 4 conjuncts: 80% | | 75% |

Alan is red. Alan is big. Alan is light. Alan is old. Alan is tall. Bob is red. Bob is small. Bob is heavy. Bob is old. Bob is tall. Charlie is blue. Charlie is big. Charlie is light. Charlie is old. Charlie is tall. David is red. David is small. David is heavy. David is young. David is tall.

Which of the following is old **and** red **and** light and big **and not** short? (A) Alan (B) Bob (C) Charlie (D) David

2019 Report Card for _Aristo_

| Subject | Grade | Teacher Comments |
|---|---|---|
| Negation | A | Nice work! |
| Conjunction | B+ | |
| Polarity | | |
| World tracking | | |
| Factivity | | |
| Counting | | |

94%

80% -98%

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# 3. More than Pattern Matching?

2019 Report Card for _____Aristo_____

| Subject | Grade | Teacher Comments | |
|---------|-------|------------------|---|
| Negation | A | Nice work! | 94% |
| Conjunction | B+ | | 80% -98% |
| **Polarity** | D+ | Could ace this with more study! | 67.1% |
| World tracking | | | |
| Factivity | | | |
| Counting | | | |

**Context:** For a given medium, sound has a slower speed at lower temperatures.

**Question:**
~~up~~
If Jim turns the thermostat ~~down~~ in his room while listening to music, what will happen to the speed of the sound waves in the room?
(A) they will speed up *[correct]* (B) they will slow down ~~*[correct]*~~

*[correct]*

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

2019 Report Card for ____Aristo____

| Subject | Grade | Teacher Comments | |
|---------|-------|------------------|---|
| Negation | A | Nice work! | 94% |
| Conjunction | B+ | | 80% -98% |
| Polarity | D+ | Could ace this with more study! | 67.1% |
| World tracking | C | | 72.5% |
| Factivity | | | |
| Counting | | | |

**Context:** If someone travels for longer, they will travel further.

**Question:** John and Rita are going for a run. Rita gets tired and takes a break on the park bench. After twenty minutes in the park, who has run farther?
(A) John *[correct]* (B) Rita

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# 3. More than Pattern Matching?

2019 Report Card for  ___*Aristo*___

| Subject | Grade | Teacher Comments | |
|---|---|---|---|
| Negation | A | *Nice work!* | *94%* |
| Conjunction | B+ | | *80% -98%* |
| Polarity | D+ | *Could ace this with more study!* | *67.1%* |
| World tracking | C | | *72.5%* |
| Factivity | D | | *66.5%* |
| Counting | | | |

If someone ***regretted*** that a particular thing happened then
    (A) that thing might or might not have happened .
    (B) that thing didn't happen .
    (C) **that thing happened *[correct]***

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# 3. More than Pattern Matching?

2019 Report Card for _____Aristo_____

| Subject | Grade | Teacher Comments | |
|---|---|---|---|
| Negation | A | *Nice work!* | 94% |
| Conjunction | B+ | | 80% -98% |
| Polarity | D+ | *Could ace this with more study!* | 67.1% |
| World tracking | C | | 72.5% |
| Factivity | D | | 66.5% |
| Counting | F | | 6% |

Daniel picked up the football. Daniel dropped the football. Daniel got the milk.

How many objects is Daniel holding? (A) zero **(B) one** (C) two (D) three

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# Outline

- Introduction
- How does Aristo work?
- What is going on behind the high scores on the exams?
- Where does Aristo fail?
- What are steps forward?

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# 4. Where is Aristo failing?

- Case study on 30 failures:

Good support for correct answer

Reading Comprehension
(IR won't help)

Good support for *in*correct answer

4 (13.3%)

8 (26.7%)

No good support

17 (56.7%)

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

Which is the best unit to measure distances between Earth and other solar systems in the universe? (A) miles (B) kilometers **(C) light years (D) astronomical units**

*In general, distances in the solar system are measured in astronomical units.*

*Distances between Earth and the stars are often measured in terms of light-years.*

Which of these objects will most likely float in water? (A) glass marble (B) steel ball **(C) hard rubber ball** **(D) table tennis ball**

- *I remember it had like a **rubber ball** in it, which would maybe **float up**…*
- *We played soccer with a giant **rubber ball that floated** like a balloon.*
- ***Rubber toys floated** on the water.*

ALLEN INSTITUTE
_for_ ARTIFICIAL INTELLIGENCE

Although they belong to the same family, an eagle and a pelican are different. What is one difference between them? (A) their preference for eating fish (B) their ability to fly **(C) their method of reproduction (D) their method of catching food**

- Need question decomposition

How are the particles in a block of iron affected when the block is melted? **(A) The particles gain mass**. (B) The particles contain less energy. **(C) The particles move more rapidly.** (D) The particles increase in volume.

- No good single supporting sentence

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

Which characteristic applies to animals in only one of these taxonomic groups: reptiles, mammals, birds, amphibians, or fishes? **(A) have hair (B) lay eggs** (C) have webbed feet (D) breathe with gills

- Boolean reasoning

Which geologic structure will most likely take the longest time to form? (A) a fault (B) a sinkhole **(C) a river meander (D) a mountain range**

- Cross-option comparative

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

- Story (experimental method)

A student wants to determine the effect of garlic on the growth of a fungus species. Several samples of fungus cultures are grown in the same amount of agar and light. Each sample is given a different amount of garlic. What is the independent variable in this investigation? (A) amount of agar (B) amount of light **(C) amount of garlic (D) amount of growth**

- Meta/sentiment

Which statement is an opinion? (A) Many plants are green. **(B) Many plants are beautiful. (C) Plants require sunlight.** (D) Plants can grow in different places.

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# Math Reasoning

About how long does it take for the Moon to complete one revolution around Earth? (A) 7 days **(B) 30 days** (C) 90 days **(D) 365 days**

- *Because it takes the moon about **27.3 days** to complete one orbit around the Earth, the moon moves a little bit further around the Earth each day.*
- *It takes **27.3 days** for the moon to complete one revolution around the earth.*
- *The moon completes one revolution of the Earth in about **29.5 days.***
- *The Moon completes one revolution around the Earth in **27.32166 days.***

# Outline

- Introduction

- How does Aristo work?

- What is going on behind the high scores on the exams?

- Where does Aristo fail?

- What are steps forward?

What virus structure is similar in function to a cell membrane?
(A) protein shell (B) internal protein...

Shin's research interests involve the structure and function of cell membrane proteins, including influenza hemaglutinin protein and an HIV virus spike protein that are responsible for cellul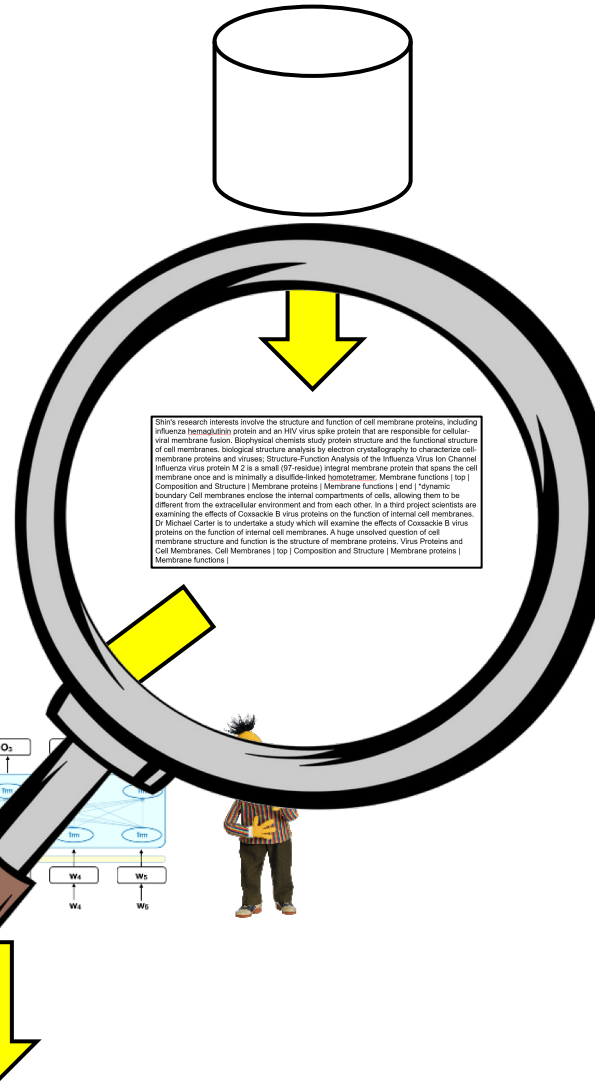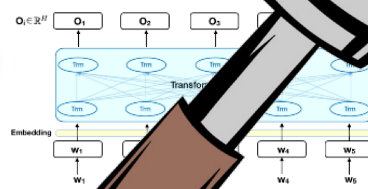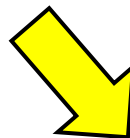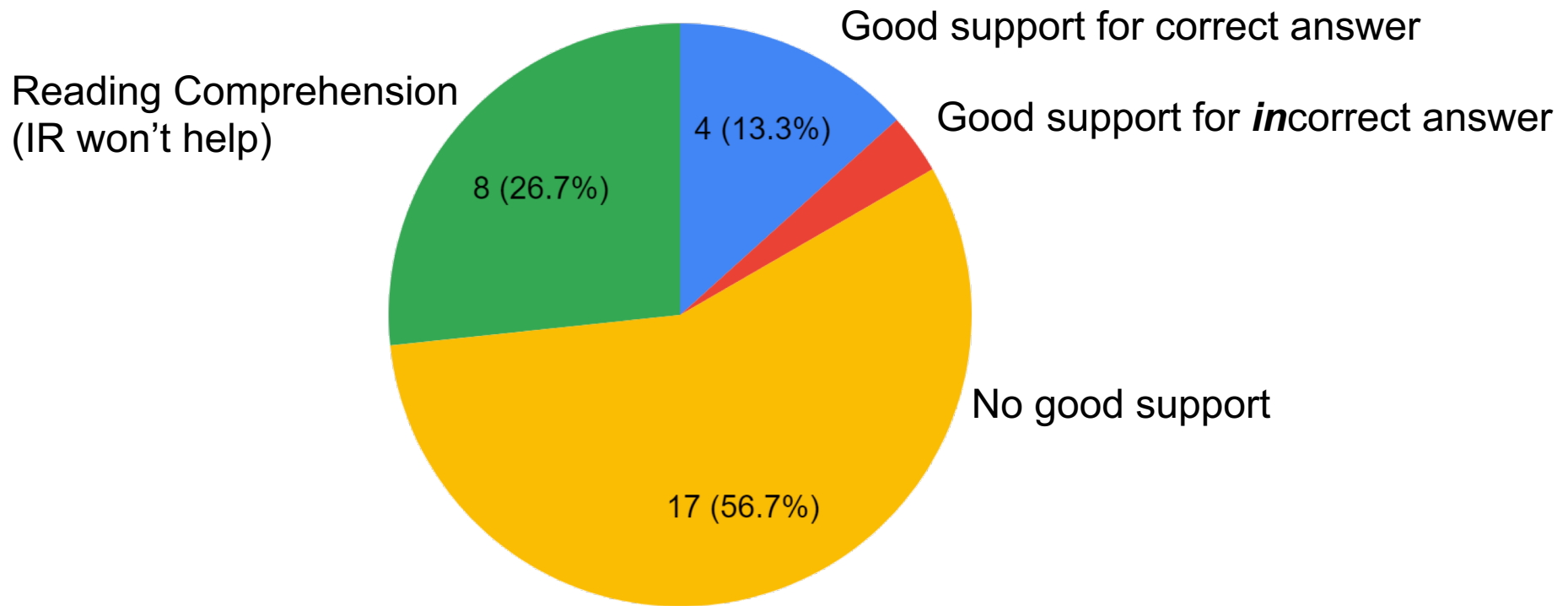ar-viral membrane fusion. Biophysical chemists study protein structure and the functional structure of cell membranes. biological structure analysis by electron crystallography to characterize cell-membrane proteins and viruses; Structure-Function Analysis of the Influenza Virus Ion Channel Influenza virus protein M 2 is a small (97-residue) integral membrane protein that spans the cell membrane once and is minimally a disulfide-linked homotetramer. Membrane functions | top | Composition and Structure | Membrane proteins | Membrane functions | end | *dynamic boundary Cell membranes enclose the internal compartments of cells, allowing them to be different from the extracellular environment and from each other. In a third project scientists are examining the effects of Coxsackie B virus proteins on the function of internal cell membranes. Dr Michael Carter is to undertake a study which will examine the effects of Coxsackie B virus proteins on the function of internal cell membranes. A huge unsolved question of cell membrane structure and function is the structure of membrane proteins. Virus Proteins and Cell Membranes. Cell Membranes | top | Composition and Structure | Membrane proteins | Membrane functions |

structure-function of membrane proteins. membrane protein structure and function; Structure and function of membrane proteins; Shin's research interests involve the structure and function of cell membrane proteins, including influenza hemaglutinin protein and an HIV virus spike protein that are responsible for cellular-viral membrane fusion. biological structure analysis by electron crystallography to characterize cell-membrane proteins and viruses; Structure-Function Analysis of the Influenza Virus Ion Channel Influenza virus protein M 2 is a small (97-residue) integral membrane protein that spans the cell membrane once and is minimally a disulfide-linked homotetramer. Biophysical chemists study protein structure and the functional structure of cell membranes. A huge unsolved question of cell membrane structure and function is the structure of membrane proteins. Virus Proteins and Cell Membranes. Cell Membranes | top | Composition and Structure | Membrane proteins | Membrane functions |

# 1. Question Decomposition

What virus structure is similar in function to a cell membrane?
(A) protein shell (B) internal protein...

➡ What is the function of a cell membrane?

⬅ Surrounds and protects, gives structure, regulates material, ….

➡ What part of the virus surrounds and protects it?

⬅ Protein shell, protein layer, …

- GapQA *(EMNLP'19)*
- New dataset coming

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

Which conducts electricity? (A) suit of armor (B) cotton candy

# 2. Multihop Reasoning

Which ==conducts electricity==? (A) suit of armor (B) cotton candy

**Retrieval 1:**
The reciprocal of the electrical resistivity is the ==electrical conductivity==.
==Electrical conductivity== is the capacity of metal to ==conduct an electric== current.
==Electrical Conductivity== Water without minerals will not ==conduct electricity.==

# 2. Multihop Reasoning

Which conducts electricity? (A) suit of armor (B) cotton candy

**Retrieval 1:**
The reciprocal of the electrical resistivity is the electrical conductivity.
Electrical conductivity is the capacity of metal to conduct an electric current.
Electrical Conductivity Water without minerals will not conduct electricity.

**Retrieval 2:**
It was not suited to be a center for extensive metal-working.
A suit of armour is a historical type of personal body armour made from metal.
Resisting arrest is a criminal charge, but civil suits can be filed.

**Form Chains:**
"suit of armor…made from metal"  AND "…metal conduct electrical current"
            => "suit of armor conducts electricity"    ☺

"Resisting arrest…suits can be filed" AND "reciprocal of resistivity is conductivity"
            => "suit of armor conducts electricity"    ☹

**Train system to recognize good chains**

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

**Photosynthesis**

Roots absorb water from the soil.

The water flows to the leaf.

Light and $CO_2$ enter leaf.

Light, water, $CO_2$ form sugar.

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

**Photosynthesis**

Roots absorb water from the soil.

The water flows to the leaf.

Light and CO2 enter leaf.

Light, water, CO2 form sugar.



Where is the sugar created? **Light, water, CO2** [BiDAF]

# 3. Modeling World States

## Paragraph

|    |                               |
|----|-------------------------------|
|    |                               |
| s1 | Roots absorb water from soil. |
|    |                               |
| s2 | The water flows to the leaf.  |
|    |                               |
| s3 | Light and CO2 enter leaf.     |
|    |                               |
| s4 | Water, light, CO2 form sugar. |
|    |                               |

## State changes: $\pi$

|    | water | light | CO2 | sugar |
|----|-------|-------|-----|-------|
|    |       |       |     |       |
| s1 |       |       |     |       |
|    |       |       |     |       |
| s2 |       |       |     |       |
|    |       |       |     |       |
| s3 |       |       |     |       |
|    |       |       |     |       |
| s4 |       |       |     |       |
|    |       |       |     |       |

ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

**Paragraph**

**State changes**: $\pi$

|  | Roots absorb water from soil. |
|---|---|
| s1 | |
| s2 | The water flows to the leaf. |
| s3 | Light and CO2 enter leaf. |
| s4 | Water, light, CO2 form sugar. |

|  | **water** | light | CO2 | sugar |
|---|---|---|---|---|
|  | **soil** | | | |
| s1 | | | | |
|  | **roots** | | | |
| s2 | | | | |
|  | | | | |
| s3 | | | | |
|  | | | | |
| s4 | | | | |
|  | | | | |

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

**Paragraph**

**State changes**: $\pi$

|     |                               |
|-----|-------------------------------|
|     |                               |
| s1  | Roots absorb water from soil. |
|     |                               |
| s2  | The water flows to the leaf.  |
|     |                               |
| s3  | Light and CO2 enter leaf.     |
|     |                               |
| s4  | Water, light, CO2 form sugar. |
|     |                               |

|    | water | light | CO2 | sugar |
|----|-------|-------|-----|-------|
|    | soil  |       |     |       |
| s1 |       |       |     |       |
|    | roots |       |     |       |
| s2 |       |       |     |       |
|    |       |       |     |       |
| s3 |       |       |     |       |
|    |       |       |     |       |
| s4 |       |       |     |       |
|    |       |       |     |       |

**Paragraph**

**State changes**: $\pi$

| | |
|---|---|
| s1 | Roots absorb water from soil. |
| s2 | The water flows to the leaf. |
| s3 | Light and CO2 enter leaf. |
| s4 | Water, light, CO2 form sugar. |

| | water | light | CO2 | sugar |
|---|---|---|---|---|
| | soil | | | |
| s1 | | | | |
| | **roots** | | | |
| s2 | | | | |
| | **leaf** | | | |
| s3 | | | | |
| | | | | |
| s4 | | | | |
| | | | | |

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# 3. Modeling World States

**Paragraph**

**State changes:** $\pi$

| | |
|---|---|
| | |
| s1 | Roots absorb water from soil. |
| | |
| s2 | The water flows to the leaf. |
| | |
| s3 | Light and CO2 enter leaf. |
| | |
| s4 | Water, light, CO2 form sugar. |
| | |

| | water | light | CO2 | sugar |
|---|---|---|---|---|
| | soil | | | |
| s1 | | | | |
| | roots | | | |
| s2 | | | | |
| | **leaf** | | | |
| s3 | | | | |
| | **leaf** | | | |
| s4 | | | | |
| | | | | |

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

# 3. Modeling World States

**Paragraph**

**State changes:** $\pi$

| | | water | light | CO2 | sugar |
|---|---|---|---|---|---|
| | | soil | sun | ? | - |
| s1 | Roots absorb water from soil. | | | | |
| | | roots | sun | ? | - |
| s2 | The water flows to the leaf. | | | | |
| | | leaf | sun | ? | - |
| s3 | Light and CO2 enter leaf. | | | | |
| | | leaf | leaf | leaf | - |
| s4 | Water, light, CO2 form sugar. | | | | |

**LEADERBOARD**

ProPara

ProPara aims to promote the research in natural language understanding in the context of procedural... (More)

Public Submissions    Getting Started    About

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

Can you pick up a penny with a magnet?

Why?

Yes

Because
- *pennies are made of metal*
- *metals are magnetic*

Actually:
*Not all metals are magnetic.*
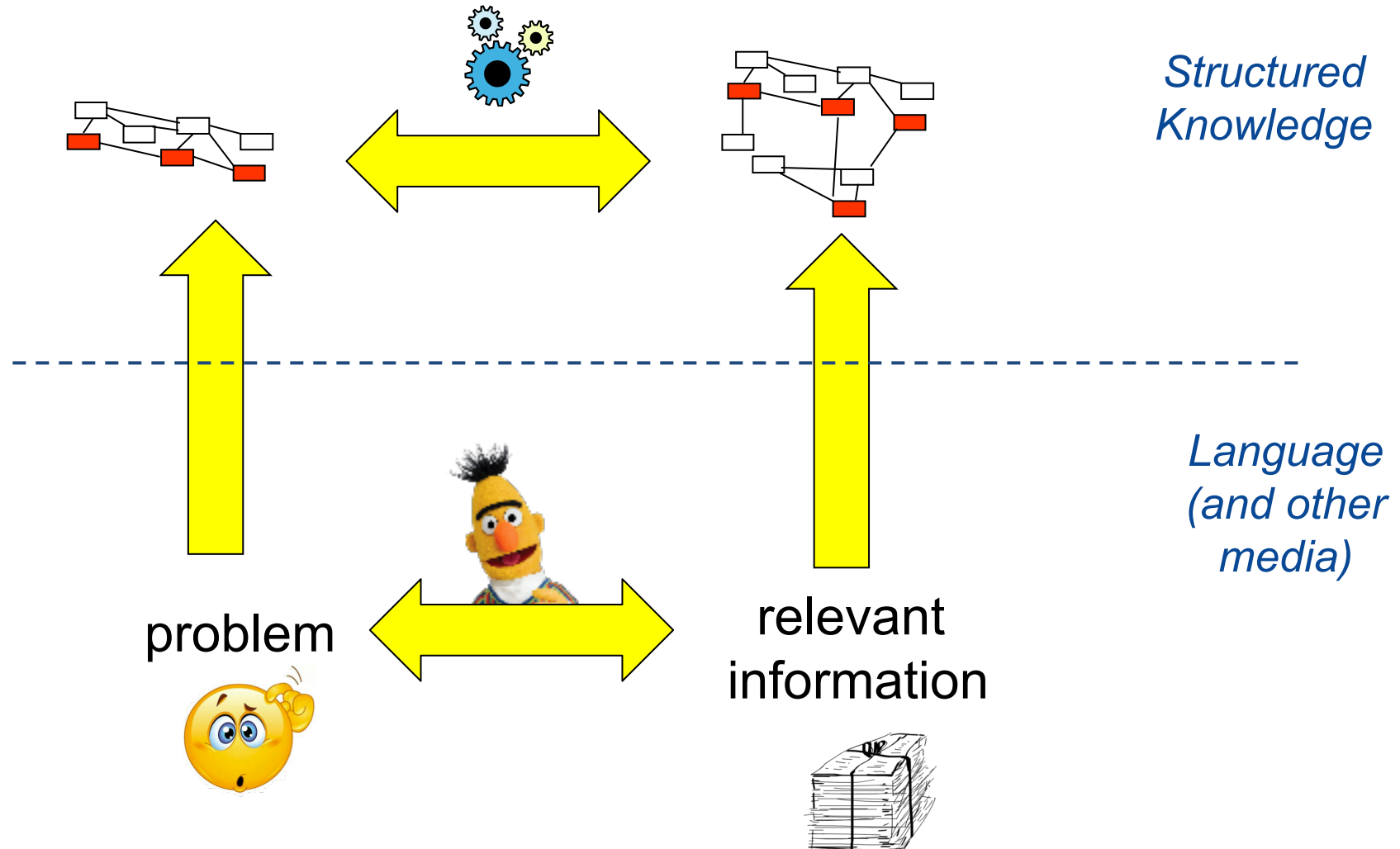*Copper is not magnetic.*

Try again!

No – because:
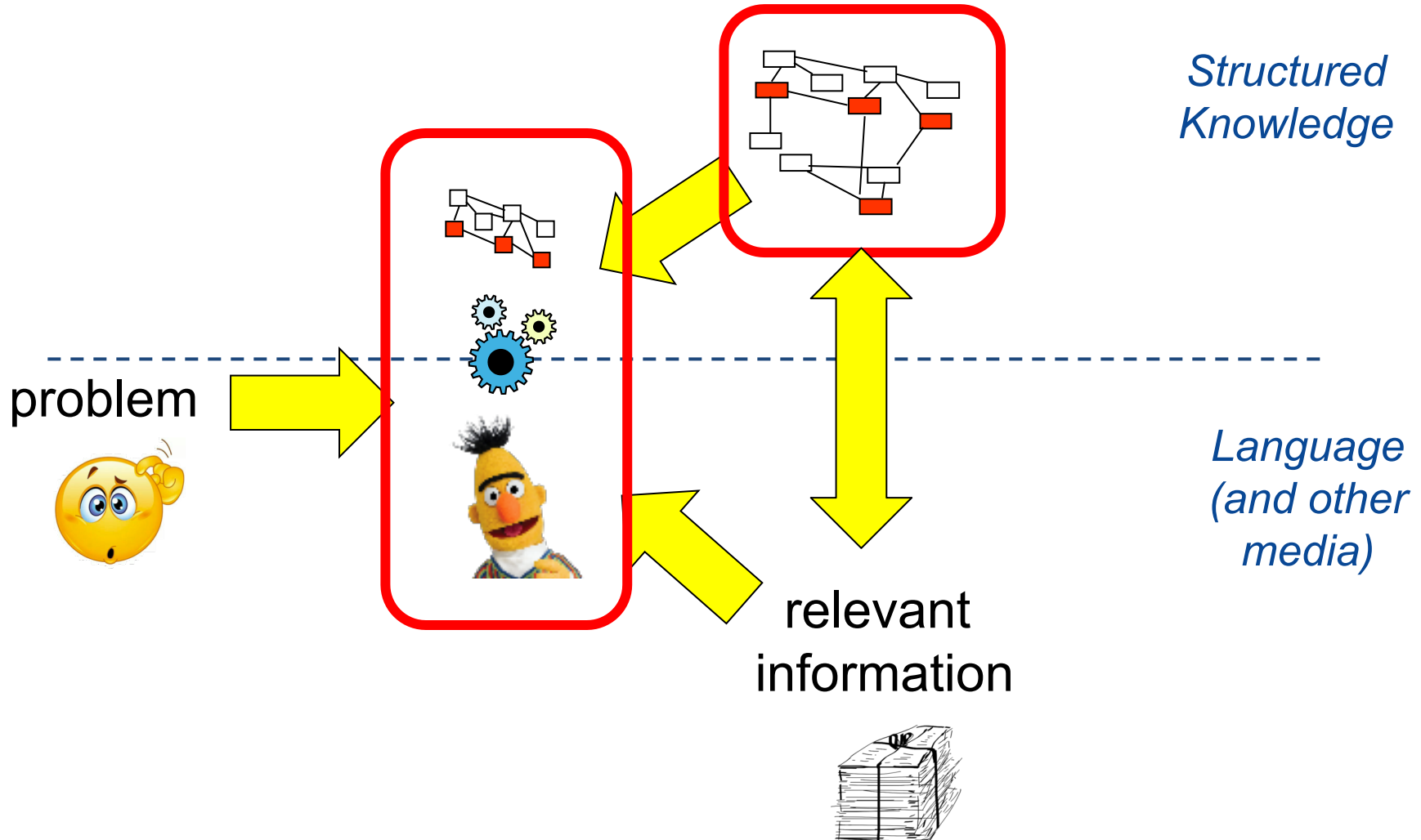- *pennies are made of copper*
- *copper is not magnetic*

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

*Structured Knowledge*

*Language (and other media)*

problem

relevant information

*Structured Knowledge*

problem

relevant information

*Language (and other media)*

ALLEN INSTITUTE *for* ARTIFICIAL INTELLIGENCE

**ARISTO**

- Surprising success!
  - LMs: Structure not essential for many tasks
  - >> "just pattern matching"
- BUT:
  - falls short with numerous types of questions
  - many other AI aspects missing

What do we need going forward?

- Structured reasoning and knowledge capture *but* with more language-like representations

**Thank you!**

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE