

# Semi-Structured Reasoning for Answering Science Questions

Daniel Khashabi, Dan Roth (UIUC)

Tushar Khot, Ashish Sabharwal, Peter Clark, Oren Etzioni (Allen Institute for Artificial Intelligence)

# Overview

**Challenge:** Build an AI system that demonstrates human-like intelligence by passing standardized science exams as written; **examples:** ----->

**Intermediate Goal: Elementary Science:** A simple embodiment of this challenge, requires question-answering significantly beyond retrieval techniques

**Approach:** A discrete optimization approach to QA for multiple-choice questions

**Results:** State of the art performance on 4<sup>th</sup> grade science (NY Regents exam)



Which physical structure would best help a bear to **survive a winter** in New York State? (A) big ears (B) black nose (C) **thick fur** (D) brown eyes



A student puts two identical plants soil. She gives them the same amount of water. She puts one of these plants near a **sunny window** and the other in a dark room. This experiment tests how the plants respond to (A) **light** (B) air (C) water (D) soil

## Main Idea

Search for the best Support Graph connecting the Question to an Answer through Tables.

Q: In New York State, the longest period of daylight occurs during which month?

The diagram illustrates the relationship between geographical data and astronomical data. It consists of three tables and a list of months.

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

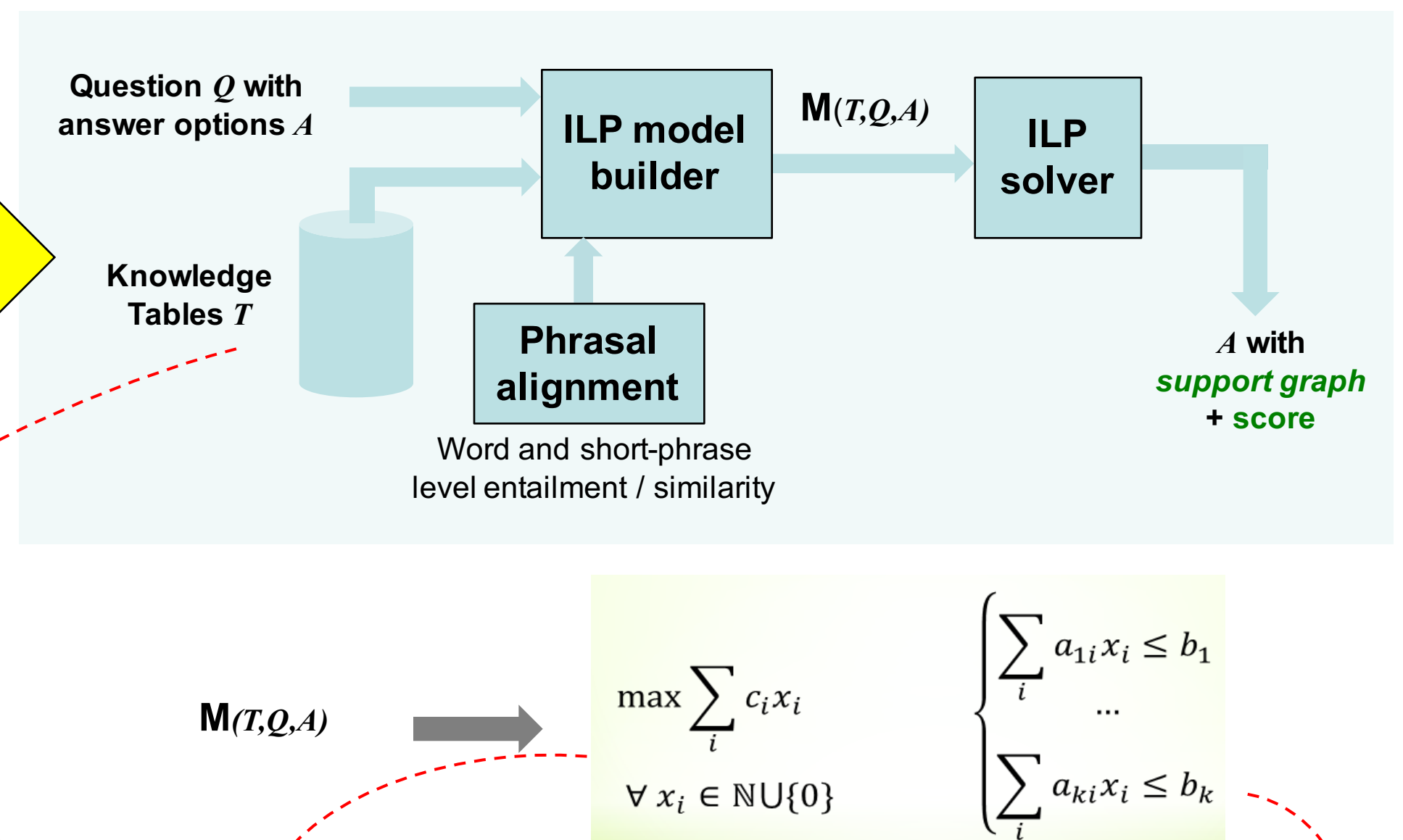
  

(A) December  
(B) June  
(C) March  
(D) September

# Architecture

A discrete **optimization** (Integer Linear Program) approach to QA for multiple-choice questions

Given as input a question, candidates answers, and a generic set of tables,  
and generates a constrained optimization problem



- **4<sup>th</sup> Grade NY Regents Science Exam**
  - 108 training set
  - 129 questions in completely unseen Test set
- **Baselines:**
  - **PMI Solver:** Statistical correlation using pointwise mutual info (280 GB of text)
  - **IR Solver:** Information Retrieval using Lucene search (280GB of text)
    - **IR Solver (table):** using only Tablestore
  - **MLN Solver:** using rules from 80K sentences

SENSE ORIGIN		PHENOMENON <i>Describe the sense detects</i>	ACTION <i>Type of sensing</i>	
		HEMISPHERE (north, south, equatorial region)	ORBITAL EVENT	MONTH OF OCCURRENCE
Yes	In	TOOL Device for	DIMENSION Parameter the	OBJECTS Type of object
No	He	IN THE PHASE CHANGE	INITIAL	FINAL
Sh	Sk			HEAT TRANSFER
		SENSE ORIGIN	PHENOMENON <i>Describe the sense detects</i>	ACTION <i>Type of sensing</i>
		Box	HEMISPHERE (north, south, equatorial region)	ORBITAL EVENT
		V4		MONTH OF OCCURRENCE
Eat	A	St	the northern hemisphere	the summer solstice occurs in June
No	C	Sr	the southern hemisphere	the summer solstice occurs in December
To	A	St	the northern hemisphere	the winter solstice occurs in December
No	Eat	Gr	the southern hemisphere	the winter solstice occurs in June
		Fr	the northern hemisphere	the spring equinox occurs in March
		Gr	the southern hemisphere	the spring equinox occurs in September
		Ce	the northern hemisphere	the fall equinox occurs in September
		Ce	the southern hemisphere	the fall equinox occurs in March

69 tables, ~7.6K rows

## Tablestore

- **Variables** define the space of “support graphs” connecting Q, A, T
  - Which nodes + edges between lexical units are active?
- **Objective** “better” support graphs = higher objective value
  - Reward active units, high lexical match links, column header match, ...
  - WH-term boost (which **form of energy**), science-term boost (**evaporation**)
  - Penalize spurious overuse of frequently occurring terms

## Objective function

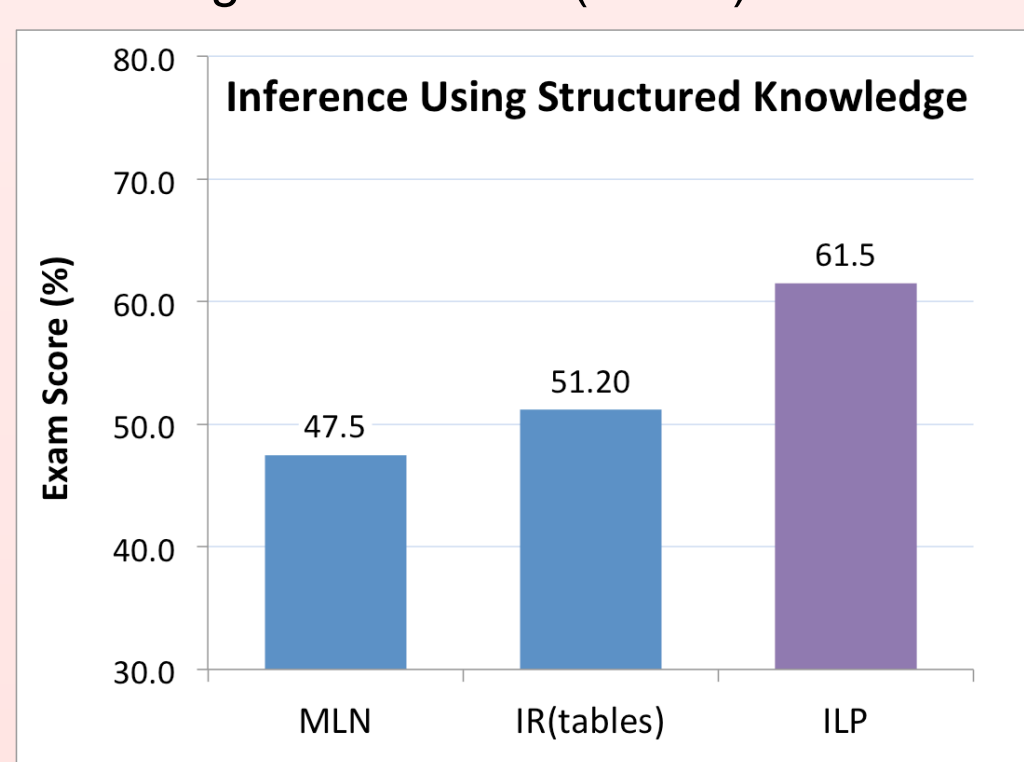
- **Structural Constraints**
  - Meaningful proof structures
    - connectedness, question coverage, appropriate table use
    - parallel evidence => identical multi-row activity signature
  - Simplicity appropriate for 4<sup>th</sup> / 8<sup>th</sup> grade
- **Semantic Constraints**
  - Chaining => table joins between semantically similar column pairs
  - Relation matching (ruler **measures** length, **change from** water **to** liquid)

## Constraints

## Evaluation

## Exploiting Structured Knowledge

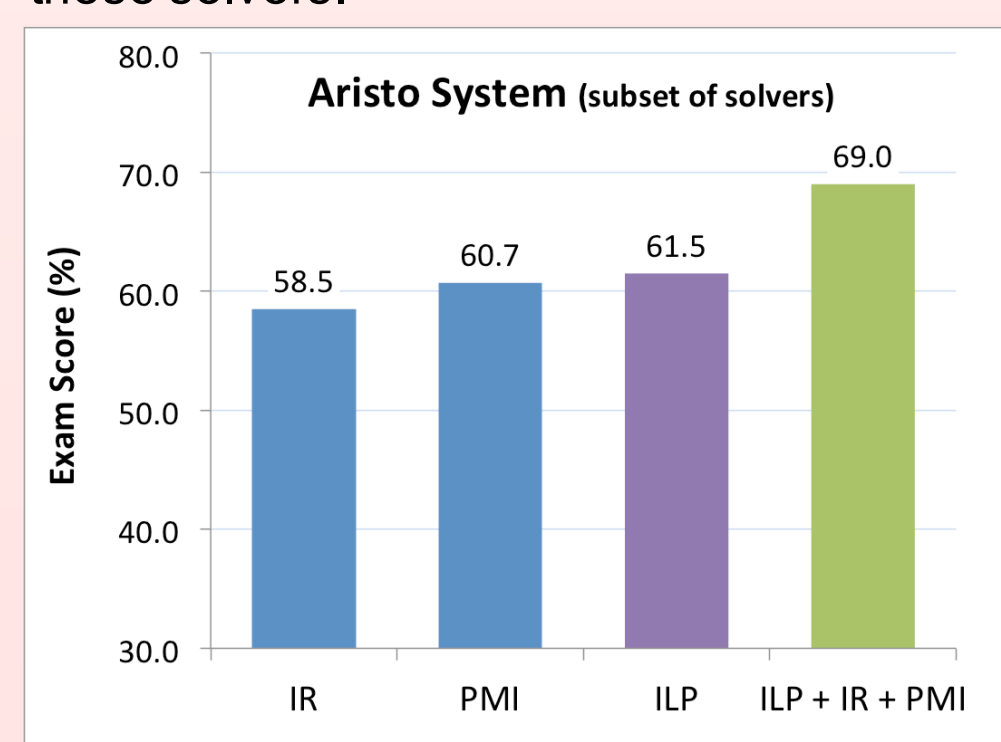
We compare the accuracy of our approach against the previous structured (MLN-based) reasoning solver and IR(tables).



1. TableLP is substantially better than IR & MLN, when given knowledge derived from the same, domain-targeted sources (6 years of exams; 95% C.I. = 9%)

## Complementary Strengths

Table I LP and IR-based methods clearly approach QA very differently. This analysis highlights the complementary strengths of these solvers.



2. Performance: % correct on 6 years of unseen questions (129 questions). The solvers ensemble performs 8-10% higher than IR baselines.

## Assessing Brittleness: Question Perturbation

How robust are approaches to simple question perturbations *that would typically make the question easier for a human*? We consider a simple, automated way to perturb each 4-way multiple-choice question using Bing:

In New York State, the longest period of daylight occurs during which month?

(A) *eastern* (B) June (C) *history* (D) *years*

Solver	Original Score (%)	% Drop with Perturbation	
		absolute	relative
IR	70.7	13.8	19.5
PMI	73.6	24.4	33.2
TableLP	85.0	<b>10.5</b>	<b>12.3</b>

3. On 1080 perturbed question of the reagents train, TableLP has the smallest drop among the solvers.

## ILP complexity, scalability

The table below summarizes various ILP and support graph statistics for TableILP, averaged across all test question

Category	Quantity	Average
ILP complexity	#variables	1043.8
	#constraints	4417.8
	#LP iterations	1348.9
Knowledge use	#rows	2.3
	#tables	1.3
Timing stats	model creation	1.9 sec
	solving the ILP	2.1 sec

4. Speed: 4 sec per question, reasoning over 140 rows across 7 tables. Contrast: 17 sec for MLN using only 1 rule per answer option!