

Question Answering via Integer Programming over Semi-Structured Knowledge

Daniel Khashabi[†], Tushar Khot[‡], Ashish Sabharwal[‡], Peter Clark[‡], Oren Etzioni[‡], Dan Roth[†]

[†]University of Illinois at Urbana-Champaign, IL, U.S.A.

{khashab2,danr}@illinois.edu

[‡]Allen Institute for Artificial Intelligence (AI2), Seattle, WA, U.S.A.

{tushark,ashishs,peterc,orene}@allenai.org

Abstract

Answering science questions posed in natural language is an important AI challenge. Answering such questions often requires non-trivial inference and knowledge that goes beyond factoid retrieval. Yet, most systems for this task are based on relatively shallow Information Retrieval (IR) and statistical correlation techniques operating on large unstructured corpora. We propose a structured inference system for this task, formulated as an Integer Linear Program (ILP), that answers natural language questions using a semi-structured knowledge base derived from text, including questions requiring multi-step inference and a combination of multiple facts. On a dataset of real, unseen science questions, our system significantly outperforms (+14%) the best previous attempt at structured reasoning for this task, which used Markov Logic Networks (MLNs). It also improves upon a previous ILP formulation by 17.7%. When combined with unstructured inference methods, the ILP system significantly boosts overall performance (+10%). Finally, we show our approach is substantially more robust to a simple answer perturbation compared to statistical correlation methods.

1 Introduction

Answering questions posed in natural language is a fundamental AI task, with a large number of impressive QA systems built over the years. Today’s Internet search engines, for instance, can successfully retrieve *factoid* style answers to many natural language queries by efficiently searching the Web. Information Retrieval (IR) systems work under the assumption that answers to many questions of interest are often explicitly stated somewhere [Kwok *et al.*, 2001], and all one needs, in principle, is access to a sufficiently large corpus. Similarly, statistical correlation based methods, such as those using Pointwise Mutual Information or PMI [Church and Hanks, 1989], work under the assumption that many questions can be answered by looking for words that tend to co-occur with the question words in a large corpus.

While both of these approaches help identify correct answers, they are not suitable for questions requiring reasoning,

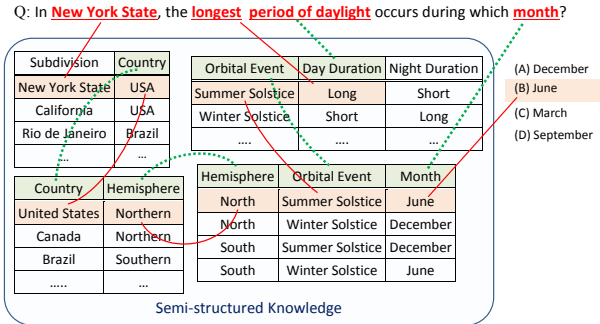


Figure 1: TableILP searches for the best support graph (chains of reasoning) connecting the question to an answer, in this case June. Constraints on the graph define what constitutes valid support and how to score it (Section 3.3).

such as chaining together multiple facts in order to arrive at a conclusion. Arguably, such reasoning is a cornerstone of human intelligence, and is a key ability evaluated by standardized science exams given to students. For example, consider a question from the NY Regents 4th Grade Science Test:

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

We would like a QA system that, even if the answer is not explicitly stated in a document, can *combine basic scientific and geographic facts* to answer the question, e.g., New York is in the north hemisphere; the longest day occurs during the summer solstice; and the summer solstice in the north hemisphere occurs in June (hence the answer is June). Figure 1 illustrates how our system approaches this, with the highlighted support graph representing its line of reasoning.

Further, we would like the system to be *robust under simple perturbations*, such as changing New York to New Zealand (in the southern hemisphere) or changing an incorrect answer option to an irrelevant word such as “last” that happens to have high co-occurrence with the question text.

To this end, we propose a structured reasoning system, called TableILP, that operates over a semi-structured knowledge base derived from text and answers questions by chaining multiple pieces of information and combining parallel

evidence.¹ The knowledge base consists of *tables*, each of which is a collection of instances of an n -ary relation defined over natural language phrases. E.g., as illustrated in Figure 1, a simple table with schema (*country, hemisphere*) might contain the instance (*United States, Northern*) while a ternary table with schema (*hemisphere, orbital event, month*) might contain (*North, Summer Solstice, June*). TableILP treats lexical constituents of the question Q , as well as cells of potentially relevant tables T , as nodes in a large graph $\mathcal{G}_{Q,T}$, and attempts to find a subgraph G of $\mathcal{G}_{Q,T}$ that “best” supports an answer option. The notion of best support is captured via a number of structural and semantic constraints and preferences, which are conveniently expressed in the Integer Linear Programming (ILP) formalism. We then use an off-the-shelf ILP optimization engine called SCIP [Achterberg, 2009] to determine the best supported answer for Q .

Following a recently proposed AI challenge [Clark, 2015], we evaluate TableILP on unseen elementary-school science questions from standardized tests. Specifically, we consider a challenge set [Clark *et al.*, 2016] consisting of all non-diagram multiple choice questions from 6 years of NY Regents 4th grade science exams. In contrast to a state-of-the-art structured inference method [Khot *et al.*, 2015] for this task, which used Markov Logic Networks (MLNs) [Richardson and Domingos, 2006], TableILP achieves a significantly (+14% absolute) higher test score. This suggests that a combination of a rich and fine-grained constraint language, namely ILP, even with a publicly available solver is more effective in practice than various MLN formulations of the task. Further, while the scalability of the MLN formulations was limited to very few (typically one or two) selected science rules at a time, our approach easily scales to hundreds of relevant scientific facts. It also complements the kind of questions amenable to IR and PMI techniques, as is evidenced by the fact that a combination (trained using simple Logistic Regression [Clark *et al.*, 2016]) of TableILP with IR and PMI results in a significant (+10% absolute) boost in the score compared to IR alone.

Our ablation study suggests that combining facts from multiple tables or multiple rows within a table plays an important role in TableILP’s performance. We also show that TableILP benefits from the table structure, by comparing it with an IR system using the same knowledge (the table rows) but expressed as simple sentences; TableILP scores significantly (+10%) higher. Finally, we demonstrate that our approach is robust to a simple perturbation of incorrect answer options: while the simple perturbation results in a relative drop of 20% and 33% in the performance of IR and PMI methods, respectively, it affects TableILP’s performance by only 12%.

2 Related Work

Clark *et al.* [2016] proposed an ensemble approach for the science QA task, demonstrating the effectiveness of a combination of information retrieval, statistical association,

¹A preliminary version of our ILP model was used in the ensemble solver of Clark *et al.* [2016]. We build upon this earlier ILP formulation, providing further details and incorporating additional syntactic and semantic constraints that improve the score by 17.7%.

rule-based reasoning, and an ILP solver operating on semi-structured knowledge. Our ILP system extends their model with additional constraints and preferences (e.g., semantic relation matching), substantially improving QA performance.

A number of systems have been developed for answering factoid questions with short answers (e.g., “What is the capital of France?”) using document collections or databases (e.g., Freebase [Bollacker *et al.*, 2008], NELL [Carlson *et al.*, 2010]), for example [Brill *et al.*, 2002; Fader *et al.*, 2014; Ferrucci *et al.*, 2010; Ko *et al.*, 2007; Yih *et al.*, 2014; Yao and Van Durme, 2014; Zou *et al.*, 2014]. However, many science questions have answers that are not explicitly stated in text, and instead require combining information together. Conversely, while there are AI systems for formal scientific reasoning (e.g., [Gunning *et al.*, 2010; Novak, 1977]), they require questions to be posed in logic or restricted English. Our goal here is a system that operates between these two extremes, able to combine information while still operating with natural language.

The task of Recognizing Textual Entailment (RTE) [Dagan *et al.*, 2010; 2013] is also closely related, as QA can be cast as entailment (Does *corpus* entail *question+answer*? [Bentivogli *et al.*, 2008]). However, RTE has primarily focused on the task of *linguistic equivalence*, and has not addressed questions where some form of scientific reasoning is required. Recent work on Natural Logic [Angeli and Manning, 2014; MacCartney, 2009] has extended RTE to account for the *logical structure* within language. Our work can be seen as going one step further, to add a layer of structured reasoning on top of this; in fact, we use an RTE engine as a basic subroutine for comparing individual table cells in our ILP formulation.

ILP based discrete optimization has been successful in several NLP tasks [Roth and Yih, 2004; Chang *et al.*, 2010; Berant *et al.*, 2010; Srikumar and Roth, 2011; Goldwasser and Roth, 2011]. While our ILP formulation also operates on natural language text, our focus is on the use of a specific semi-structured table representation for QA. Cohen [2000] studied tables with natural language text requiring soft matching, with a focus on efficiently computing the top few candidates given a database query. In contrast, our system, given a natural language question, (implicitly) seeks to generate a query that produces the most supported answer.

3 QA as Subgraph Optimization

We begin with our knowledge representation formalism, followed by our treatment of QA as an optimal subgraph selection problem over such knowledge, and then briefly describe our ILP model for subgraph selection.

3.1 Semi-Structured Knowledge as Tables

We use semi-structured knowledge represented in the form of n -ary predicates over natural language text [Clark *et al.*, 2016]. Formally, a k -column table in the knowledge base is a predicate $r(x_1, x_2, \dots, x_k)$ over strings, where each string is a (typically short) natural language phrase. The column headers capture the table schema, akin to a relational database. Each row in the table corresponds to an instance of this predicate. For example, a simple country-hemisphere table represents the binary predicate $r_{\text{ctry-hems}}(c, h)$ with instances such

as (Australia, Southern) and (Canada, Northern). Since table content is specified in natural language, the same entity is often represented differently in different tables, posing an additional inference challenge.

Although techniques for constructing this knowledge base are outside the scope of this paper, we briefly mention them. Tables were constructed using a mixture of manual and semi-automatic techniques. First, the table schemas were manually defined based on the syllabus, study guides, and training questions. Tables were then populated both manually and semi-automatically using IKE [Dalvi *et al.*, 2016], a table-building tool that performs interactive, bootstrapped relation extraction over a corpus of science text. In addition, to augment these tables with the broad knowledge present in study guides that doesn't always fit the manually defined table schemas, we ran an Open IE [Banko *et al.*, 2007] pattern-based *subject-verb-object* (SVO) extractor from Clark *et al.* [2014] over several science texts to populate three-column Open IE tables. Methods for further automating table construction are under development.

3.2 QA as a Search for Desirable Support Graphs

We treat question answering as the task of pairing the question with an answer such that this pair has the best support in the knowledge base, measured in terms of the strength of a “support graph” defined as follows.

Given a multiple choice question Q and tables T , we can define a labeled undirected graph $\mathcal{G}_{Q,T}$ over nodes \mathcal{V} and edges \mathcal{E} as follows. We first split Q into lexical constituents (e.g., non-stopword tokens, or chunks) $\mathbf{q} = \{q_\ell\}$ and answer options $\mathbf{a} = \{a_m\}$. For each table T_i , we consider its cells $\mathbf{t} = \{t_{ijk}\}$ as well as column headers $\mathbf{h} = \{h_{ik}\}$. The nodes of $\mathcal{G}_{Q,T}$ are then $\mathcal{V} = \mathbf{q} \cup \mathbf{a} \cup \mathbf{t} \cup \mathbf{h}$. For presentation purposes, we will equate a graph node with the lexical entity it represents (such as a table cell or a question constituent). The undirected edges of $\mathcal{G}_{Q,T}$ are $\mathcal{E} = ((\mathbf{q} \cup \mathbf{a}) \times (\mathbf{t} \cup \mathbf{h})) \cup (\mathbf{t} \times \mathbf{t}) \cup (\mathbf{h} \times \mathbf{h})$ excluding edges both whose endpoints are within a single table.

Informally, an edge denotes (soft) equality between a question or answer node and a table node, or between two table nodes. To account for lexical variability (e.g., that *tool* and *instrument* are essentially equivalent) and generalization (e.g., that a *dog* is an *animal*), we replace string equality with a phrase-level entailment or similarity function $w : \mathcal{E} \rightarrow [0, 1]$ that labels each edge $e \in \mathcal{E}$ with an associated score $w(e)$. We use entailment scores (directional) from \mathbf{q} to $\mathbf{t} \cup \mathbf{h}$ and from $\mathbf{t} \cup \mathbf{h}$ to \mathbf{a} , and similarity scores (symmetric) between two nodes in \mathbf{t} .² In the special case of column headers across two tables, the score is (manually) set to either 0 or 1, indicating whether this corresponds to a meaningful join.

Intuitively, we would like the support graph for an answer option to be connected, and to include nodes from the ques-

²In our evaluations, w for entailment is a simple WordNet-based [Miller, 1995] function that computes the best word-to-word alignment between phrases, scores these alignments using WordNet's hypernym and synonym relations normalized using relevant word-sense frequency, and returns the weighted sum of the scores. w for similarity is the maximum of the entailment score in both directions. Alternative definitions for these functions may also be used.

ELEMENT	DESCRIPTION
T_i	table i
h_{ik}	header of the k -th column of i -th table
t_{ijk}	cell in row j and column k of i -th table
r_{ij}	row j of i -th table
ℓ_{ik}	column k of i -th table
q_ℓ	ℓ -th lexical constituent of the question Q
a_m	m -th answer option

Table 1: Notation for the ILP formulation.

tion, the answer option, and at least one table. Since each table row represents a coherent piece of information but cells within a row do not have any edges in $\mathcal{G}_{Q,T}$ (the same holds also for cells and the corresponding column headers), we use the notion of an augmented subgraph to capture the underlying table structure. Let $G = (V, E)$ be a subgraph of $\mathcal{G}_{Q,T}$. The *augmented subgraph* G^+ is formed by adding to G edges (v_1, v_2) such that v_1 and v_2 are in V and they correspond to either the same row (possibly the header row) of a table in T or to a cell and the corresponding column header.

Definition 1. A *support graph* $G = G(Q, T, a_m)$ for a question Q , tables T , and an answer option a_m is a subgraph (V, E) of $\mathcal{G}_{Q,T}$ with the following basic properties:

1. $V \cap \mathbf{a} = \{a_m\}$, $V \cap \mathbf{q} \neq \phi$, $V \cap \mathbf{t} \neq \phi$;
2. $w(e) > 0$ for all $e \in E$;
3. if $e \in E \cap (\mathbf{t} \times \mathbf{t})$ then there exists a corresponding $e' \in E \cap (\mathbf{h} \times \mathbf{h})$ involving the same columns; and
4. the augmented subgraph G^+ is connected.

A support graph thus connects the question constituents to a unique answer option through table cells and (optionally) table headers corresponding to the aligned cells. A given question and tables give rise to a large number of possible support graphs, and the role of the inference process will be to choose the “best” one under a notion of *desirable* support graphs developed next. We do this through a number of additional structural and semantic properties; the more properties the support graph satisfies, the more desirable it is.

3.3 ILP Formulation

We model the above support graph search for QA as an ILP optimization problem, i.e., as maximizing a linear objective function over a finite set of variables, subject to a set of linear inequality constraints. A summary of the model is given below.³ We note that the ILP objective and constraints aren't tied to the particular domain of evaluation; they represent general properties that capture what constitutes a well-supported answer for a given question.

Table 1 summarizes the notation for various elements of the problem, such as t_{ijk} for cell (j, k) of table i . All core variables in the ILP model are binary, i.e., have domain $\{0, 1\}$. For each element, the model has a unary variable capturing whether this element is part of the support graph G , i.e., it is “active”. For instance, row r_{ij} is active if at least one cell in

³Details of the ILP model are omitted and may be found in the extended version of this paper [Khashabi *et al.*, 2016].

BASIC PAIRWISE ACTIVITY VARIABLES	
$y(t_{ijk}, t_{i'j'k'})$	cell to cell
$y(t_{ijk}, q_\ell)$	cell to question constituent
$y(h_{ik}, q_\ell)$	header to question constituent
$y(t_{ijk}, a_m)$	cell to answer option
$y(h_{ik}, a_m)$	header to answer option
$y(\ell_{ik}, a_m)$	column to answer option
$y(T_i, a_m)$	table to answer option
$y(\ell_{ik}, \ell_{ik'})$	column to column relation
HIGH-LEVEL UNARY VARIABLES	
$x(T_i)$	active table
$x(r_{ij})$	active row
$x(\ell_{ik})$	active column
$x(h_{ik})$	active column header
$x(q_\ell)$	active question constituent
$x(a_m)$	active answer option

Table 2: Variables used for defining the optimization problem for TableILP solver. All variables have domain $\{0, 1\}$.

row j of table i is in G . The model also has pairwise ‘‘alignment’’ variables, capturing edges of $\mathcal{G}_{Q,T}$. The alignment variable for an edge e in $\mathcal{G}_{Q,T}$ is associated with the corresponding weight $w(e)$, and captures whether e is included in G . To improve efficiency, we create a pairwise variable for e only if $w(e)$ is larger than a certain threshold. These unary and pairwise variables are then used to define various types of constraints and preferences, as discussed next.

To make the definitions clearer, we introduce all basic variables used in our optimization in Table 2, and will use them later to define constraints explicitly. We use the notation $x(\cdot)$ to refer to a unary variable parameterized by a single element of the optimization, and $y(\cdot, \cdot)$ to refer to a pairwise variable parameterized by a pair of elements. Unary variables represent the presence of a specific element as a node in the support graph G . For example $x(T_i) = 1$ if and only if the table T_i is active in G . Similarly, $y(t_{ijk}, q_\ell) = 1$ if and only if the corresponding edge is present in G , which we alternatively refer to as an *alignment* between cell (j, k) of table i and the ℓ -th constituent of the question.

As previously mentioned, in practice we do not create all possible pairwise variables. Instead we choose the pairs alignment score $w(e)$ exceeds a pre-set threshold. For example, we create $y(t_{ijk}, t_{i'j'k'})$ only if $w(t_{ijk}, t_{i'j'k'}) \geq \text{MINCELLALIGNMENT}$.⁴

The objective function is a weighted linear sum over all variables instantiated for a given question answering problem.⁵ A small set of auxiliary variables is defined for linearizing complicated constraints.

Constraints are a significant part of our model, used for imposing the desired behavior on the support graph. Due to lack of space, we discuss only a representative subset here.⁶

⁴An exhaustive list of the minimum alignment thresholds for creating pairwise variables is in Table 10 of [Khashabi *et al.*, 2016].

⁵The complete list of weights for the pairwise and unary variables are included in Table 9 of [Khashabi *et al.*, 2016].

⁶The complete list of the constraints is explained in Table 13 of [Khashabi *et al.*, 2016].

Some constraints relate variables to each other. For example, unary variables are defined through constraints that relate them to the corresponding pairwise variables. For instance, for active row variable $x(r_{ij})$, we ensure that it is 1 if and only if at least one cell in row j is active:

$$x(r_{ij}) \geq y(t_{ijk}, *) \quad \forall (t_{ijk}, *) \in \mathcal{R}_{ij}, \forall i, j, k,$$

where \mathcal{R}_{ij} is collection of pairwise variables with one end in row j of table i .

In the remainder of this section, we outline some of the important characteristics we expect in our model, and provide details of a few illustrative constraints.

Basic Lookup

Consider the following question:

Which characteristic helps a fox find food? (A) sense of smell (B) thick fur (C) long tail (D) pointed teeth

In order to answer such lookup-style questions, we generally seek a row with the highest aggregate alignment to question constituents. We achieve this by incorporating the question-table alignment variables with the alignment scores, $w(e)$, as coefficients and the active question constituents variable with a constant coefficient in the objective function. Since any additional question-table edge with a positive entailment score (even to irrelevant tables) in the support graph would result in an increase in the score, we disallow tables with alignments only to the question (or only to a choice) and add a small penalty for every table used in order to reduce noise in the support graph. We also limit the maximum number of alignments of a question constituent and table cells, in order to prevent one constituent or cell from having a large influence on the objective function and thereby the solution:

$$\sum_{(*, q_\ell) \in \mathcal{Q}_l} y(*, q_\ell) \leq \text{MAXALIGNMENTSPERQCONS}, \forall l$$

where \mathcal{Q}_l is the set of all pairwise variables with one end in the question constituent ℓ .

Parallel Evidence

For certain questions, evidence needs to be combined from multiple rows of a table. For example,

Sleet, rain, snow, and hail are forms of (A) erosion (B) evaporation (C) groundwater (D) precipitation

To answer this question, we need to combine evidence from multiple table entries from the weather terms table, (*term*, *type*), namely (sleet, precipitation), (rain, precipitation), (snow, precipitation), and (hail, precipitation). To achieve this, we allow multiple active rows in the support graph. Similar to the basic constraints, we limit the maximum number of active rows per table and add a penalty for every active row to ensure only relevant rows are considered for reasoning:

$$\sum_j x(r_{ij}) \leq \text{MAXROWSPERTABLE}, \forall i$$

To encourage only coherent parallel evidence within a single table, we limit our support graph to always use the same columns across multiple rows within a table, i.e., every active row has the active cells corresponding to the same set of columns.

Evidence Chaining

Questions requiring chaining of evidence from multiple tables, such as the example in Figure 1, are typically the most challenging in this domain. Chaining can be viewed as performing a *join* between two tables. We introduce alignments between cells across columns in pairs of tables to allow for chaining of evidence. To help minimize potential noise introduced by chaining irrelevant facts, we add a penalty for every inter-table alignment and also rely on the 0/1 weights of header-to-header edges to ensure only semantically meaningful table joins are considered.

Semantic Relation Matching

Our constraints so far have only looked at the content of the table cells, or the structure of the support graph, without explicitly considering the *semantics* of the table schema. By using alignments between the question and column headers (i.e., type information), we exploit the table schema to prefer alignments to columns relevant to the “topic” of the question. In particular, for questions of the form “which X . . .”, we prefer answers that directly entail X or are connected to cells that entail X. However, this is not sufficient for questions such as:

What is one way to change water from a liquid to a solid? (A) decrease the temperature (B) increase the temperature (C) decrease the mass (D) increase the mass

Even if we select the correct table, say $r_{\text{change-init-fin}}(c, i, f)$ that describes the initial and final states for a phase change event, both choice (A) and choice (B) would have the exact same score in the presence of table rows (increase temperature, solid, liquid) and (decrease temperature, liquid, solid). The table, however, does have the initial vs. final state structure. To capture this semantic structure, we annotate pairs of columns within certain tables with the semantic relationship present between them. In this example, we would annotate the phase change table with the relations: $\text{changeFrom}(c, i)$, $\text{changeTo}(c, f)$, and $\text{fromTo}(i, f)$.

Given such semantic relations for table schemas, we can now impose a preference towards question-table alignments that respect these relations. We associate each semantic relation with a set of linguistic patterns describing how it might be expressed in natural language. TableILP then uses these patterns to spot possible mentions of the relations in the question Q . We then add the soft constraint that for every pair of active columns in a table (with an annotated semantic relation) aligned to a pair of question constituents, there should be a valid expression of that relation in Q between those constituents. In our example, we would match the relation $\text{fromTo}(\text{liquid}, \text{solid})$ in the table to “liquid to a solid” in the question via the pattern “X to a Y” associated with $\text{fromTo}(X, Y)$, and thereby prefer aligning with the correct row (decrease temperature, liquid, solid).

4 Evaluation

We compare our approach to three existing methods, demonstrating that it outperforms the best previous structured approach [Khot *et al.*, 2015] and produces a statistically significant improvement when used in combination with IR-based

methods [Clark *et al.*, 2016]. For evaluations, we use a 2-core 2.5 GHz Amazon EC2 linux machine with 16 GB RAM.

Question Set. We use the same question set as Clark *et al.* [2016], which consists of all non-diagram multiple-choice questions from 12 years of the NY Regents 4th Grade Science exams.⁷ The set is split into 108 development questions and 129 hidden test questions based on the year they appeared in (6 years each). All numbers reported below are for the hidden test set, except for question perturbation experiments which relied on the 108 development questions.

Test scores are reported as percentages. For each question, a solver gets a score of 1 if it chooses the correct answer and $1/k$ if it reports a k -way tie that includes the correct answer. On the 129 test questions, a score difference of 9% (or 7%) is statistically significant at the 95% (or 90%, resp.) confidence interval based on the binomial exact test [Howell, 2012].

Corpora. We work with three knowledge corpora:

1. Web Corpus: This corpus contains 5×10^{10} tokens (280 GB of plain text) extracted from Web pages. It was collected by Charles Clarke at the University of Waterloo, and has been used previously by Turney [2013] and Clark *et al.* [2016]. We use it here to compute statistical co-occurrence values for the PMI solver.
2. Sentence Corpus [Clark *et al.*, 2016]: This includes sentences from the Web corpus above, as well as around 80,000 sentences from various domain-targeted sources for elementary science: a Regents study guide, CK12 textbooks (www.ck12.org), and web sentences with similar content as the course material.
3. Table Corpus (cf. Section 3.1): This includes 65 tables totaling around 5,000 rows, designed based on the development set and study guides, as well as 4 Open IE-style [Banko *et al.*, 2007] automatically generated tables totaling around 2,600 rows.⁸

4.1 Solvers

TableILP (our approach). Given a question Q , we select the top 7 tables from the Table Corpus using the the standard TF-IDF score of Q with tables treated as bag-of-words documents. For each selected table, we choose the 20 rows that overlap with Q the most. This filtering improves efficiency and reduces noise. We then generate an ILP and solve it using the open source SCIP engine [Achterberg, 2009], returning the active answer option a_m from the optimal solution. To check for ties, we disable a_m , re-solve the ILP, and compare the score of the second-best answer, if any, with that of a_m .

MLN Solver (structured inference baseline). We consider the current state-of-the-art structured reasoning method developed for this specific task by Khot *et al.* [2015]. We compare against their best performing system, namely Praline, which uses Markov Logic Networks [Richardson and Domingos, 2006] to (a) align lexical elements of the question with

⁷These are the only publicly available state-level science exams. <http://www.nysedregents.org/Grade4/Science/home.html>

⁸Table Corpus and the ILP model are available at allenai.org.

probabilistic first-order science rules and (b) to control inference. We use the entire set of 47,000 science rules from their original work, which were also derived from same domain-targeted sources as the ones used in our Sentence Corpus.

IR Solver (information retrieval baseline). We use the IR baseline by Clark *et al.* [2016], which selects the answer option that has the best matching sentence in a corpus. Specifically, for each answer option a_i , the IR solver sends $q + a_i$ as a query to a search engine (we use Lucene) on the Sentence Corpus, and returns the search engine’s score for the top retrieved sentence s , where s must have at least one non-stopword overlap with q , and at least one with a_i . The option with the highest Lucene score is returned as the answer.

PMI Solver (statistical co-occurrence baseline). We use the PMI-based approach by Clark *et al.* [2016], which selects the answer option that most frequently co-occurs with the question words in a corpus. Specifically, it extracts unigrams, bigrams, trigrams, and skip-bigrams from the question and each answer option. For a pair (x, y) of n -grams, their pointwise mutual information (PMI) [Church and Hanks, 1989] in the corpus is defined as $\log \frac{p(x,y)}{p(x)p(y)}$ where $p(x, y)$ is the co-occurrence frequency of x and y (within some window) in the corpus. The solver returns the answer option that has the largest average PMI in the Web Corpus, calculated over all pairs of question n -grams and answer option n -grams.

4.2 Results

We first compare the accuracy of our approach against the previous structured (MLN-based) reasoning solver. We also compare against IR(tables), an IR solver using table rows expressed as sentences, thus embodying an unstructured approach operating on the same knowledge as TableILP.

Solver	Test Score (%)
MLN	47.5
IR(tables)	51.2
TableILP	61.5

Table 3: TableILP significantly outperforms both the prior MLN reasoner, and IR using identical knowledge as TableILP

As Table 3 shows, among the two structured inference approaches, TableILP outperforms the MLN baseline by 14%. The preliminary ILP system reported by Clark *et al.* [2016] achieves only a score of 43.8% on this question set. Further, given the same semi-structured knowledge (i.e., the Table Corpus), TableILP is substantially (+10%) better at exploiting the structure than the IR(tables) baseline, which, as mentioned above, uses the same data expressed as sentences.

Complementary Strengths

While their overall score is similar, TableILP and IR-based methods clearly approach QA very differently. To assess whether TableILP adds any new capabilities, we considered the 50 (out of 129) questions incorrectly answered by PMI solver (ignoring tied scores). On these unseen but arguably more difficult questions, TableILP answered 27 questions correctly, achieving a score of 54% compared to the random

chance of 25% for 4-way multiple-choice questions. Results with IR solver were similar: TableILP scored 24.75 on the 52 questions incorrectly answered by IR (i.e., 47.6% accuracy).

Solver	Test Score (%)
IR	58.5
PMI	60.7
TableILP	61.5
TableILP + IR	66.1
TableILP + PMI	67.6
TableILP + IR+ PMI	69.0

Table 4: Solver combination results

This analysis highlights the complementary strengths of these solvers. Following Clark *et al.* [2016], we create an ensemble of TableILP, IR, and PMI solvers, combining their answer predictions using a simple Logistic Regression model trained on the development set. This model uses 4 features derived from each solver’s score for each answer option, and 11 features derived from TableILP’s support graphs.⁹ Table 4 shows the results, with the final combination at 69% representing a significant improvement over individual solvers.

ILP Solution Properties

Table 5 summarizes various ILP and support graph statistics for TableILP, averaged across all test questions.

The optimization model has around 50 high-level constraints, which result, on average, in around 4000 inequalities over 1000 variables. Model creation, which includes computing pairwise entailment scores using WordNet, takes 1.9 seconds on average per question, and the resulting ILP is solved by the SCIP engine in 2.1 seconds (total for all four options), using around 1,300 LP iterations for each option.¹⁰ Thus, TableILP takes only 4 seconds to answer a question using multiple rows across multiple tables (typically 140 rows in total), as compared to 17 seconds needed by the MLN solver for reasoning with four rules (one per answer option).

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

Table 5: TableILP statistics averaged across questions

While the final support graph on this question set relies mostly on a single table to answer the question, it generally combines information from more than two rows (2.3 on average) for reasoning. This suggests parallel evidence is more frequently used on this dataset than evidence chaining.

⁹Details of the 11 features may be found in the extended version [Khashabi *et al.*, 2016].

¹⁰Commercial ILP solvers (e.g., CPLEX, Gurobi) are much faster than the open-source SCIP solver we used for evaluations.

4.3 Ablation Study

To quantify the importance of various components of our system, we performed several ablation experiments, summarized in Table 6 and described next.

Solver	Test Score (%)
TableILP	61.5
No Multiple Row Inference	51.0
No Relation Matching	55.6
No Open IE Tables	52.3
No Lexical Entailment	50.5

Table 6: Ablation results for TableILP

No Multiple Row Inference: We modify the ILP constraints to limit inference to a single row (and hence a single table), thereby disallowing parallel evidence and evidence chaining (Section 3.3). This drops the performance by 10.5%, highlighting the importance of being able to combine evidence from multiple rows (which would correspond to multiple sentences in a corpus) from one or more tables.

No Relation matching: To assess the importance of considering the semantics of the table, we remove the requirement of matching the semantic relation present between columns of a table with its lexicalization in the question (Section 3.3). The 6% drop indicates TableILP relies strongly on the table semantics to ensure creating meaningful inferential chains.

No Open IE tables: To evaluate the impact of relatively unstructured knowledge from a large corpus, we removed the tables containing Open IE extractions (Section 3.2). The 9% drop in the score shows that this knowledge is important and TableILP is able to exploit it even though it has a very simple triple structure. This opens up the possibility of extending our approach to triples extracted from larger knowledge bases.

No Lexical Entailment: Finally, we test the effect of changing the alignment metric w (Section 3.2) from WordNet based scores to a simple asymmetric word-overlap measured as $score(T, H) = \frac{|T \cap H|}{|H|}$. Relying on just word-matching results in an 11% drop, which is consistent with our knowledge often being defined in terms of generalities.

4.4 Question Perturbation

One desirable property of QA systems is robustness to simple variations of a question, *especially when a variation would make the question arguably easier for humans*.

To assess this, we consider a simple, automated way to perturb each 4-way multiple-choice question: (1) query Microsoft’s Bing search engine (www.bing.com) with the question text and obtain the text snippet of the top 2,000 hits; (2) create a list of strings by chunking and tokenizing the results; (3) remove stop words and special characters, as well as any words (or their lemma) appearing in the question; (4) sort the remaining strings based on their frequency; and (5) replace the three incorrect answer options in the question with the most frequently occurring strings, thereby generating a new question. For instance:

In New York State, the longest period of daylight occurs during which month? (A) *eastern* (B) June (C) *history* (D) *years*

As in this example, the perturbations (italicized) are often not even of the correct “type”, typically making them much easier for humans. They, however, still remain difficult for solvers.

Solver	Original Score (%)	% Drop with Perturbation	
		absolute	relative
IR	70.7	13.8	19.5
PMI	73.6	24.4	33.2
TableILP	85.0	10.5	12.3

Table 7: Drop in solver scores (on the development set, rather than the hidden test set) when questions are perturbed

For each of the 108 development questions, we generate 10 new perturbed questions, using the 30 most frequently occurring words in step (5) above. While this approach can introduce new answer options that should be considered correct as well, only 3% of the questions in a random sample exhibited this behavior. Table 7 shows the performance of various solvers on the resulting 1,080 perturbed questions. As one might expect, the PMI approach suffers the most at a 33% relative drop. TableILP’s score drops as well (since answer type matching isn’t perfect), but only by 12%, attesting to its higher resilience to simple question variation.

5 Conclusion

Answering real science questions is a challenging task because they are posed in natural language, require extensive domain knowledge, and often require combining multiple facts together. We presented TableILP, a system that can answer such questions, using a semi-structured knowledge base. We treat QA as a subgraph selection problem and then formulate this as an ILP optimization. Most importantly, this formulation allows multiple, semi-formally expressed facts to be combined to answer questions, a capability outside the scope of IR-based QA systems. In our experiments, this approach significantly outperforms both the previous best attempt at structured reasoning for this task, and an IR engine provided with the same knowledge. It also significantly boosts performance when combined with unstructured methods (IR and PMI). These results suggest that the approach is both viable and promising for natural language question answering.

Acknowledgments

D.K. is in part supported by AI2 and Google. The authors would like to thank Christos Christodoulopoulos, Sujay Jauhar, Sam Skjonsberg, and the Aristo Team at AI2 for invaluable discussions and insights.

References

[Achterberg, 2009] Tobias Achterberg. SCIP: solving constraint integer programs. *Math. Prog. Computation*, 1(1):1–41, 2009.

- [Angeli and Manning, 2014] Gabor Angeli and Christopher D. Manning. NaturalLI: Natural logic inference for common sense reasoning. In *EMNLP*, 2014.
- [Banko *et al.*, 2007] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. Open information extraction from the web. In *IJCAI*, 2007.
- [Bentivogli *et al.*, 2008] Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The sixth pascal recognizing textual entailment challenge. In *TAC*, 2008.
- [Berant *et al.*, 2010] Jonathan Berant, Ido Dagan, and Jacob Goldberger. Global learning of focused entailment graphs. In *ACL*, pages 1220–1229, 2010.
- [Bollacker *et al.*, 2008] Kurt D. Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD*, 2008.
- [Brill *et al.*, 2002] Eric Brill, Susan Dumais, and Michele Banko. An analysis of the AskMSR question-answering system. In *Proceedings of EMNLP*, pages 257–264, 2002.
- [Carlson *et al.*, 2010] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam R. Hruschka Jr., and Tom M. Mitchell. Toward an architecture for never-ending language learning. In *AAAI*, 2010.
- [Chang *et al.*, 2010] Ming-Wei Chang, Dan Goldwasser, Dan Roth, and Vivek Srikumar. Discriminative learning over constrained latent representations. In *2010 NAACL*, pages 429–437, 2010.
- [Church and Hanks, 1989] Kenneth Ward Church and Patrick Hanks. Word association norms, mutual information and lexicography. In *27th ACL*, pages 76–83, 1989.
- [Clark *et al.*, 2014] Peter Clark, Niranjan Balasubramanian, Sumithra Bhakthavatsalam, Kevin Humphreys, Jesse Kinkade, Ashish Sabharwal, and Oyvind Tafjord. Automatic construction of inference-supporting knowledge bases. In *4th AKBC Workshop*, Montreal, Canada, December 2014.
- [Clark *et al.*, 2016] Peter Clark, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Turney, and Daniel Khashabi. Combining retrieval, statistics, and inference to answer elementary science questions. In *30th AAAI*, 2016.
- [Clark, 2015] Peter Clark. Elementary school science and math tests as a driver for AI: take the Aristo challenge! In *29th AAAI/IAAI*, pages 4019–4021, Austin, TX, 2015.
- [Cohen, 2000] William W. Cohen. Data integration using similarity joins and a word-based information representation language. *ACM Transactions on Information Systems*, 18(3):288–321, 2000.
- [Dagan *et al.*, 2010] Ido Dagan, Bill Dolan, Bernardo Magnini, and Dan Roth. Recognizing textual entailment: Rational, evaluation and approaches—erratum. *Natural Language Engineering*, 16(01):105–105, 2010.
- [Dagan *et al.*, 2013] Ido Dagan, Dan Roth, Mark Sammons, and Fabio Massimo Zanzotto. Recognizing textual entailment: Models and applications. *Synthesis Lectures on Human Language Technologies*, 6(4):1–220, 2013.
- [Dalvi *et al.*, 2016] Bhavana Dalvi, Sumithra Bhakthavatsalam, and Peter Clark. IKE - an interactive tool for knowledge extraction. In *5th AKBC Workshop*, 2016.
- [Fader *et al.*, 2014] Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. Open question answering over curated and extracted knowledge bases. In *Proceedings of SIGKDD*, pages 1156–1165, 2014.
- [Ferrucci *et al.*, 2010] David Ferrucci, Eric Brown, Jennifer Chu-Carroll, James Fan, David Gondek, Aditya A Kalyanpur, Adam Lally, J William Murdock, Eric Nyberg, John Prager, et al. Building Watson: An overview of the DeepQA project. *AI Magazine*, 31(3):59–79, 2010.
- [Goldwasser and Roth, 2011] D. Goldwasser and D. Roth. Learning from natural instructions. In *IJCAI*, 2011.
- [Gunning *et al.*, 2010] D. Gunning, V. Chaudhri, P. Clark, K. Barker, J. Chaw, and M. Greaves. Project Halo update - progress toward digital Aristotle. *AI Magazine*, 31(3), 2010.
- [Howell, 2012] David Howell. *Statistical methods for psychology*. Cengage Learning, 2012.
- [Khashabi *et al.*, 2016] Daniel Khashabi, Tushar Khot, Ashish Sabharwal, Peter Clark, Oren Etzioni, and Dan Roth. Question answering via integer programming over semi-structured knowledge (extended version). Available at <http://allenai.org/papers.html>, 2016.
- [Khot *et al.*, 2015] Tushar Khot, Niranjan Balasubramanian, Eric Gribkoff, Ashish Sabharwal, Peter Clark, and Oren Etzioni. Exploring Markov logic networks for question answering. In *2015 EMNLP*, Lisbon, Portugal, September 2015.
- [Ko *et al.*, 2007] Jeongwoo Ko, Eric Nyberg, and Luo Si. A probabilistic graphical model for joint answer ranking in question answering. In *Proceedings of SIGIR*, pages 343–350, 2007.
- [Kwok *et al.*, 2001] Cody C. T. Kwok, Oren Etzioni, and Daniel S. Weld. Scaling question answering to the web. In *WWW*, 2001.
- [MacCartney, 2009] Bill MacCartney. *Natural Language Inference*. PhD thesis, 2009.
- [Miller, 1995] George Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [Novak, 1977] G. Novak. Representations of knowledge in a program for solving physics problems. In *IJCAI-77*, 1977.
- [Richardson and Domingos, 2006] Matthew Richardson and Pedro Domingos. Markov logic networks. *Machine learning*, 62(1–2):107–136, 2006.
- [Roth and Yih, 2004] D. Roth and W. Yih. A linear programming formulation for global inference in natural language tasks. In *CoNLL*, pages 1–8. *ACL*, 2004.
- [Srikumar and Roth, 2011] V. Srikumar and D. Roth. A joint model for extended semantic role labeling. In *EMNLP*, Edinburgh, Scotland, 2011.
- [Turney, 2013] Peter D Turney. Distributional semantics beyond words: Supervised learning of analogy and paraphrase. *TACL*, 1:353–366, 2013.
- [Yao and Van Durme, 2014] Xuchen Yao and Benjamin Van Durme. Information extraction over structured data: Question answering with Freebase. In *52nd ACL*, 2014.
- [Yih *et al.*, 2014] Wen-tau Yih, Xiaodong He, and Christopher Meek. Semantic parsing for single-relation question answering. In *ACL (2)*, pages 643–648, 2014.
- [Zou *et al.*, 2014] Lei Zou, Ruizhe Huang, Haixun Wang, Jeffer Xu Yu, Wenqiang He, and Dongyan Zhao. Natural language question answering over RDF: a graph data driven approach. In *SIGMOD*, pages 313–324, 2014.