

COMP 790.139 (Fall 2017)

Natural Language Processing

Lecture 6: Semantic Parsing 2; Question Answering



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

Mohit Bansal

(various slides adapted/borrowed from courses by Dan Klein, JurafskyMartin-SLP3, others)

Announcements

- ▶ Coding-HW1 (on word vector training+evaluation_+visualization) was due Oct5 midnight!
- ▶ Midterm project presentation next week (look for details in email).
- ▶ Coding-HW2 will be announced soon!

SRL and Semantic Parsing 2

(AMR, Neural Models, etc.)

SRL Features

Features

Headword of constituent

Examiner

Headword POS

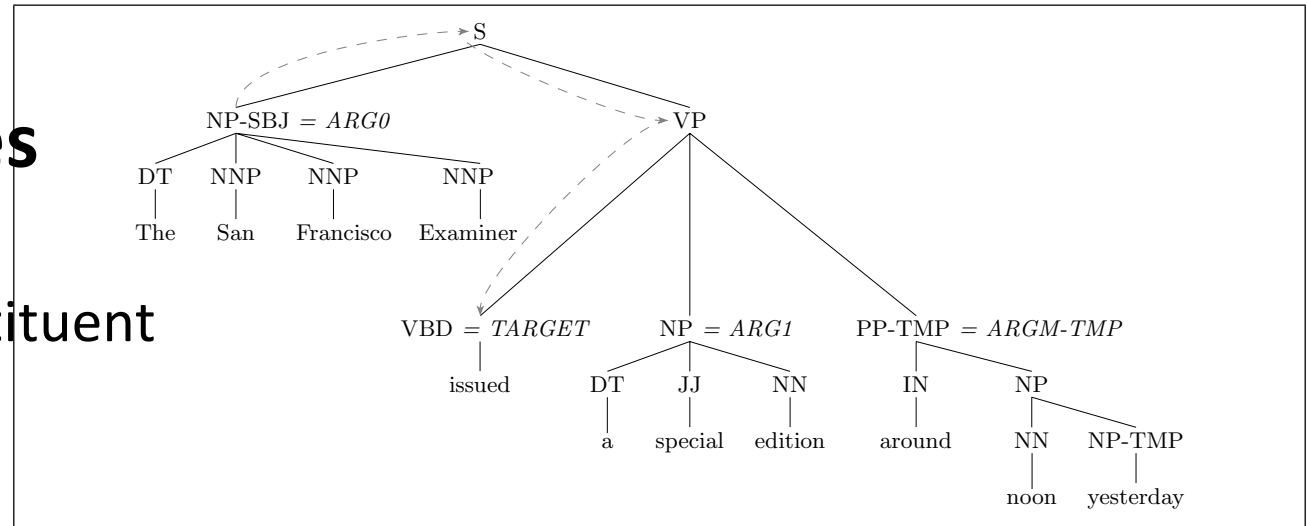
NNP

Voice of the clause

Active

Subcategorization of pred

VP -> VBD NP PP



Named Entity type of constit

ORGANIZATION

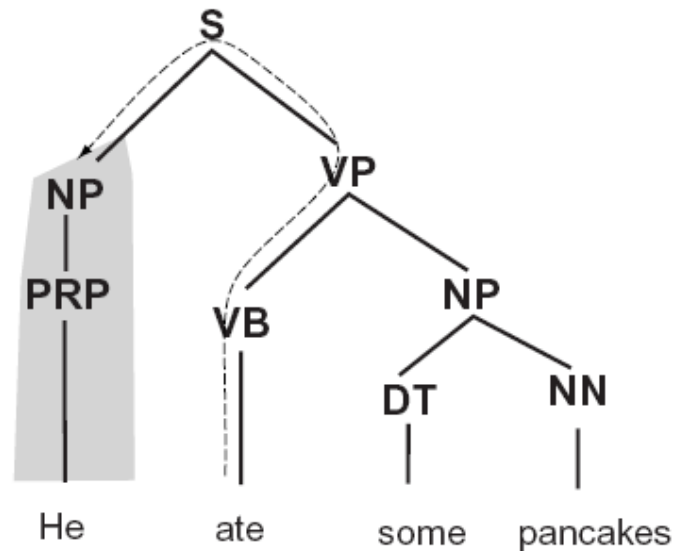
First and last words of constit

The, Examiner

Linear position, clause re: predicate

before

Path-based Features for SRL



<i>Path</i>	<i>Description</i>
VB↑VP↓PP	PP argument/adjunct
VB↑VP↑S↓NP	subject
VB↑VP↓NP	object
VB↑VP↑VP↑S↓NP	subject (embedded VP)
VB↑VP↓ADVP	adverbial adjunct
NN↑NP↑NP↓PP	prepositional complement of noun

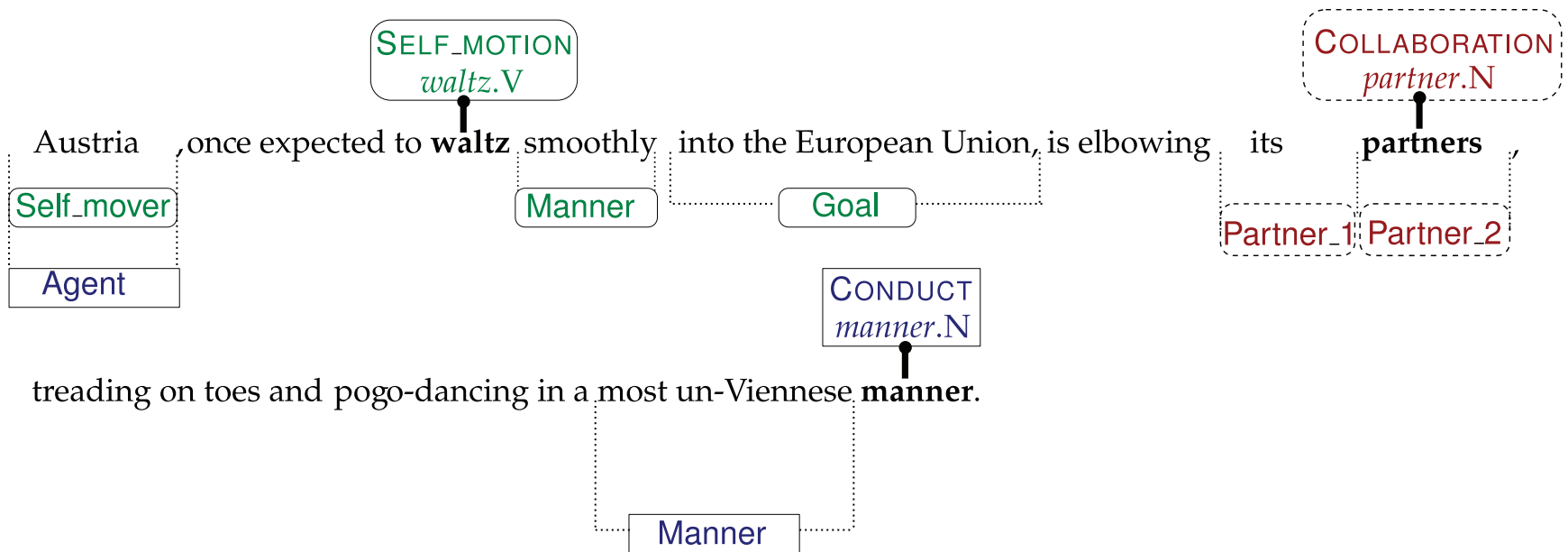
Some SRL Results

- ▶ So major feature categories in traditional feature-based SRL models were:
 - ▶ Headword, syntactic type, case, etc. of candidate node/constituent
 - ▶ Linear and tree path from predicate target to node
 - ▶ Active vs. passive voice
 - ▶ Second order and higher order features
- ▶ Accuracy for such feature-based SRL models then highly depends on accuracy of underlying parse tree!
 - ▶ So quite high SRL results when using ground-truth parses
 - ▶ Much lower results with automatically-predicted parses!

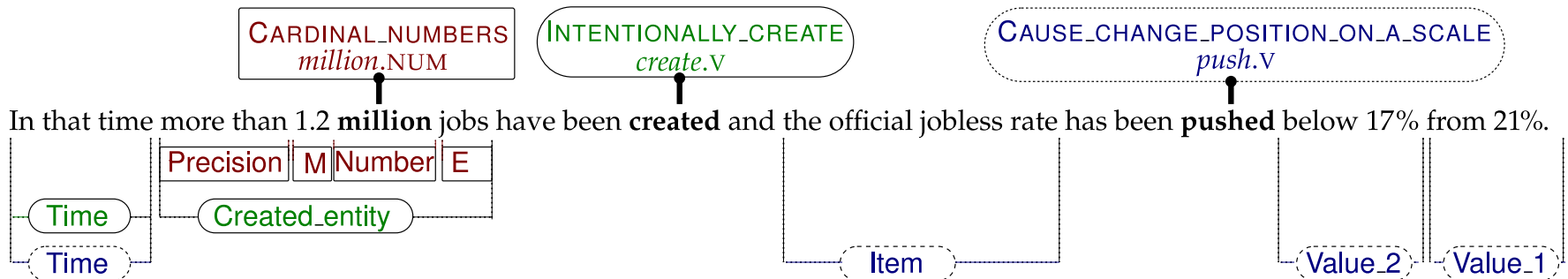
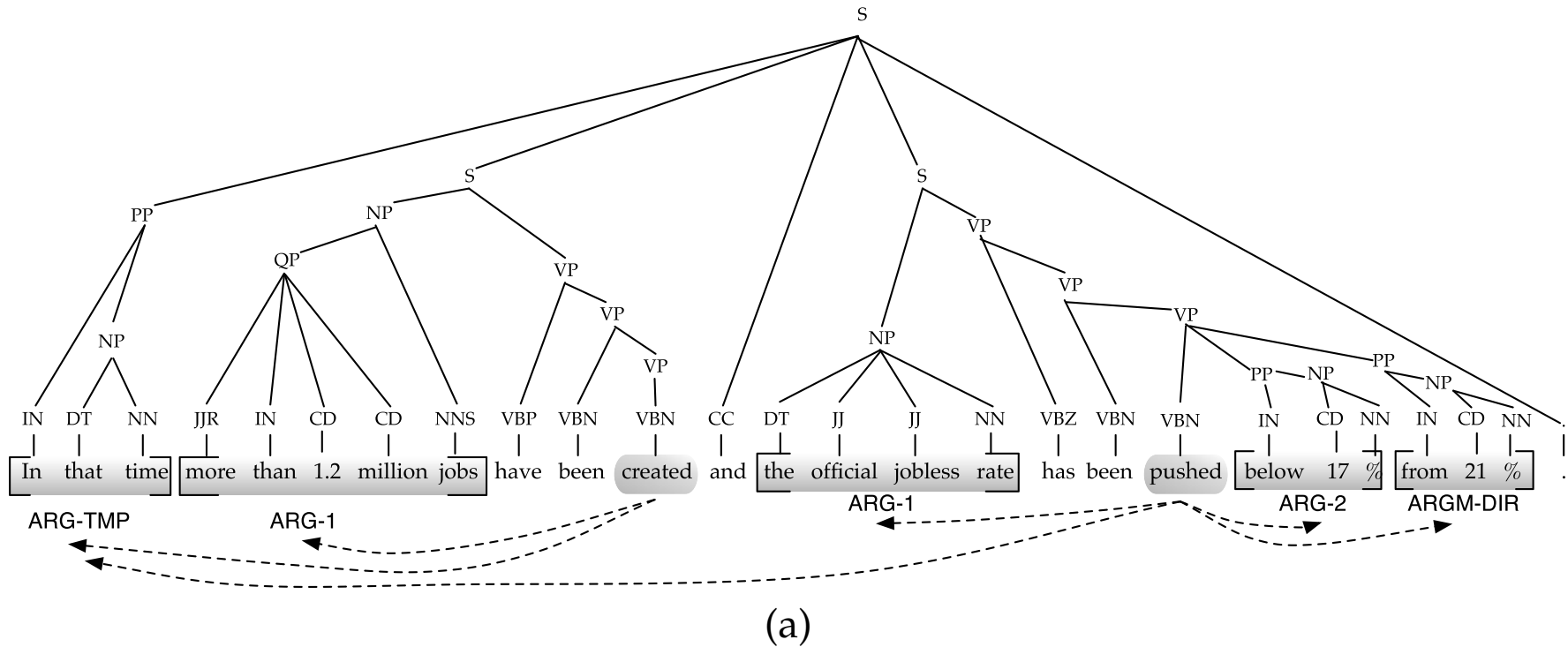
CORE		ARGM	
F1	Acc.	F1	Acc.
92.2	80.7	89.9	71.8

CORE		ARGM	
F1	Acc.	F1	Acc.
84.1	66.5	81.4	55.6

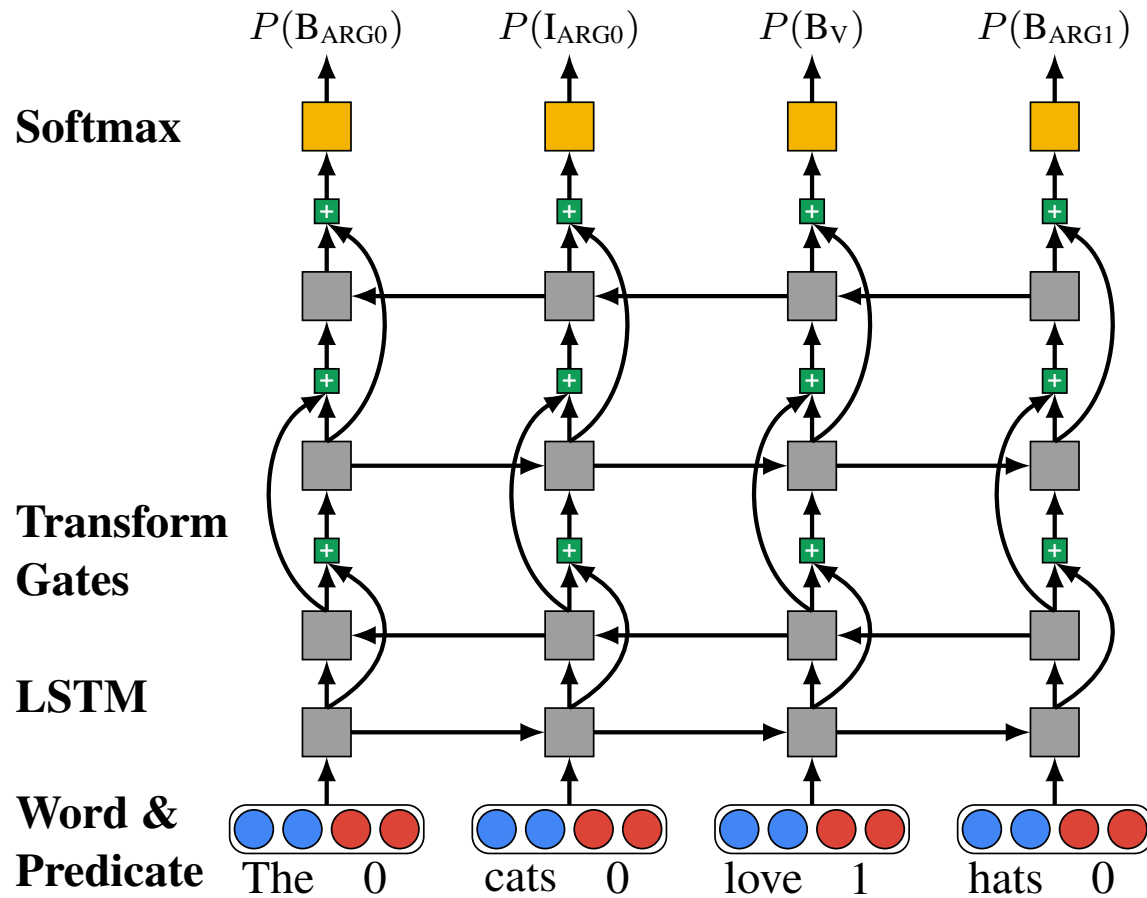
Schematic of Frame Semantics (FrameNet)



PropBank vs. FrameNet Representations



Neural SRL



Neural SRL

► CoNLL 2005 dataset:

Method	Development				WSJ Test				Brown Test				Combined
	P	R	F1	Comp.	P	R	F1	Comp.	P	R	F1	Comp.	F1
Ours (PoE)	83.1	82.4	82.7	64.1	85.0	84.3	84.6	66.5	74.9	72.4	73.6	46.5	83.2
Ours	81.6	81.6	81.6	62.3	83.1	83.0	83.1	64.3	72.9	71.4	72.1	44.8	81.6
Zhou	79.7	79.4	79.6	-	82.9	82.8	82.8	-	70.7	68.2	69.4	-	81.1
FitzGerald (Struct.,PoE)	81.2	76.7	78.9	55.1	82.5	78.2	80.3	57.3	74.5	70.0	72.2	41.3	-
Täckström (Struct.)	81.2	76.2	78.6	54.4	82.3	77.6	79.9	56.0	74.3	68.6	71.3	39.8	-
Toutanova (Ensemble)	-	-	78.6	58.7	81.9	78.8	80.3	60.1	-	-	68.8	40.8	-
Punyakanok (Ensemble)	80.1	74.8	77.4	50.7	82.3	76.8	79.4	53.8	73.4	62.9	67.8	32.3	77.9

► CoNLL 2012 dataset:

Method	Development				Test			
	P	R	F1	Comp.	P	R	F1	Comp.
Ours (PoE)	83.5	83.2	83.4	67.5	83.5	83.3	83.4	68.5
Ours	81.8	81.4	81.5	64.6	81.7	81.6	81.7	66.0
Zhou	-	-	81.1	-	-	-	81.3	-
FitzGerald (Struct.,PoE)	81.0	78.5	79.7	60.9	81.2	79.0	80.1	62.6
Täckström (Struct.)	80.5	77.8	79.1	60.1	80.6	78.2	79.4	61.8
Pradhan (revised)	-	-	-	-	78.5	76.6	77.5	55.8

Neural SRL

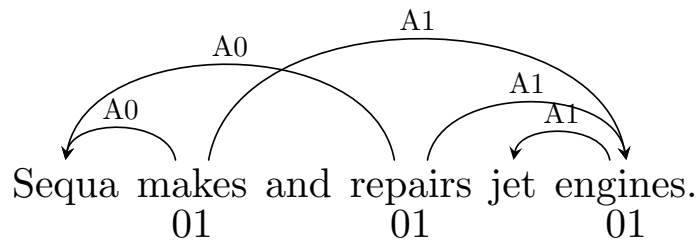


Figure 1: A semantic dependency graph.

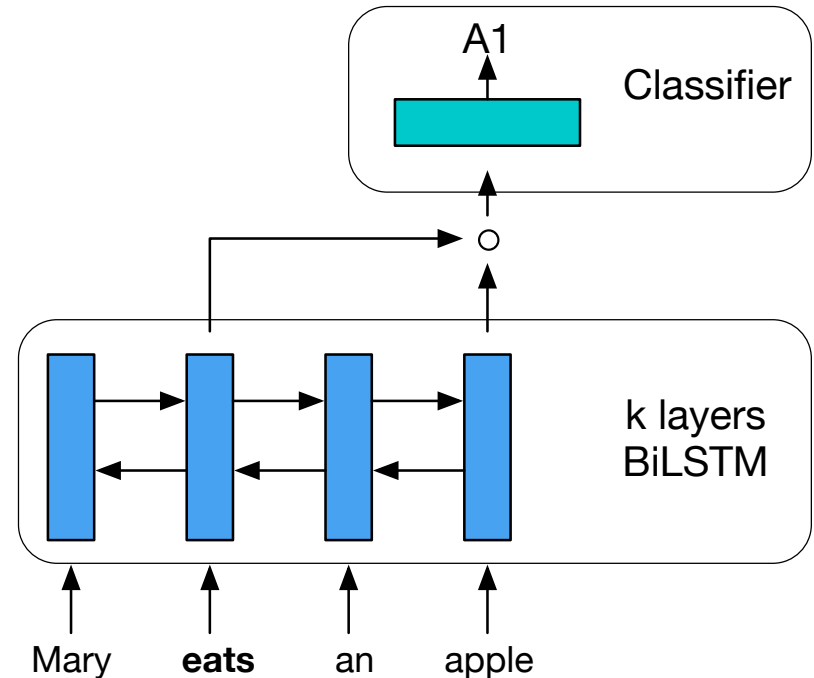


Figure 2: Predicting an argument and its label with an LSTM encoder.

Neural Semantic Parsing

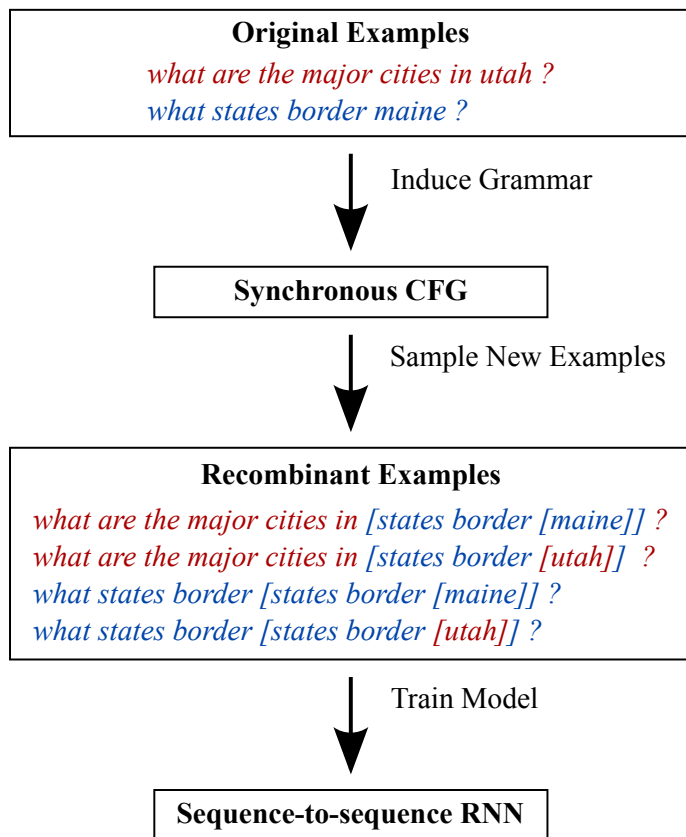


Figure 1: An overview of our system. Given a dataset, we induce a high-precision synchronous context-free grammar. We then sample from this grammar to generate new “recombinant” examples, which we use to train a sequence-to-sequence RNN.

GEO

x: “what is the population of iowa ?”

```
y: _answer ( NV , (
  _population ( NV , V1 ) , _const (
    V0 , _stateid ( iowa ) ) ) )
```

ATIS

x: “can you list all flights from chicago to milwaukee”

```
y: ( _lambda $0 e ( _and
  ( _flight $0 )
  ( _from $0 chicago : _ci )
  ( _to $0 milwaukee : _ci ) ) )
```

Overnight

x: “when is the weekly standup”

```
y: ( call listValue ( call
  getProperty meeting.weekly_standup
  ( string start_time ) ) )
```

Figure 2: One example from each of our domains. We tokenize logical forms as shown, thereby casting semantic parsing as a sequence-to-sequence task.

Neural Semantic Parsing

Examples

```
("what states border texas ?",  
answer(NV, (state(V0), next_to(V0, NV), const(V0, stateid(texas)))))  
("what is the highest mountain in ohio ?",  
answer(NV, highest(V0, (mountain(V0), loc(V0, NV), const(V0, stateid(ohio))))))
```

Rules created by ABSENTITIES

```
ROOT → ⟨ "what states border STATEID ?",  
    answer(NV, (state(V0), next_to(V0, NV), const(V0, stateid(STATEID)))) ⟩  
STATEID → ⟨ "texas", texas ⟩  
ROOT → ⟨ "what is the highest mountain in STATEID ?",  
    answer(NV, highest(V0, (mountain(V0), loc(V0, NV),  
        const(V0, stateid(STATEID)))) ⟩  
STATEID → ⟨ "ohio", ohio ⟩
```

Rules created by ABSWHOLEPHRASES

```
ROOT → ⟨ "what states border STATE ?", answer(NV, (state(V0), next_to(V0, NV), STATE)) ⟩  
STATE → ⟨ "states border texas", state(V0), next_to(V0, NV), const(V0, stateid(texas)) ⟩  
ROOT → ⟨ "what is the highest mountain in STATE ?",  
    answer(NV, highest(V0, (mountain(V0), loc(V0, NV), STATE)) ⟩
```

Rules created by CONCAT-2

```
ROOT → ⟨ SENT1 </s> SENT2, SENT1 </s> SENT2 ⟩  
SENT → ⟨ "what states border texas ?",  
    answer(NV, (state(V0), next_to(V0, NV), const(V0, stateid(texas)))) ⟩  
SENT → ⟨ "what is the highest mountain in ohio ?",  
    answer(NV, highest(V0, (mountain(V0), loc(V0, NV), const(V0, stateid(ohio)))) ⟩
```

Figure 3: Various grammar induction strategies illustrated on GEO. Each strategy converts the rules of an input grammar into rules of an output grammar. This figure shows the base case where the input grammar has rules $ROOT \rightarrow \langle x, y \rangle$ for each (x, y) pair in the training dataset.

Neural Semantic Parsing

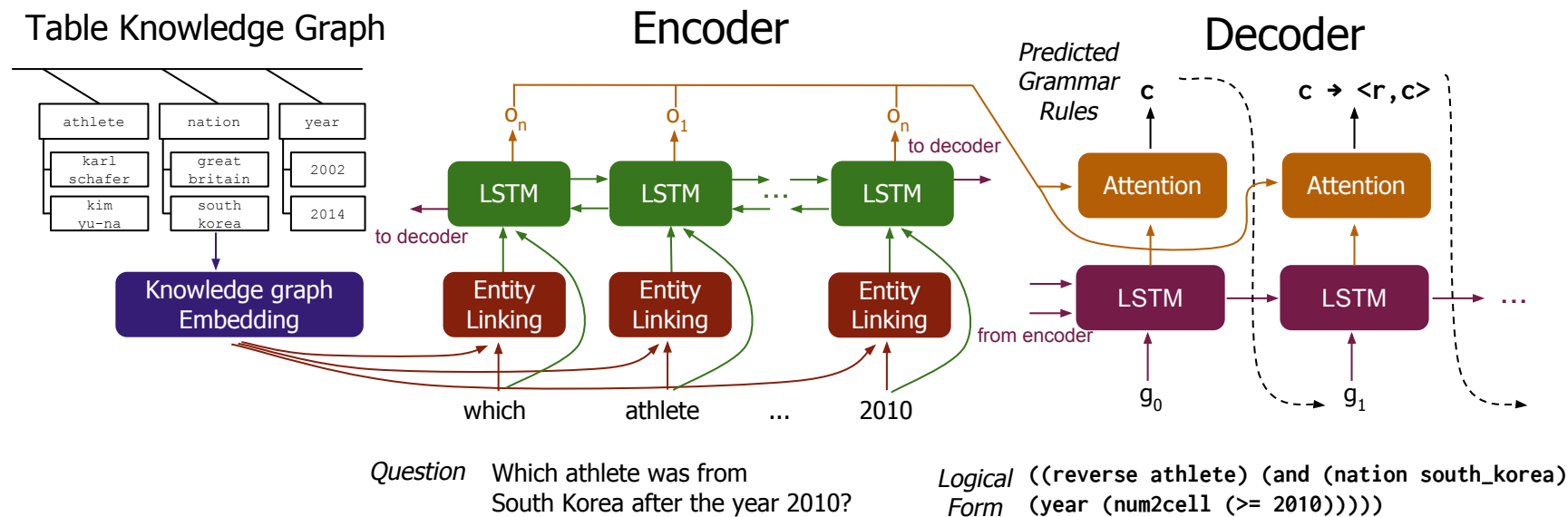


Figure 1: Overview of our semantic parsing model. The encoder performs entity embedding and linking before encoding the question with a bidirectional LSTM. The decoder predicts a sequence of grammar rules that generate a well-typed logical form.

Neural Semantic Parsing

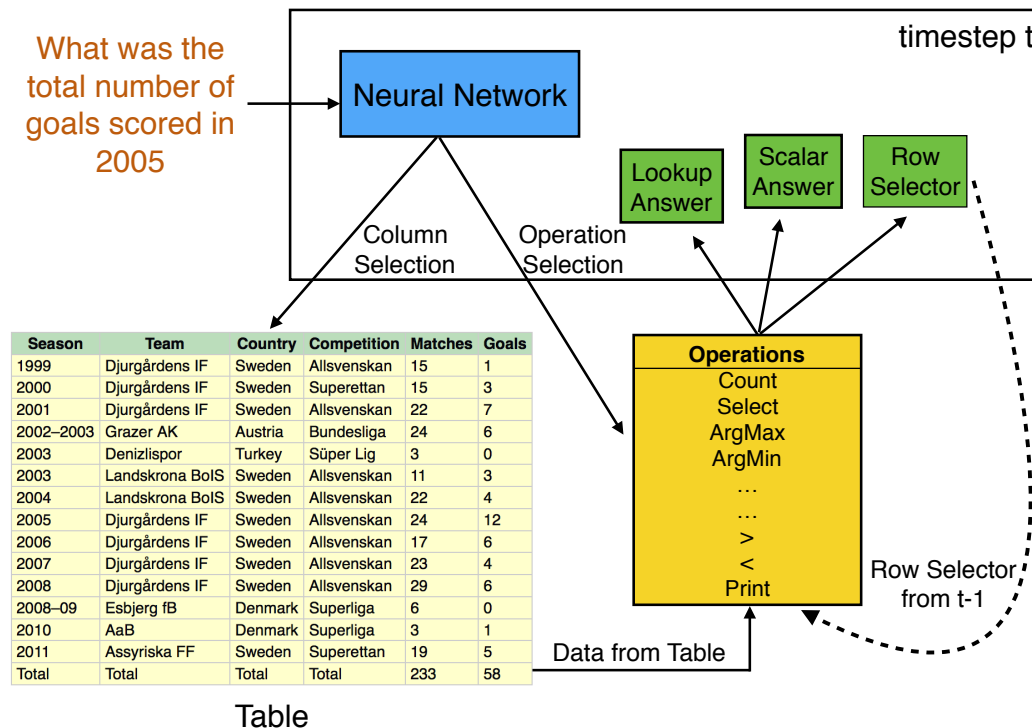


Figure 1: Neural Programmer is a neural network augmented with a set of discrete operations. The model runs for a fixed number of time steps, selecting an operation and a column from the table at every time step. The induced program transfers information across timesteps using the *row selector* variable while the output of the model is stored in the *scalar answer* and *lookup answer* variables.

Neural AMR Parsing

Obama was elected and his voters celebrated

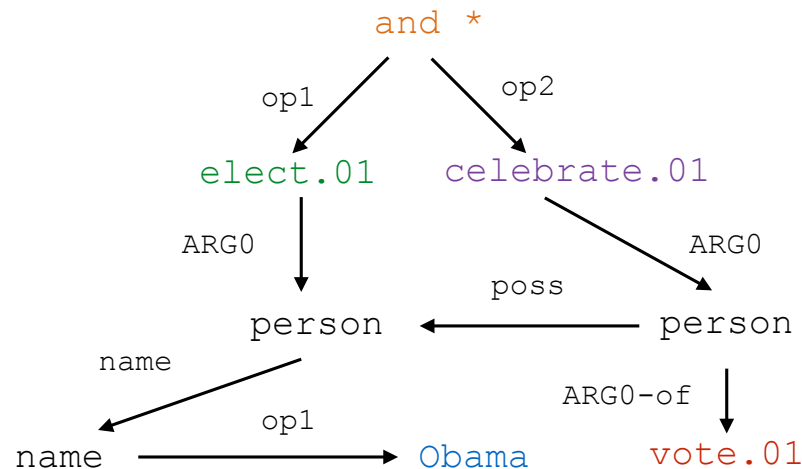
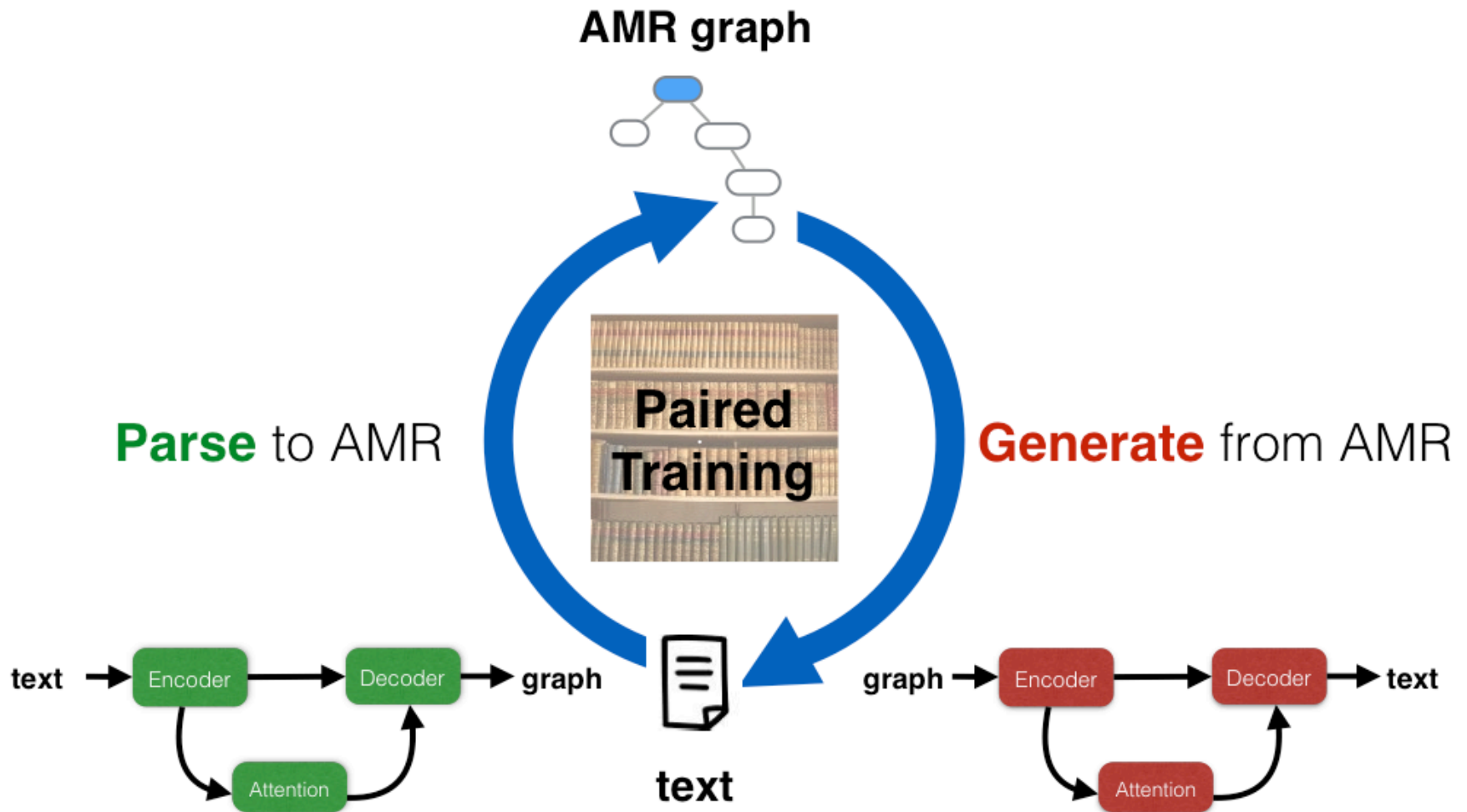


Figure 1: An example sentence and its corresponding Abstract Meaning Representation (AMR). AMR encodes semantic dependencies between entities mentioned in the sentence, such as “Obama” being the “arg0” of the verb “elected”.

Neural AMR Parsing



Question Answering

IR-based Question Answering

- ▶ Initial approaches to Q&A: pattern matching, pattern learning, query rewriting, information extraction

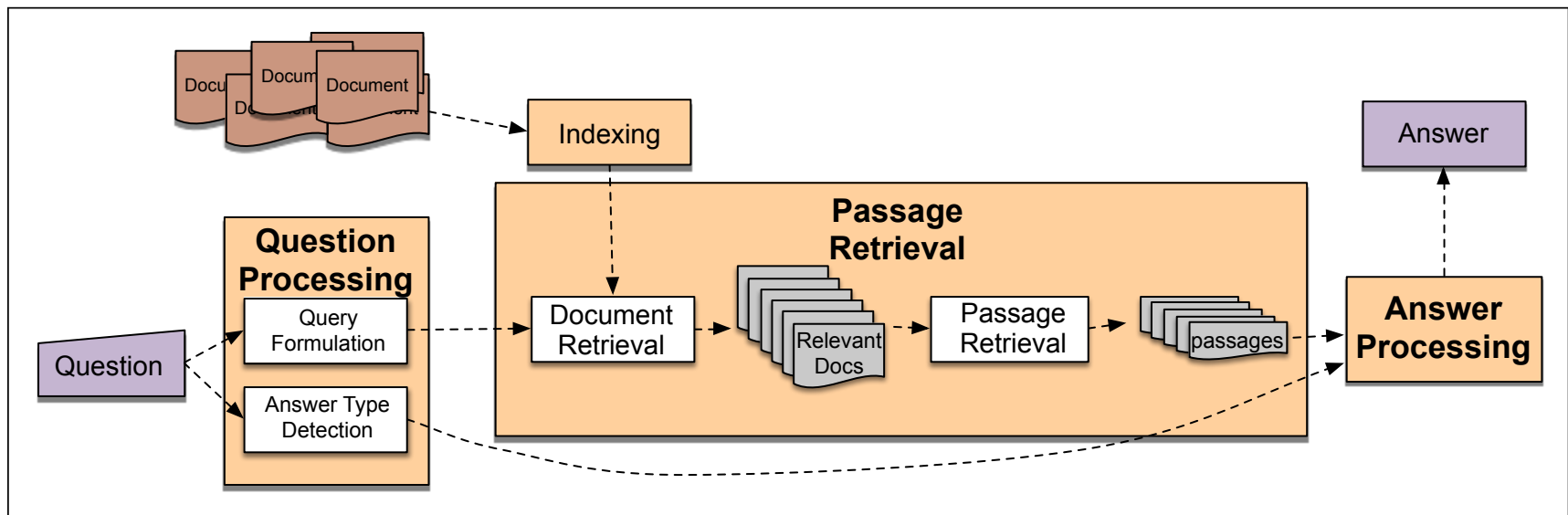


Figure 28.2 IR-based factoid question answering has three stages: question processing, passage retrieval, and answer processing.

IR-based Question Answering

- ▶ Initial approaches to Q&A: pattern matching, pattern learning, query rewriting, information extraction

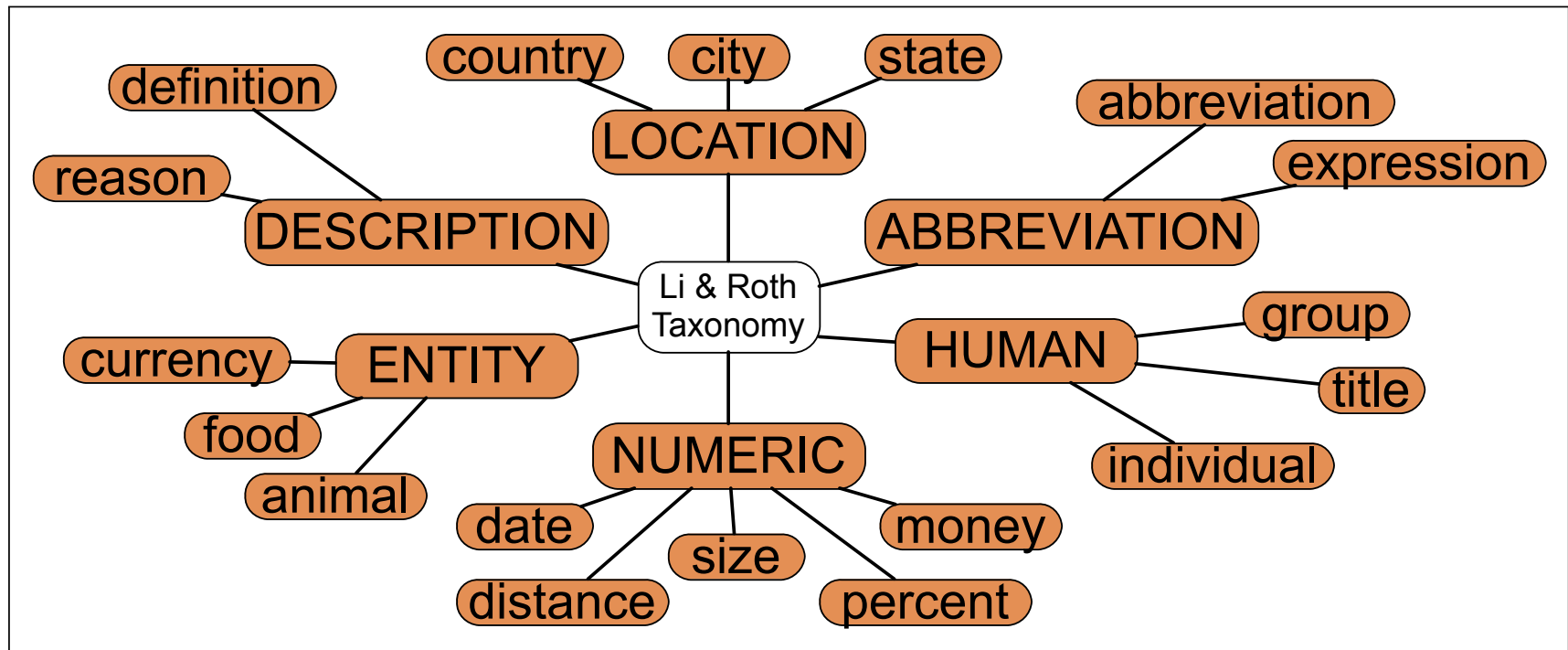


Figure 28.3 A subset of the Li and Roth (2005) answer types.

IR-based Question Answering

- ▶ Next came a large-scale, open-domain IE system like IBM Watson

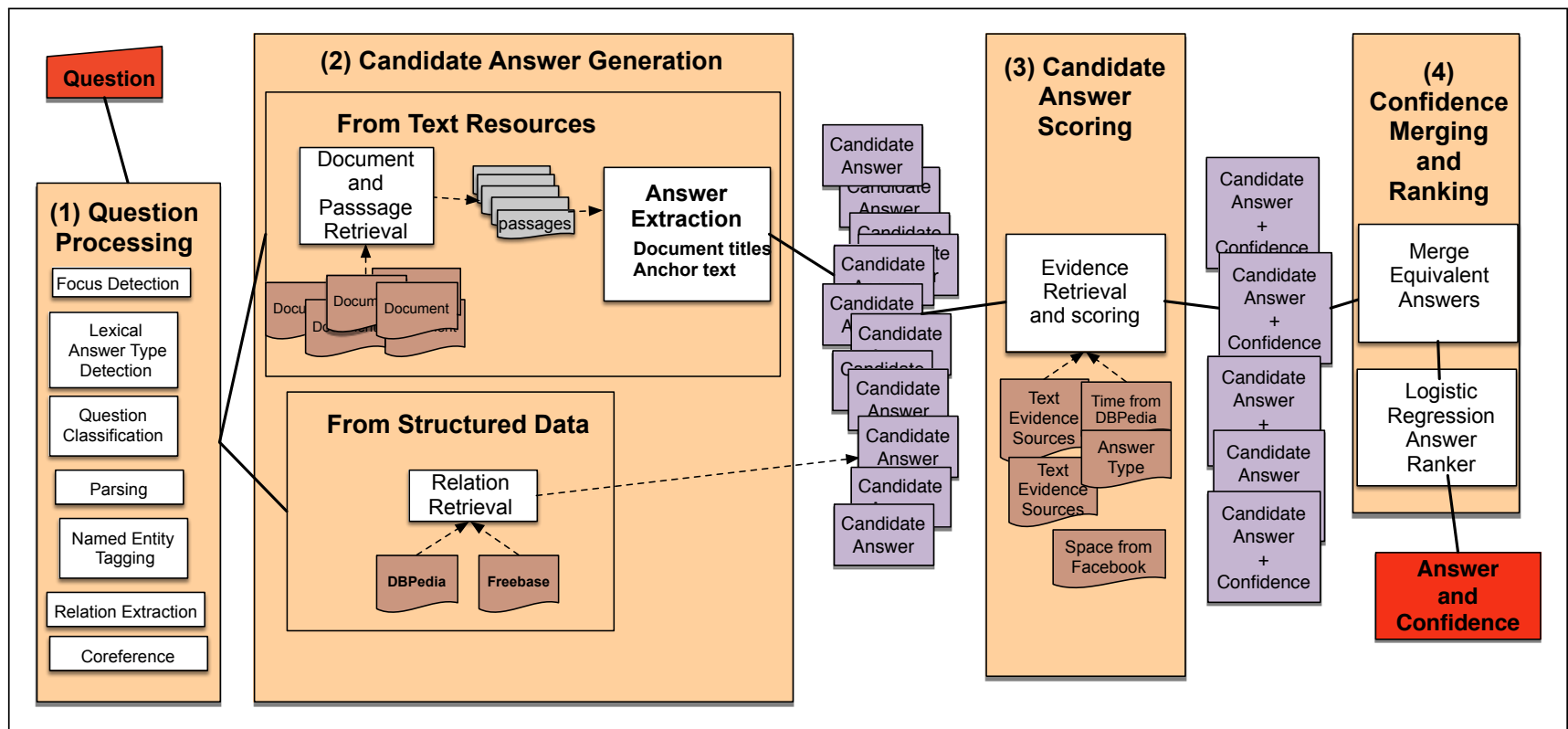
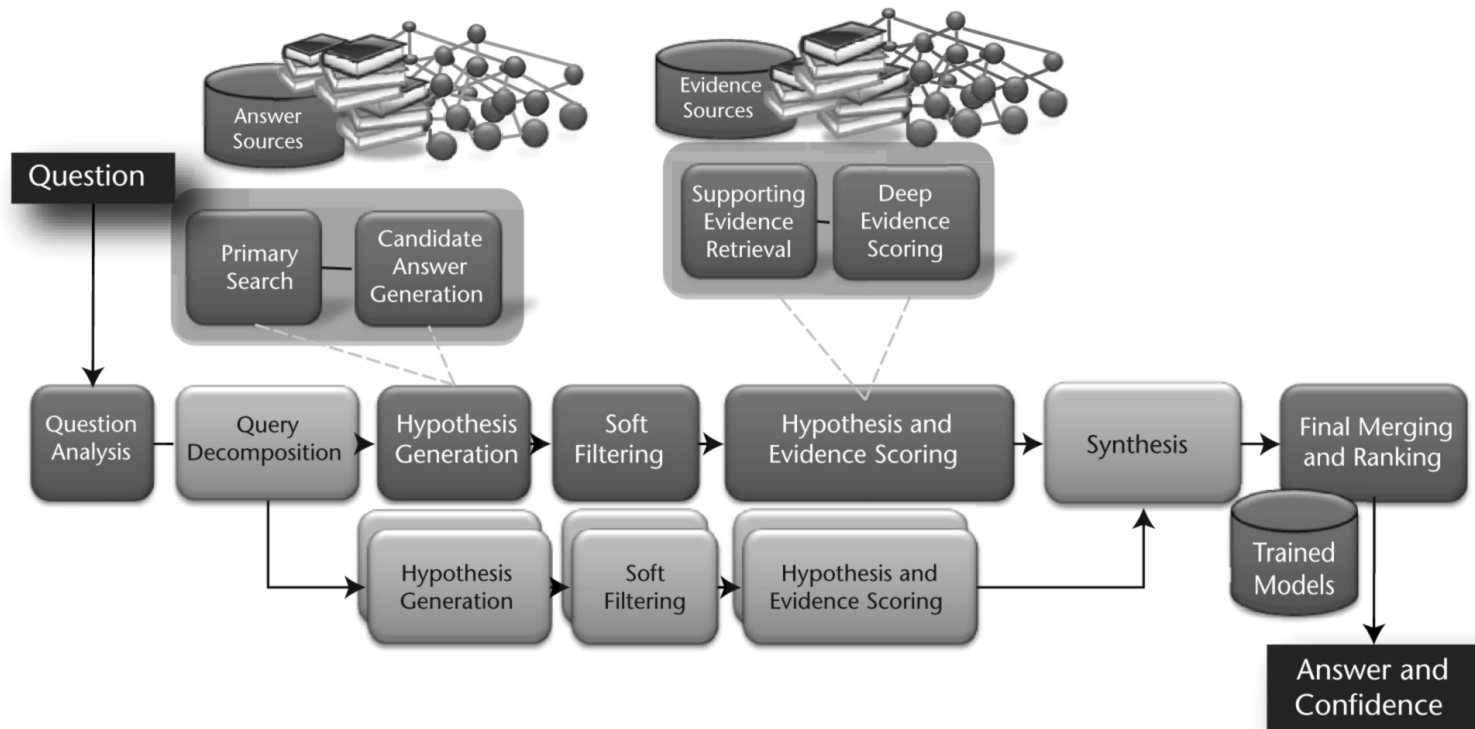


Figure 28.9 The 4 broad stages of Watson QA: (1) Question Processing, (2) Candidate Answer Generation, (3) Candidate Answer Scoring, and (4) Answer Merging and Confidence Scoring.

IR-based Question Answering

- ▶ Next came a large-scale, open-domain IE system like IBM Watson



Knowledge Base Q&A (Semantic Parsing)

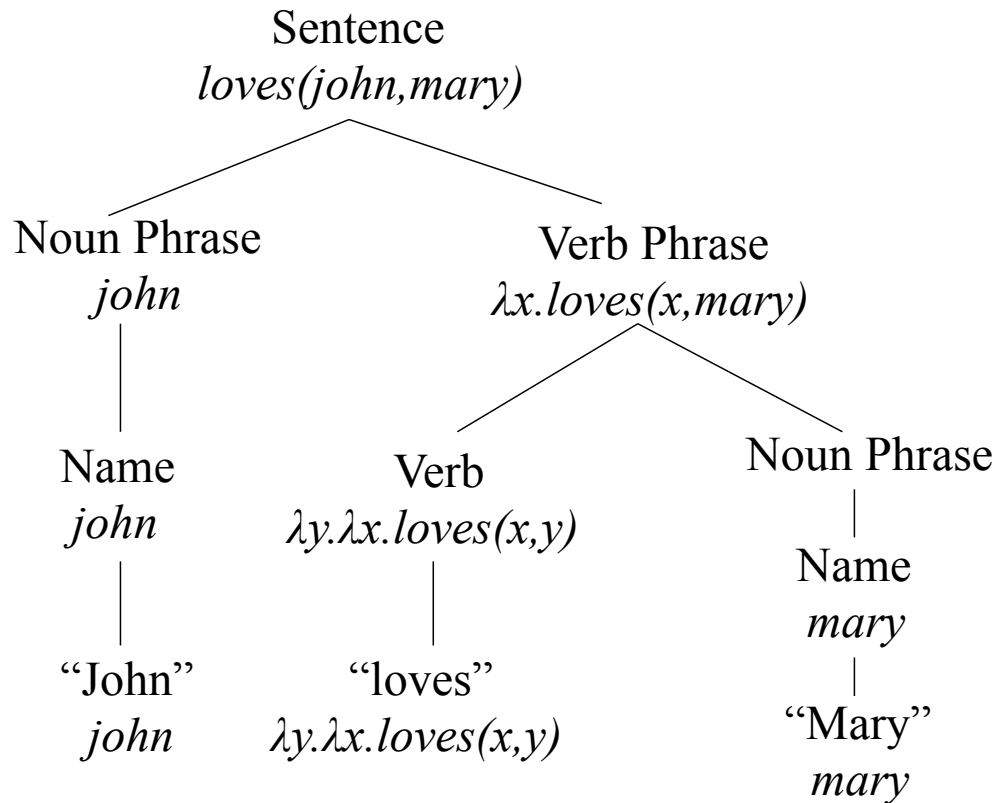
- ▶ Answering question by mapping it to a query (e.g., based on logical forms) executable on a structured database

Question	Logical form
When was Ada Lovelace born?	<code>birth-year (Ada Lovelace, ?x)</code>
What states border Texas?	<code>$\lambda x.state(x) \wedge borders(x,texas)$</code>
What is the largest state	<code>$argmax(\lambda x.state(x), \lambda x.size(x))$</code>
How many people survived the sinking of the Titanic	<code>(count (!fb:event.disaster.survivors fb:en.sinking_of_the_titanic))</code>

Figure 28.7 Sample logical forms produced by a semantic parser for question answering. These range from simple relations like `birth-year`, or relations normalized to databases like Freebase, to full predicate calculus.

Semantic Parsing Recap

- ▶ Parsing with logic (booleans, individuals, functions) and lambda forms



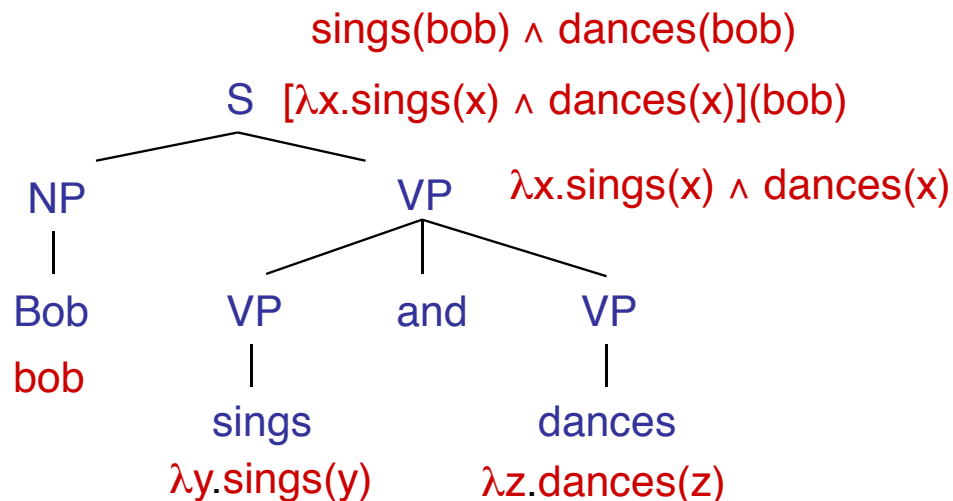
[Wong and Mooney, 2007; Zettlemoyer and Collins, 2007; Poon and Domingos, 2009; Artzi and Zettlemoyer, 2011, 2013; Kwiatkowski et al., 2013; Cai and Yates, 2013; Berant et al., 2013; Poon 2013; Berant and Liang, 2014; Iyer et al., 2014]

Compositional Semantics

- ▶ Now after we have these meanings for words, we want to combine them into meaning for phrases and sentences
- ▶ For this, we associate a combination rule with each grammar rule of the parse tree, e.g.:

S: $\beta(\alpha) \rightarrow$ NP: α VP: β *(function application)*

VP: $\lambda x . \alpha(x) \wedge \beta(x) \rightarrow$ VP: α and: \emptyset VP: β *(intersection)*



CCG Parsing Example

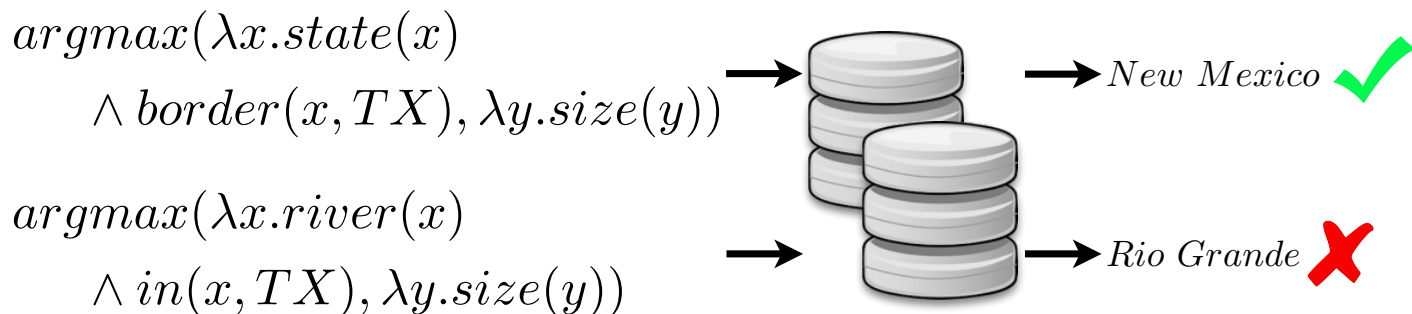
Show me	flights	to	Prague
S/N $\lambda f.f$	N $\lambda x.flight(x)$	(N\N) / NP $\lambda y.\lambda f.\lambda x.f(y) \wedge to(x,y)$	NP PRG
		N\N $\lambda f.\lambda x.f(x) \wedge to(x, PRG)$	
		N $\lambda x.flight(x) \wedge to(x, PRG)$	
		S $\lambda x.flight(x) \wedge to(x, PRG)$	

Weak Supervision

- ▶ Instead of relying on sentence-logicform pairs as training data, we can learn from query-answer pairs
- ▶ Logical forms are latent, and we can check which one gets the correct answer on being executed against a knowledge base (KB)

What is the largest state that borders Texas?

New Mexico



Weak Supervision

► Learning from Instruction-Demonstration Pairs

at the chair, move forward three steps past the sofa



Some examples from other domains:

- Sentences and labeled game states [Goldwasser and Roth 2011]
- Sentences and sets of physical objects [Matuszek et al. 2012]

Weak Supervision

► Learning from Conversation Logs

SYSTEM how can I help you ? (**OPEN_TASK**)

USER i ' d like to fly to new york

SYSTEM flying to new york . (**CONFIRM:** $from(fl, ATL)$) leaving what city ?
(**ASK:** $\lambda x.from(fl, x)$)

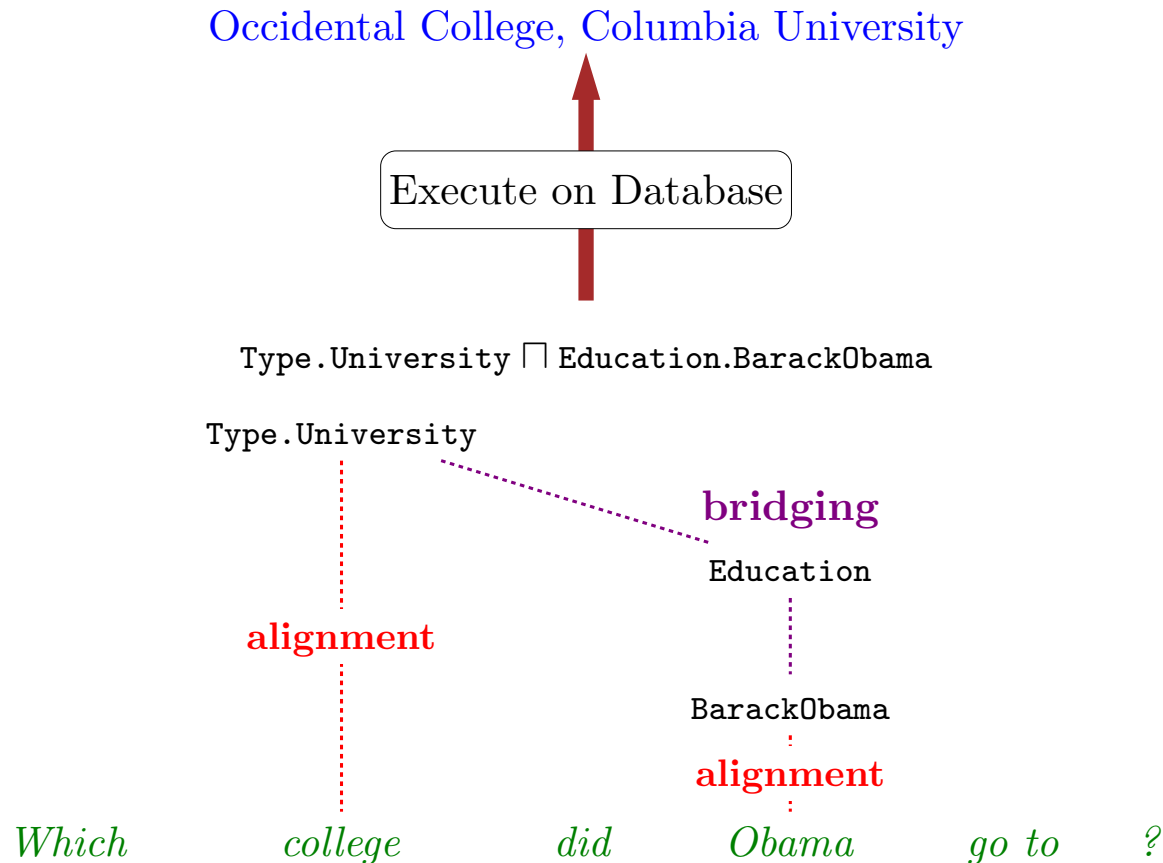
USER from boston on june seven with american airlines

SYSTEM flying to new york . (**CONFIRM:** $to(fl, NYC)$) what date would you like to depart boston ? (**ASK:** $\lambda x.date(fl, x) \wedge to(fl, BOS)$)

USER june seventh

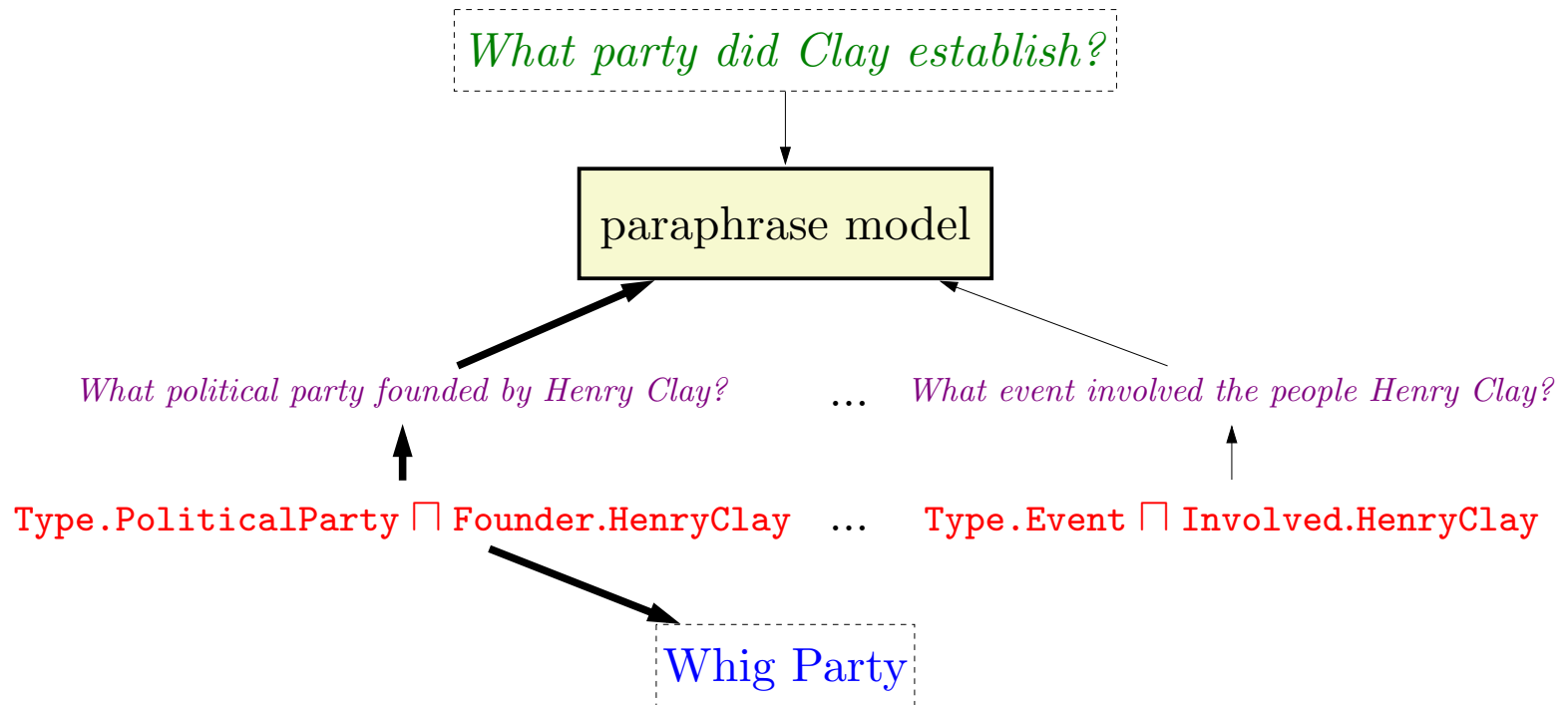
[CONVERSATION CONTINUES]

Semantic Parsing on Freebase



Mapping questions to answers via latent logical forms. To narrow down the logical predicate space, they use a (i) coarse *alignment* based on Freebase and a text corpus and (ii) a *bridging* operation that generates predicates compatible with neighboring predicates.

Semantic Parsing via Paraphrasing



For each candidate logical form (red), they generate canonical utterances (purple). The model is trained to paraphrase the input utterance (green) into the canonical utterances associated with the correct denotation (blue).

Reading Comprehension or Passage-based Q&A

Passage-based Q&A

James the Turtle was always getting in trouble. Sometimes he'd reach into the freezer and empty out all the food. Other times he'd sled on the deck and get a splinter. ... He went to the grocery store and pulled all the pudding off the shelves and ate two jars. Then he walked to the fast food restaurant and ordered 15 bags of fries.

Q. What did James pull off of the shelves in the grocery store?

(A) pudding, (B) fries, (C) food, (D) splinters

Q. Where did James go after eating two jars of pudding?

(A) grocery, (B) restaurant, (C) freezer, (D) home

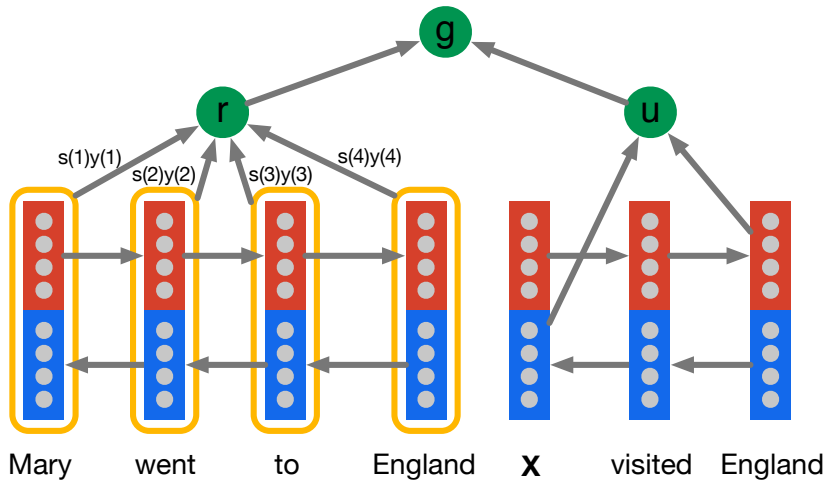
CNN/DailyMail RC Datasets

Original Version	Anonymised Version
Context The BBC producer allegedly struck by Jeremy Clarkson will not press charges against the “Top Gear” host, his lawyer said Friday. Clarkson, who hosted one of the most-watched television shows in the world, was dropped by the BBC Wednesday after an internal investigation by the British broadcaster found he had subjected producer Oisin Tymon “to an unprovoked physical and verbal attack.” ...	the <i>ent381</i> producer allegedly struck by <i>ent212</i> will not press charges against the “ <i>ent153</i> ” host , his lawyer said friday . <i>ent212</i> , who hosted one of the most - watched television shows in the world , was dropped by the <i>ent381</i> wednesday after an internal investigation by the <i>ent180</i> broadcaster found he had subjected producer <i>ent193</i> “ to an unprovoked physical and verbal attack . ” ...
Query Producer X will not press charges against Jeremy Clarkson, his lawyer says.	producer X will not press charges against <i>ent212</i> , his lawyer says .
Answer Oisin Tymon	<i>ent193</i>

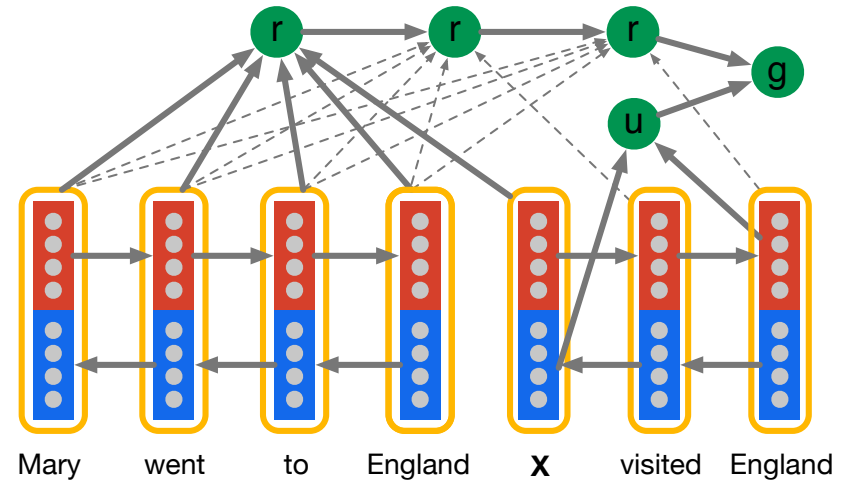
CNN/DailyMail RC Datasets

	CNN			Daily Mail		
	train	valid	test	train	valid	test
# months	95	1	1	56	1	1
# documents	90,266	1,220	1,093	196,961	12,148	10,397
# queries	380,298	3,924	3,198	879,450	64,835	53,182
Max # entities	527	187	396	371	232	245
Avg # entities	26.4	26.5	24.5	26.5	25.5	26.0
Avg # tokens	762	763	716	813	774	780
Vocab size	118,497			208,045		

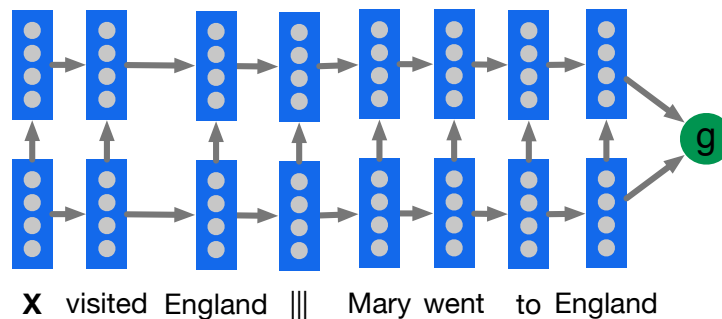
Attentive/Impatient Readers



(a) Attentive Reader.



(b) Impatient Reader.



(c) A two layer Deep LSTM Reader with the question encoded before the document.

Attentive/Impatient Readers

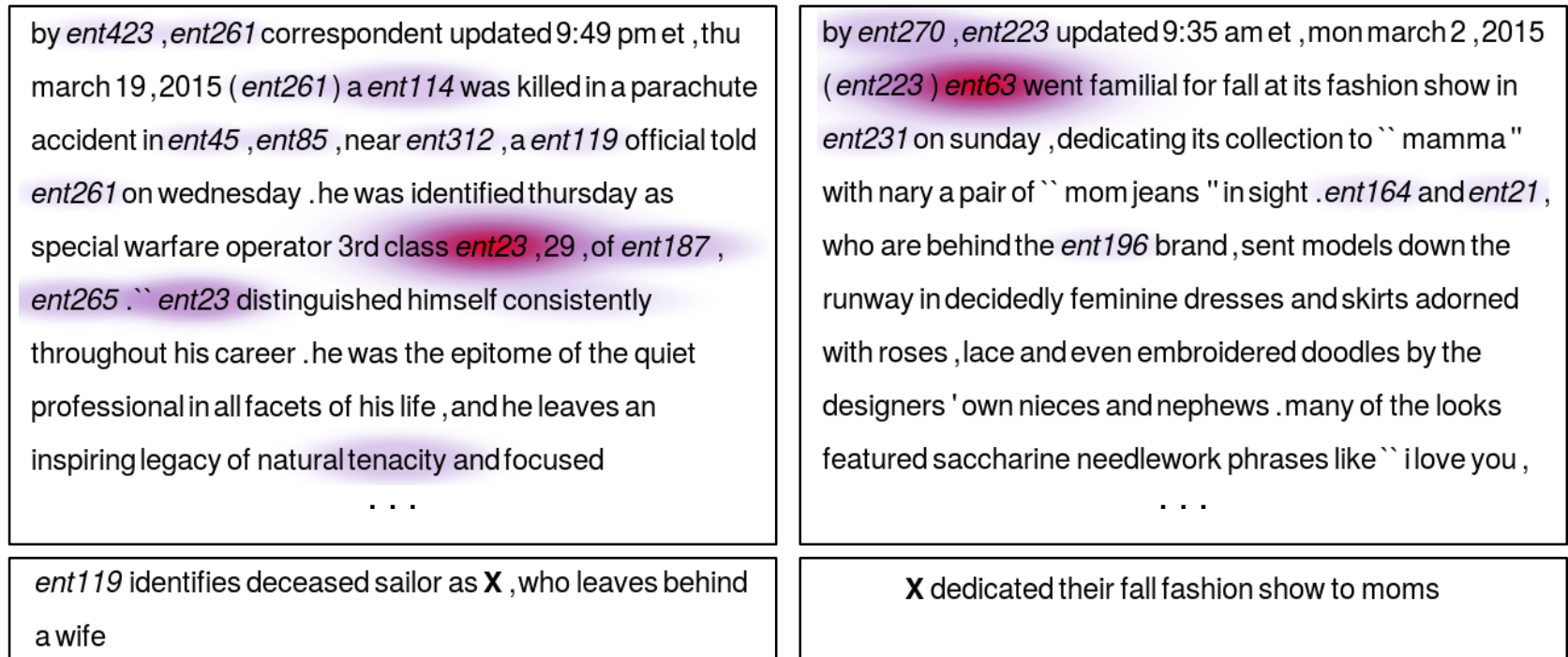


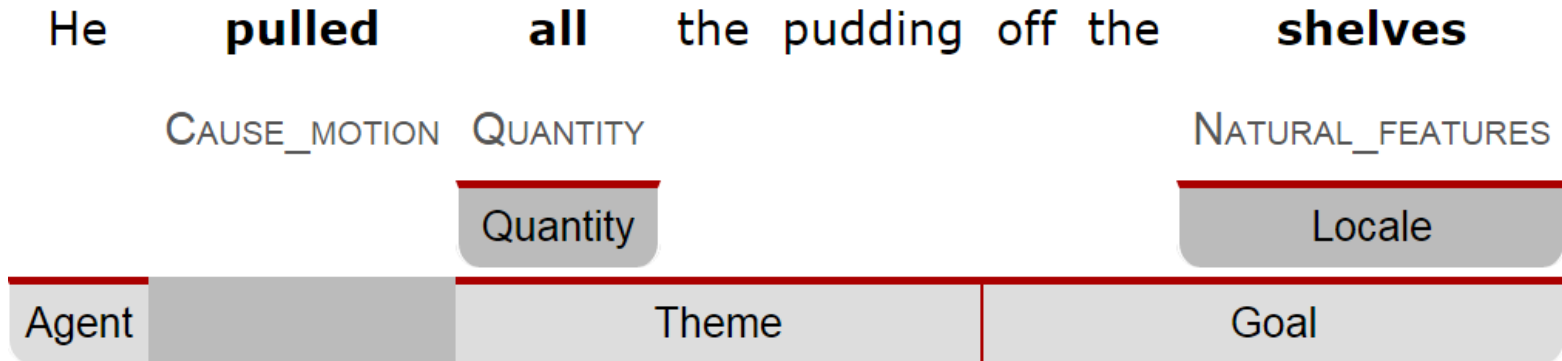
Figure 3: Attention heat maps from the Attentive Reader for two correctly answered validation set queries (the correct answers are *ent23* and *ent63*, respectively). Both examples require significant lexical generalisation and co-reference resolution in order to be answered correctly by a given model.

Attentive/Impatient Readers

	CNN		Daily Mail	
	valid	test	valid	test
Maximum frequency	30.5	33.2	25.6	25.5
Exclusive frequency	36.6	39.3	32.7	32.8
Frame-semantic model	36.3	40.2	35.5	35.5
Word distance model	50.5	50.9	56.4	55.5
Deep LSTM Reader	55.0	57.0	63.3	62.2
Uniform Reader	39.0	39.4	34.6	34.4
Attentive Reader	61.6	63.0	70.5	69.0
Impatient Reader	61.8	63.8	69.0	68.0

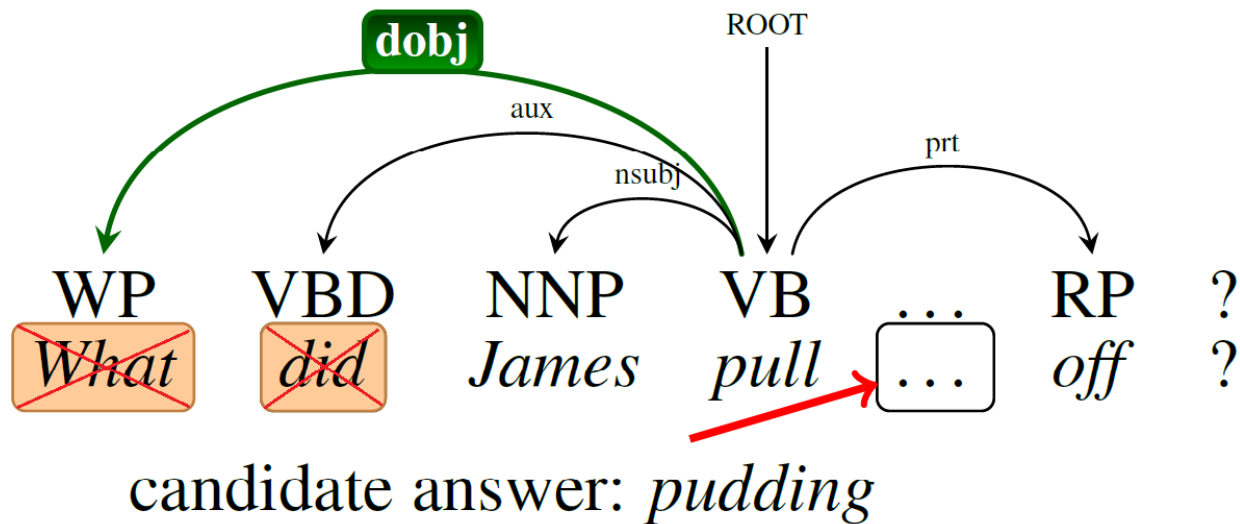
Feature-based Model

- ▶ Weighted word overlap between the bag of words constructed from the question/answer and in the window (and their word embedding versions)
- ▶ Minimal distance between two word occurrences in the passage that are also contained in the question/answer pair
- ▶ Frame semantics (predicates, frames evoked, and predicted argument labels) match between passage sentence and question+answer



Feature-based Model

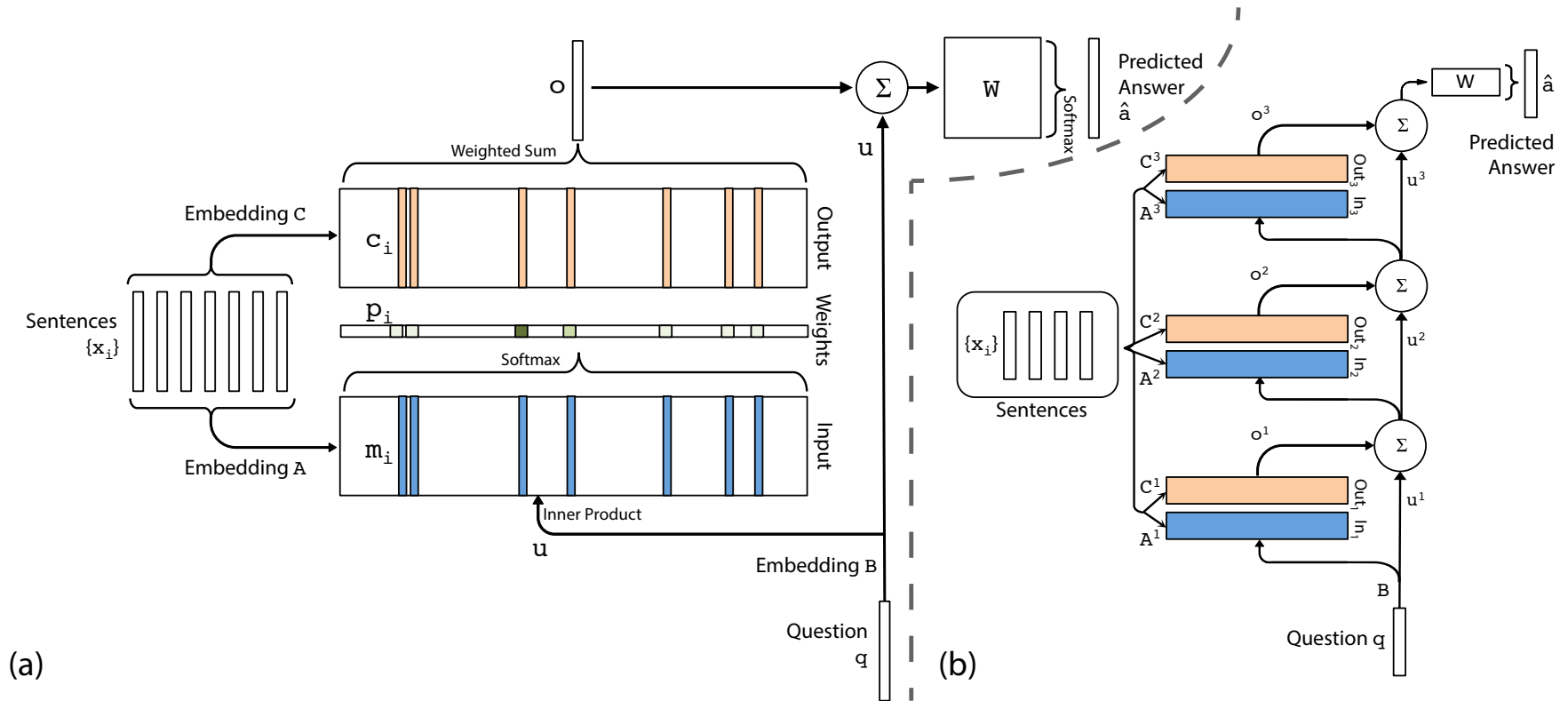
- ▶ Syntactic dependencies match between passage sentence and ques+ans converted to statement
- ▶ Extra features computed after coreference resolution of pronouns/nominals to map to their entity clusters



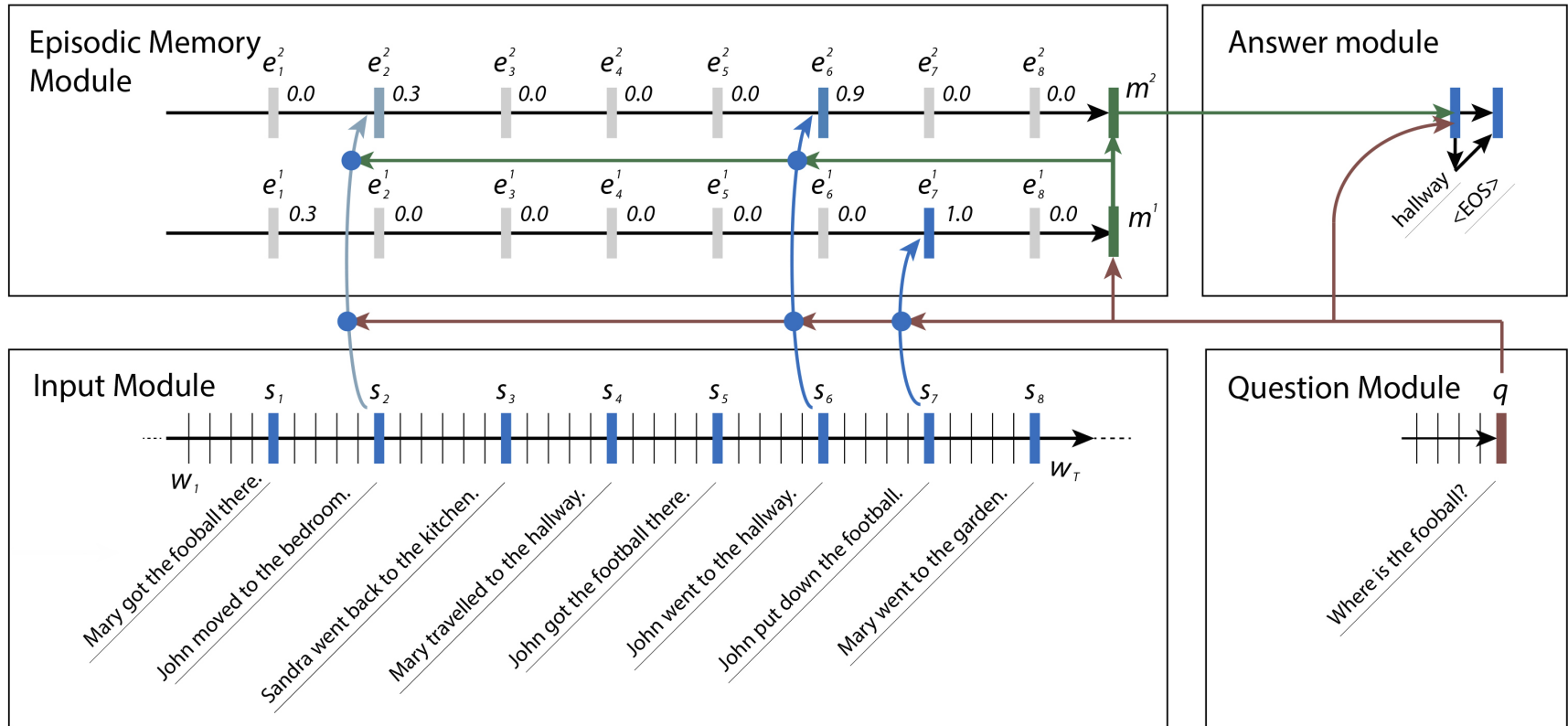
Multi-Hop Memory Models

- ▶ Several questions need multi-hop (e.g., path or count-based) reasoning to answer
- ▶ Memory models perform multiple passes over the text to collect the multiple evidence pieces
- ▶ Some example models:
 - ▶ End-to-End Memory Networks
 - ▶ Dynamic Memory Networks
 - ▶ Gated Attention Readers

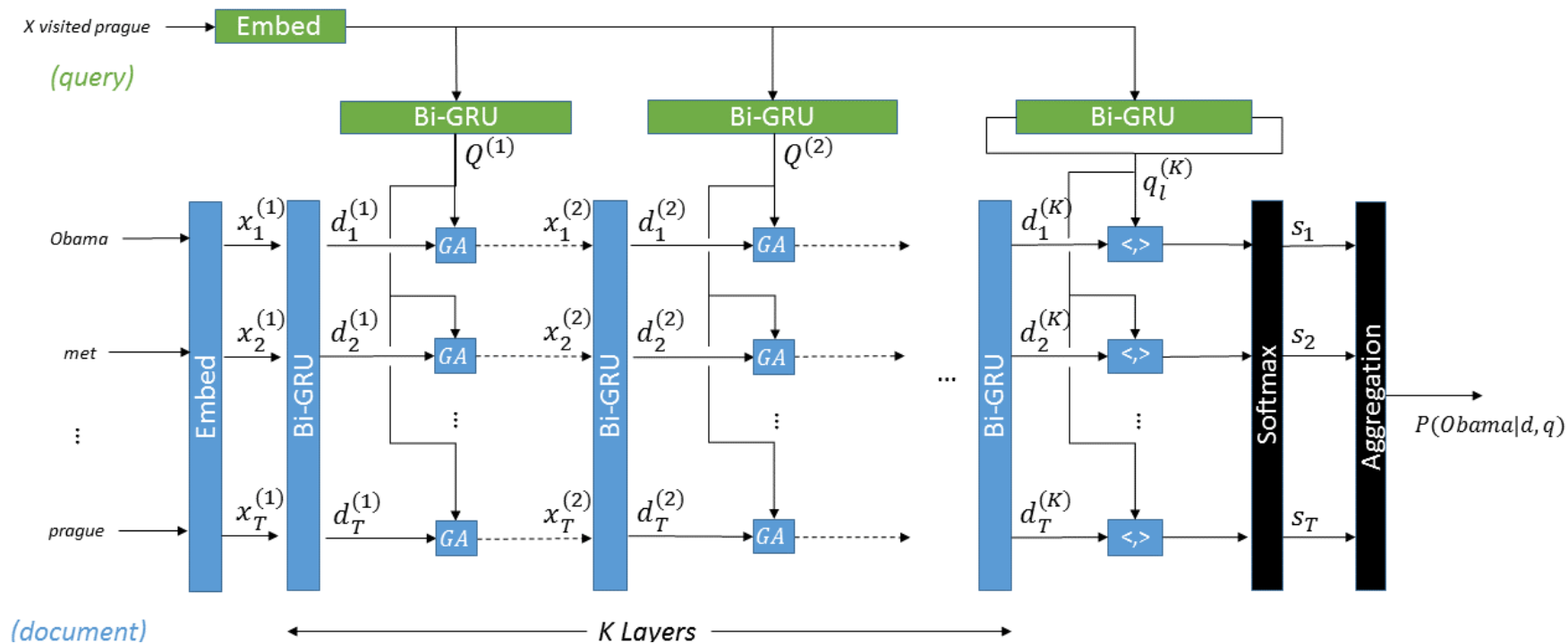
End-to-End Memory Networks



Dynamic Memory Networks

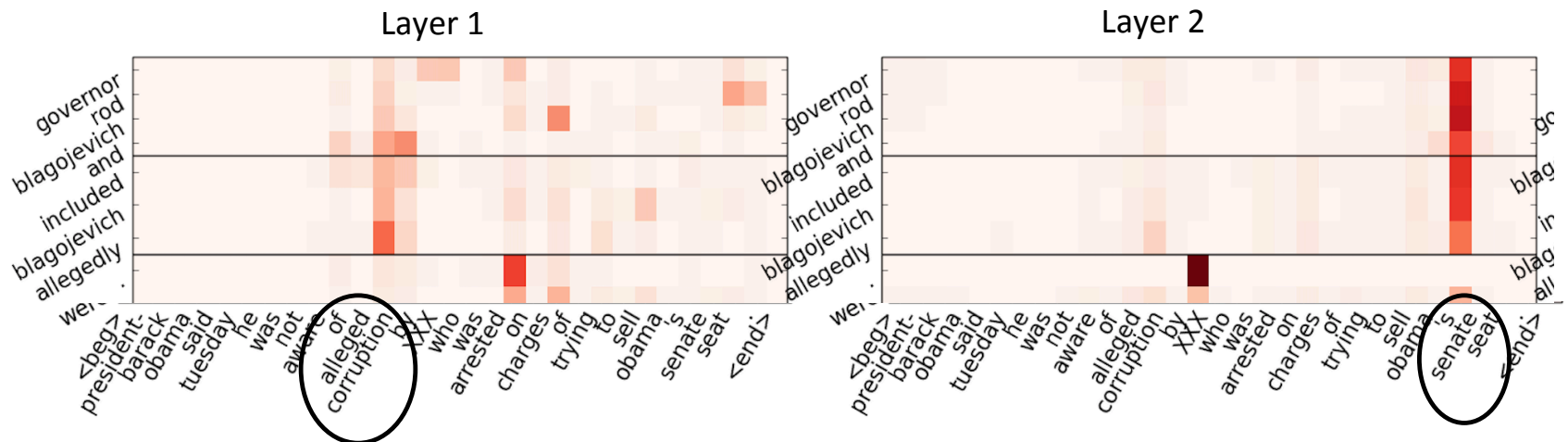


Gated Attention Readers



Gated Attention Readers

- **Context:** “...arrested Illinois **governor Rod Blagojevich** and his chief of staff John Harris on corruption charges ... **included Blagojevich** allegedly conspiring to sell or trade the **senate seat** left vacant by President-elect Barack Obama...”
- **Query:** “President-elect Barack Obama said Tuesday he was not aware of **alleged corruption** by **X** who was arrested on charges of trying to sell Obama’s **senate seat**.”
- **Answer:** Rod Blagojevich



Code + Data: <https://github.com/bdhingra/ga-reader>

Gated Attention Readers

Table 1: Validation/Test accuracy (%) on WDW dataset for both “Strict” and “Relaxed” settings. Results with “†” are of previously published works.

Model	Strict		Relaxed	
	Val	Test	Val	Test
Human †	–	84	–	–
Attentive Reader †	–	53	–	55
AS Reader †	–	57	–	59
Stanford AR †	–	64	–	65
NSE †	66.5	66.2	67.0	66.7
GA-- †	–	57	–	60.0
GA (update $L(w)$)	67.8	67.0	67.0	66.6
GA (fix $L(w)$)	68.3	68.0	69.6	69.1
GA (+feature, update $L(w)$)	70.1	69.5	70.9	71.0
GA (+feature, fix $L(w)$)	71.6	71.2	72.6	72.6

Table 2: **Top:** Performance of different gating functions. **Bottom:** Effect of varying the number of hops K . Results on WDW without using the qe-comm feature and with fixed $L(w)$.

Gating Function	Accuracy	
	Val	Test
Sum	64.9	64.5
Concatenate	64.4	63.7
Multiply	68.3	68.0
K		
1 (AS) †	–	57
2	65.6	65.6
3	68.3	68.0
4	68.3	68.2

Gated Attention Readers

Table 3: Validation/Test accuracy (%) on CNN, Daily Mail and CBT. Results marked with “†” are of previously published works. Results marked with “‡” were obtained by training on a larger training set. Best performance on standard training sets is in bold, and on larger training sets in italics.

Model	CNN		Daily Mail		CBT-NE		CBT-CN	
	Val	Test	Val	Test	Val	Test	Val	Test
Humans (query) †	–	–	–	–	–	52.0	–	64.4
Humans (context + query) †	–	–	–	–	–	81.6	–	81.6
LSTMs (context + query) †	–	–	–	–	51.2	41.8	62.6	56.0
Deep LSTM Reader †	55.0	57.0	63.3	62.2	–	–	–	–
Attentive Reader †	61.6	63.0	70.5	69.0	–	–	–	–
Impatient Reader †	61.8	63.8	69.0	68.0	–	–	–	–
MemNets †	63.4	66.8	–	–	70.4	66.6	64.2	63.0
AS Reader †	68.6	69.5	75.0	73.9	73.8	68.6	68.8	63.4
DER Network †	71.3	72.9	–	–	–	–	–	–
Stanford AR (relabeling) †	73.8	73.6	77.6	76.6	–	–	–	–
Iterative Attentive Reader †	72.6	73.3	–	–	75.2	68.6	72.1	69.2
EpiReader †	73.4	74.0	–	–	75.3	69.7	71.5	67.4
AoA Reader †	73.1	74.4	–	–	77.8	72.0	72.2	69.4
ReasoNet †	72.9	74.7	77.6	76.6	–	–	–	–
NSE †	–	–	–	–	78.2	73.2	74.3	71.9
BiDAF †	76.3	76.9	80.3	79.6	–	–	–	–
MemNets (ensemble) †	66.2	69.4	–	–	–	–	–	–
AS Reader (ensemble) †	73.9	75.4	78.7	77.7	76.2	71.0	71.1	68.9
Stanford AR (relabeling,ensemble) †	77.2	77.6	80.2	79.2	–	–	–	–
Iterative Attentive Reader (ensemble) †	75.2	76.1	–	–	76.9	72.0	74.1	71.0
EpiReader (ensemble) †	–	–	–	–	76.6	71.8	73.6	70.6
AS Reader (+BookTest) † ‡	–	–	–	–	80.5	76.2	83.2	80.8
AS Reader (+BookTest,ensemble) † ‡	–	–	–	–	82.3	78.4	85.7	83.7
GA--	73.0	73.8	76.7	75.7	74.9	69.0	69.0	63.9
GA (update $L(w)$)	77.9	77.9	81.5	80.9	76.7	70.1	69.8	67.3
GA (fix $L(w)$)	77.9	77.8	80.4	79.6	77.2	71.4	71.6	68.0
GA (+feature, update $L(w)$)	77.3	76.9	80.7	80.0	77.2	73.3	73.0	69.8
GA (+feature, fix $L(w)$)	76.7	77.4	80.0	79.3	78.5	74.9	74.4	70.7

Facebook bAbI Tasks (Synthetic)

Task 1: Single Supporting Fact

Mary went to the bathroom.
John moved to the hallway.
Mary travelled to the office.
Where is Mary? A:office

Task 2: Two Supporting Facts

John is in the playground.
John picked up the football.
Bob went to the kitchen.
Where is the football? A:playground

Task 3: Three Supporting Facts

John picked up the apple.
John went to the office.
John went to the kitchen.
John dropped the apple.
Where was the apple before the kitchen? A:office

Task 4: Two Argument Relations

The office is north of the bedroom.
The bedroom is north of the bathroom.
The kitchen is west of the garden.
What is north of the bedroom? A: office
What is the bedroom north of? A: bathroom

Task 5: Three Argument Relations

Mary gave the cake to Fred.
Fred gave the cake to Bill.
Jeff was given the milk by Bill.
Who gave the cake to Fred? A: Mary
Who did Fred give the cake to? A: Bill

Task 6: Yes/No Questions

John moved to the playground.
Daniel went to the bathroom.
John went back to the hallway.
Is John in the playground? A:no
Is Daniel in the bathroom? A:yes

Task 7: Counting

Daniel picked up the football.
Daniel dropped the football.
Daniel got the milk.
Daniel took the apple.
How many objects is Daniel holding? A: two

Task 8: Lists/Sets

Daniel picks up the football.
Daniel drops the newspaper.
Daniel picks up the milk.
John took the apple.
What is Daniel holding? milk, football

Task 9: Simple Negation

Sandra travelled to the office.
Fred is no longer in the office.
Is Fred in the office? A:no
Is Sandra in the office? A:yes

Task 10: Indefinite Knowledge

John is either in the classroom or the playground.
Sandra is in the garden.
Is John in the classroom? A:maybe
Is John in the office? A:no

Facebook bAbI Tasks (Synthetic)

Task 11: Basic Coreference

Daniel was in the kitchen.
Then he went to the studio.
Sandra was in the office.
Where is Daniel? A: studio

Task 12: Conjunction

Mary and Jeff went to the kitchen.
Then Jeff went to the park.
Where is Mary? A: kitchen
Where is Jeff? A: park

Task 13: Compound Coreference

Daniel and Sandra journeyed to the office.
Then they went to the garden.
Sandra and John travelled to the kitchen.
After that they moved to the hallway.
Where is Daniel? A: garden

Task 14: Time Reasoning

In the afternoon Julie went to the park.
Yesterday Julie was at school.
Julie went to the cinema this evening.
Where did Julie go after the park? A: cinema
Where was Julie before the park? A: school

Task 15: Basic Deduction

Sheep are afraid of wolves.
Cats are afraid of dogs.
Mice are afraid of cats.
Gertrude is a sheep.
What is Gertrude afraid of? A: wolves

Task 16: Basic Induction

Lily is a swan.
Lily is white.
Bernhard is green.
Greg is a swan.
What color is Greg? A: white

Task 17: Positional Reasoning

The triangle is to the right of the blue square.
The red square is on top of the blue square.
The red sphere is to the right of the blue square.
Is the red sphere to the right of the blue square? A: yes
Is the red square to the left of the triangle? A: yes

Task 18: Size Reasoning

The football fits in the suitcase.
The suitcase fits in the cupboard.
The box is smaller than the football.
Will the box fit in the suitcase? A: yes
Will the cupboard fit in the box? A: no

Task 19: Path Finding

The kitchen is north of the hallway.
The bathroom is west of the bedroom.
The den is east of the hallway.
The office is south of the bedroom.
How do you go from den to kitchen? A: west, north
How do you go from office to bathroom? A: north, west

Task 20: Agent's Motivations

John is hungry.
John goes to the kitchen.
John grabbed the apple there.
Daniel is hungry.
Where does Daniel go? A: kitchen
Why did John go to the kitchen? A: hungry

Facebook bAbI Tasks (Synthetic)

TASK	Weakly Supervised		Uses External Resources	Strong Supervision (using supporting facts)						
	<i>N-gram Classifier</i>	<i>LSTM</i>	<i>Structured SVM COREF+SRL features</i>	<i>MemNN Weston et al. (2014)</i>	<i>MemNN ADAPTIVE-MEMORY</i>	<i>MemNN AM + N-GRAMS</i>	<i>MemNN AM + NONLINEAR</i>	<i>MemNN AM + KG + NL</i>	<i>No. of ex. req. ≥ 95</i>	<i>MultiTask Training</i>
1 - Single Supporting Fact	36	50	99	100	100	100	100	100	250 ex.	100
2 - Two Supporting Facts	2	20	74	100	100	100	100	100	500 ex.	100
3 - Three Supporting Facts	7	20	17	20	100	99	100	100	500 ex.	98
4 - Two Arg. Relations	50	61	98	71	69	100	73	100	500 ex.	80
5 - Three Arg. Relations	20	70	83	83	83	86	86	98	1000 ex.	99
6 - Yes/No Questions	49	48	99	47	52	53	100	100	500 ex.	100
7 - Counting	52	49	69	68	78	86	83	85	FAIL	86
8 - Lists/Sets	40	45	70	77	90	88	94	91	FAIL	93
9 - Simple Negation	62	64	100	65	71	63	100	100	500 ex.	100
10 - Indefinite Knowledge	45	44	99	59	57	54	97	98	1000 ex.	98
11 - Basic Coreference	29	72	100	100	100	100	100	100	250 ex.	100
12 - Conjunction	9	74	96	100	100	100	100	100	250 ex.	100
13 - Compound Coref.	26	94	99	100	100	100	100	100	250 ex.	100
14 - Time Reasoning	19	27	99	99	100	99	100	99	500 ex.	99
15 - Basic Deduction	20	21	96	74	73	100	77	100	100 ex.	100
16 - Basic Induction	43	23	24	27	100	100	100	100	100 ex.	94
17 - Positional Reasoning	46	51	61	54	46	49	57	65	FAIL	72
18 - Size Reasoning	52	52	62	57	50	74	54	95	1000 ex.	93
19 - Path Finding	0	8	49	0	9	3	15	36	FAIL	19
20 - Agent's Motivations	76	91	95	100	100	100	100	100	250 ex.	100
Mean Performance	34	49	79	75	79	83	87	93		92

Who-did-What (WDW) Dataset

- ▶ Solves several issues with CNN/DM dataset:
 - ▶ Starts with the selection of a question article from Gigaword corpus
 - ▶ Question is formed by deleting a person named entity from the first sentence of the question article
 - ▶ An information retrieval system is then used to select a passage with high overlap with the first sentence of the question article, and an answer choice list is generated from the person named entities in the passage
- ▶ Forms questions from two distinct articles rather than summary points
- ▶ Allows using documents that don't contain manually-written summaries
- ▶ Reduces syntactic similarity between question & relevant passage sentences
- ▶ Selectively remove problems so as to suppress four simple baselines — selecting the most mentioned person, the first mentioned person, and two language model baselines
- ▶ The resulting dataset yields a larger gap between human and machine performance than existing ones, i.e., humans can answer more questions, while existing state-of-the-art models perform worse!

Who-did-What (WDW) Dataset

Passage: Britain's decision on Thursday to drop extradition proceedings against Gen. Augusto Pinochet and allow him to return to Chile is understandably frustrating ... Jack Straw, the home secretary, said the 84-year-old former dictator's ability to understand the charges against him and to direct his defense had been seriously impaired by a series of strokes. ... Chile's president-elect, Ricardo Lagos, has wisely pledged to let justice run its course. But the outgoing government of President Eduardo Frei is pushing a constitutional reform that would allow Pinochet to step down from the Senate and retain parliamentary immunity from prosecution. ...

Question: Sources close to the presidential palace said that Fujimori declined at the last moment to leave the country and instead he will send a high level delegation to the ceremony, at which Chilean President Eduardo Frei will pass the mandate to XXX.

Choices: (1) Augusto Pinochet (2) Jack Straw (3) Ricardo Lagos

Passage: Tottenham won 2-0 at Hapoel Tel Aviv in UEFA Cup action on Thursday night in a defensive display which impressed Spurs skipper Robbie Keane. ... Keane scored the first goal at the Bloomfield Stadium with Dimitar Berbatov, who insisted earlier on Thursday he was happy at the London club, heading a second. The 26-year-old Berbatov admitted the reports linking him with a move had affected his performances ... Spurs manager Juande Ramos has won the UEFA Cup in the last two seasons ...

Question: Tottenham manager Juande Ramos has hinted he will allow XXX to leave if the Bulgaria striker makes it clear he is unhappy.

Choices: (1) Robbie Keane (2) Dimitar Berbatov

Who-did-What (WDW) Dataset

Baseline	Accuracy	
	Before	After
First person in passage	0.60	0.32
Most frequent person	0.61	0.33
<i>n</i> -gram	0.53	0.33
Unigram	0.43	0.32
Random*	0.32	0.32

Table 2: Performance of suppressed baselines. *Random performance is computed as a deterministic function of the number of times each choice set size appears. Many questions have only two choices and there are about three choices on average.

	relaxed train	train	valid	test
# queries	185,978	127,786	10,000	10,000
Avg # choices	3.5	3.5	3.4	3.4
Avg # tokens	378	365	325	326
Vocab size	347,406		308,602	

Table 3: Dataset statistics.

Who-did-What (WDW) Dataset

System	WDW	CNN
Word overlap	0.47	—
Sliding window	0.48	—
Distance	0.46	—
Sliding window + Distance	0.51	—
Semantic features	0.52	—
Attentive Reader	0.53	0.63 ^I
Attentive Reader (relaxed train)	0.55	
Stanford Reader	0.64	0.73 ^{II}
Stanford Reader (relaxed train)	0.65	
AS Reader	0.57	0.70 ^{III}
AS Reader (relaxed train)	0.59	
GA Reader	0.57	0.74 ^{IV}
GA Reader (relaxed train)	0.60	
Human Performance	84/100	0.75+ ^{II}

SQuAd Dataset (100K Manually-Labeled)

- ▶ Based on manual annotation from Mturk on Wiki articles, as opposed to cloze/fill-in-the-blank on summaries, etc.; large size (100K+)
- ▶ Answer is a span in the document:

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called “showers”.

What causes precipitation to fall?

gravity

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

graupel

Where do water droplets collide with ice crystals to form precipitation?

within a cloud

SQuAd Dataset (100K Manually-Labeled)

Dataset	Question source	Formulation	Size
SQuAD	crowdsourced	RC, spans in passage	100K
MCTest (Richardson et al., 2013)	crowdsourced	RC, multiple choice	2640
Algebra (Kushman et al., 2014)	standardized tests	computation	514
Science (Clark and Etzioni, 2016)	standardized tests	reasoning, multiple choice	855
WikiQA (Yang et al., 2015)	query logs	IR, sentence selection	3047
TREC-QA (Voorhees and Tice, 2000)	query logs + human editor	IR, free form	1479
CNN/Daily Mail (Hermann et al., 2015)	summary + cloze	RC, fill in single entity	1.4M
CBT (Hill et al., 2015)	cloze	RC, fill in single word	688K

Table 1: A survey of several reading comprehension and question answering datasets. SQuAD is much larger than all datasets except the semi-synthetic cloze-style datasets, and it is similar to TREC-QA in the open-endedness of the answers.

SQuAd Dataset (100K Manually-Labeled)

Paragraph 1 of 43

Spend around 4 minutes on the following paragraph to ask 5 questions! If you can't ask 5 questions, ask 4 or 3 (worse), but do your best to ask 5. Select the answer from the paragraph by clicking on 'Select Answer', and then highlight the smallest segment of the paragraph that answers the question.

Oxygen is a chemical element with symbol O and atomic number 8. It is a member of the chalcogen group on the periodic table and is a highly reactive nonmetal and oxidizing agent that readily forms compounds (notably oxides) with most elements. By mass, oxygen is the third-most abundant element in the universe, after hydrogen and helium. At standard temperature and pressure, two atoms of the element bind to form dioxygen, a colorless and odorless diatomic gas with the formula O

2. Diatomic oxygen gas constitutes 20.8% of the Earth's atmosphere. However, monitoring of atmospheric oxygen levels show a global downward trend, because of fossil-fuel burning. Oxygen is the most abundant element by mass in the Earth's crust as part of oxide compounds such as silicon dioxide, making up almost half of the crust's mass.

When asking questions, **avoid using** the same words/phrases as in the paragraph. Also, you are encouraged to pose **hard questions**.

Ask a question here. Try using your own words

Select Answer

Ask a question here. Try using your own words

Select Answer

SQuAd Dataset (100K Manually-Labeled)

Reasoning	Description	Example	Percentage
Lexical variation (synonymy)	Major correspondences between the question and the answer sentence are synonyms.	Q: What is the Rankine cycle sometimes called ? Sentence: The Rankine cycle is sometimes referred to as a <u>practical Carnot cycle</u> .	33.3%
Lexical variation (world knowledge)	Major correspondences between the question and the answer sentence require world knowledge to resolve.	Q: Which governing bodies have veto power? Sen.: The European Parliament and the Council of the European Union have powers of amendment and veto during the legislative process.	9.1%
Syntactic variation	After the question is paraphrased into declarative form, its syntactic dependency structure does not match that of the answer sentence even after local modifications.	Q: What Shakespeare scholar is currently on the faculty ? Sen.: Current faculty include the anthropologist Marshall Sahllins, ..., Shakespeare scholar <u>David Bevington</u> .	64.1%
Multiple sentence reasoning	There is anaphora, or higher-level fusion of multiple sentences is required.	Q: What collection does the V&A Theatre & Performance galleries hold? Sen.: The V&A Theatre & Performance galleries opened in March 2009. ... They hold the UK's biggest national collection of material about live performance.	13.6%
Ambiguous	We don't agree with the crowdworkers' answer, or the question does not have a unique answer.	Q: What is the main goal of criminal punishment? Sen.: Achieving crime control via incapacitation and deterrence is a major goal of criminal punishment.	6.1%

Table 3: We manually labeled 192 examples into one or more of the above categories. Words relevant to the corresponding reasoning type are bolded, and the crowdsourced answer is underlined.

Adversarial Examples for Evaluating RC Systems

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

Adversarial Examples for Evaluating RC Systems



	Image Classification	Reading Comprehension
Possible Input		Tesla moved to the city of Chicago in 1880.
Similar Input		Tadakatsu moved to the city of Chicago in 1881.
Semantics	Same	Different
Model's Mistake	Considers the two to be different	Considers the two to be the same
Model Weakness	Overly sensitive	Overly stable

Table 1: Adversarial examples in computer vision exploit model oversensitivity to small perturbations. In contrast, our adversarial examples work because models do not realize that a small perturbation can completely change the meaning of a sentence. Images from [Szegedy et al. \(2014\)](#).

Adversarial Examples for Evaluating RC Systems

Article: **Nikola Tesla**

Paragraph: "In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for *Prague* where he was to study. Unfortunately, he arrived too late to enroll at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses."

Question: "What city did Tesla move to in 1880?"

Answer: *Prague*

Model Predicts: *Prague*

AddAny

Randomly initialize d words:

spring attention income getting reached

↓ Greedily change one word

spring attention income other reached

↓ Repeat many times

Adversary Adds: **tesla move move other george**

Model Predicts: *george*

AddSent

What city did *Tesla* move to in *1880*?

Prague

(Step 1)
Mutate
question

(Step 2)
Generate
fake answer

What city did *Tadakatsu* move to in *1881*?

Chicago

(Step 3)
Convert into
statement

Tadakatsu moved the city of *Chicago* to in *1881*.

(Step 4)
Fix errors with
crowdworkers,
verify resulting
sentences with
other crowdworkers

Adversary Adds: **Tadakatsu moved to the city of Chicago in 1881.**

Model Predicts: *Chicago*

Adversarial Examples for Evaluating RC Systems

Model	Original	ADDSSENT	ADDONESSENT
ReasoNet-E	81.1	39.4	49.8
SEDT-E	80.1	35.0	46.5
BiDAF-E	80.0	34.2	46.9
Mnemonic-E	79.1	46.2	55.3
Ruminating	78.8	37.4	47.7
jNet	78.6	37.9	47.0
Mnemonic-S	78.5	46.6	56.0
ReasoNet-S	78.2	39.4	50.3
MPCM-S	77.0	40.3	50.0
SEDT-S	76.9	33.9	44.8
RaSOR	76.2	39.5	49.5
BiDAF-S	75.5	34.3	45.7
Match-E	75.4	29.4	41.8
Match-S	71.4	27.3	39.0
DCR	69.3	37.8	45.1
Logistic	50.4	23.2	30.4