

Question Answering

Le Thanh Huong

School of Information and Communication Technology

huonglt@soict.hust.edu.vn

How can you solve the following problem?

1. Build a QA system based on a set of FAQs/user manual
2. Build a chatbot for a computer store
3. Build a QA system with the knowledge base from Wikipedia
4. Build a QA system for advising education enrollment
5. Build a QA system for study counseling

Question Answering

- Target: Build a system that can answer human questions automatically in natural language



- Sources:
 - Text passages, web documents, knowledge bases, databases, available Q&A sets
- Question types: return a value or not, open/closed domain, simple/complex questions, etc.
- Answer types: a few words, a paragraph, a list, yes/no, etc.

Sample TREC questions

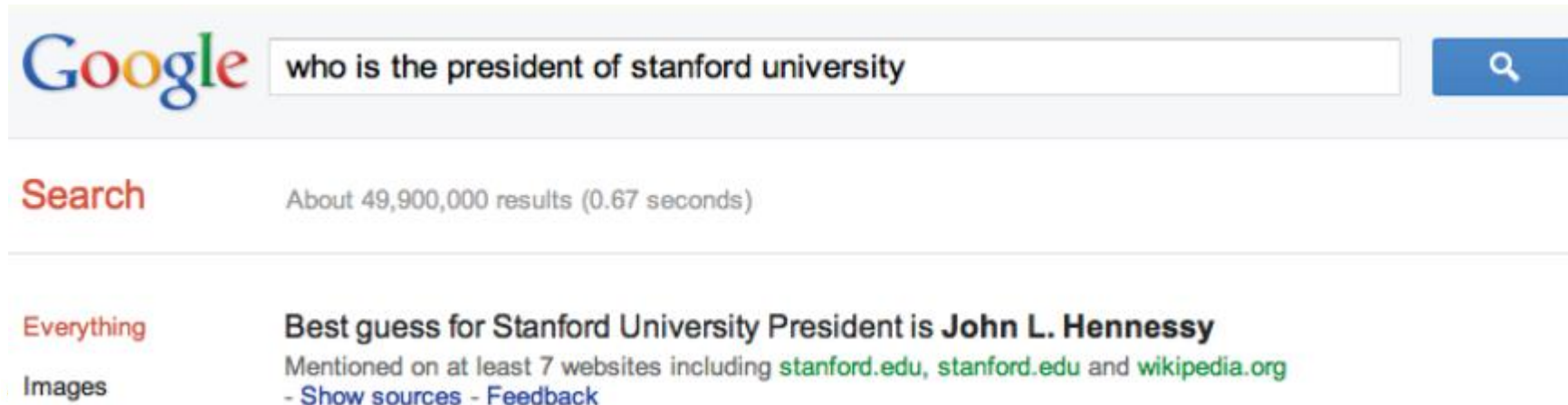
- Who is the author of the book “The Iron Lady: A Biography of Margaret Thatcher”?
- What was the monetary value of the Nobel Peace Prize in 1989?
- What does the Peugeot company manufacture?
- How much did Mercury spend on advertising in 1993?
- Why did David Koresh ask the FBI for a word processor?

People want to ask questions

- Examples from AltaVista query log (late 1990s)
 - Who invented surf music?
 - How to make stink bombs
 - Which english translation of the bible is used in official catholic liturgies?
- Examples from Excite query log (12/1999)
 - How can i find someone in Texas
 - Where can i find information on puritan religion?
 - What vacuum cleaner does Consumers Guide recommend

Online QA Examples

- LCC: http://www.languagecomputer.com/demos/question_answering/index.html
- AnswerBus is an open-domain question answering system: www.answerbus.com
- EasyAsk, AnswerLogic, AnswerFriend, Start, Quasm, Mulder, Webclopedia, TextMap, etc.
- Google



AskJeeves

- ...is most hyped example of QA
- ...does pattern matching to match your question to their own knowledge base of questions
 - If that works, you get the human-curated answers to that known question
 - If that fails, return regular web search
- A potentially interested middle ground, but a weak shadow of real QA

The screenshot shows the AskJeeves website interface. At the top, there's a search bar with the URL `uk.ask.com/web?qsrc=1&o=0&l=dir&q=who+is+the+president+of+The+United+States++2012&dm=all`. Below the search bar, there's a navigation bar with links for "Answers", "Advanced Search", "Settings", and "Your Cookie Choices". The main search bar contains the text "who is the president of The United States 2012" and a blue "Find Answers" button. To the left of the search bar, there's a sidebar with the AskJeeves logo and a list of categories: "Everything", "Images", "Video", "Reference", and "Q&A". Below the sidebar, there's a section titled "Explore Answers About" with links to "United States History Timeline", "United States Atlas", "United States Road Map", "United States Facts", "Visitors Visa Requirements United States America", and "A List of Presidents in Order". To the right of the search bar, there's a section titled "Popular Q&A" with two questions and answers. The first question is "Q: Who the president of the united states 2012?" and the answer is "A: President Barack Obama has won re-election and will serve 4 more years a... Read More »" with a source link to www.chacha.com. The second question is "Q: Who is vice president of united states 2012?" and the answer is "A: The vice president of The United States of America in 2012 is Joe Biden. Read More »" with a source link to wiki.answers.com.

Answers

Advanced Search Settings Your Cookie Choices

AskJeeves

who is the president of The United States 2012

Find Answers

The We
UK Onl

Explore Answers About

Everything ▶

Images

Video

Reference

Q&A

United States History Timeline

United States Atlas

United States Road Map

United States Facts

Visitors Visa Requirements United States America

A List of Presidents in Order

Popular Q&A

Q: Who the president of the united states 2012?

A: President Barack Obama has won re-election and will serve 4 more years a... [Read More »](#)

Source: www.chacha.com

Q: Who is vice president of united states 2012?

A: The vice president of The United States of America in 2012 is Joe Biden. [Read More »](#)

Source: wiki.answers.com



TEXTMAP
THE ENTITY SEARCH ENGINE

Monitoring the World So You Don't Have To ...



ENTITIES

SOL

Search!

[TextMap](#) : [TextMed](#) : [Textblq](#) : [TextBiz](#) : [Make homepage!](#) : [Link to us](#) : [Help?](#)

Question Answering

Wednesda

in what year did John Lennon die?

Answer: 1980

[[The Beatles Anthology](#) 02/28/2006 [wiki](#)]



TEXTMAP

THE ENTITY SEARCH ENGINE

Monitoring the World So You Don't Have To ...

SOURCES

CONTACT

who is the Prime Minister of vietnam

Search!

☒ TextMap ☐ All Sources

[TextMap](#) : [TextMed](#) : [TextBlg](#) : [TextBiz](#) : [Make homepage!](#) : [Link to us](#) : [Help?](#)

Search Results 1-25 of about 330,000

[Next >>](#)

Rank	Entity	Score	Type	Popularity	Top Month for Query
1	Vietnam	<div><div></div></div>	COUNTRY	<div><div></div></div>	November 2006
2	Iraq	<div><div></div></div>	COUNTRY	<div><div></div></div>	November 2006
3	Tony Blair	<div><div></div></div>	PERSON	<div><div></div></div>	May 2007

Search!

☒ TextMap ☐ All Sources[TextMap](#) : [TextMed](#) : [Textblq](#) : [TextBiz](#) : [Make homepage!](#) : [Link to us](#) : [Help?](#)

Vietnam COUNTRY

Sentiment Score: 67.3 + 21.9

Articles Referencing Vietnam [\[More Articles\]](#) [\(What is this?\)](#)

Title[VA hospital honors veterans with carnival](#)[Lead-tainted toys recalled](#)[Homemade explosives found in Fife](#)[Thompson is ho-hum in debate debut](#)[Central America faces new test in Asia](#)[Bush's fear factor](#)[Two doctors blame boot camp death on sickle cell](#)**Relational Network:** [\(What is this?\)](#)

Referen

News S

Sentim

Top Performing Systems

- ...can answer ~70% of the questions
- Approaches:
 - Knowledge-rich approaches, using many NLP techniques (Harabagiu, Moldovan et al.-SMU/UTD/LCC)
 - AskMRS: shallow approach
 - Middle ground use large collection of surface matching patterns (ISI)
 - ...

AskMRS: shallow approach

- In what year did Abraham Lincoln die?
- Ignore hard documents and find easy ones

Abraham Lincoln, 1809-1865

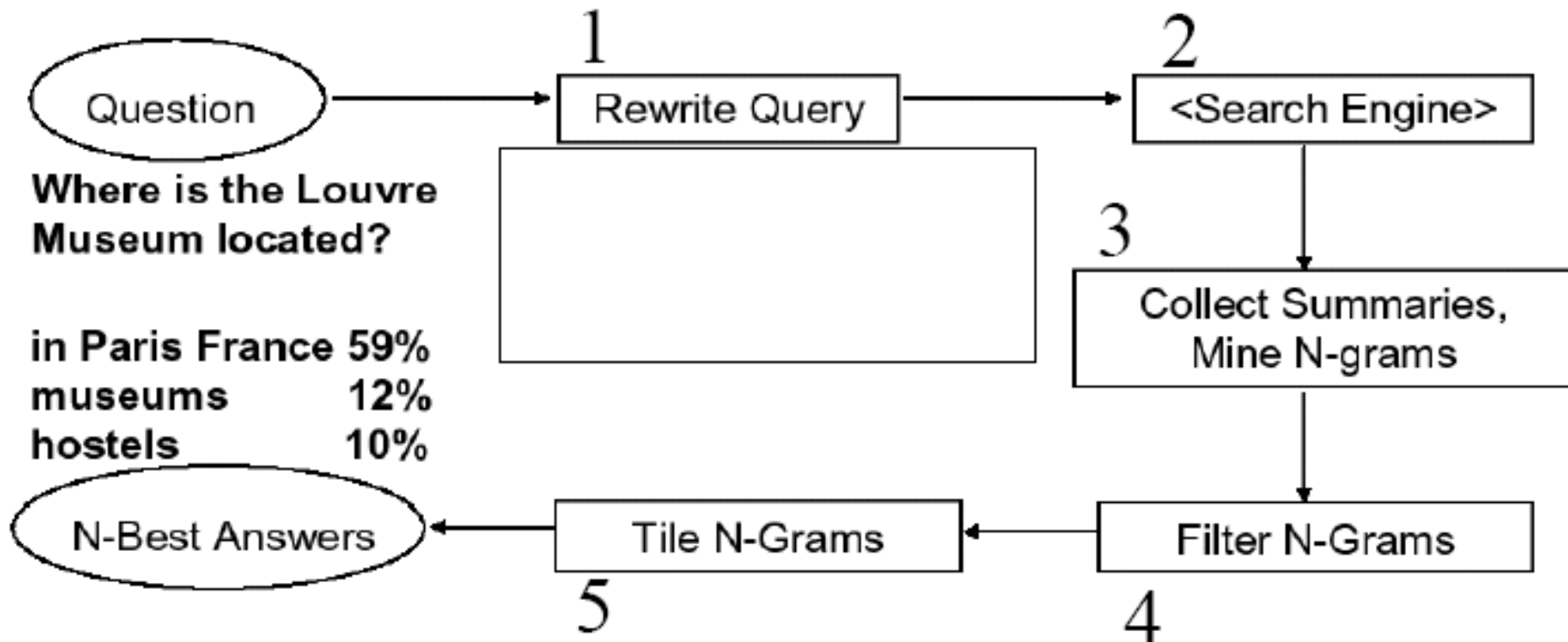
"LINCOLN, ABRAHAM was born near Hodgenville, Kentucky, on February 12, 1809. In 1816, the Lincoln family moved to Pigeon Creek in Perry (now Spencer) County. Two years later, Abraham Lincoln's mother died and his father married a woman his 'angel' mother. Lincoln attended a formal school for only a few months but acquired knowledge through the reading of books. In 1830 where he obtained a job as a store clerk and the local postmaster. He served without distinction in the Black Hawk War, lost his attempt at the state legislature, but two years later he tried again, was successful, and Lincoln was admitted to the bar and became noteworthy as a witty, honest, competent lawyer. In 1846, at which time he opposed the war with Mexico. By 1858, he had gained national attention for his series of debates with Stephen A. Douglas. In 1860, he lost the election but became a significant figure in his party. On March 4, seven southern states had seceded. Lincoln called for 75,000 volunteers (approximately 11,000 responded). Lincoln immediately took action. The Emancipation Proclamation which expanded the purpose of the war. The dedication of a national cemetery in Gettysburg, Lincoln explained the meaning of the war. He died of disease during the war. He was General of the Army. He was President at the time of the war.

Sixteenth President
1861-1865
Married to Mary Todd Lincoln

Abraham Lincoln
16th President of the United States (March 4, 1861 to April 15, 1865)
Born: February 12, 1809, in Hardin County, Kentucky
Died: April 15, 1865, at Petersen's Boarding House in Washington, D.C.

"I was born February 12, 1809, in Hardin County, Kentucky. My parents were both born in Virginia, of undistinguished families, perhaps I should say. My mother, who died in my tenth year, was of a family of the name of Lincoln."

AskMSR: Details



Step 1: Rewrite queries

- Intuition: The user's question is often syntactically quite close to sentences that contain the answer
 - Where is the Louvre Museum located?
 - The Louvre Museum is located in *Paris*
 - Who created the character of Scroogle?
 - *Charles Dickens* created the character of Scrooge.

Query rewriting

- Classify question into 7 categories
 - Who is/was/are/were...?
 - When is/did/will/are/were...?
 - Where is/are/were...?
- a) Category-specific transformation rules
 - E.g., For Where question, move “is” to all possible locations
 - Where **is** the Louvre Museum located?
 - **is** the Louvre Museum located?
 - the **is** Louvre Museum located?
 - the Louvre **is** Museum located?
 - the Louvre Museum **is** located?
 - the Louvre Museum located **is**?
- b) Expected answer “Datatype” (eg, Date, Person, Location,...)
 - When was the French Revolution? → DATE
- Hand-crafted classification/rewrite/datatype rules

Query Rewriting - weights

- Some query rewrites are more reliable than others

Where is the Louvre Museum located?

Weight 1

Lots of non-answers
could come back too

Weight 5

if we get a match,
it's probably right

+“the Louvre Museum is located”

+Louvre +Museum +located

Step 2: Query search engine

- Send all rewrites to a Web search engine
- Retrieve top N answers (100?)
- Rely just on search engine's words/phrases, not the full text of the actual document

Step 3: Mining N-Grams

- Unigram, bigram, trigram, ..., N-gram: list of N adjacent term in a sequence
 - Eg. “Web Question Answering: Is More Always Better”
 - Unigram: Web, Question, Answering, Is, More, Always, Better
 - Bigram: Web Question, Question Answering, Answering Is, Is More, More Always, Always Better
 - Trigram: ...

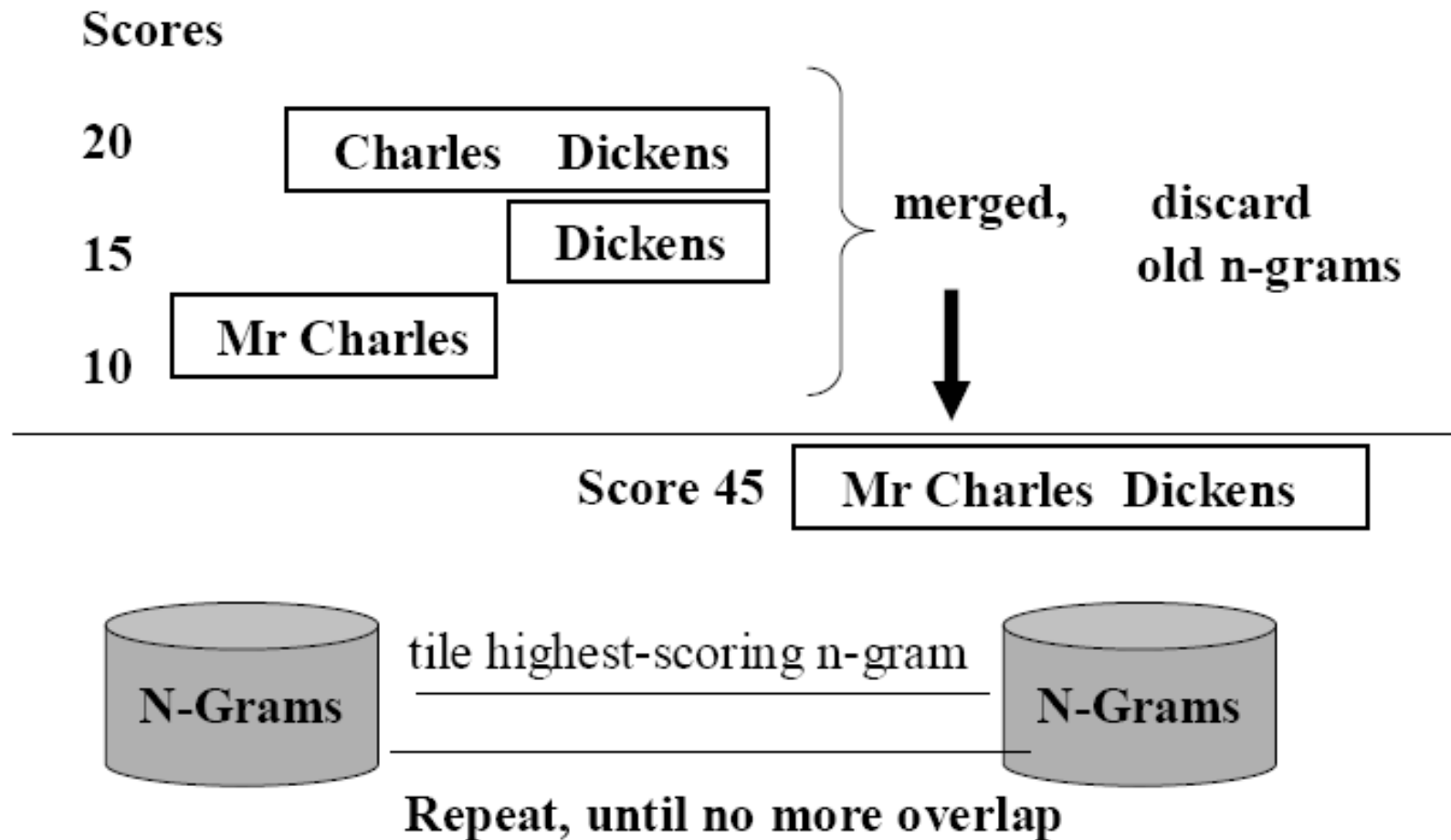
Mining N-grams

- Simple: Enumerate all N-grams ($N=1,2,3\dots$) in all retrieved phrases
 - Use hash table and other tools to make this efficient
- Weight of an n-gram: occurrence count
 - Eg, “Who created the character of Scrooge?”
 - Dickens – 117
 - Christmas Carol – 78
 - Charles Dickens – 75
 - Disney – 72
 - Carl Banks – 54
 - A Christmas – 41
 - Christmas Carol - 45

Step 4: Filtering N-Grams

- Each question type is associated with one or more “data-type filters” = regular expression
- When... Date
- Where... Location
- What... Person
- Who...
- Boost score of n-grams that do match regexp
- Lower score of n-grams that don't match regexp

Step 5: Tiling the Answers



Results

- Standard TREC contest test-bed:
 - ~1M documents; 900 questions
 - Technique doesn't do well (but rank in top 9/30 participants!)
- Limitation:
 - Works best only for fact-based questions
 - Limited range of
 - Question categories
 - Answer data types
 - Query rewriting rules

Surface matching patterns (Ravichandran and Hovy, ISI)

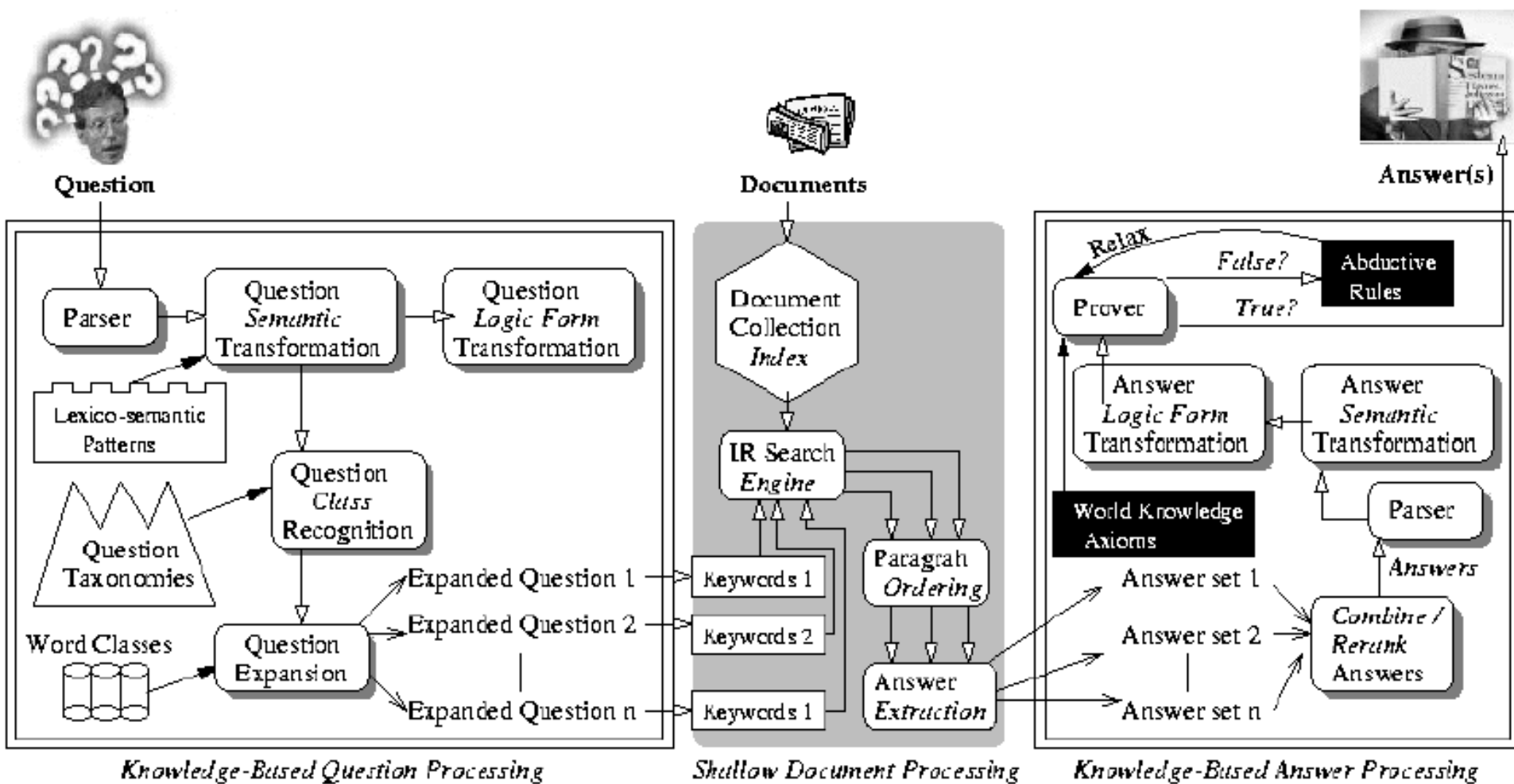
- When was X born?
 - Mozart was born in 1756
 - Gandhi (1869—1948)
- <NAME> was born in <BIRTHDATE>
- <NAME> (<BIRTHDATE>-
- Use a Q-A pair to query a search engine
- Extract patterns and compute their accuracy

Example: INVENTOR

- <ANSWER> invents <NAME>
 - the <NAME> was invented by <ANSWER>
 - <ANSWER> invented the <NAME> in
 - <ANSWER>'s invention of the <NAME>
 - ...
-
- Many of these patterns have high accuracy
 - But still some mistakes

Full NLP QA

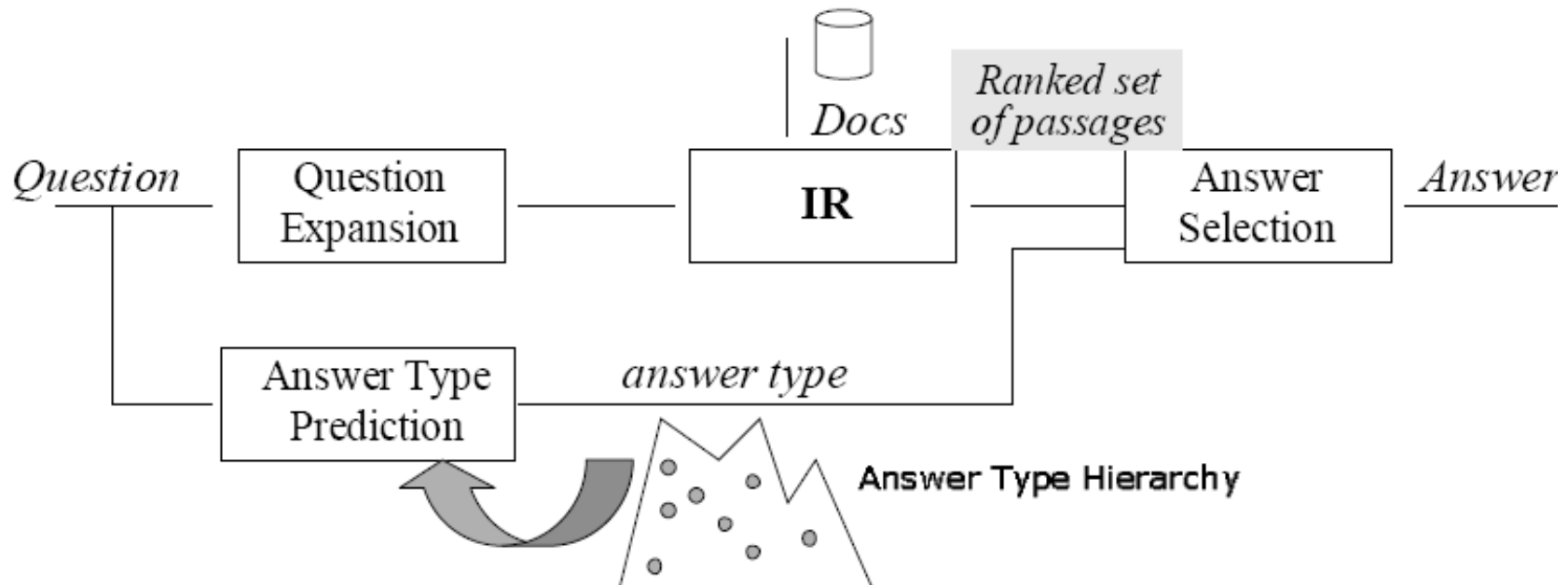
LCC: Harabagiu, Moldovan et al.



Value from sophisticated NLP – Pasca & Harabagiu (2001)

- Good IR is needed: SMART paragraph retrieval
- Large taxonomy of question types and expected answer types is crucial
- Statistical parser is used to parse questions and relevant text for answers, and to build KB
- Query expansion loops (morphological, lexical synonyms, and semantic relations) important
- Answer ranking by simply ML method

Answer types in State-of-the-art QA systems



Features:

- Answer type:
 - Labels questions with answer type based on a taxonomy
 - Classifies questions (eg., by using a maximum entropy model)

Answer Types

- “Who” questions can have organizations as answers
 - Who sells the most hybrid cars?
- “Which” questions can have people as answers
 - Which president went to war with Mexico?

Keyword Selection Algorithm

Select all...

- Non-stopwords in quotations
- NNP words in recognized named entities
- Complex nominals with their adjectival modifiers
- Other complex nominals
- Nouns with adjectival modifiers
- Other nouns
- Verbs
- The answer type word

Passage Extraction Loop

- Passage Extraction Component
 - Extracts passages that contain all selected keywords
 - Passage size/start position dynamic
- Passage quality and keyword adjustment
 - 1st iteration: use the first 6 keyword selection heuristics
 - If $\# \text{passages} < \theta \rightarrow$ query is too strict \rightarrow drop a keyword
 - If $\# \text{passages} > \theta \rightarrow$ query is too relaxed \rightarrow add a keyword

Passage Scoring

Involve 3 scores:

- #words from the question that are recognized in the same sequence in the window
- #words that separate the most distant keywords in the window
- #unmatched keywords in the window

Rank candidate answers in the retrieved passages

- Name the first private citizen to fly in space
- Answer type: Person
- Text passage:

“Among them was Christa McAuliffe, the first private citizen to fly in space. Karen Ailen, best known for her starring role in “Raiders of the Lost Ark”, plays McAuliffe. Brian Kerwin is featured as shuttle pilot Mike Smith...”
- Best candidate answer: Christa McAuliffe

Name Entity Recognition

- Several QA systems are determined by the recognition of name entities. E.g.,
 - Give me information about some Dell laptops with prices in the range from 18 million to 22 million VND
 - Can you recommend me some tourist attraction in Hanoi?
 - Give me some information about houses for rent near HUST with the price under 5 million VND
 - Which city has the largest population in Vietnam?
- Important features:
 - Precision of recognition
 - Coverage of name classes (Producer - Branch name)
 - Mapping into concept hierarchies (computer, laptop, iPad,...)
 - Participation into semantic relations (eg, predicate-argument structures or frame semantics)

Semantics and Reasoning for QA: Predicate-argument structure

- *When was Microsoft established?*

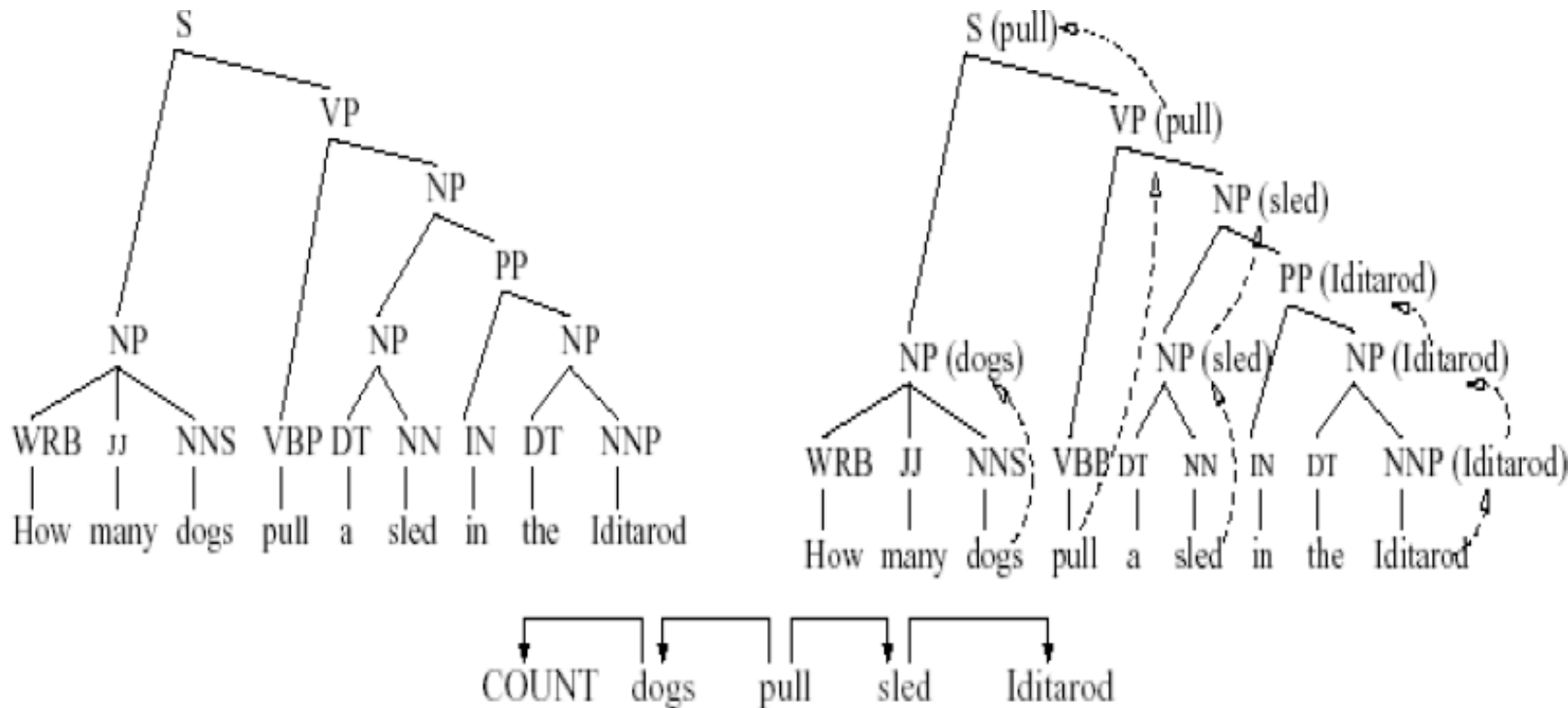
Microsoft plans to establish manufacturing partnerships in Brazil and Mexico in May.

- Need to be able to detect sentences in which ‘Microsoft’ is object of ‘establish’ or close synonym.
- Matching sentence:

Microsoft Corp was founded in the US in 1975, incorporated in 1981, and established in the UK in 1982.

- Require analysis of sentence syntax/ semantics

Semantics and Reasoning for QA: Syntax to Logical Forms

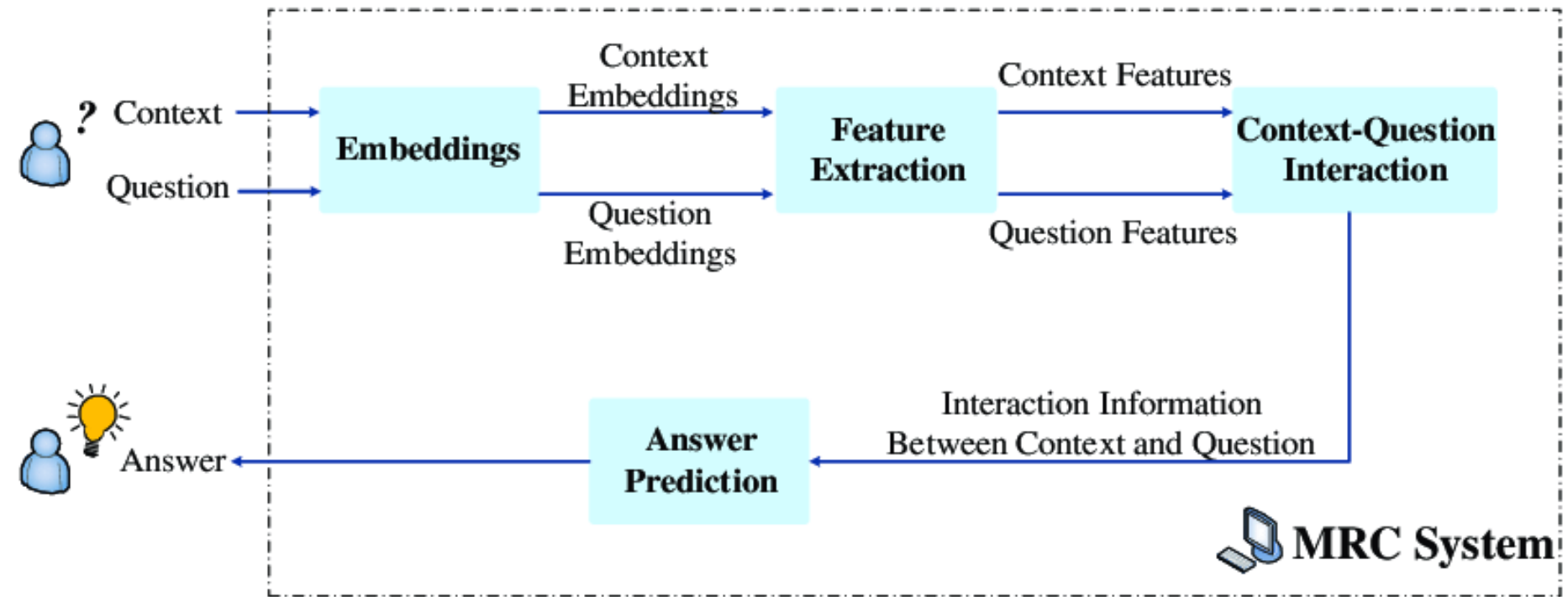


- Syntactic analysis plus semantic → logical form
- Mapping of question and potential answer LFs to find the best match

Inference

- System attempts inference to justify an answer (often following lexical chains)
- Their inference is a middle ground between logic and pattern matching
- But very effective: 30% improvement
- Q: When was the internal combustion engine invented?
- A: The first internal-combustion engine was built in 1867.
- Invent → create_mentally → create → build

Neural models for reading comprehension



Stanford question answering dataset (SQuAD)

- 100k annotated (passage, question, answer) triples 16 Large-scale supervised datasets are also a key ingredient for training effective neural models for reading comprehension! This is a limitation— not all the questions can be answered in this way!
- Passages are selected from English Wikipedia, usually 100~150 words.
- Questions are crowd-sourced.
- Each answer is a short segment of text (or span) in the passage.
- SQuAD was for years the most popular reading comprehension dataset; it is “almost solved” today (though the underlying task is not,) and the state-of-the-art exceeds the estimated human performance.

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

(Rajpurkar et al., 2016): SQuAD: 100,000+ Questions for Machine Comprehension

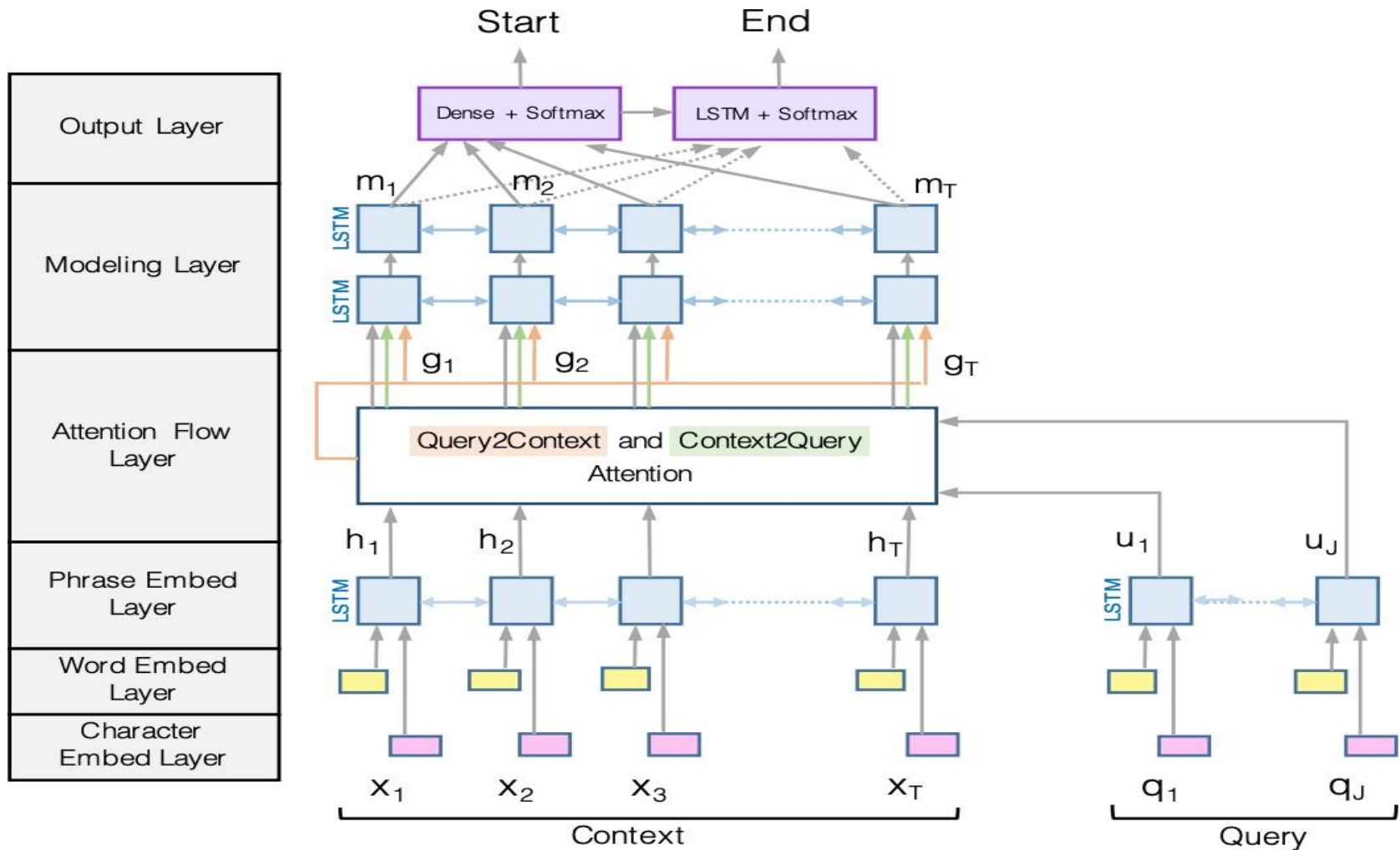
Stanford question answering dataset (SQuAD)

- Evaluation: exact match (0 or 1) and F1 (partial credit).
- For development and testing sets, 3 gold answers are collected, because there could be multiple plausible answers.
- We compare the predicted answer to each gold answer (a, an, the, punctuations are removed) and take max scores. Finally, we take the average of all the examples for both exact match and F1.
- Estimated human performance: EM = 82.3, F1 = 91.2
 - Q: What did Tesla do in December 1878?
 - A: {left Graz, left Graz, left Graz and severed all relations with his family}
 - Prediction: {left Graz and served}
 - Exact match: $\max\{0, 0, 0\} = 0$
 - F1: $\max\{0.67, 0.67, 0.61\} = 0.6$

Other question answering datasets

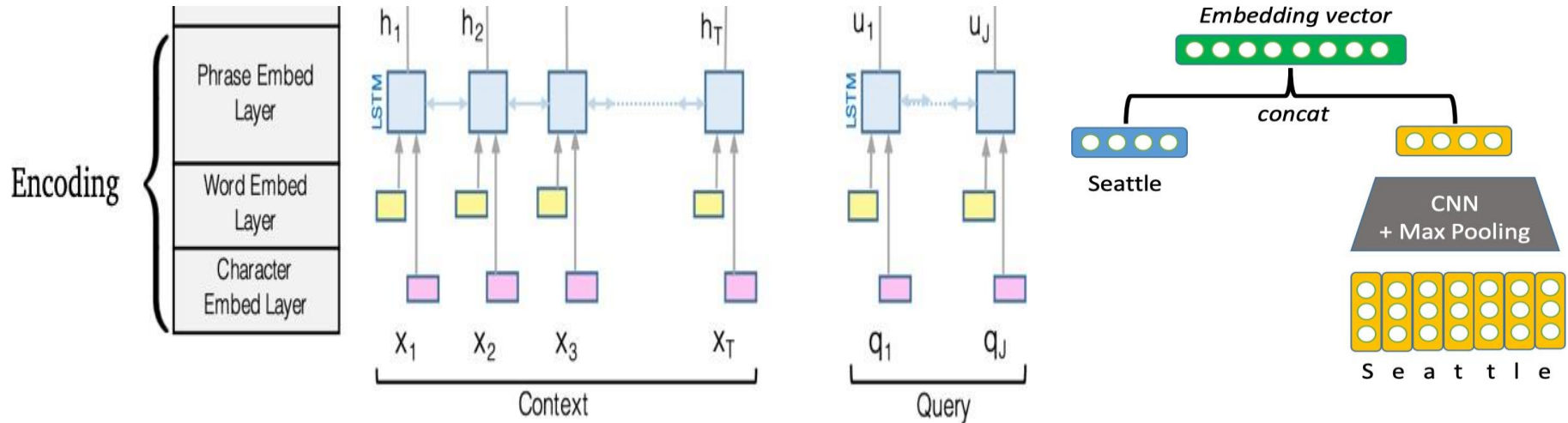
- **TriviaQA:** Questions and answers by trivia enthusiasts. Independently collected web paragraphs that contain the answer and seem to discuss question, but no human verification that paragraph supports answer to question
- **Natural Questions:** Question drawn from frequently asked Google search questions. Answers from Wikipedia paragraphs. Answer can be substring, yes, no, or NOT_PRESENT. Verified by human annotation.
- **HotpotQA.** Constructed questions to be answered from the whole of Wikipedia which involve getting information from two pages to answer a multistep query:
 - Q: Which novel by the author of “Armada” will be adapted as a feature film by Steven Spielberg?
 - A: Ready Player One

BiDAF: the Bidirectional Attention Flow model



(Seo et al., 2017): Bidirectional Attention Flow for Machine Comprehension

BiDAF: Encoding



- Use a concatenation of word embedding (GloVe) and character embedding (CNNs over character embeddings) for each word in context and query

$$e(c_i) = f([\text{GloVe}(c_i); \text{charEmb}(c_i)]) \quad e(q_i) = f([\text{GloVe}(q_i); \text{charEmb}(q_i)])$$

- Then, use two bidirectional LSTMs separately to produce contextual embeddings for both context and query

$$\vec{c}_i = \text{LSTM}(\vec{c}_{i-1}, e(c_i)) \in \mathbb{R}^H$$

$$\overleftarrow{c}_i = \text{LSTM}(\overleftarrow{c}_{i+1}, e(c_i)) \in \mathbb{R}^H$$

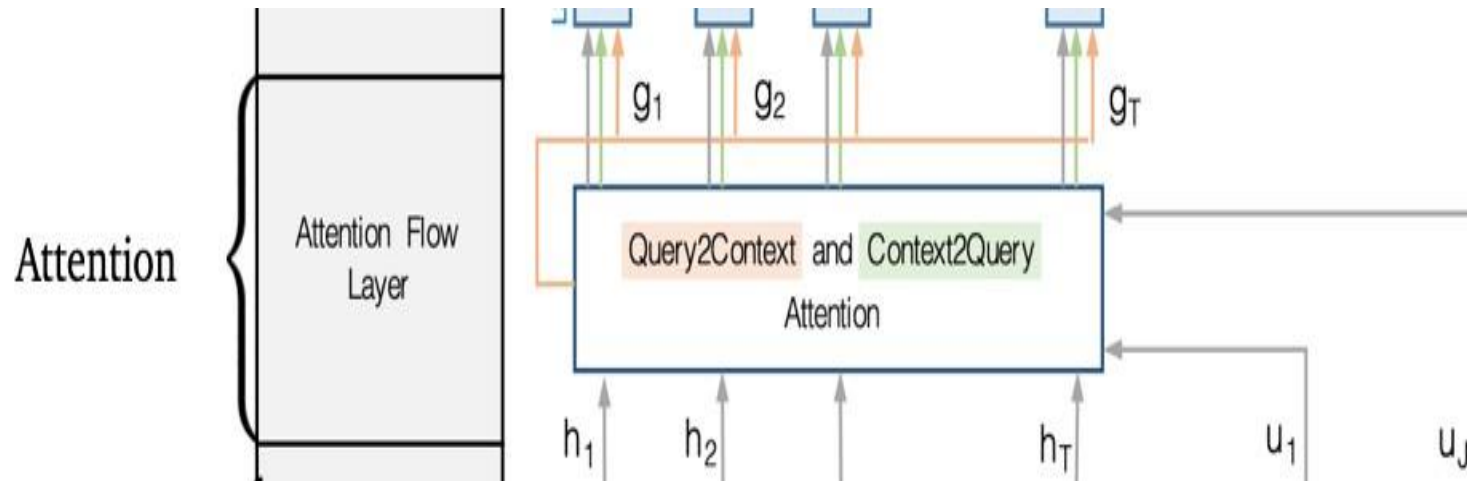
$$\mathbf{c}_i = [\vec{c}_i; \overleftarrow{c}_i] \in \mathbb{R}^{2H}$$

$$\vec{q}_i = \text{LSTM}(\vec{q}_{i-1}, e(q_i)) \in \mathbb{R}^H$$

$$\overleftarrow{q}_i = \text{LSTM}(\overleftarrow{q}_{i+1}, e(q_i)) \in \mathbb{R}^H$$

$$\mathbf{q}_i = [\vec{q}_i; \overleftarrow{q}_i] \in \mathbb{R}^{2H}$$

BiDAF: Attention



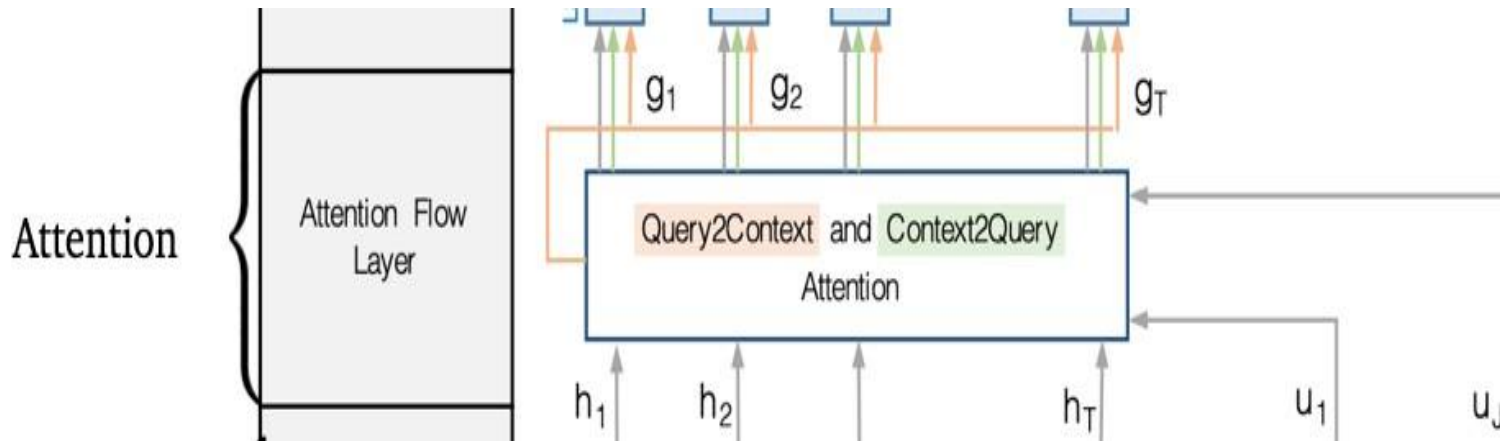
- Context-to-query attention: For each context word, choose the most relevant words from the query words

Q: *Who leads the United States?*

C: *Barak Obama is the president of the USA.*

(Slides adapted from Minjoon Seo)

BiDAF: Attention



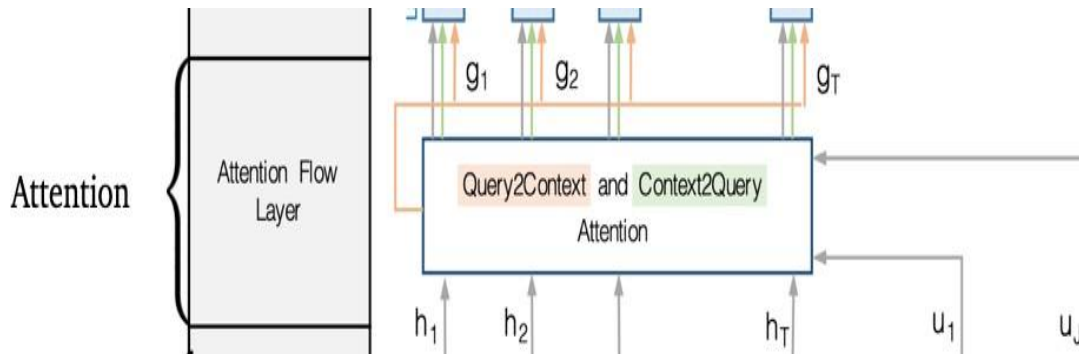
- Query-to-context attention: choose the context words that are most relevant to one of query words

While **Seattle**'s weather is very nice in summer, its weather is very rainy in winter, making it one of the most gloomy cities in the U.S. LA is ...

Q: Which city is gloomy in winter?

(Slides adapted from Minjoon Seo)

BiDAF: Attention



The final output is

$$\mathbf{g}_i = [\mathbf{c}_i; \mathbf{a}_i; \mathbf{c}_i \odot \mathbf{a}_i; \mathbf{c}_i \odot \mathbf{b}] \in \mathbb{R}^{8H}$$

- First, compute a similarity score for every pair of $(\mathbf{c}_i, \mathbf{q}_j)$:

$$S_{i,j} = \mathbf{w}_{\text{sim}}^T [\mathbf{c}_i; \mathbf{q}_j; \mathbf{c}_i \odot \mathbf{q}_j] \in \mathbb{R} \quad \mathbf{w}_{\text{sim}} \in \mathbb{R}^{6H}$$

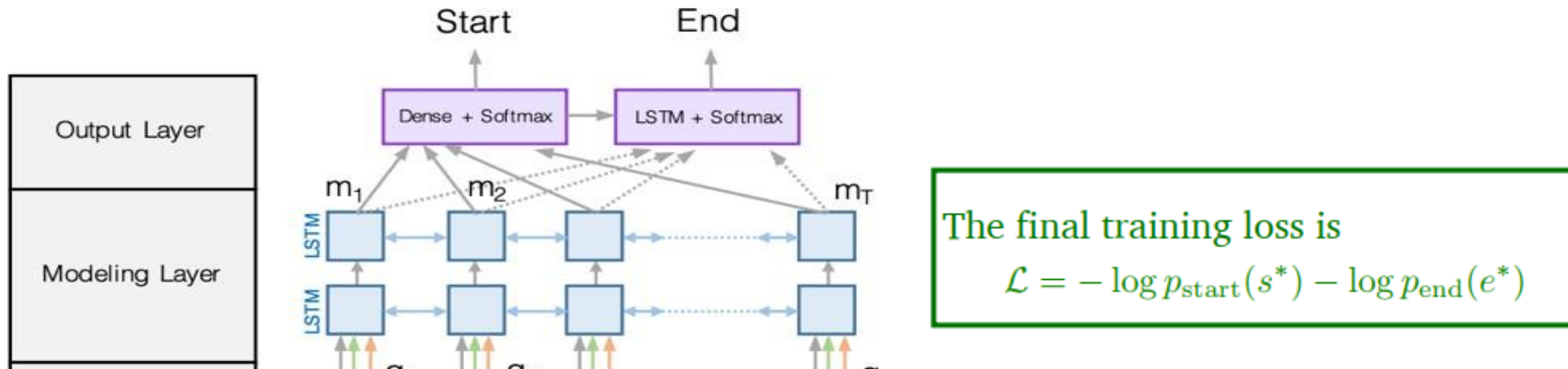
- Context-to-query attention (which question words are more relevant to \mathbf{c}_i)

$$\alpha_{i,j} = \text{softmax}_j(S_{i,j}) \in \mathbb{R} \quad \mathbf{a}_i = \sum_{j=1}^M \alpha_{i,j} \mathbf{q}_j \in \mathbb{R}^{2H}$$

- Query-to-context attention (which context words are relevant to some question words):

$$\beta_i = \text{softmax}_i(\max_{j=1}^M (S_{i,j})) \in \mathbb{R}^N \quad \mathbf{b} = \sum_{i=1}^N \beta_i \mathbf{c}_i \in \mathbb{R}^{2H}$$

BiDAF: Modeling and output layers



Modeling layer: pass to another two layers of bi-directional LSTMs

- Attention layer is modeling interactions between query and context
- Modeling layer is modeling interactions within context words

$$\mathbf{m}_i = \text{BiLSTM}(\mathbf{g}_i) \in \mathbb{R}^{2H}$$

Output layer: two classifiers predicting the start and end positions:

$$p_{\text{start}} = \text{softmax}(\mathbf{w}_{\text{start}}^T [\mathbf{g}_i; \mathbf{m}_i]) \quad p_{\text{end}} = \text{softmax}(\mathbf{w}_{\text{end}}^T [\mathbf{g}_i; \mathbf{m}'_i])$$

$$\mathbf{m}'_i = \text{BiLSTM}(\mathbf{m}_i) \in \mathbb{R}^{2H} \quad \mathbf{w}_{\text{start}}, \mathbf{w}_{\text{end}} \in \mathbb{R}^{10H}$$

BiDAF: Performance on SQuAD

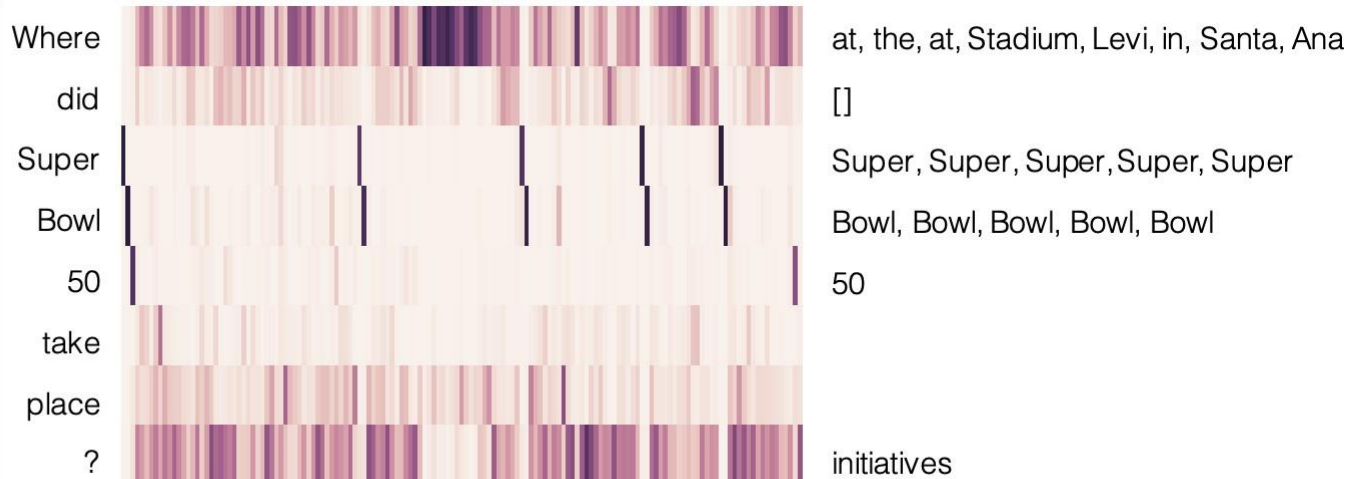
- F1=77.3 % on SQuAD v1.1.
- Without context-to-query 67.7 F1 attention
- Without query-to-context 73.7 F1 attention
- Without character embedding 75.4 F1

Single Model	Published ¹² EM / F1	LeaderBoard ¹³ EM / F1
LR Baseline (Rajpurkar et al., 2016)	40.4 / 51.0	40.4 / 51.0
Dynamic Chunk Reader (Yu et al., 2016)	62.5 / 71.0	62.5 / 71.0
Match-LSTM with Ans-Ptr (Wang & Jiang, 2016)	64.7 / 73.7	64.7 / 73.7
Multi-Perspective Matching (Wang et al., 2016)	65.5 / 75.1	70.4 / 78.8
Dynamic Coattention Networks (Xiong et al., 2016)	66.2 / 75.9	66.2 / 75.9
FastQA (Weissenborn et al., 2017)	68.4 / 77.1	68.4 / 77.1
BiDAF (Seo et al., 2016)	68.0 / 77.3	68.0 / 77.3
SEDt (Liu et al., 2017a)	68.1 / 77.5	68.5 / 78.0
RaSoR (Lee et al., 2016)	70.8 / 78.7	69.6 / 77.7
FastQAExt (Weissenborn et al., 2017)	70.8 / 78.9	70.8 / 78.9
ReasonNet (Shen et al., 2017b)	69.1 / 78.9	70.6 / 79.4
Document Reader (Chen et al., 2017)	70.0 / 79.0	70.7 / 79.4
Ruminating Reader (Gong & Bowman, 2017)	70.6 / 79.5	70.6 / 79.5
jNet (Zhang et al., 2017)	70.6 / 79.8	70.6 / 79.8
Conductor-net	N/A	72.6 / 81.4
Interactive AoA Reader (Cui et al., 2017)	N/A	73.6 / 81.9
Reg-RaSoR	N/A	75.8 / 83.3
DCN+	N/A	74.9 / 82.8
AIR-FusionNet	N/A	76.0 / 83.9
R-Net (Wang et al., 2017)	72.3 / 80.7	76.5 / 84.3
BiDAF + Self Attention + ELMo	N/A	77.9 / 85.3
Reinforced Mnemonic Reader (Hu et al., 2017)	73.2 / 81.8	73.2 / 81.8

(Seo et al., 2017): Bidirectional Attention Flow for Machine Comprehension

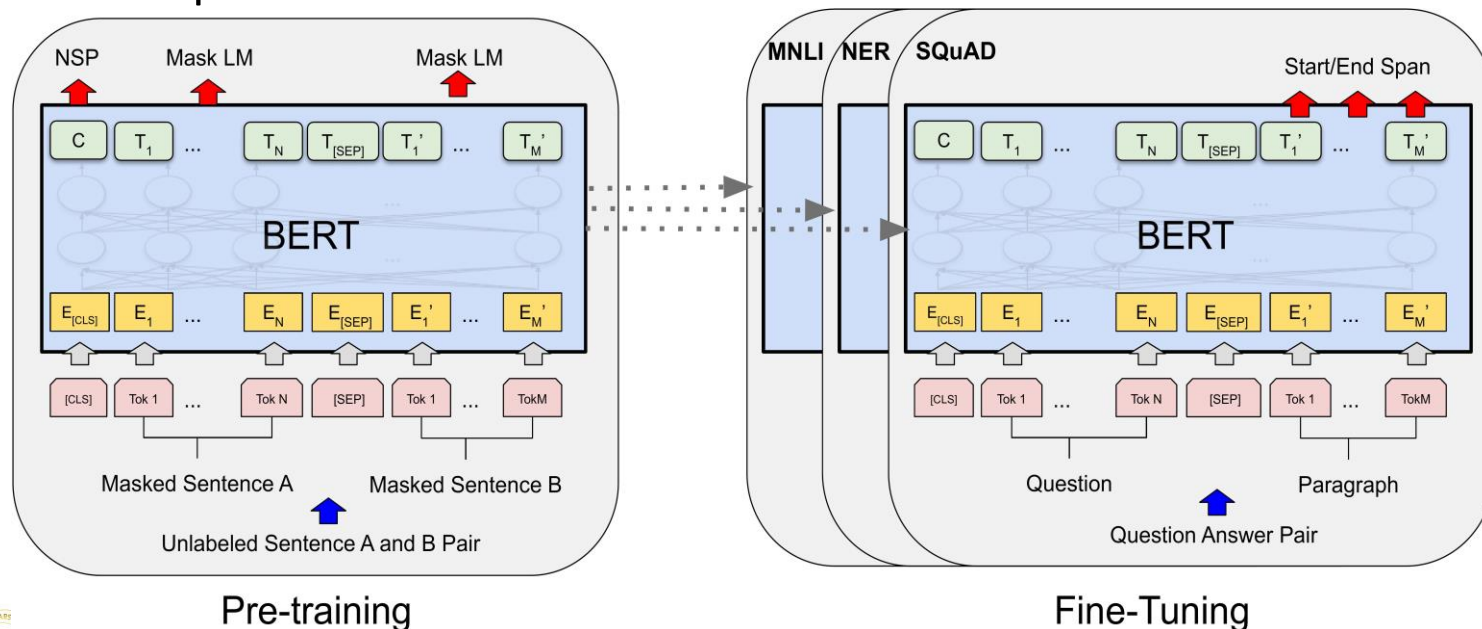
Attention visualization

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California. As this was the 50th Super Bowl, the league emphasized the "golden anniversary" with various gold-themed initiatives, as well as temporarily suspending the tradition of naming each Super Bowl game with Roman numerals (under which the game would have been known as "Super Bowl L"), so that the logo could prominently feature the Arabic numerals 50.



BERT for reading comprehension

- BERT is a deep bidirectional Transformer encoder pre-trained on large amounts of text (Wikipedia + BooksCorpus)
- BERT is pre-trained on two training objectives:
 - Masked language model (MLM)
 - Next sentence prediction (NSP)
- BERTbase has 12 layers and 110M parameters, BERTlarge has 24 layers and 330M parameters



Transformer

- A deep learning model that transforms an input sequence to an output sequence using the encoder – decoder architecture
- Multi-head attention
- Feed forward layers
- Positional embeddings

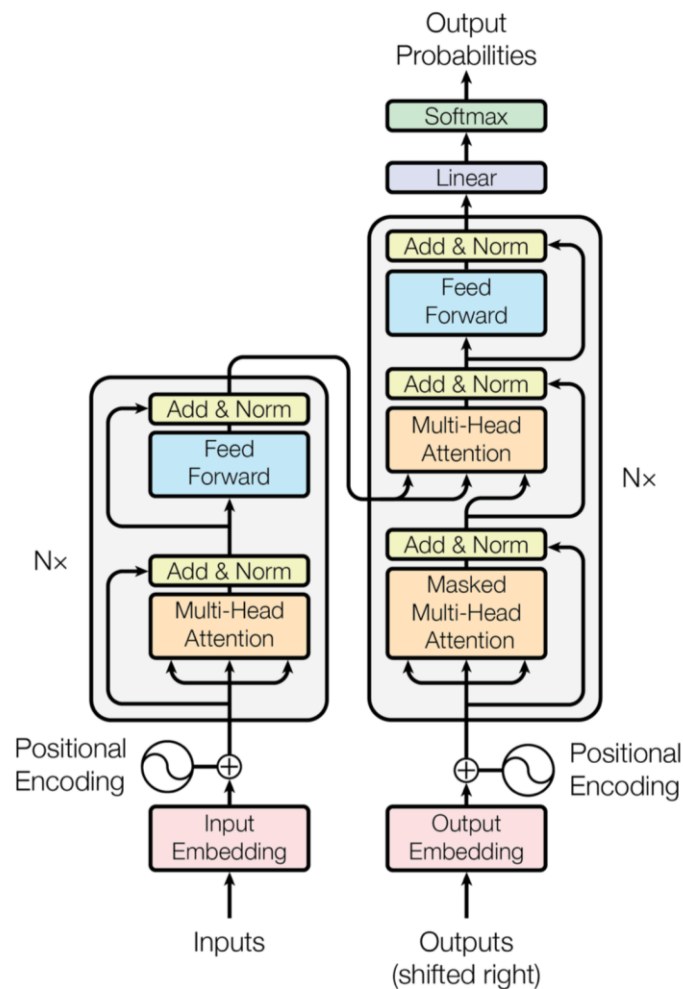
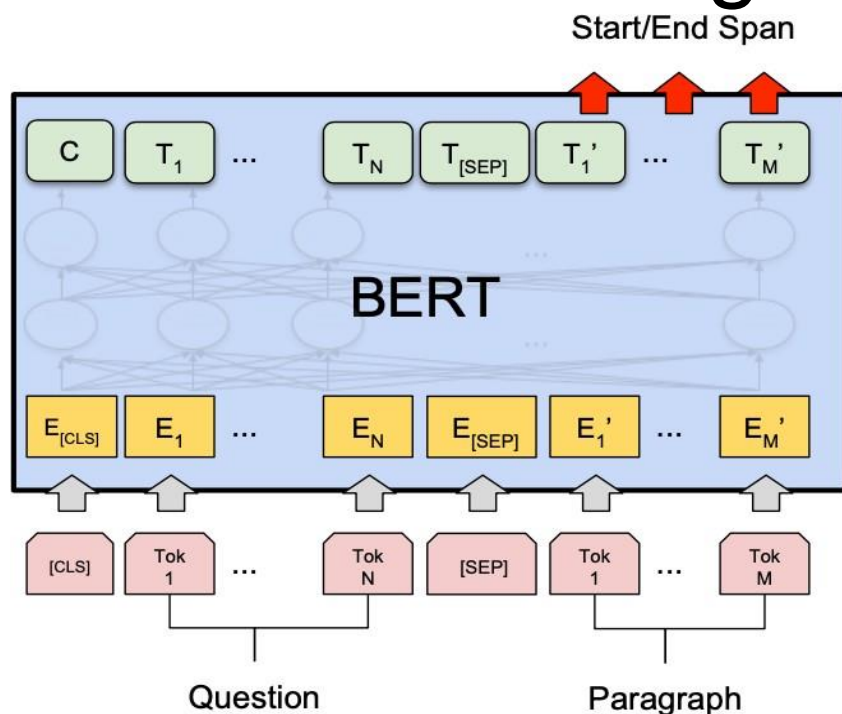


Figure 1: The Transformer - model architecture.

<https://medium.com/inside-machine-learning/what-is-a-transformer-d07dd1fbec04>

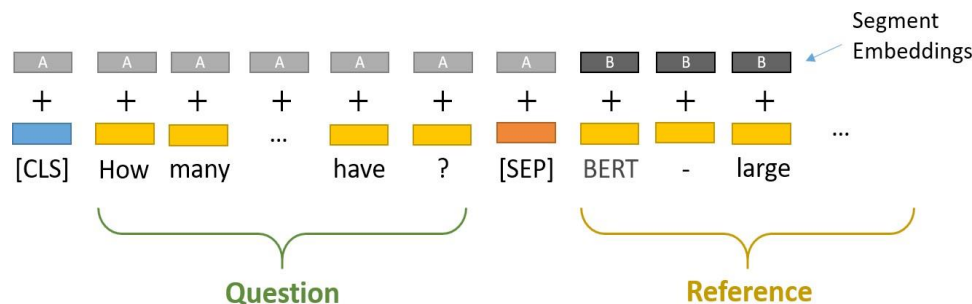
BERT for reading comprehension



Question = Segment A

Passage = Segment B

Answer = predicting two endpoints in segment B



Question: How many parameters does BERT-large have?

Reference Text: BERT-large is really big... it has 24 layers and an embedding size of 1,024, for a total of 340M parameters! Altogether it is 1.34GB, so expect it to take a couple minutes to download to your Colab instance.

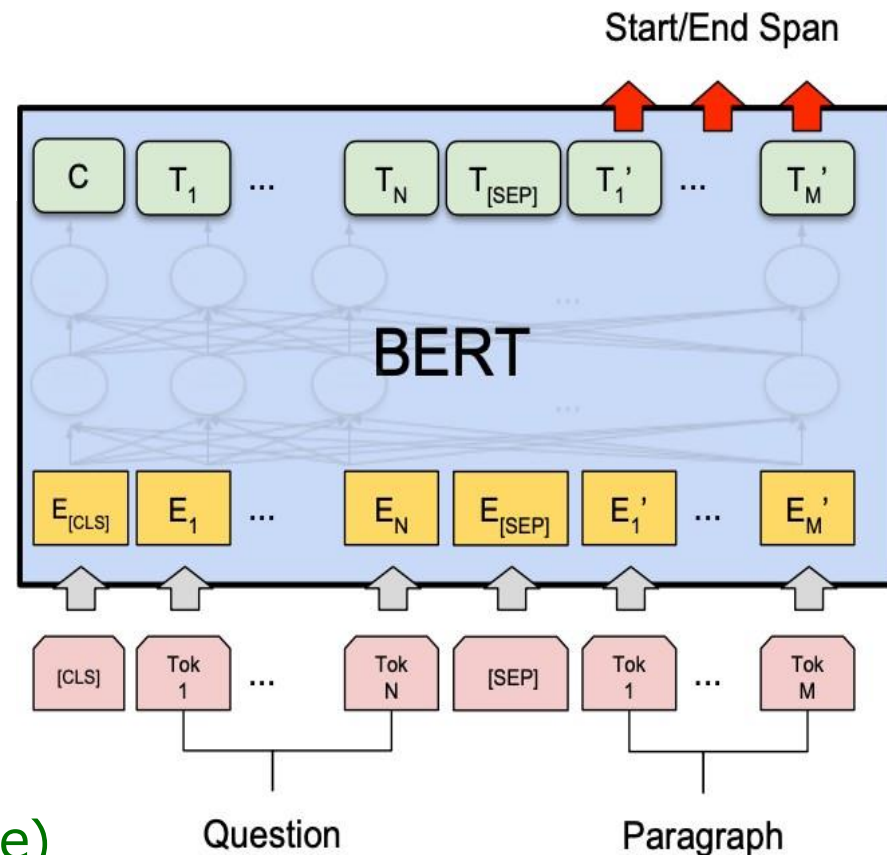
Image credit: <https://mccormickml.com/>

BERT for reading comprehension

$$\mathcal{L} = -\log p_{\text{start}}(s^*) - \log p_{\text{end}}(e^*)$$

	F1	EM
Human performance	91.2	82.3
BiDAF	77.3	67.7
BERT-base	88.5	80.8
BERT-large	90.9	84.1
XLNet	94.5	89.0
RoBERTa	94.6	88.9
ALBERT	94.8	89.3

(dev set, except for human performance)



ChatGPT

- Launched in 11/2022 by OpenAI
- ChatGPT stands for Generative Pre-Trained Transformer. ChatGPT is powered by LLM(Large Language Model) → it can resemble human-like responses
- Data: from online databases, containing books, webtexts, Wikipedia, articles, and other online literature.
- Usages: chatbots, virtual assistants, other applications requiring a high level of natural language processing ability.
- Anyone can use Chat GPT by integrating it into their own apps or by using one of the many prebuilt chatbot platforms that incorporate the technology.

COMPARING GENERATIVE AI

Google Bard AI

- Answers real-time queries
- Regular Google search results
- Based on LaMDA
- No plagiarism detector
- Free for now

ChatGPT

- Answers are based on data recorded up to 2021
- Text-only responses
- Based on GPT
- Has plagiarism detector
- ChatGPT Plus is a paid plan

Exercise

How can you solve the following problem?

1. Build a chatbot for a computer store.
2. Build a QA system with the knowledge base from Wikipedia

Chatbot for a computer store

- Give me information about some **Dell** laptops with prices in the range from **18 million** to **22 million VND**
 - Laptop's attributes:
 - Branch name: **Dell**
 - Processor: Core i5 1135G7 2.4GHz
 - RAM: 8Gb DDR4 3200
 - Hard Drive: 256GB SSD
 - Video card: VGA onboard - Intel Iris Xe Graphics
 - Screen Size: About 14 inches
 - Operating system: NoOS
 - Price: **20 million VND**