

### **CS4248: Natural Language Processing**

Lecture 11 — Classification Applications

### **Recap of Week 10**





Linea



#### GPT — RLHF (Reinforcement Learning from Human Feedback)

- RLHF two common setups
  - Use human-generated responses to prompts to fine-tune the pretrained model
  - Generate multiple response for same prompt; human ranks response; use ranking for fine-tuning



74

### Announcements

A2 grades should be out soon.

A2 is a representation of real interest in SoC, SG and worldwide about current research in fake news.



#### Abstract

social context representation and learning framework for fake news detection. Unlike previous contextual models that have targeted performance, our focus is on representation learning. Compared to transductive models, FANG is scalable in training as it does not have to maintain the and is efficient at inference time, without the need to reprocess the entire graph. Our experimental results show that FANG is better at capturing the social context into a highfidelity representation, compared to recent graphical and nongraphical models. In particular, FANG yields significant improvements for the task of fake news detection and is robust in the case of limited training data. We further demonstrate that the representations learned by FANG generalize to related tasks, such as predicting the factuality of reporting of a news medium.

after its publication. These are mainly verbatim recircula-We propose Factual News Graph (FANG), a novel graphical tions with negative sentiment of the original post explained by the typically appalling content of fake news. After that short time window, we see denial posts questioning the validity of the news, and the stance distribution stabilizes afterwards with virtually no support. In contrast, the real news example in Table 1 leads to moderate engagement, social entities involved in the propagation of other news mainly comprised of supportive posts with neutral sentiment that stabilize quickly. Such temporal shifts in user perception serve as important signals to distinguish fake from real news.

> Previous work proposed partial representations of social context with (i) news, sources, and users as major entities and (ii) stances, friendship, and publication as major interactions.5, 16, 17, 22 However, they did not put much emphasis on the quality of the representation, on modeling the entities and their interactions, and on minimally supervised settings.



# Announcements

24th STePS is on next Wed 15:00–18:00

Come down to SoC COM3 MPH and support your fellow CS4248 teams and check out what other teams have done on their projects!



Given a collection of English sentences denoted as X, where each sentence may or may not contain

approach A linguistic approach to a robust defence

OOPSI

# Outline

### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

### • Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

### **Text Summarization**

- Text Summarization basic goal
  - Generate a condensed version of a (large) document or multiple documents
  - Summarization should convey the main idea of the original document(s) to the reader
- Wide range of applications
  - Outlines or abstracts of any document, article, etc.
  - Summaries of email threads
  - Action items from a meeting
  - Simplifying text by compressing sentences



Google's cloud unit looked into using artificial intelligence to help a financial firm decide whom to lend money to. It turned down the client's idea after weeks of internal discussions, deeming the project too ethically dicey. Google has also blocked new AI features analysing emotions, fearing cultural insensitivity. Microsoft restricted software mimicking voices and IBM rejected a client request for an advanced facial-recognition system.

wide drive to balance the pursuit of lucrative AI system with a greater consideration of social responsibility.

### **Text Summarization — Dimensions**

Input / Source

### **Single Document**

- Input: single document (e.g., news article, web page, blog post, etc.)
- Common outputs:
  - abstract
  - outline
  - headline

### **Multiple Documents**

- Input: group/set of documents
- Case 1: documents are about similar topic (e.g., multiple news stories about the same event)
  - → Output: "proper" text summary
- Case 2: documents are about diverse topics (e.g., all news stories over the course of a day)
  - → Output: clusters / categories of document (potentially with a text summary for each cluster/category)

### **Text Summarization — Dimensions**

Trigger

#### **Generic Summarization**

- "Just" summarize the content
- No additional factor/requirement/etc. driving the summarization process

### **Query-focused** Summarization

- Summarize a document with respect to an information need expressed in a **user query**
- Kind of a complex Q&A task

# **Query-Focused Summarization — Example**



https://singapur.diplo.de > sg-en > service > 15-Covid19

General Information on the Covid-19-Situation in Singapore ... Travellers from Singapore to Germany aged 6 years and above are required to have a vaccination certificate, a proof of recovery after an infection or a negative ... Vaccinated Travel Lane · Covid-19 · Vaccinated Travel Lane (VTL)

https://www.mfa.gov.sg > Germany > Travel-Page

#### Germany - Ministry of Foreign Affairs Singapore

COVID-19 Travel Restrictions. As of 3 March 2022, Singapore is no longer classified by Germany as a "high risk area", according to the list of designated ...

https://www.singaporeair.com > en\_UK > travel-info

#### Vaccinated Travel Lanes (VTL) | Singapore Airlines

Country. Germany. Covid-19 test. -. Quarantine. -. Other information. Travellers departing from Singapore must complete the Digital Entry Application before ...

#### Online search

- Summary = sentence snippets
   from the search result page
- Heuristics pick snippets that
  - Include many search terms
  - Appear early in the document
  - Have special markups (e.g., bold)

■ ...

### **Text Summarization — Dimensions**

Summarization Approach

### **Extractive Summarization**

- Summary = selected phrases or sentences from source document(s)
- No "true" text generation task
- Challenge: risk of incoherent summaries

### **Abstractive Summarization**

- Summary = newly generated text (potentially using completely different words)
- Advantage: generally much more coherent
- Challenge: generally more difficult (compared to extractive summarization)

**Note:** Both approaches can be combined, e.g: use extractive summarization to find subset of important sentences and that apply abstractive summarization over this subset.

# **Abstractive Summarization — Example**

Dow C Sand Conventio	NTOWN ORE s Expo & Control Con	arina Bay Inds Singap
Marina Bay Sands S	ingapore	
Website     Directions     Save       4.6 ★★★★     8,982 Google review	ews	
Reviews @	Write a review	Add a photo
Rooms · 4.3 ★ ★ ★ ★ Rooms had views · Guests liked t some said they were dated & main improved · Guests liked the large they could be improved	he large, clean roon ntenance could be bathrooms, though	ns, though some said
Location · 4.6 * * * * * Shopping, sightseeing, restaurant accessible by car · Near public tra	s & bars nearby · Ea	asily
Service & facilities · 4.1 * * * * Guests enjoyed the pool & fitness	centre · Guests spo	oke highly of

### Google hotel review summary

- Identification of frequent phrases (with either positive or negative sentiment)
- Display of most common phrases (potentially a canonical version of similar phrases)
- Generation of very simple sentences (e.g.: "Guest liked [...] but some said [...]")
- Sentence generation based on templates (disclaimer: my personal opinion; might be wrong!)
- Advantages
  - Simple but still appropriate results
  - "Safe" results (no risk of weird reviews)

# Outline

### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

### • Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts





Create your own baseline summarization system by specifying a simple method for each of the three steps

#### (1) Content Selection

Choose sentence (or phrases) to extract

#### (2) Information Reordering

Choose and order to place them

#### (3) Sentence Realization

Clean up sentences; finalize summary

In-Lecture Activity (7 mins)







Clean up sentences; finalize summary

(i.e., no rewriting, simplification, generation)

### **Content Selection — Baseline Algorithm**

### • Naive approach: Pick the first sentence(s)



#### → Summary:

"Singapore, officially the Republic of Singapore, is a sovereign island city-state in maritime Southeast Asia."

# **Unsupervised Content Selection**

- Core idea: Finding keywords
  - Choose sentences with many [ important / informative / salient / etc. ] words
- Various strategies proposed, e.g.:
  - tf-idf (we already know how to do this)
  - Log-likelihood ratio (LLR)
  - TextRank graph-based approach (supports keyword & sentence extractions)

# Log-Likelihood Ratio (LLR)

- Step 1: Identify salient words
  - Assign words with a minimum LLR with a positive weight
  - Option: assign words that are in the query/question with a positive weight

$$weight(w_i) = \begin{cases} 1 & \text{if } -2\log\lambda(w_i) > 10 \\ 1 & \text{if } w_i \in \text{query/question} \end{cases} \xrightarrow{\text{1. Important Words}} 2. \text{ Overlap with the query} \\ \text{In case of query-focused summarization} \end{cases}$$

- Step 2: Score sentences
  - Score of a sentence = average weight over all words in the sentence

$$weight(\mathbf{S}) = \frac{1}{|S|} \sum_{w \in S} weight(w)$$

# Log-Likelihood Ratio (LLR)

- Underlying assumption
  - **Binomial** distribution for generating *w* in a text

P(word w appears k times in a text) = 
$$b(p, k, n) = \binom{n}{k} p^k (1-p)^{n-k}$$
  
probability of w; estimate via MLE:  $p = \frac{k}{k}$ 

Log-Likelihood Ratio

probability of observing *w* in document *d* and corpus *c* assuming **equal** probabilities *p* in both *d* and *c* 

n

number of

worde in text

 $\lambda(w_i) = \frac{b(p, k_c, n_c) \cdot b(p, k_c, n_c)}{b(p_d, k_d, n_d) \cdot b(p_c, k_c, n_c)}$ 

probability of observing *w* in document *d* and corpus *c* assuming **different** probabilities  $p_d$  and  $p_c$  in *d* and *c* 



#### **Core algorithm**

- Identify meaningful text units → set of vertices V (either words or sentences depending on task)
- Identify meaningful relations between text units → set edges E (e.g.: co-occurrence of text units ot similarity between text units)
- 3) Apply graph-based ranking algorithm over G(V, E) (proposed in original paper: Weighted PageRank)
- 4) Sort vertices based on their final score

Represent text as a graph Important vertex in graph ⇔ Important text unit in document

### **TextRank**

### • Identification of keyword

- Text units = words → vertices = words
- Unweighted edge = "binary" co-occurrence

(there exists an edge between to vertices if the two corresponding words appear together within a window)

- Apply PageRank over resulting Graph
- Choose keywords with highest scores

**Note:** PageRank is defined over direct graphs, but an indirect edge can be represented as 2 directed edges.

Source: TextRank: Bringing Order into Texts

Compatibility of systems of linear constraints over the set of natural numbers. Criteria of compatibility of a system of linear Diophantine equations, strict inequations, and nonstrict inequations are considered. Upper bounds for components of a minimal set of solutions and algorithms of construction of minimal generating sets of solutions for all types of systems are given. These criteria and the corresponding algorithms for constructing a minimal supporting set of solutions can be used in solving all the considered types systems and systems of mixed types.



#### Keywords assigned by TextRank:

linear constraints; linear diophantine equations; natural numbers; nonstrict inequations; strict inequations; upper bounds

#### Keywords assigned by human annotators:

linear constraints; linear diophantine equations; minimal generating sets; nonstrict inequations; set of natural numbers; strict inequations; upper bounds

### **TextRank**

#### • Sentence extraction

- Text units = sentences → vertices = sentences
- Weighted edge = sentence similarity (e.g., Jaccar, cosine between tf-idf / embedding vectors)
- Apply PageRank over resulting Graph
- Choose sentences with highest scores

**Note:** PageRank is defined over unweighted graphs, but can trivially extended to weighted graphs.



# Quick Side Note — PageRank

- PageRank centrality measure
  - Quantifies importance of a node in a graph (pages in the Web Graph connected by links)
  - Recursive definition: A node is important if many other important nodes point to it
  - Computing PageRank = Finding the largest
     Eigenvector of a matrix derived from graph



In-Lecture Activity (2 mins)



What is the **interpretation** of the text unit (word or sentence) with the **highest** TextRank score?



# Outline

### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

### • Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

# **Evaluating Summaries — ROUGE**

- ROUGE ("Roo J" Recall Oriented Understudy for Gisting Evaluation)
  - Measure similarity between 2 texts based on n-gram overlap
  - Not as good as human evaluation shown to be a convenient proxy
- Basic procedure: Given a document d and a generated summary  $\hat{y}$ 
  - Have N humans produce a set of reference summaries  $S_H$
  - What percentage of the n-grams from the reference summaries appear in  $\hat{y}$ ?

$$\label{eq:ROUGEN} \text{ROUGE-N} = \frac{\sum\limits_{s \in S_H} \sum\limits_{g_N \in \hat{y}} min(Count(g_N, \hat{y}), Count(g_N, s))}{\sum\limits_{s \in S_H} \sum\limits_{g_n \in \hat{y}} Count(g_N, s)}$$
 specifies of the size of the n-grams to be considered

In-Lecture Activity (3 mins)



Let's practice ROUGE ("bigram style")! Calculate R-2, given the 4 summaries below:



"water spinach is a leaf vegetable commonly eaten in tropical areas of asia"

3 human-generated summaries (reference)

\*\*\* "water spinach is a semi-aquatic tropical plant grown as a vegetable"

2. "water spinach is a semi-aquatic tropical plant grown as a vegetable"

: "water spinach is a commonly eaten leaf vegetable of asia"

$$\operatorname{ROUGE-N} = \frac{\sum\limits_{s \in S_H} \sum\limits_{g_N \in \widehat{\boldsymbol{y}}} \min(Count(g_N, \hat{\boldsymbol{y}}), Count(g_N, s))}{\sum\limits_{s \in S_H} \sum\limits_{g_n \in \widehat{\boldsymbol{y}}} Count(g_N, s)}$$



In-Lecture Activity (3 mins)



#### System-generated summary

"water spinach is a leaf vegetable commonly eaten in tropical areas of asia"

#### 3 human-generated summaries (reference)

"water spinach is a semi-aquatic tropical plant grown as a vegetable" → 10 bigrams
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "
 "

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

2 water spinach is a semi-aquatic tropical plant grown as a vegetable → 10 bigrams



# **Outline**

### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

### • Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

# **Query-Focused Multidocument Summarization**



# **Sentence Simplification**

- Unsupervised approach
  - Sentence simplification by sentence trimming
  - Input: parse tree of sentence → trimmed parse tree (remove "less important" subtrees based on linguistically-motivated rules)

appositives	Rajam <del>, 28, an artist who was living at the time in Philadelphia,</del> found the inspiration in the back of city magazines.
attribution clauses	Rebels agreed to talks with government officials, international observers said Tuesday.
Prepositional phrases without named entities	The commercial fishing restrictions in Washington will not be lifted unless the salmon population increases to a sustainable number.
initial adverbials	"For example, []", "On the other hand, []", "As a matter of fact, []", "At this point, []"

### **MDS Sentence Extraction — Maximal Marginal Relevance (MMR)**

- Maximal Marginal Relevance (MMR)
  - Iteratively, greedily pick the best sentence to add to existing summary (stop when desired length of summary is reached)
  - 2 selection criteria
    - (1) Relevance
      - Sentence s<sub>i</sub> is maximally relevant to user's query q
      - Example: high cosine similarity between s; and q

#### (2) Novelty

• Sentence is minimally redundant to existing summary S so far

**Note:** Sim1 and Sim2 can be the same similarity measure

$$MMR = \underset{s_i \in C \setminus S}{\operatorname{argmax}} \left[ \alpha \cdot Sim_1(s_i, q) - (1 - \alpha) \cdot \underset{s_j \in S}{\max} Sim_2(s_i, s_j) \right]$$
  
all sentence not selected so far similarity between  $s_i$  and query q max. similarity between  $s_i$  and an sentence in current summary

# **Information Ordering**

- Chronological ordering:
  - Order sentences by the date of the document, e.g., for summarizing news (Source: Inferring Strategies for Sentence Ordering in Multidocument News Summarization, 2002)
- Coherence:
  - Choose orderings that make neighboring sentences similar (by cosine).
  - Choose orderings in which neighboring sentences discuss the same entity (Source: Modeling Local Coherence: An Entity-Based Approach, 2007)
- Topical ordering
  - Learn the ordering of topics in the source documents

# **Domain-Specific Information Extraction**

#### • Domain: definitions

- a word's hypernym/genus, synonyms, etc.
- Domain: biographies
  - a person's birth/death, fame factor, education, nationality and so on

### • Domain: drugs / drug use

- **Problem** (the medical condition)
- Intervention (the drug or procedure)
- **Comparison** (e.g., control group)
- Outcome (the result of the study)

PICO

# **Definitional Templates**

• Domain: definitions

hypernym	The Hajj is a type of ritual
synonym	The Hajj, or Pilgrimage of Mecca, is the central duty of Islam
subtype	Qiran, Tamattu's, and Ifrad are three different types of Hajj

• Domain: biographies

dates	was assassinated on April 4, 1968
nationality	was born in Atlanta, Georgia
education	entered Boston University as a doctoral student

• Domain: drugs / drug use

population	37 otherwise healthy children aged 2 to 12 years
intervention	acetaminophen (10 mg/kg)
outcome	ibuprofen provided greater temperature decrement and longer duration of antipyresis than acetaminophen when the two drugs were administered in approximately equal dose

### **2011 Rise of the Machines**

**A VANTAGE POINT OR A BELIEF** 

**IBM Watson won** Jeopardy! on February 16, 2011





# **Outline**

### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

### Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts
## **Pre-Lecture Activity from Last Week**

#### • Assigned Task

- Do a web search and for the question stated below
- Post your answer(s) to the question into your Tutorial's Discussion in Canvas (please cite or quote your sources)

*"What is the relationship between information retrieval and natural language processing?"* 

#### Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers, but this won't help you learn better

### **Pre-Lecture Activity from Last Week**



NLP techniques are crucial to the performance of IR systems. To be able to surface what users intend to find based on a string query requires more than simple pattern matching.

More often than not, IR systems would need more complex algorithms to be useful, and this is where NLP comes in. For example, the concept of word embeddings can be used to deal with related queries by enabling meaningful comparison of semantic similarity. NLP is a branch of AI that deals with natural languages. It can be used to analyze or generate data. Information retrieval aims to perform analysis or search data from a large set of data collection. It uses some of the NLP techniques to achieve its tasks.

4



Information retrieval is concerned with the indexing and retrieval of relevant information according to a search term. NLP can make use of IR techniques to obtain better representations of inputs and to retrieve relevant information not present in the input

Retrieval-augmented generation

### **Question Answering via Web Search**



## **Question Answering via Web Search**

	are the highest mount	ains?	8 =
k NA	TURAL LANGUAGE ∫Σ MATH	I INPUT	🌐 EXTENDED KEYBOARD : EXAMPLES 👲 UPLOAD 🔀 RANDOM
Ass	uming "highest" refers to	elevation   l	prominence or peak distance to earth center instead
nput	interpretation		
m	ountains with highes	t elevation	
lesu	lt		More
		8849 m	
1	Mount Everest		
1 2	Mount Everest K2	8612 m	
1 2 3	Mount Everest K2 Kangchenjunga	8612 m 8586 m	
1 2 3 4	Mount Everest K2 Kangchenjunga Lhotse	8612 m 8586 m 8516 m	

Related to Retrieval Augmented Generation (RAG) and other use of tools by LLMs common in the last ½ year.

#### More next week.

## The Latest King in Town: GPT-x / ChatGPT



You are playing Jeopardy, and the answer is "4-letter word for a vantage point or a belief". What is the correct questions?



#### What is the word "View"?



## **Question Answering — Dimensions**

#### **Context / Source**

- Passage, document, corpus, ..., the Web
- Knowledge base
- Semi-structured tables
- Images / Video
- ...combination of sources

#### **Question Types**

• Factoid questions (typically direct and clear answers)

"How many calories does a tub of Ben & Jerry's have?"

• Open-ended questions (narratives, opinions, descriptions, etc.)

> "What is the healthiest way to quickly lose weight?"

#### **Answer Types**

- Yes/No
- Short text span/paragraph (extracted or generated)
- Database entry
- List of alternatives

#### • SQuAD (2016)

- Stanford Question Answering Dataset
- Over 100k and questions & answers generated by crowdworkers

Source: SQuAD: 100.000+ Questions for Machine Comprehension of Text

#### • SQuAD 2.0 (2018)

- Over 50k+ question & answers (crowdsourced)
- Twist: <u>unanswerable</u> question & <u>plausible</u> answers

Source: Know What You Don't Know: Unanswerable Questions for SQuAD

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

Article: Endangered Species Act

**Paragraph:** "... Other legislation followed, including the Migratory Bird Conservation Act of 1929, a 1937 treaty prohibiting the hunting of right and gray whales, and the Bald Eagle Protection Act of 1940. These later laws had a low cost to society—the species were relatively rare—and little opposition was raised."

Question 1: "Which laws faced significant opposition?" Plausible Answer: later laws

**Question 2:** "What was the name of the 1937 treaty?" Plausible Answer: Bald Eagle Protection Act

- MCTest (2013)
  - MCT: machine comprehension of text
  - Generation of dataset done by crowdworkers

(short stories + factoid questions with 4 multiple choice answers + opened-ended (more challenging) questions including answers)

Source: MCTest: A Challenge Dataset for the Open-Domain Machine Comprehension of Text

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. His aunt Jane tried as hard as she could to keep him out of trouble, but he was sneaky and got into lots of trouble behind her back.

One day, James thought he would go into town and see what kind of trouble he could get into. 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. He didn't pay, and instead headed home.

His aunt was waiting for him in his room. She told James that she loved him, but he would have to start acting like a well-behaved turtle.

After about a month, and after getting into lots of trouble, James finally made up his mind to be a better turtle.

1) What is the name of the trouble making turtle? A) Fries B) Pudding C) James D) Jane 2) What did James pull off of the shelves in the grocery store? A) pudding B) fries C) food D) splinters 3) Where did James go after he went to the grocery store? A) his deck B) his freezer C) a fast food restaurant D) his room

#### • CoQA (2019)

- CoQA: Conversational Question Answering
- Dataset generation by pairs of crowdworkers (one asking the questions, one answering the questions)
- 127k questions with answers from 8k conversations

Source: CoQA: A Conversational Question Answering Challenge

The Virginia governor's race, billed as the marquee battle of an otherwise anticlimactic 2013 election cycle, is shaping up to be a foregone conclusion. Democrat Terry McAuliffe, the longtime political fixer and moneyman, hasn't trailed in a poll since May. Barring a political miracle, Republican Ken Cuccinelli will be delivering a concession speech on Tuesday evening in Richmond. In recent ...

Q<sub>1</sub>: What are the candidates **running** for? A<sub>1</sub>: Governor R<sub>1</sub>: The Virginia governor's race

#### Q<sub>2</sub>: Where?

A<sub>2</sub>: Virginia R<sub>2</sub>: The Virginia governor's race

Q3: Who is the democratic candidate?

A<sub>3</sub>: Terry McAuliffe R<sub>3</sub>: Democrat Terry McAuliffe

Q4: Who is his opponent?

A<sub>4</sub>: Ken Cuccinelli R<sub>4</sub> Republican Ken Cuccinelli

Q<sub>5</sub>: What party does he belong to? A<sub>5</sub>: Republican

R5: Republican Ken Cuccinelli

Q6: Which of them is winning?

A6: Terry McAuliffe

R<sub>6</sub>: Democrat Terry McAuliffe, the longtime political fixer and moneyman, hasn't trailed in a poll since May

Dataset	Conversational	Answer Type	Domain
MCTest (Richardson et al., 2013)	×	Multiple choice	Children's stories
CNN/Daily Mail (Hermann et al., 2015)	×	Spans	News
Children's book test (Hill et al., 2016)	×	Multiple choice	Children's stories
SQuAD (Rajpurkar et al., 2016)	×	Spans	Wikipedia
MS MARCO (Nguyen et al., 2016)	×	Free-form text, Unanswerable	Web Search
NewsQA (Trischler et al., 2017)	×	Spans	News
SearchQA (Dunn et al., 2017)	×	Spans	Jeopardy
TriviaQA (Joshi et al., 2017)	×	Spans	Trivia
RACE (Lai et al., 2017)	×	Multiple choice	Mid/High School Exams
Narrative QA (Kočiský et al., 2018)	×	Free-form text	Movie Scripts, Literature
SQuAD 2.0 (Rajpurkar et al., 2018)	×	Spans, Unanswerable	Wikipedia

#### • Long Form Question Answering

- Question & answers collected from Subreddit ELI5 (Explain Like I'm 5: <u>https://www.reddit.com/r/explainlikeimfive/</u>)
- Collecting of source document from Common Crawl (<u>https://commoncrawl.org/</u>)
- 272k questions with at least 1 answer each

Source: ELI5: Long Form Question Answering

**Question:** How do Jellyfish function without brains or nervous systems? [...] (60 words)

**Answer:** Jellyfish may not have a brain, but they have a rough nervous system and innate behaviours. However, they are very simple creatures. They're invertebrate: creatures without a backbone. Most jellyfish have really short life spans. Sometimes just a couple of hours. [...] As their name implies, they are largely composed of basically jelly inside a thin membrane. They're over 95% water. (327 words)

**Documents:** [...] Jellyfish do not have brains, and most barely have nervous systems. They have primitive nerve cells that help them orient themselves in the water and sense light and touch. [...] While they dont possess brains, the animals still have neurons that send all sorts of signals throughout their body. [...] They may accomplish this through the assistance of their nerve rings. Jellyfish don't have brains, and that's just where things begin. They don't have many of the body parts that are typical in other animals. [...] (1070 words)

### Outline

#### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

#### Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

# QA Systems — Main Paradigms

- Information retrieval-based QA systems
  - Built on top of large text corpora (unstructured data)
  - Use IR techniques find relevant passages (or documents)
  - Apply reading comprehension algorithms over passages and draw answer (algorithms can be feature-based, neural-based, or both)

#### • Knowledge-based QA systems

- Built on top of semantic representations (structured data, e.g., knowledge graphs)
- Parse question to predicate calculus (e.g., FOL) or a query language (e.g., SQL, SPARQL)
- Optional: Generate "nice" answer from results

#### Hybrid Q&A systems

### **IR-Based Factoid QA Systems**



#### Basic components and architecture

## **IR-Based Factoid QA Systems**

- Question Processing
  - Detect question type, answer type, focus, relations
  - Formulate queries to send to a search engine / database
- Passage Retrieval
  - Retrieve ranked documents
  - Break into suitable passages and rerank
- Answer Processing
  - Extract candidate answers
  - Rank candidates using evidence from the text and external sources



## **Knowledge-Based Factoid QA Systems**

- Build semantic representation of question
  - times, dates, locations, entities, numeric quantities, etc.
- Use representations to query structured data or resources
  - Geospatial databases
  - Ontologies (Wikipedia Infoboxes, dbPedia, WordNet, Yago)
  - Scientific databases
  - etc.



## **Hybrid QA Systems**

- Example: IBM Watson
  - Build a shallow semantic representation of the query
  - Generate answer candidates using IR methods (Augmented with ontologies and semi-structured data)
  - Score each candidate using richer knowledge sources (geospatial databases, temporal reasoning, taxonomical classification)



53

# QA using Large Language Models: GPT (Generative Pretrained Transformer)

#### • GPT

- Uses only the Transformer Decoder without the encoder attention block (alternatively: encoder with "do not look ahead" masking)
- Self-supervised training
- Learning objectives
  - Auto-regressive Language Model
- (Very) oversimplified history of GPT
  - GPT-1/2/3: text only, "just" making it larger; GPT-4: multimodal
  - GPT-3+: reinforcement learning from human feedback (RLHF)



### **Outline**

#### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

#### Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

## **Question Processing**

- Things to extract from the question:
  - Answer Type Detection (decide the named entity type (e.g., person, place) of the answer)
  - Query Formulation (choose query keywords for the IR system)
  - Question Type classification (factoid question? definition question? math question? etc?)
  - Focus Detection (find the question words that are replaced by the answer)
  - Relation Extraction (find relations between entities in the question)

*"Who was the first president of Singapore?"* Question word: "who" important keywords Answer is a person (name)

*"who"*  $\rightarrow$  factoid questions

Relation extraction  $\rightarrow$  FOL PresidentOf(x, Singapore)

# Answer Type Taxonomy (Li & Roth, 2002)

- 2- layered taxonomy
  - 6 coarse classes

     (ABBREVIATION, ENTITY, DESCRIPTION, HUMAN, LOCATION and NUMERIC VALUE)
  - 50 fine classes
  - On the right: distribution of 500 questions in <u>TREC-10 Question Classification</u> test dataset

Class	# Class		#
ABBREV.	9	description	7
abb	1	manner	2
exp	8	reason	6
ENTITY	94	HUMAN	65
animal	16	group	6
body	2	individual	55
color	10	title	1
creative	0	description	3
currency	6	LOCATION	81
dis.med.	2	city	18
event	2	country	3
food	4	mountain	3
instrument	1	other	50
lang	2	state	7
letter	0	NUMERIC	113
other	12	code	0
plant	5	count	9
product	4	date	47
religion	0	distance	16
sport	1	money	3
substance	15	order	0
symbol	0	other	12
technique	1	period	8
term	7	percent	3
vehicle	4	speed	6
word	0	temp	5
DESCRIPTION	138	size	0
definition	123	weight	4

# **Answer Type Detection**

- Hand-written rules, e.g.:
  - Regular Expressions
  - Question headword

(first noun phrase after the wh-word)

(Enough data)

- Machine Learning
  - Requires annotated question datasets
  - Train classifier(s) over annotated dataset (feature-based, neural-based, or both)
  - Wide range of relevant features
     (question words, POS tags, parse features, named entities, etc.)

#### • Hybrid Methods

"Which city in Asia is also called the Garden City?"

"What is the official mascot of Singapore."



# **Query Formulation — Keyword Selection**

- Keyword heuristics (ordered list!)
  - (1) Select all non-stop words in quotations
    - (2) Select all NNP words in recognized named entities
    - (3) Select all complex nominals with their adjectival modifiers
    - (4) Select all other complex nominals
    - (5) Select all nouns with their adjectival modifiers
    - (6) Select all other nouns
    - (7) Select all verbs
    - (8) Select all adverbs
    - (9) Select the question focus word(s)
       (skipped in all previous steps)
    - (10) Select all other words

Question Processing				
Query Formulation				
Answer Type Detection				

### **Query Formulation — Keyword Selection**





### Outline

#### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

#### Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts



- IR engine retrieves documents using query terms
- Segment the documents into shorter units
  - Typically paragraphs, sentences, text spans
- Passage ranking
  - Use answer type to help re-rank passages



How do you rank passages?

Write a sample web query of your choice, or choose between learning about nigritude ultramarine or the relationship between telicity and aspect.

What "features" do you use to judge a passage's goodness?

# **Passage Ranking**

#### • Common features

- Number of Named Entities of the right type in passage
- Number of query words in passage
- Number of question N-grams also in passage
- Proximity of query keywords to each other in passage
- Longest sequence of question words
- Rank of the document containing passage

### **Answer Extractions**

- Answer extraction core task
  - Extract a specific answer from the passage (typically multiple answer candidates)
  - Span labeling: given a passage, identifying span of text which constitutes an answer
- Different strategies
  - Simple baseline: Run NER tagger on passage and return span in the passage is the correct answer type
  - Feature-based answer extraction
  - Neural-based answer extraction

```
"Who was the first president of Singapore?" (PERSON)
"Yusof bin Ishak was a Singaporean politician who was the
first president of Singapore, serving from 1965 to 1970."
```

detected answer type

### **Feature-Based Answer Extraction**

#### Common features

- Answer type match (candidate contains a phrase with the correct answer type)
- Pattern match (regular expression pattern matches the candidate)
- Question keywords (number of question keywords in the candidate)
- Keyword distance (distance in words between the candidate and query keywords)
- Novelty factor (a word in the candidate is not in the query)
- Apposition features (candidate is an appositive to question terms)
- Punctuation location (candidate is immediately followed by a comma, period, quotation marks, semicolon, or exclamation mark)
- Sequences of question terms (the length of the longest sequence of question terms that occurs in the candidate answer)

### **Neural-Based Answer Extraction**

• Example: DrQA



#### **Document Retriever**

- Basic IR-based approach
- Articles and questions are compared as TF-IDF weighted BoW vectors

#### **Document Reader**

- Vector representations of questions and paragraphs using RNN encoder
- Train 2 independent classifiers over encoded question and paragraphs
  - (1) Predict the start of answer span
  - (2) Predict the end of answer span

### **Outline**

#### • Text Summarization

- Overview & Categorization
- Basic Architecture
- Evaluation
- Query-Focused Summarization

#### Question Answering

- Overview & Categorization
- Factoid QA (Basic Architecture)
- Core Components
- Extended Concepts

# **Knowledge-Based QA Systems**

- Information source: knowledge graphs
  - Structured representation of knowledge
  - e.g.: DBpedia, Wikidata, YAGO, NELL, etc.





## **Knowledge-Based Factoid QA Systems**

- Knowledge graph: database of relations
  - (Semi-)automatic extractions from public data sources (often manually curated sources, e.g., Wikipedia infoboxes)
  - Relation extraction from unstructured text corpora (tricky task, many research papers)

Spouse(Barack\_Obama, Michelle\_Robinson) Occupation(Barack\_Obama, Politician) Occupation(Barack\_Obama, Lawyer) GraduatedFrom(Barack\_Obama, Columbia\_University)

#### Extraction relations in questions

 e.g., meaning representation with FOL (question → FOL expression typically contains variables)

"What college did Obama go to?"

 $GraduatedFrom(Barack_Obama, x)$ 



## **Geospatial Knowledge**

- Knowledge about containment, overlap, directionality, borders, e.g.:
  - "Singapore" a possible answer for "Asian city"
  - "Woodlands" is an area/zone/region in "Singapore"

#### GeoNames knowledge graph

	Name	Country	Feature class	Latitude	Longitude
1 🖗	Singapore SIN,Sin-ka-po,Singapore,Singapore City,Singapour,Singapur,Singapura,Sinkapoure,Sin-ka-po,Tumasik,cin	Singapore, SG.01	capital of a political entity population 3,547,809	N 1° 17' 22''	E 103° 51' 0"
2 🥊	Singapore Cingapura,Republic of Singapore,Sigapoa,Singaboor,Singaepuru,Singapo,Singapoa,Singapoer,Singapoo,Sin	Singapore,	independent political entity population 5,638,676	N 1° 22' 0''	E 103° 48' 0''
з 🖲	Singapore Changi Airport Aerodrom Singapur,Aeroport Changi,Aeroport Internacional de Singapur-Changi,Aeroport de Singapour Ch	Singapore, SG.02	airport elevation 6m	N 1° 21' 18''	E 103° 59' 24''
4 🕓	Central Singapore Community Development Council	Singapore, SG.01	region	N 1° 17' 55''	E 103° 51' 13"
5 🕅	Woodlands Woodlands New Town	Singapore, SG.03	populated place population 252,530	N 1° 26' 16''	E 103° 47' 19''
6 🧐	National University of Singapore <sup>(3)</sup> Gjal hoc Quoc gia Singapore,Nacional'nij universitet Singapuru,Nacional'nyj universitet Singapura,Na	Singapore,	college	N 1° 17' 46''	E 103° 46' 47''
7 🖲	Singapore River Rio Singapur,Riviere Singapour,Rivière Singapour,Río Singapur,Sin-ka-pho Ho,Sin-ka-pho Hô,Singapore,	<u>Singapore</u> ,	stream	N 1° 17' 12''	E 103° 51' 9''
8 🕓	Universal Studios Singapore	<u>Singapore</u>	amusement park	N 1° 15' 20"	E 103° 49' 15''
9 🕓	Singapore Botanic Gardens	<u>Singapore</u> ,	nature reserve	N 1° 18' 37''	E 103° 48' 59"

# **Temporal Reasoning**

- Common observation
  - Answers depend on current time or time frame
  - Common attribute in many knowledge graphs (also interesting: biographical dictionaries, obituaries, etc.)

• Example from IBM Watson

"In 1594 he took a job as a tax collector in Andalusia"

#### Candidate answers

- "Thoreau" is a bad answer (born in 1817)
- "Cervantes" is possible (was alive in 1594)
# **Context and Conversation in Virtual Assistants**

- Coreference helps resolve ambiguities
  - Question focus outside the actual question



#### • Clarification questions

- Insufficient information to find answer
- Too many possible answer candidates



## **Common Evaluation Metrics**

• Accuracy

. . .

- Does the answer match gold-standard answer? Exact Match (1/0 binary score)
- Often too "harsh", since an answer might be partially correct
- Mean Reciprocal Rank (MRR)
  - For each questions, return a ranked list of *m* candidate answers.
  - Question score is 1/Rank of the first correct answer

if the 1st answer is correct: 1.0 else if the 2nd answer is correct: 1/2 else if the 3rd answer is correct: 1/3

else if the *m*-th answer is correct: 1/m else: 0 (none of the *m* answers is correct)

■ Take mean over scores for all *n* questions

#### Summary

- Classification as core task of "higher-level" NLP applications
  - Often in combination with different core tasks (e.g., information retrieval, document ranking, etc.)
    - (1) Fake News Detection (Assignment 2)
      - Predict if a document (e.g., news article, tweet) is fake
    - (2) Text Summarization
      - Predict if a sentence is relevant to be part of a summary
    - (3) **Question Answering** 
      - Predict question and answer type
      - Feature-based answer extraction

### **Pre-Lecture Activity**

- Assigned Task
  - Do a web search and for the question stated below
  - Post your answer(s) to the question into your Tutorial's Discussion in Canvas (please cite or quote your sources)

# *"What are current limitations and challenges of LLMs (and using LLMs)?"*

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers, but this won't help you learn better

