

# Sequencing Legal DNA

## NLP for Law and Political Economy

### 11. Global Semantics

## Q&A Page

[bit.ly/NLP-QA11](https://bit.ly/NLP-QA11)

What is the endpoint of NLP?

# What is the endpoint of NLP?

Machine understanding of text **discourse**.

## What is the endpoint of NLP?

Machine understanding of text **discourse**.

- ▶ good summaries of long texts: extraction of relevant information, discarding of irrelevant information.

## What is the endpoint of NLP?

Machine understanding of text **discourse**.

- ▶ good summaries of long texts: extraction of relevant information, discarding of irrelevant information.
- ▶ question answering: retrieving evidence and answers from large corpora

## What is the endpoint of NLP?

Machine understanding of text **discourse**.

- ▶ good summaries of long texts: extraction of relevant information, discarding of irrelevant information.
- ▶ question answering: retrieving evidence and answers from large corpora
- ▶ storytelling: enriched text generation that understands characters and their relationships.

# Outline

Information Extraction and Knowledge Graphs

Long-Context Transformers

Text Summarization

Question Answering and Claim Verification

Narratives

## Knowledge Graphs

- ▶ A structured graph representing facts (and assertions?) as tuples.
- ▶ Entities are nodes, relations are edges:
  - ▶ (head entity, relation, tail entity)

# Knowledge Graphs

- ▶ A structured graph representing facts (and assertions?) as tuples.
- ▶ Entities are nodes, relations are edges:
  - ▶ (head entity, relation, tail entity)
- ▶ E.g., DBpedia: crowd-sourced effort to extract structured information from Wikipedia and make it available as linked open data.

GENERATING FACTS FOR THE ENTITY BILLIE HOLIDAY

“Facts” as RDF Triples

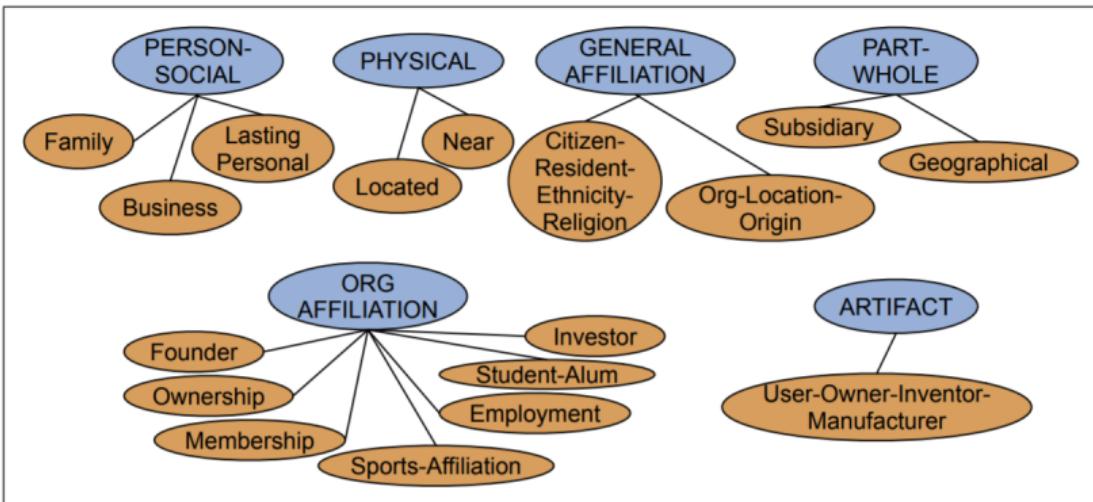


Subject ..... Predicate ..... Object  
(Thing)



S <[http://dbpedia.org/resource/Billie\\_Holiday](http://dbpedia.org/resource/Billie_Holiday)>  
P <<http://xmlns.com/foaf/0.1/name>>  
O "Billie Holiday"

# Relation Extraction



**Figure 17.1** The 17 relations used in the ACE relation extraction task.

Relations	Types	Examples
Physical-Located	PER-GPE	He was in Tennessee
Part-Whole-Subsidiary	ORG-ORG	XYZ, the parent company of ABC
Person-Social-Family	PER-PER	Yoko's husband John
Org-AFF-Founder	PER-ORG	Steve Jobs, co-founder of Apple...

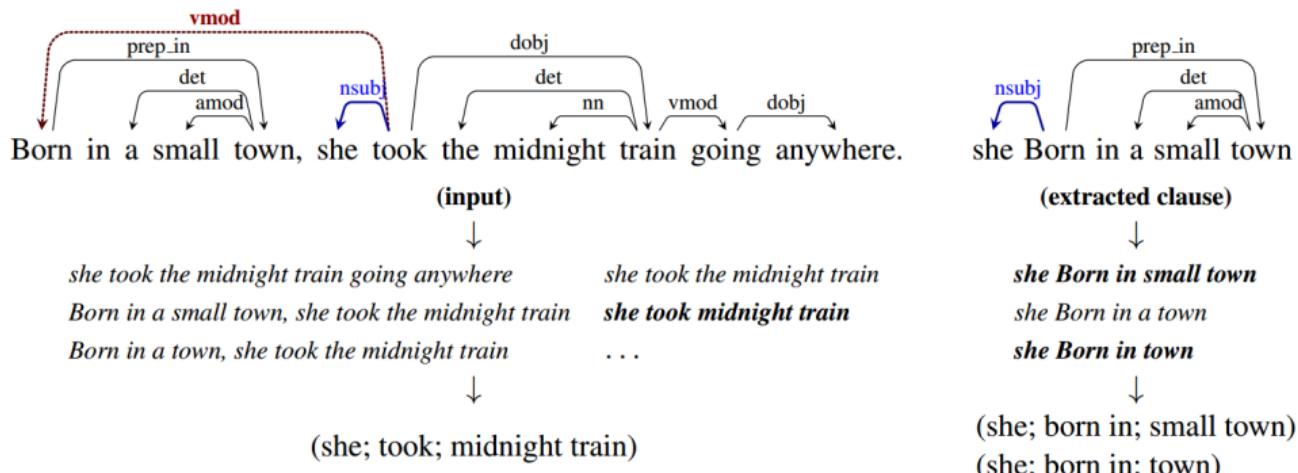
**Figure 17.2** Semantic relations with examples and the named entity types they involve.

## Information Extraction

- ▶ Manual information extraction is too costly.
- ▶ Can use relation extraction from text to fill in the knowledge graph.

# Information Extraction

- ▶ Manual information extraction is too costly.
- ▶ Can use relation extraction from text to fill in the knowledge graph.
- ▶ Rule-based approach:
  - ▶ use dependency parse to extract relations (e.g. Bank et al 2007, Fader et al 2011, Angeli et al 2015):



- ▶ filters: relation must contain a predicate; subject and object must be noun phrases.
- ▶ thresholds: aggregate over large corpora and keep only frequent (therefore potentially reliable) relations.

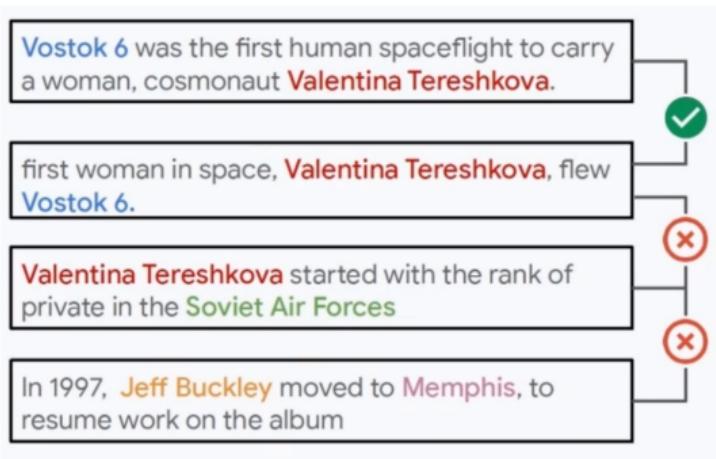
## Soares et al (2019): Contrastive Learning for Relations

- ▶ English Wikipedia:
  - ▶ use named entity recognition and co-reference resolution to assign unique ID's.
  - ▶ take sentences with at least 2 entities.

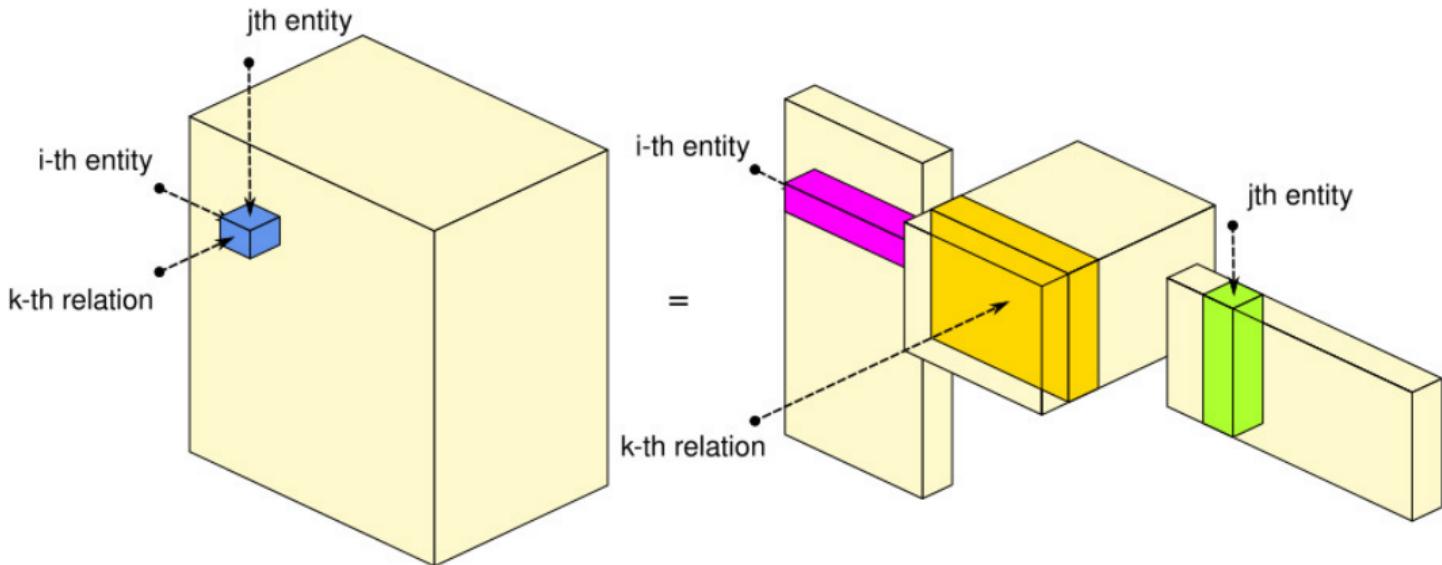
- ▶ English Wikipedia:
  - ▶ use named entity recognition and co-reference resolution to assign unique ID's.
  - ▶ take sentences with at least 2 entities.

## Negative sampling objective:

- ▶ Positive samples:
  - ▶ statements with the same two entities
- ▶ Negative samples:
  - ▶ same entity, one other entity
  - ▶ two different entities



# Knowledge Graph Embeddings



- ▶ Many proposed approaches for embedding KG entities, for example by predicting the presence of a bilateral relation.

## Learning from Knowledge Graph Relations

- ▶ Knowledge graph embeddings can be used to fill in missing relations.

e.g., parents of a person are often married, so

$$(\text{John, parent of, Anne}) + (\text{Mary, parent of, Anne}) \rightarrow (\text{John, married to, Mary})$$

# Learning from Knowledge Graph Relations

- ▶ Knowledge graph embeddings can be used to fill in missing relations.

e.g., parents of a person are often married, so

$$(\text{John, parent of, Anne}) + (\text{Mary, parent of, Anne}) \rightarrow (\text{John, married to, Mary})$$

EXAMPLES OF PATHS LEARNED BY PRA ON FREEBASE TO PREDICT WHICH COLLEGE A PERSON ATTENDED

<b>Relation Path</b>	<b>F1</b>	<b>Prec</b>	<b>Rec</b>	<b>Weight</b>
<i>(draftedBy, school)</i>	0.03	1.0	0.01	2.62
<i>(sibling(s), sibling, education, institution)</i>	0.05	0.55	0.02	1.88
<i>(spouse(s), spouse, education, institution)</i>	0.06	0.41	0.02	1.87
<i>(parents, education, institution)</i>	0.04	0.29	0.02	1.37
<i>(children, education, institution)</i>	0.05	0.21	0.02	1.85
<i>(placeOfBirth, peopleBornHere, education)</i>	0.13	0.1	0.38	6.4
<i>(type, instance, education, institution)</i>	0.05	0.04	0.34	1.74
<i>(profession, peopleWithProf., edu., inst.)</i>	0.04	0.03	0.33	2.19

# Language models and knowledge graphs

Hayashi et al (2019):

Topic: Barack Obama

Article **Barack Hussein Obama II** (...; born August 4, 1961) is an American[nationality] attorney[occupation] and politician[occupation] who served as the 44th president of the United States[position held] from 2009 to 2017. ...

Knowledge Graph

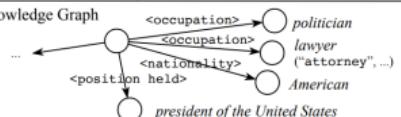


Figure 1: Overview of our task of language modeling conditioned on structured knowledge. For a given topic, we want to learn an LM that leverages the knowledge graph through relations when modeling the text.

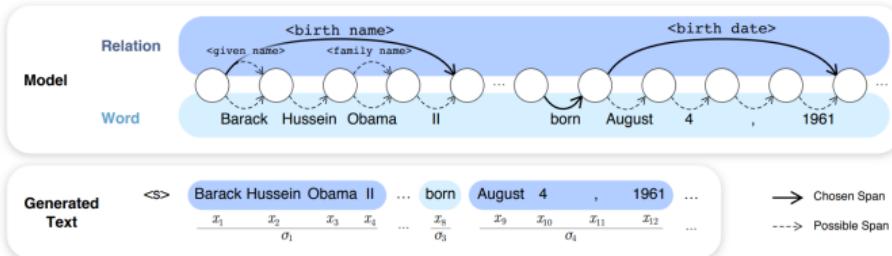


Figure 2: While generating, our model switches between the two sources, namely “Relation” and “Word”. Nodes represent hidden states up to each token, and edges represent possible span matches, i.e., choice of latent variables. In this example, we show one choice of latent variables with solid lines, and other options as dashed lines. We also show an “annotation” of the token sequence by the spans and sources we choose.

# Language models and knowledge graphs

Hayashi et al (2019):

Topic: Barack Obama

Article Barack Hussein Obama II (...; born August 4, 1961) is an American[nationality] attorney[occupation] and politician[occupation] who served as the 44th president of the United States[position held] from 2009 to 2017. ...

Knowledge Graph

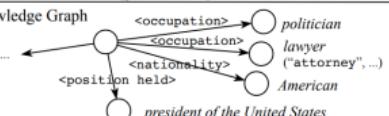


Figure 1: Overview of our task of language modeling conditioned on structured knowledge. For a given topic, we want to learn an LM that leverages the knowledge graph through relations when modeling the text.

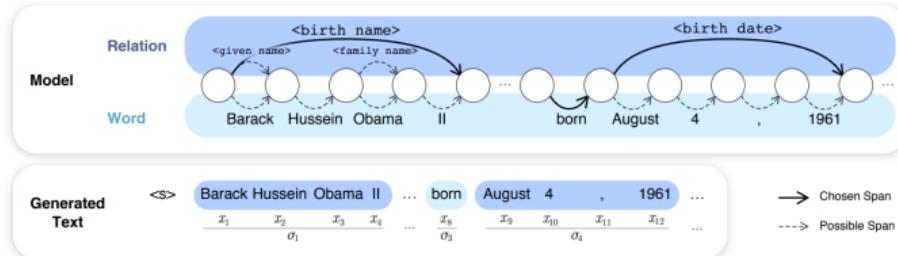


Figure 2: While generating, our model switches between the two sources, namely “Relation” and “Word”. Nodes represent hidden states up to each token, and edges represent possible span matches, i.e., choice of latent variables. In this example, we show one choice of latent variables with solid lines, and other options as dashed lines. We also show an “annotation” of the token sequence by the spans and sources we choose.

Wang et al (2020):

- ▶ Language models are knowledge graphs
- ▶ iterate over entity pairs, search LM attention weights for highest-probability phrase that connect them.
- ▶ construct / complete knowledge graph with these LM outputs.

# Outline

Information Extraction and Knowledge Graphs

Long-Context Transformers

Text Summarization

Question Answering and Claim Verification

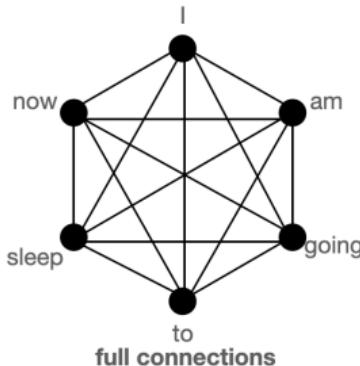
Narratives

## The Problem with BERT/GPT

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to  $n = 512$ ).

## The Problem with BERT/GPT

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to  $n = 512$ ).
- ▶ BERT's attention heads take as input the embeddings for each pair-wise interaction between tokens.

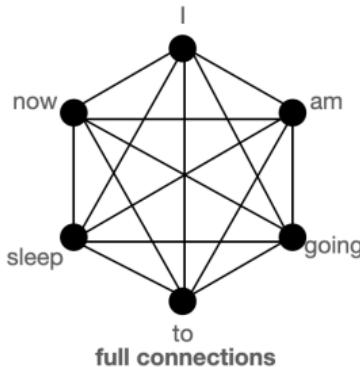


$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

- ▶  $n^2$  computations are needed at each step, so computation time is convex in sequence length.

## The Problem with BERT/GPT

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to  $n = 512$ ).
- ▶ BERT's attention heads take as input the embeddings for each pair-wise interaction between tokens.



$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

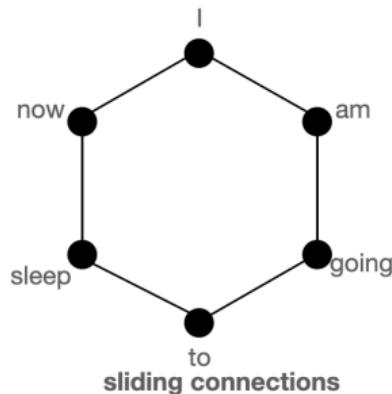
- ▶  $n^2$  computations are needed at each step, so computation time is convex in sequence length.
- ▶ Long-document transformers like BigBird try to approximate fully connected attention while enforcing sparsity between some/most tokens.

## Three alternatives to full pairwise attention

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

## Three alternatives to full pairwise attention

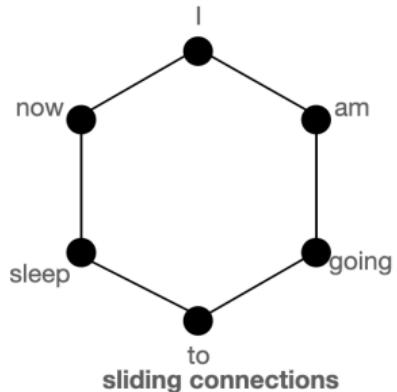
$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$



- (1)  $a(\cdot)$  always includes the tokens  $j$  before and after  $i$ .

## Three alternatives to full pairwise attention

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j)x_j, \forall i \in \{1, \dots, n_L\}$$



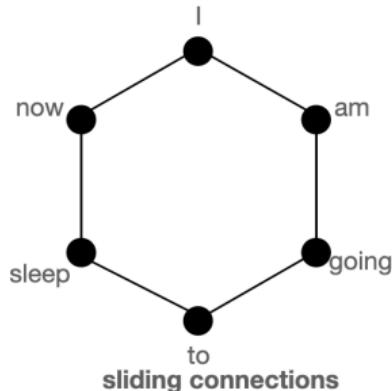
(1)  $a(\cdot)$  always includes the tokens  $j$  before and after  $i$ .



(2) Pick some important tokens (eg the first and last), and always include them in  $a(\cdot)$ .

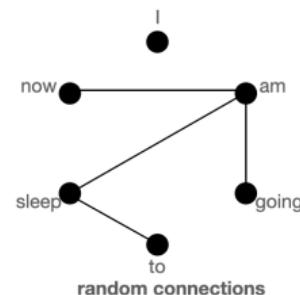
## Three alternatives to full pairwise attention

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j)x_j, \forall i \in \{1, \dots, n_L\}$$



(1)  $a(\cdot)$  always includes the tokens  $j$  before and after  $i$ .

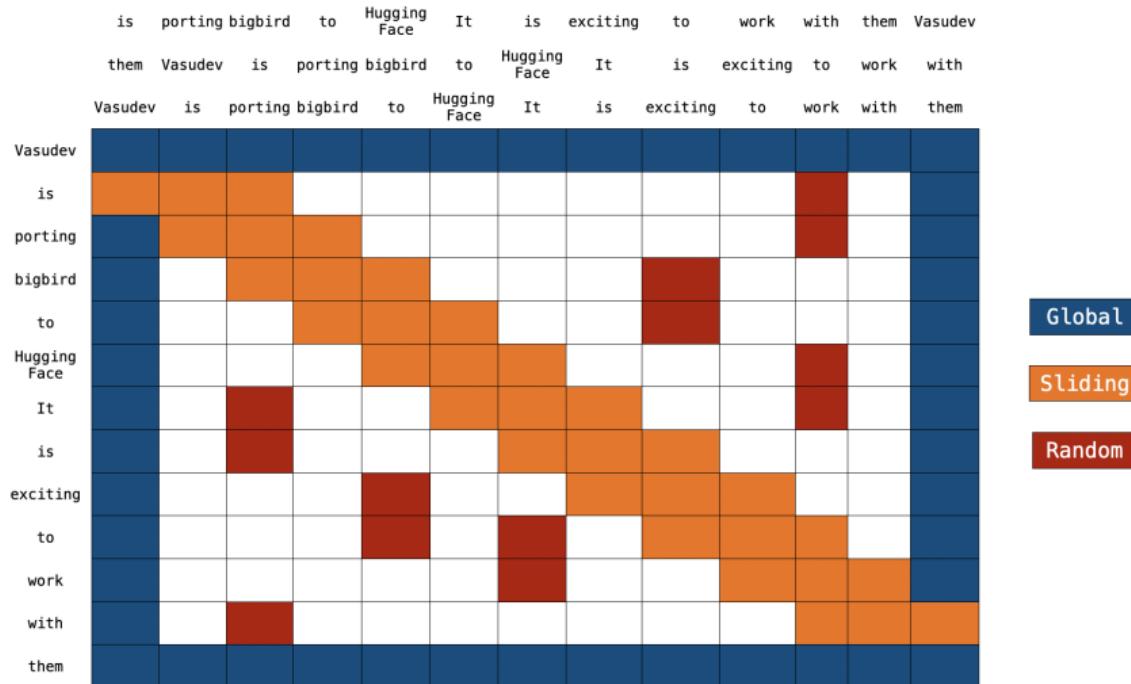
(2) Pick some important tokens (eg the first and last), and always include them in  $a(\cdot)$ .



(3) Pick some random tokens  $j$  to be included in  $a(\cdot)$ .

## BigBird's Block Sparse Attention

# BigBird's Block Sparse Attention



- ▶ Block sparse attention is an efficient implementation of these three alternative attention mechanisms: Each token attends to sliding tokens, some global tokens, & some random tokens.
- ▶ Extends the context window from 512 tokens to 4096 tokens, superior on long-document tasks.

# Outline

Information Extraction and Knowledge Graphs

Long-Context Transformers

**Text Summarization**

Question Answering and Claim Verification

Narratives

## Text Summarization

Goal: produce a shorter version of a text that contains the most relevant or important information.

- ▶ obvious applications in law / legal practice.
- ▶ not proven: dimension reduction or information extraction for social science measurement.

# Single Document Summarization

## Document

Cambodian leader Hun Sen on Friday rejected opposition parties ' demands for talks outside the country , accusing them of trying to "internationalize " the political crisis .

Government and opposition parties have asked King Norodom Sihanouk to host a summit meeting after a series of post-election negotiations between the two opposition groups and Hun Sen 's party to form a new government failed .

Opposition leaders Prince Norodom Ranariddh and Sam Rainsy , citing Hun Sen 's threats to arrest opposition figures after two alleged attempts on his life , said they could not negotiate freely in Cambodia and called for talks at Sihanouk 's residence in Beijing .Hun Sen , however , rejected that ."

I would like to make it clear that all meetings related to Cambodian affairs must be conducted in the Kingdom of Cambodia , " Hun Sen told reporters after a Cabinet meeting on Friday ." No-one should internationalize Cambodian affairs .

It is detrimental to the sovereignty of Cambodia , " he said .Hun Sen 's Cambodian People 's Party won 64 of the 122 parliamentary seats in July 's elections , short of the two-thirds majority needed to form a government on its own .Ranariddh and Sam Rainsy have charged that Hun Sen 's victory in the elections was achieved through widespread fraud .They have demanded a thorough investigation into their election complaints as a precondition for their cooperation in getting the national assembly moving and a new government formed .....



## Summary

Cambodian government rejects opposition's call for talks abroad

- ▶ **Extractive summarization:**
  - ▶ create the summary from phrases or sentences in the source document(s)
  - ▶ SOTA: Fine-tune BERT on corpora to predict start and end tokens of annotated summary sentences.

- ▶ **Extractive summarization:**
  - ▶ create the summary from phrases or sentences in the source document(s)
  - ▶ SOTA: Fine-tune BERT on corpora to predict predict start and end tokens of annotated summary sentences.
- ▶ **Abstractive summarization:**
  - ▶ express the ideas in the source documents using (at least in part) different words
  - ▶ SOTA: fine-tune **Big Bird Pegasus** to reconstruct provided summaries.

## Summarization with Human Feedback (Stiennon et al 2020)

Collect a large, high-quality dataset of human comparisons between summaries, train a model to predict the human-preferred summary, and use that model as a reward function to fine-tune a summarization policy using reinforcement learning

# Summarization with Human Feedback (Stiennon et al 2020)

Collect a large, high-quality dataset of human comparisons between summaries, train a model to predict the human-preferred summary, and use that model as a reward function to fine-tune a summarization policy using reinforcement learning

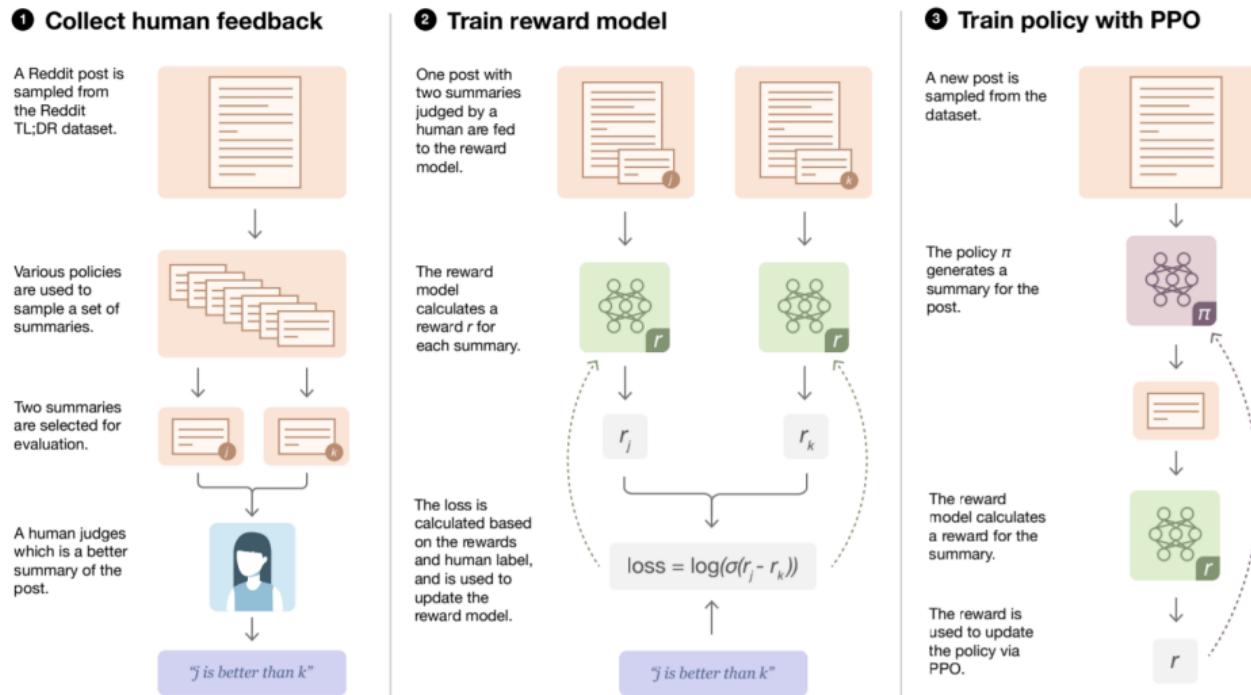
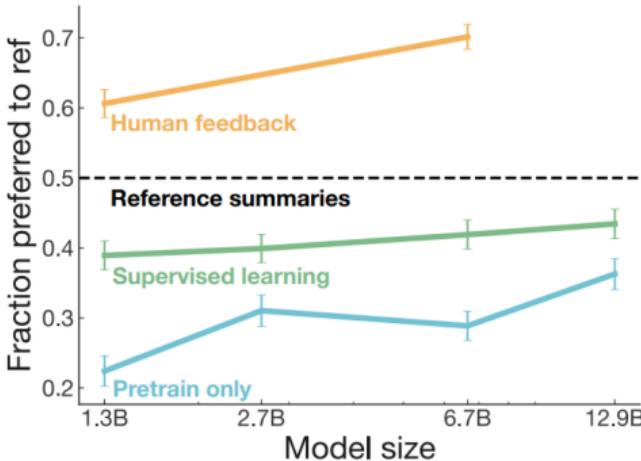


Figure 2: Diagram of our human feedback, reward model training, and policy training procedure.




---

[r/dating\_advice] First date ever, going to the beach. Would like some tips

Hey Reddit! I (20M) would like some tips, because I have my first ever date tomorrow (although I've had a gf for 3 years, but no actual dating happened), and we're going to the beach.

I met this girl, we have mutual friends, at a festival a few days ago. We didn't kiss, but we talked, held hands, danced a bit. I asked her to go on a date with me, which was super hard as it is the first time I've asked this to anybody. What I mean to say is, it's not like a standard \*first\* date because we already spent some time together.

I'm really nervous and excited. I'm going to pick her up tomorrow, we're cycling to the beach which will take 30 minutes, and then what? I'm a bit scared. Should I bring something (the weather, although no rain and sunny, is not super so no swimming), should we do something. I'd like all the tips I can get. Thanks!

---

**Human written reference TL;DR**

First date after 3 years in a relationship, going to the beach, terrified. What to bring with me, what to do?

**6.7B supervised model**

Going on a date with a girl I met a few days ago, going to the beach. What should I bring, what should we do?

**6.7B human feedback model**

Going on my first ever date tomorrow, cycling to the beach. Would like some tips on what to do and bring. I'm a bit nervous and excited. Thanks!

Table 1: Example of post and samples on the TL;DR dataset, chosen to be particularly short. For random samples (along with posts), see Appendix H and our website.

# Outline

Information Extraction and Knowledge Graphs

Long-Context Transformers

Text Summarization

Question Answering and Claim Verification

Narratives

# Open Question Answering and Claim Verification

Perhaps the most difficult global semantics tasks:

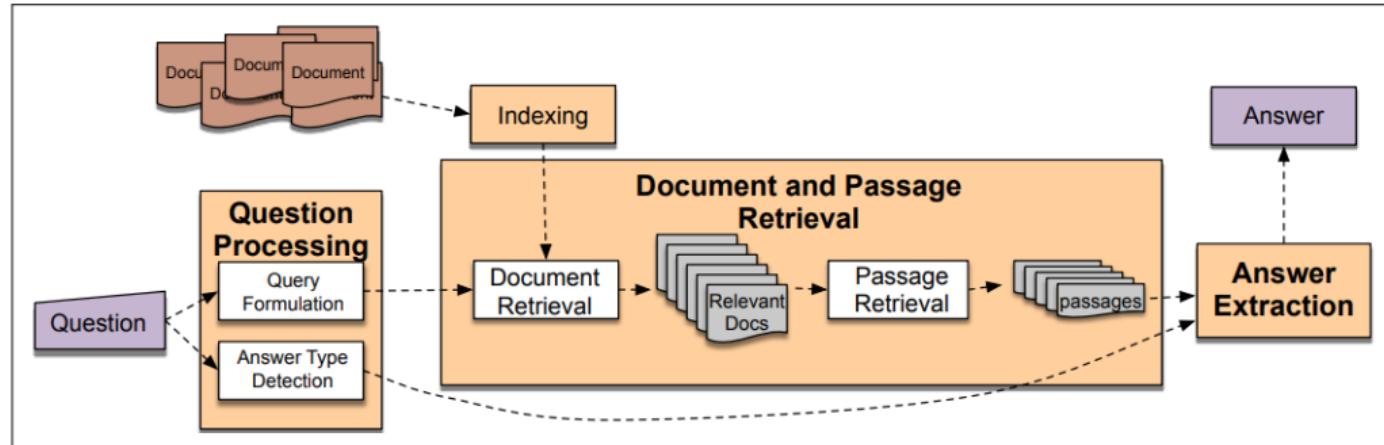
- ▶ Open question answering:
  - ▶ Answer any question.
  - ▶ “What are the responsibilities of the mayor of Zurich?”
- ▶ Open claim verification:
  - ▶ Check whether a plain-text claim is true or false.
  - ▶ “Zurich has the highest per-capita income of any city in Europe.”

# Open Question Answering and Claim Verification

Perhaps the most difficult global semantics tasks:

- ▶ Open question answering:
  - ▶ Answer any question.
  - ▶ “What are the responsibilities of the mayor of Zurich?”
- ▶ Open claim verification:
  - ▶ Check whether a plain-text claim is true or false.
  - ▶ “Zurich has the highest per-capita income of any city in Europe.”
- ▶ Both problems are solved using information retrieval pipelines:
  - ▶ search large corpora or knowledge graphs for evidence
  - ▶ use evidence to answer the question or check the claim

# Information Retrieval for Question Answering



**Figure 25.2** IR-based factoid question answering has three stages: question processing, passage retrieval, and answer processing.

- ▶ e.g., IBM Watson is a fast search engine over a knowledge base.

# Automated Claim Verification

**Claim** (*by Minister Shailesh Vara*)

“The average criminal bar barrister working full-time is earning some £84,000.”

**Verdict:** FALSE (*by Channel 4 Fact Check*)

The figures the Ministry of Justice have stressed this week seem decidedly dodgy. Even if you do want to use the figures, once you take away the many overheads self-employed advocates have to pay you are left with a middling sum of money.

1. Claim spotting (what to fact check – facts vs opinions, etc)
2. Evidence retrieval
3. Evidence filtering
4. Fact-check claim given evidence (textual entailment)

## Information Retrieval

- ▶ Input is a plain text query (a question or a claim)

## Information Retrieval

- ▶ Input is a plain text query (a question or a claim)
- ▶ The standard approach is **BM25**:
  - ▶ roughly, rank all documents by their TF-IDF cosine similarity to the query.
  - ▶ index every document by which words it contains, to quickly filter out irrelevant documents.

- ▶ Input is a plain text query (a question or a claim)
- ▶ The standard approach is **BM25**:
  - ▶ roughly, rank all documents by their TF-IDF cosine similarity to the query.
  - ▶ index every document by which words it contains, to quickly filter out irrelevant documents.
- ▶ Problem: requires exact word overlap (not synonyms).
  - ▶ alternatives use embeddings, e.g. S-BERT (but separate problem of how to encode long documents).
  - ▶ then can do fast approximate search over dense vectors (e.g. Faiss, Johnson et al 2017)
  - ▶ BM25 still usually does better, but this is being actively researched

## Inference Step

- ▶ Question answering:
  - ▶ take retrieved evidence as the context passage, and do local question answering (as in Week 10 lecture)

## Inference Step

- ▶ Question answering:
  - ▶ take retrieved evidence as the context passage, and do local question answering (as in Week 10 lecture)
- ▶ Claim verification:
  - ▶ take retrieved evidence as the premise, and the claim as the hypothesis, and do textual entailment (also like in Week 10).

How can we use global semantics in social science?

***Group Activity***

# Outline

Information Extraction and Knowledge Graphs

Long-Context Transformers

Text Summarization

Question Answering and Claim Verification

Narratives

- ▶ At the beginning I mentioned that long-range narrative comprehension and storytelling is another endpoint of NLP.
- ▶ But relative to the other tasks, it is not well-researched.

- ▶ At the beginning I mentioned that long-range narrative comprehension and storytelling is another endpoint of NLP.
- ▶ But relative to the other tasks, it is not well-researched.
- ▶ Meanwhile, “narratives” are also of major interest in social science, e.g. political economy:

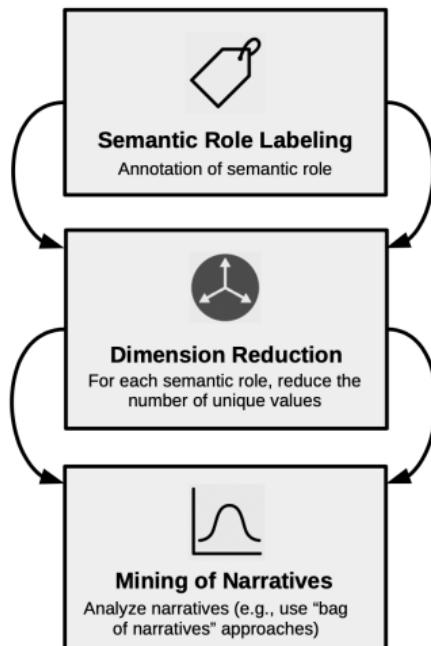
- ▶ At the beginning I mentioned that long-range narrative comprehension and storytelling is another endpoint of NLP.
- ▶ But relative to the other tasks, it is not well-researched.
- ▶ Meanwhile, “narratives” are also of major interest in social science, e.g. political economy:

*“Higher taxes will hurt the economy.”*

*“Immigrants steal our jobs.”*

*“Murderers deserve the death penalty.”*

# Ash, Gauthier, and Widmer (2021), "Mining Narratives from Large Text Corpora"



## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabrugge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabrugge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.
- ▶ **TF-IDF with clustering:** e.g., Hoberg and Phillips (JPE 2016), categorizing industries in SEC filings to measure concentration.

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabrugge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.
- ▶ **TF-IDF with clustering:** e.g., Hoberg and Phillips (JPE 2016), categorizing industries in SEC filings to measure concentration.
- ▶ **Topic models:** e.g., Hansen, McMahon, and Prat (QJE 2018), measuring allocation of attention in Central Bank discussion transcripts.

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabrugge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.
- ▶ **TF-IDF with clustering:** e.g., Hoberg and Phillips (JPE 2016), categorizing industries in SEC filings to measure concentration.
- ▶ **Topic models:** e.g., Hansen, McMahon, and Prat (QJE 2018), measuring allocation of attention in Central Bank discussion transcripts.
- ▶ **Syntactic parsing:** e.g. Vannoni, Ash, and Morelli (PA 2020) extracting modal verb structures ("governor shall do X").

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabruegge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.
- ▶ **TF-IDF with clustering:** e.g., Hoberg and Phillips (JPE 2016), categorizing industries in SEC filings to measure concentration.
- ▶ **Topic models:** e.g., Hansen, McMahon, and Prat (QJE 2018), measuring allocation of attention in Central Bank discussion transcripts.
- ▶ **Syntactic parsing:** e.g. Vannoni, Ash, and Morelli (PA 2020) extracting modal verb structures ("governor shall do X").
- ▶ **Word embeddings:** e.g. Ash, Chen, and Ornaghi (2021), measuring gender stereotypes in judicial opinions.

## Methods Review: Text Data in Social Science

- ▶ **Dictionary pattern matching:** e.g., Baker, Bloom, and Davis (QJE 2017), measuring economic policy uncertainty from newspaper text.
- ▶ **Supervised learning to scale a dimension in text:** e.g. Gentzkow and Shapiro (Ema 2010), Gentzkow, Shapiro, and Taddy (Ema 2019), measuring partisanship of language in congressional speeches / newspapers; Osabruegge, Ash, and Morelli (PA 2021) doing cross-domain topic classification in political text.
- ▶ **TF-IDF with clustering:** e.g., Hoberg and Phillips (JPE 2016), categorizing industries in SEC filings to measure concentration.
- ▶ **Topic models:** e.g., Hansen, McMahon, and Prat (QJE 2018), measuring allocation of attention in Central Bank discussion transcripts.
- ▶ **Syntactic parsing:** e.g. Vannoni, Ash, and Morelli (PA 2020) extracting modal verb structures ("governor shall do X").
- ▶ **Word embeddings:** e.g. Ash, Chen, and Ornaghi (2021), measuring gender stereotypes in judicial opinions.
- ▶ **Semantic role labeling**, extracting agents, actions, patients, and attributes (**Ash, Gauthier, and Widmer 2021**).

## Narratives as Semantic Roles

- ▶ Let  $S$  be a text corpus: a sequence of statements (or sentences)  $S_j$

## Narratives as Semantic Roles

- ▶ Let  $S$  be a text corpus: a sequence of statements (or sentences)  $S_j$
- ▶ **Semantic role:** functional component of a sentence, articulating “who” is doing “what” to “whom”
  - ▶ Agent (who): (health insurance)
  - ▶ Verb (what action): (saves)
  - ▶ Patient (to whom): (lives)

## Narratives as Semantic Roles

- ▶ Let  $S$  be a text corpus: a sequence of statements (or sentences)  $S_j$
- ▶ **Semantic role:** functional component of a sentence, articulating “who” is doing “what” to “whom”
  - ▶ Agent (who): (health insurance)
  - ▶ Verb (what action): (saves)
  - ▶ Patient (to whom): (lives)
- ▶ We decompose each statement  $S_j$  into a role sequence:

$$S_j = \{\text{Agent, Negation, Verb, Patient}\}$$

## Main Challenge: High Dimensionality

Summary statistics from U.S. Congressional Record

	All narratives	Complete narratives
Sentences	6,405,907	1,456,308
Statements	27,157,787	2,045,488
Agents, raw	4,000,807	2,045,481
Agents, raw unique	706,332	412,061
<hr/>		
Patients, raw	13,111,368	2,045,487
Patients, raw unique	3,543,212	758,804
<hr/>		
Verbs, raw (cf. statements)	27,157,787	2,045,488
Verbs, raw unique	24,807	10,788
<hr/>		

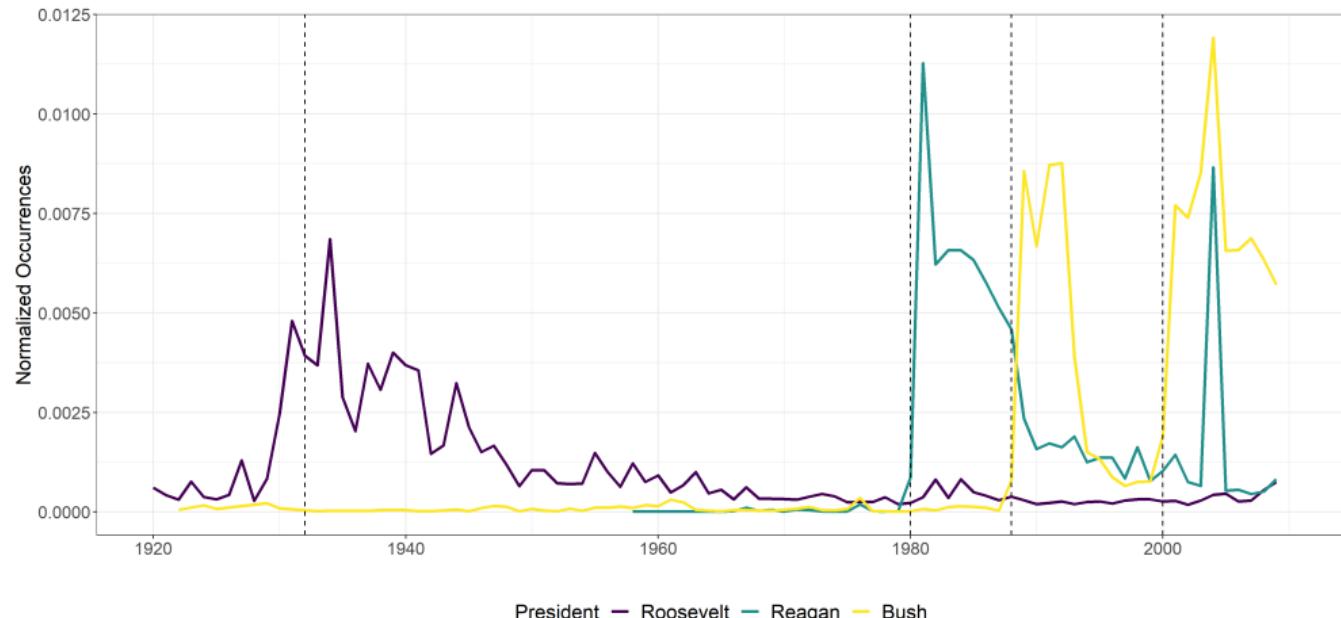
## Dimension Reduction: Verbs (Dictionary Approach)

- ▶ Replace verbs by their most common synonym (from WordNet)
  - ▶ “**make**”: make, create, produce, establish, build, ...
  - ▶ “**provide**”: provide, offer, furnish, cater, ...
  - ▶ “**use**”: use, employ, utilize, expend, habituate, ...
- ▶ Where possible, replace negated verbs by their (non-negated) most common antonym
  - ▶ **Example:** “can't find job” becomes “lose job”

Verb	Count
make	2374401
take	1055386
see	482254
provide	462867
give	377449
go	299348
have	287331
get	215879
do	206180
use	204300

## Dimension Reduction: Agents and Patients (NER)

- ▶ Identify named entities using spaCy
- ▶ Replace the 500 most frequent named entities (personalities, places, ...) by a unique token



## Dimension Reduction: Agents and Patients (Embedding Clusters)

Take agent/patient phrases that are not captured by NER:

- ▶ Encode each (e.g., "health care", "income tax credit") as phrase embedding
- ▶ Train k-means clustering on phrase embeddings (batch of 50,000 speeches)
  - ▶ → 500 additional entity clusters
- ▶ Labeled by the most frequent token in the cluster

## Dimension Reduction: Agents and Patients (Embedding Clusters)

Take agent/patient phrases that are not captured by NER:

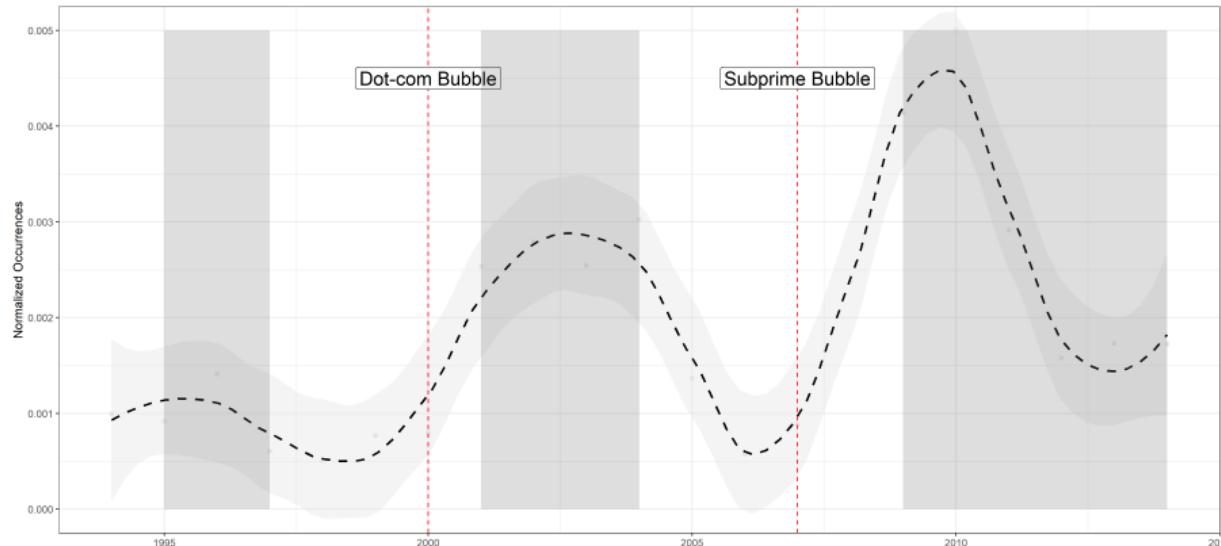
- ▶ Encode each (e.g., "health care", "income tax credit") as phrase embedding
- ▶ Train k-means clustering on phrase embeddings (batch of 50,000 speeches)
  - ▶ → 500 additional entity clusters
- ▶ Labeled by the most frequent token in the cluster
- ▶ "**american**": american people, million american, american taxpayer, many american, american public, every american, american citizen, american family, american worker
- ▶ "**government**": government, local government, federal government, government official, american government, government employee, government operation, branch government
- ▶ "**job**": job, good job, million job, american job, outstanding job, great job, job do, excellent job, job creation, do job
- ▶ "**money**": money, great money, percent money, tax money, public money, monies, good money, people money, make money, use money

Agent	Count	Patient	Count
president	310228	job	144314
people	298532	do	115236
american	218049	be	112563
man	215616	thing	90334
government	174977	money	87056
anyone	161031	service	86167
administration	124774	call	84745
company	121730	american	84220
law	105464	matter	83270
program	101246	family	77324

## Raw sentences and their mined narratives

- ▶ "President, I think the administration has begun to address the overseas basing issue."  
→ (administration, address, foreign policy)
- ▶ "As always, God bless and protect our troops and their families."  
→ (god, bless, troop)  
→ (god, protect, troop)
- ▶ "We need to pay attention to agriculture and the survival of the family farm as other countries protect and subsidize their farmers."  
→ (country, protect, farmer)  
→ (country, subsidize, farmer)

## Narratives Capture Major Events: "people lose job" and Economic Crises



# Narratives Capture Political Partisanship

Top Democrat Narrative Statements				Top Republican Narrative Statements			
Agent	Verb	Patient	Odds Ratio	Agent	Verb	Patient	Odds Ratio
people	make	billions	8.92	government	borrow	money	0.471
company	denies	benefits	7.33	president	balanced	budget	0.464
medicare	benefits	a lot	5.80	americans	pay	tax	0.461
president bush	took	government	4.42	god	bless	country	0.457
worker	lost	job	4.40	president	made	hearings	0.456
agency	receives	grant	3.18	citizens	abide	law	0.448
children	have	benefit	3.06	government	sponsors	economy	0.435
people	lost	job	2.95	people	buy	benefits	0.429
people	make	year	2.850	act	appropriates	fund	0.416
millions of americans	lost	job	2.785	doctor	takes	care	0.406

## Narratives as Worldviews

- ▶ At the document or speaker level, we obtain a list of relationships between entities.
- ▶ This list may also be represented as an adjacency matrix (network representation).
- ▶ The resulting graphs uncover the broader connections between narratives.

<https://sites.google.com/view/trump-narratives/accueil>