



# Modeling the Spread of Information within Novels

David Bamman  
School of Information, UC Berkeley  
[dbamman@berkeley.edu](mailto:dbamman@berkeley.edu)

In collaboration with Matt Sims,  
Olivia Lewke, Anya Mansoor, Sejal  
Popat and Sheng Shen

# Computational Humanities

Ted Underwood (2018), “Why Literary **Time** is Measured in Minutes”

Algee-Hewitt et al. (2016), “Canon/Archive: Large-Scale Dynamics in the Information Field”

Richard Jean So and Hoyt Long (2015), “Literary Pattern Recognition”

Ted Underwood, David Bamman and Sabrina Lee, The Transformation of **Gender** in English-Language Fiction (2018)

Holst Katsma (2014), **Loudness** in the Novel

So et al (2014), “**Cents** and Sensibility”

Matt Wilkens (2013), “The **Geographic** Imagination of Civil War Era American Fiction”

Jockers and Mimno (2013), “Significant **Themes** in 19th-Century Literature,”

Ted Underwood and Jordan Sellers (2012). “The Emergence of **Literary Diction**.” JDH

## Fiction as data

- BookCorpus (self-publishing)  
Zhu et al. 2015
- NarrativeQA  
Kočiský et al. 2017
- Commonsense stories  
Mostafazadeh et al. 2016
- Google Books  
Michel et al. 2010; Goldberg and Orwant 2013

## Modeling literary phenomena

- Character types  
Bamman et al. 2013, 2014
- Relationships  
Iyer et al. 2016, Chaturvedi et al. 2017
- Sentiment/plot  
Elsner 2012, Mohammad 2011,  
Jockers 2015, Reagan et al. 2018
- Character psychology  
Rashkin et al. 2018

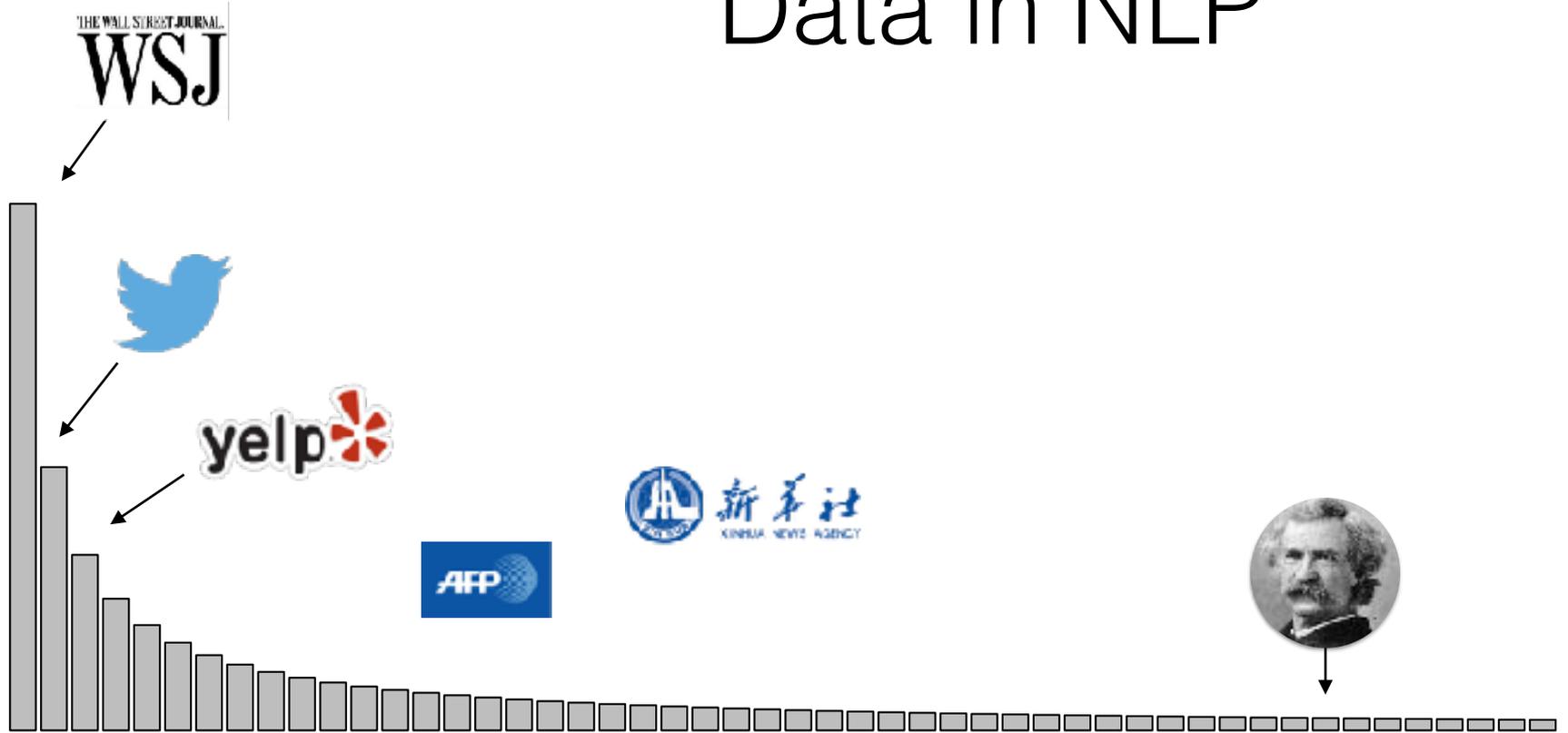
How do we design NLP to drive *insight* into literary texts?

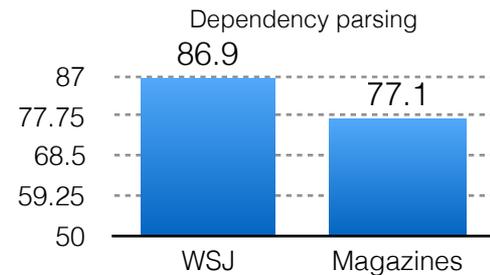
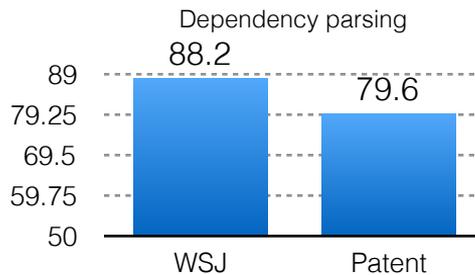
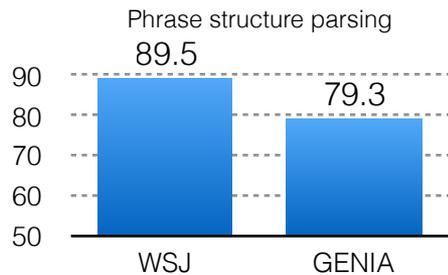
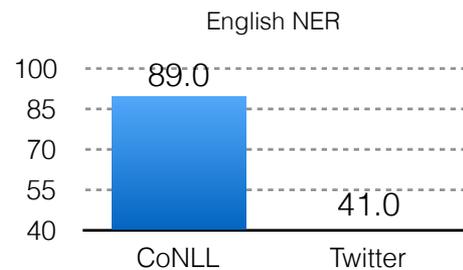
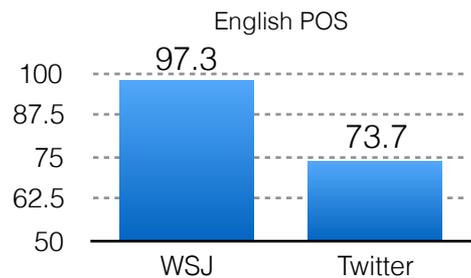
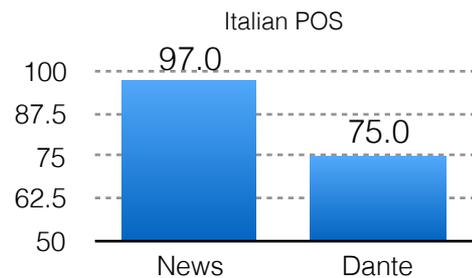
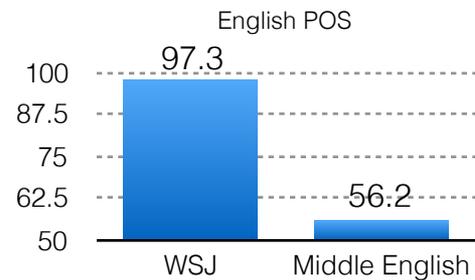
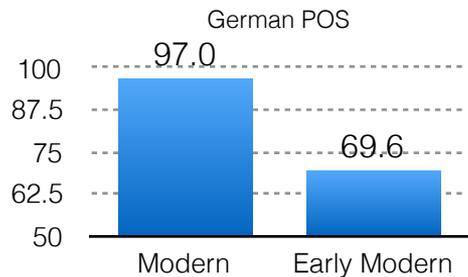
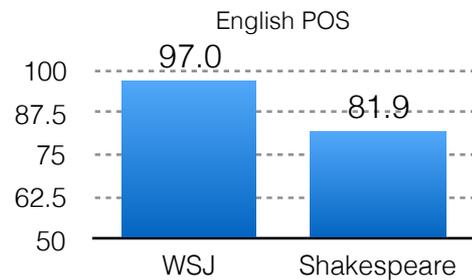
# NLP Pipeline



NLP Task	Accuracy
Tokenization	100%
Part-of-speech tagging	98.0% [Bohnet et al. 2018]
Named entity recognition	93.1 [Akbik et al. 2018]
Syntactic parsing	95.1 F [Kitaev and Klein 2018]
Coreference resolution	73.0 F [Lee et al. 2018]

# Data in NLP





# Active work

- Domain adaptation

[Chelba and Acero, 2006; Daumé and Marcu, 2006; Daumé 2009; Duong et al. 2015; Glorot et al. 2011, Chen et al. 2012, Yang and Eisenstein 2014, Schnabel and Schütz 2014]

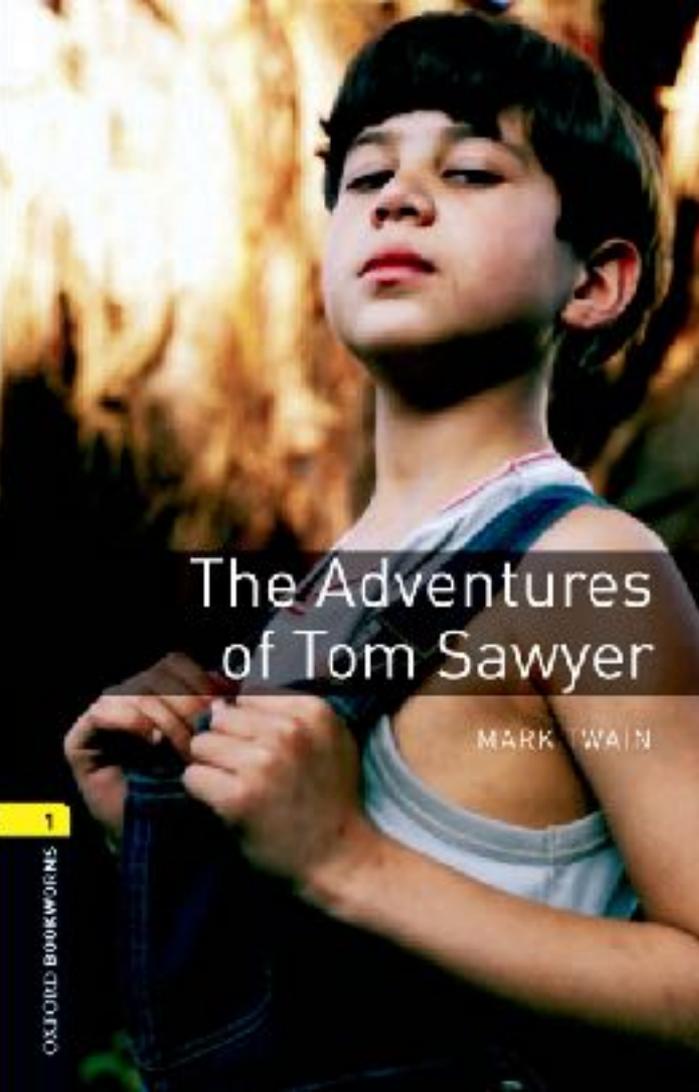
- Contextualized word representations

[Peters et al. 2018; Devlin et al. 2018; Howard and Ruder 2018; Radford et al. 2019]

- Data annotation. 210,532 tokens from 100 different novels, annotated for:

- Entities (person/place, etc.)
- Events
- Coreference

} available on Github now



# Literary entities

"TOM!"

No answer.

"TOM!"

No answer.

"What's gone with **that boy**, I wonder? You **TOM!**"

No answer.

**The old lady** pulled her spectacles down and looked over them about **the room**.

# Literary entities

Most work in NLP focuses on *named* entity recognition — mentions of specific categories (person, place, organization) that are explicitly named.

Mr. Knightley, a sensible man about seven or eight-and-thirty, was not only a very old and intimate friend of the family, but particularly connected with it, as the elder brother of Isabella's husband.

Austen, *Emma*

# Entity recognition

- Mr. Knightley
- a sensible man about seven or eight-and-thirty
- a very old and intimate friend of the family
- the family
- Isabella
- Isabella's husband
- the elder brother of Isabella's husband

Mr. Knightley, a sensible man about seven or eight-and-thirty, was not only a very old and intimate friend of the family, but particularly connected with it, as the elder brother of Isabella's husband.

# Entity recognition

- Mr. Knightley
- a sensible man about seven or eight-and-thirty
- a very old and intimate friend of the family
- the elder brother of Isabella's husband

- the family

- Isabella

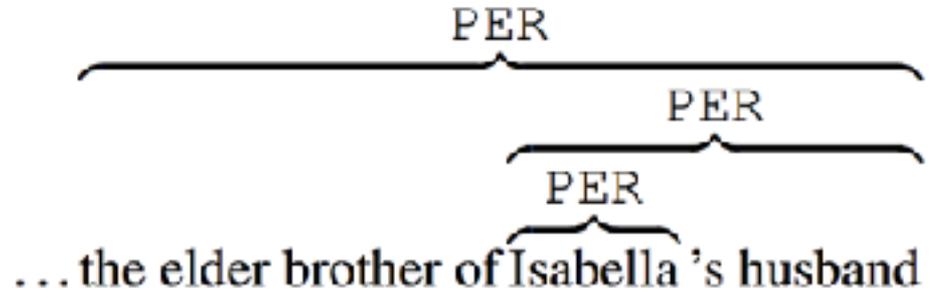
- Isabella's husband

Mr. Knightley, a sensible man about seven or eight-and-thirty, was not only a very old and intimate friend of the family, but particularly connected with it, as the elder brother of Isabella's husband.

Austen, *Emma*

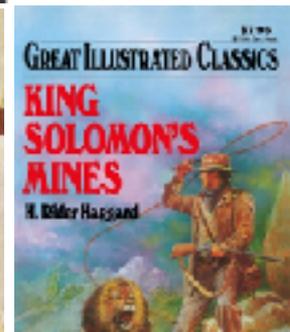
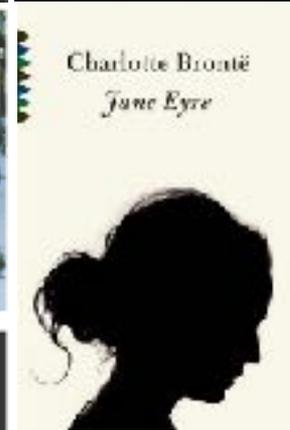
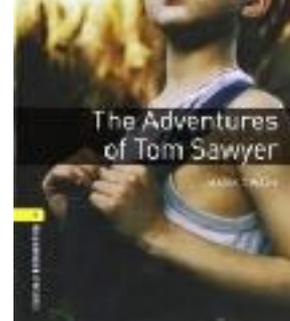
# Nested entity recognition

- Recognize spans of text that correspond to categories of entities (whether named or not).



# Dataset

- 100 books from Project Gutenberg
- Mix of high literary style (e.g., Edith Wharton's *Age of Innocence*, James Joyce's *Ulysses*) and popular pulp (Haggard's *King Solomon's Mines*, Alger's *Ragged Dick*).
- Select first 2000 words from each text



# Entity classes

- **Person**. Single person with proper name (Tom Sawyer) or common entity (the boy); set of people (her daughters).
- **Organization**. Formal association (the army, the Church as an administrative entity).
- **Vehicle**. Devices primarily designed to move an object from one location to another (ships, trains, carriages).

# Entity classes

- **GPE.** Entities that contain a population, government, physical location and political boundaries (New York, the village)
- **Location.** Entities with physicality but w/o political status (New England, the South, Mars), including natural settings (the country, the valley, the forest)
- **Facility.** Functional, primarily built structure designed for habitation (buildings), storage (barns), transportation (streets) and maintained outdoor space (gardens).

# Metaphor

- Only annotate phrases whose types denotes an entity class.

PER PER

John is a doctor

PER

PER

???

the young man was not really a poet; but surely he was a poem

Chesterton, *The Man  
Who Was Thursday*

# Personification

- **Person** includes characters who engage in dialogue or have reported internal monologue, regardless of human status (includes aliens and robots as well).

As soon as I was old enough to eat grass **my mother** used to go out to work in the daytime, and come back in the evening.

Sewell, *Black Beauty*

# Data

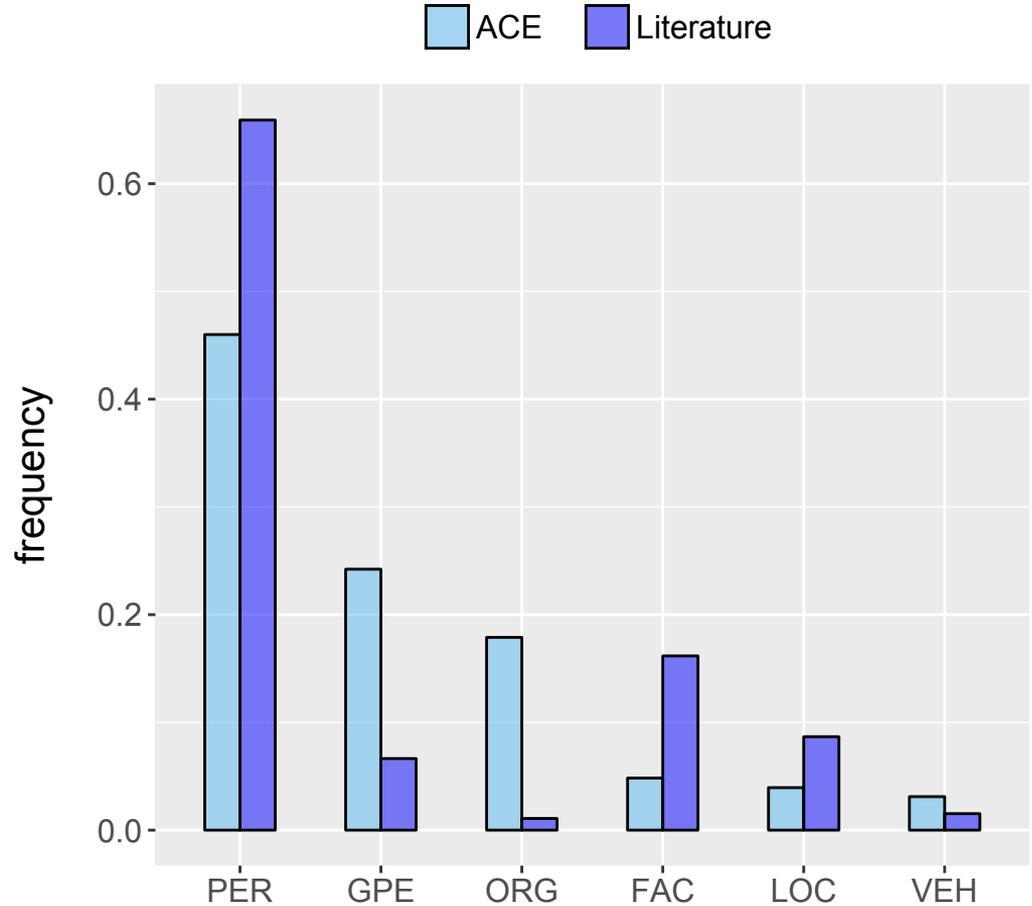
Cat	Count	Examples
PER	9,383	my mother, Jarndyce, the doctor, a fool, his companion
FAC	2,154	the house, the room, the gardne, the drawing-room, the library
LOC	1,170	the sea, the river, the country, the woods, the forest
GPE	878	London, England, the town, New York, the village
VEH	197	the ship, the car, the train, the boat, the carriage
ORG	130	the army, the Order of Elks, the Church, Blodgett College

# Prediction

How well can find these entity mentions in text as a function of **the training domain**?

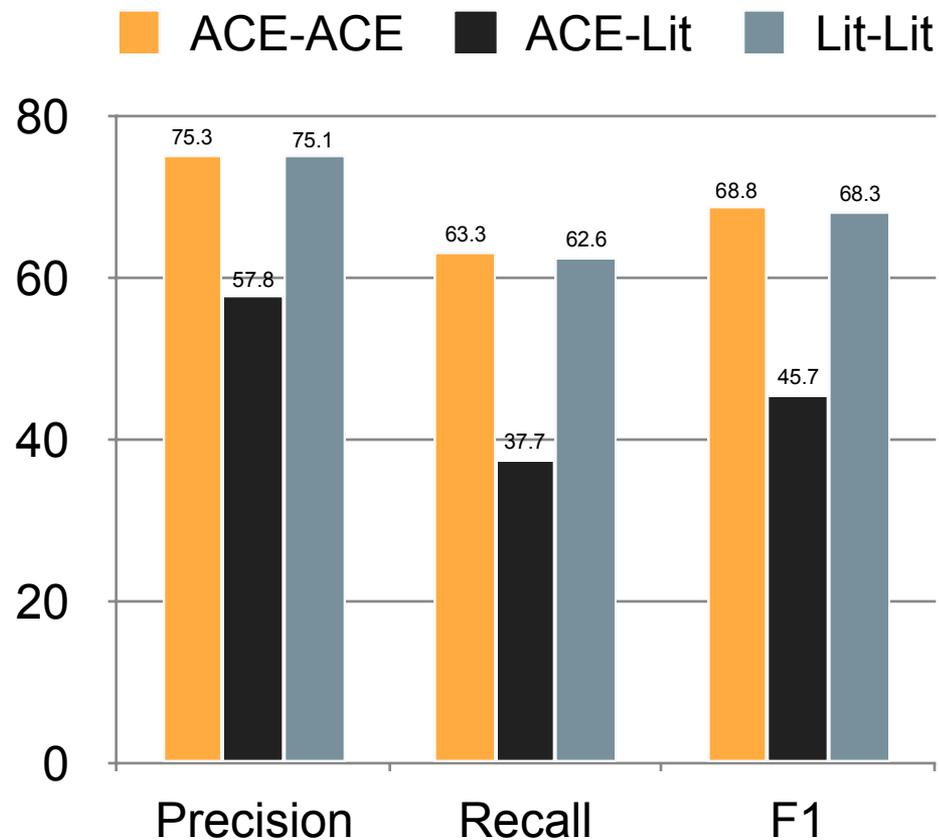
# Data

- ACE (2005) data from newswire, broadcast news, broadcast conversation, weblogs



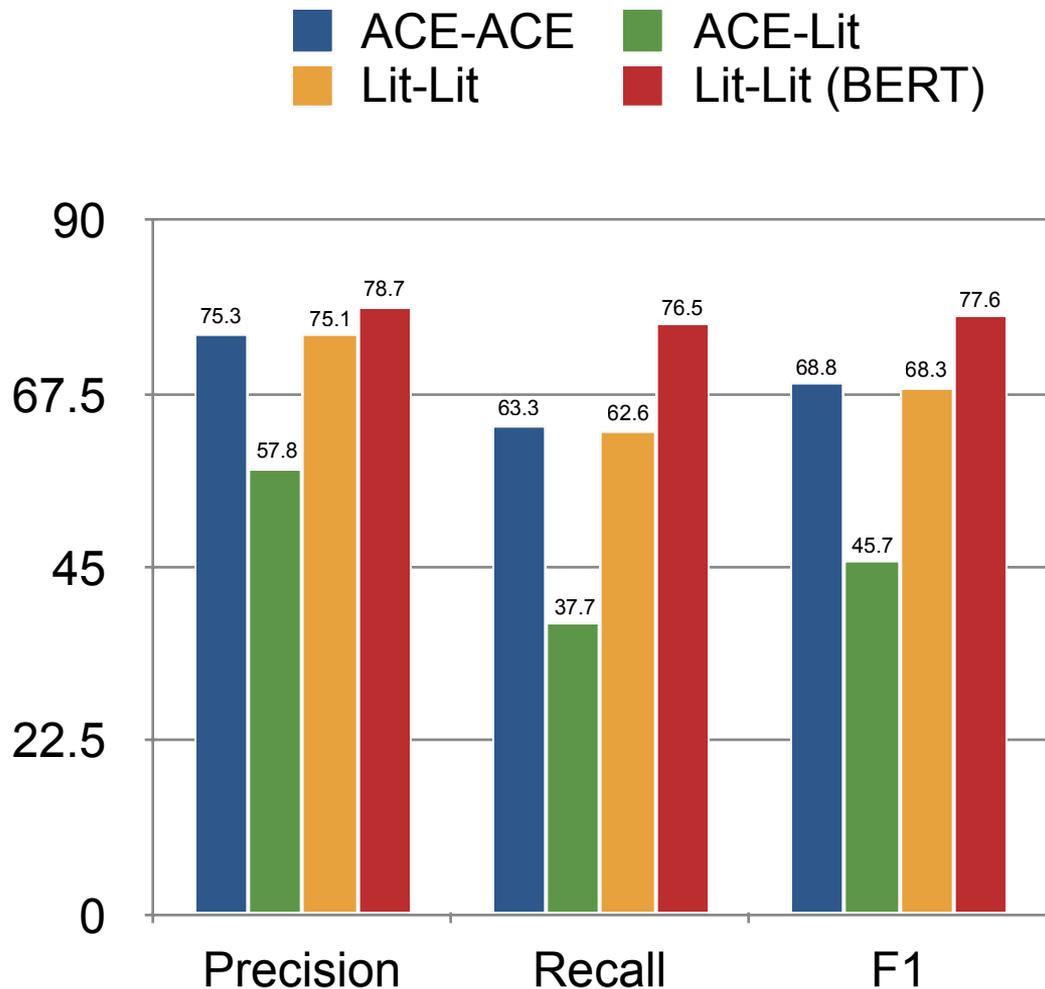
# Prediction

- Ju et al. (2018): layered BiLSTM-CRF; state-of-the-art on ACE 2005.
- Evaluate performance difference when altering the training/test domain.



# Prediction

- Ju et al. (2018): layered BiLSTM-CRF; state-of-the-art on ACE 2005.
- Evaluate performance difference when altering the training/test domain.
- Adding BERT contextual embeddings (Devlin et al. 2019) yields +9.3 F1 score



# Analysis

- Tag entities in 1000 new Gutenberg texts (78M tokens) using the two models (ACE vs. LIT) and analyze the difference in frequencies with which a given string is tagged as **PER** under both models.

Mrs.
Miss
Lady
Aunt

MOSCOW, April 17 (AFP)

Silence is golden -- especially when your hand is weak -- top Moscow policy analysts said in an assessment of the fallout from Russia's vocal opposition to what turned out to be a swift US-led campaign in Iraq.

Several top diplomacy experts told a Kremlin-run forum that countries like China and India that said little about the conflict before its March 20 launch were already reaping the benefits.

Some suggested that Russian President **Vladimir Putin** will now be scrambling to contain the damage to his once-budding friendship with US President **George W. Bush** because he was poorly advised by his intelligence and defense aides.

AFP\_ENG\_20030417.0307

## Chapter I: The Bertolini

“**The Signora** had no business to do it,” said **Miss Bartlett**, “no business at all. She promised us south rooms with a view close together, instead of which here are north rooms, looking into a courtyard, and a long way apart. Oh, **Lucy!**”

“And a Cockney, besides!” said **Lucy**, who had been further saddened by **the Signora**’s unexpected accent. “It might be London.”

Forster, *A Room with a View*

# Analysis

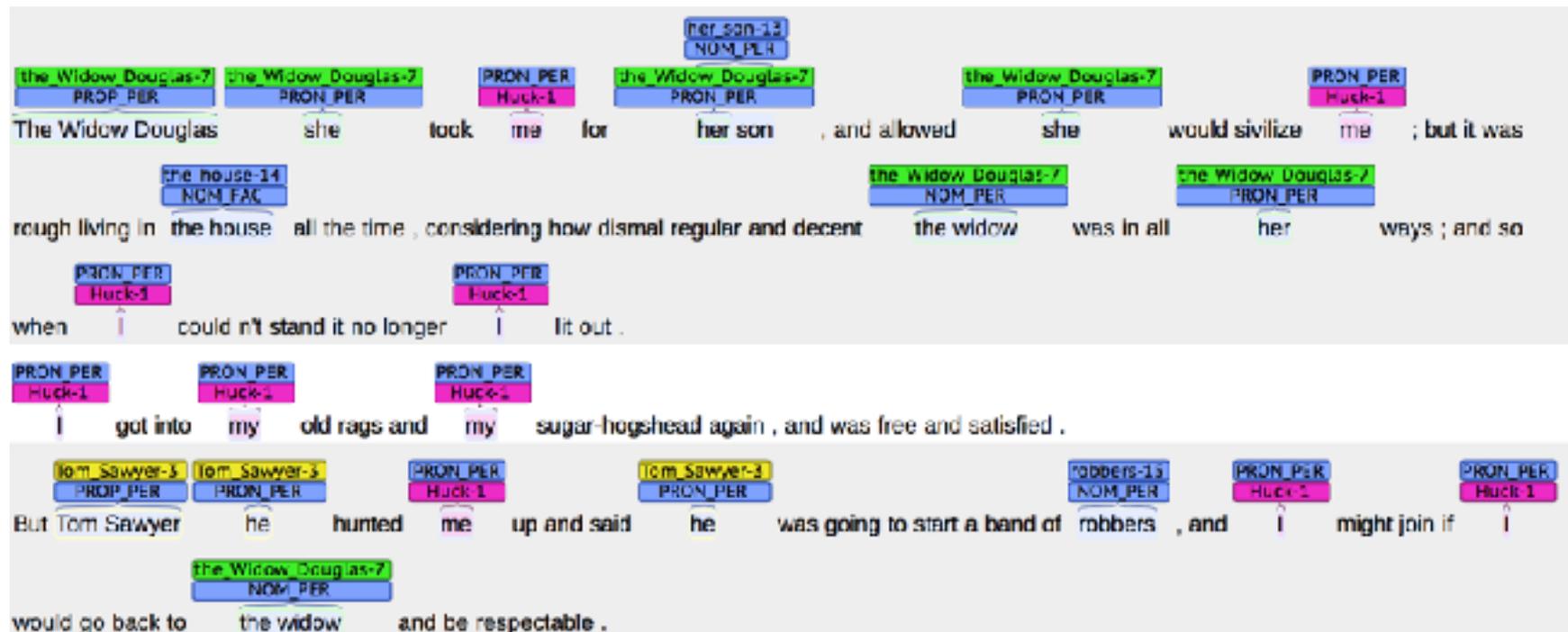
- How well does each model identify entities who are men and women?
- We annotate the gender for all **PER** entities in the literary test data and measure the recall of each model with respect to those entities.

Training	Women	Men	Diff
ACE	38.0	49.6	-11.6
Literary	69.3	68.2	1.1

# Coreference

Challenge: cluster together all mentions in a text that co-refer to the **same** entity.

# Coreference



# Coreference

[Elco Industries Inc.] said [it] expects net income in the year ending June 30, 1990, to fall below a recent analyst's estimate of \$ 1.65 a share. [The Rockford, Il. maker of fasteners] also said...

# Literary conference

- Coreference in literary texts raises a number of difficult complications that require special attention.
- Many of these complications involve the nature of **identity**.





# Identity and near-identity

- Entities in literary texts *change* (over the course of a lifetime or other long periods).
- Classical works in coreference determine whether two mentions refer to the same entity in the real world, which becomes entangled in deep metaphysical complexities on the nature of identity.
- We follow Recasens et al. (2010) in judging the similarity between two entities **constructed in discourse**.

- On homecoming night [Postville] feels like Hometown, USA, but a look around [this town of 2,000] shows it's become a miniature Ellis Island. This was an all-white, all-Christian community . . . For those who prefer [the old Postville], Mayor John Hyman has a simple answer. (Recasens et al., 2011, 10)

Halfway down a by-street of one of our New England towns stands [a rusty wooden house, with seven acutely peaked gables, facing towards various points of the compass, and a huge, clustered chimney in the midst]. The street is Pyncheon Street; [the house] is [the old Pyncheon House]<sub>cop</sub>

- [The House of the Seven Gables] and [the old Pyncheon House] could refer to different things (the same structure but different temporal extents); but here that are equated in the discourse

# Specificity

specific



My son just watched *Frozen*

generic



Kids like *Frozen*

Value	Meaning
specific	singular occurrence at a particular place and time
general	claim about groups, abstractions

Whereas with respect to Turkey, I had much ado to keep him from being a reproach to me. His clothes were apt to look oily and smell of eatinghouses. He wore his pantaloons very loose and baggy in summer. His coats were execrable; his hat not to be handled. But while the hat was a thing of indifference to me, inasmuch as his natural civility and deference, as a dependent Englishman, always led him to doff it the moment he entered the room, yet his coat was another matter. Concerning his coats, I reasoned with him; but with no effect. The truth was, I suppose, that [a man of so small an income] could not afford to sport such a lustrous face and a lustrous coat at one and the same time. (Melville, *Bartleby, The Scrivener*)

Class *near*-identity (Recasens et al. 2010)

# Near-identity

- A proper noun appears first, and a subsequent noun phrase refers to some aspect of the discourse entity (Recasens et al. 2010)
  - Role
  - Location
  - Organization
  - Information realization
  - Representation
  - Class

# Class near-identity

Two noun phrases share an is-a relationship, but they stand in a different position in the categorical hierarchy so that one can be viewed as more general or specific to the other.

Diego looked for information about **his character** in the novel forgetting that Saramago does not usually describe **them**.

# Revelations of identity



# Revelations of identity

“There will call upon you to-night, at a quarter to eight o’clock,” it said, “[a gentleman who desires to consult you upon a matter of the very deepest moment]<sub>x</sub>. ...Be in your chamber then at that hour, and do not take it amiss if [your visitor]<sub>x</sub> wear a mask.”

...

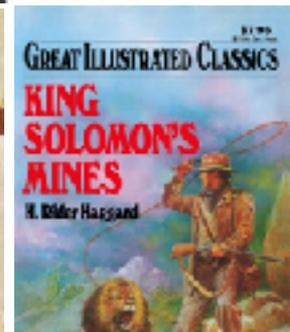
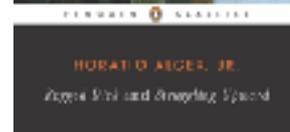
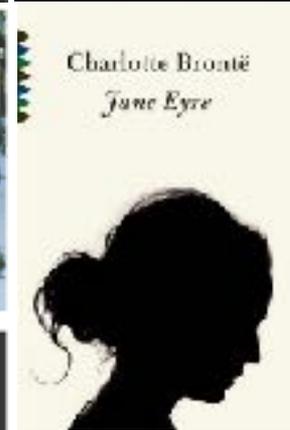
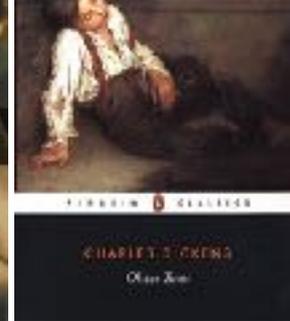
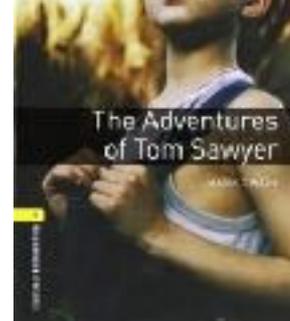
“It only remains, therefore, to discover what is wanted by [this German who writes upon Bohemian paper and prefers wearing a mask to showing [his]<sub>x</sub> face]<sub>x</sub>. And here [he]<sub>x</sub> comes, if I am not mistaken, to resolve all our doubts.” (Conan Doyle, *Adventures of Sherlock Holmes*)

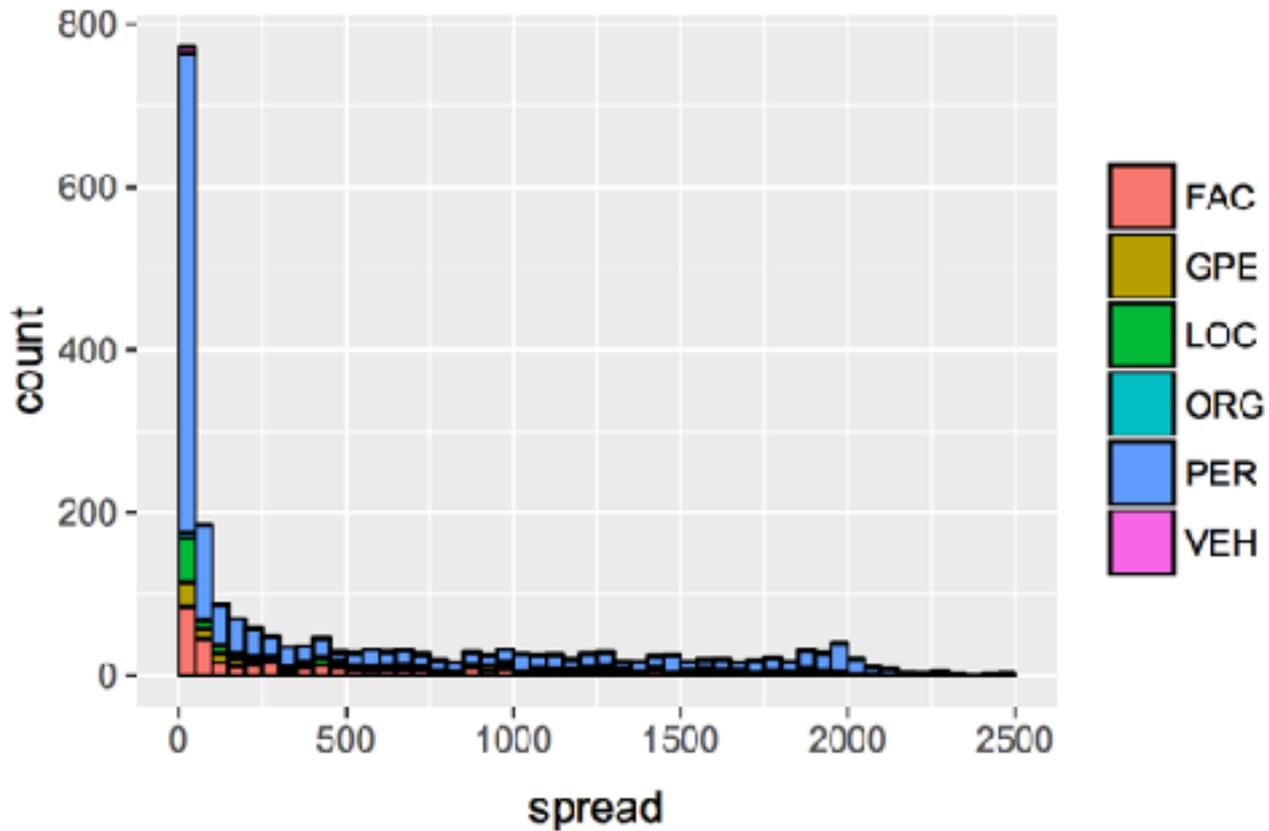
# Revelations of identity

We annotate identity from the reader's point of view (as distinct from the characters'): all mentions are coreferent if a reader can determine that they are identical at any point in the narrative.

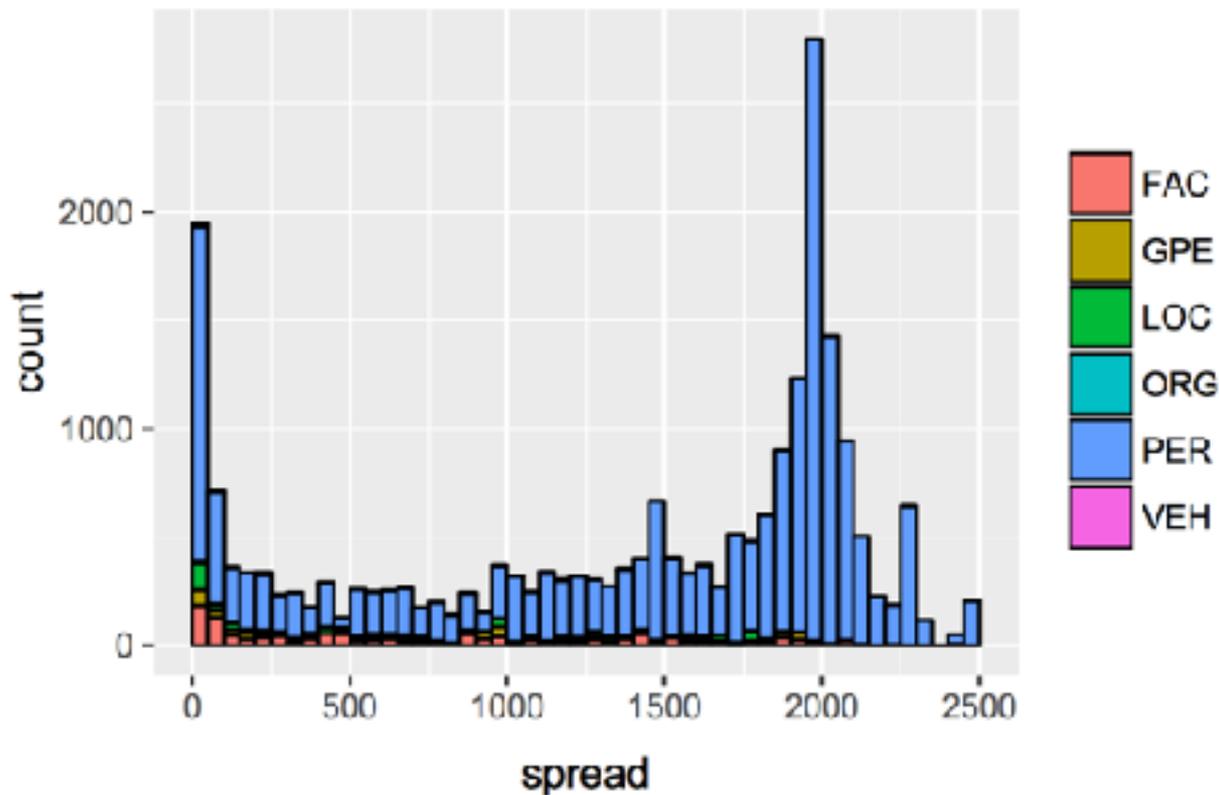
# Dataset

- 100 books from Project Gutenberg (same as for entity annotations)
- First 2000 words from each text
- 29,104 mentions



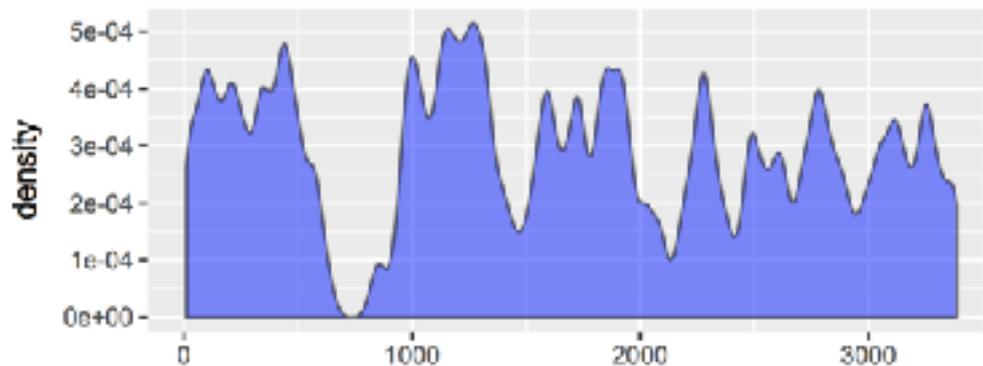
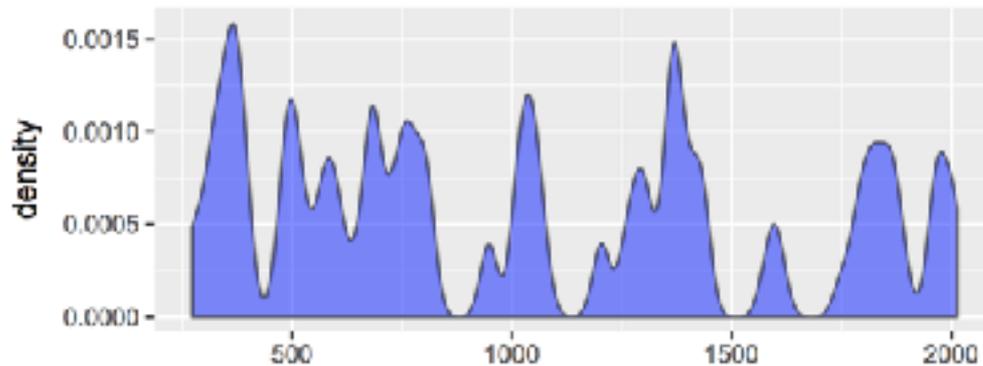


Spread between first and last mention of an entity, measured in tokens



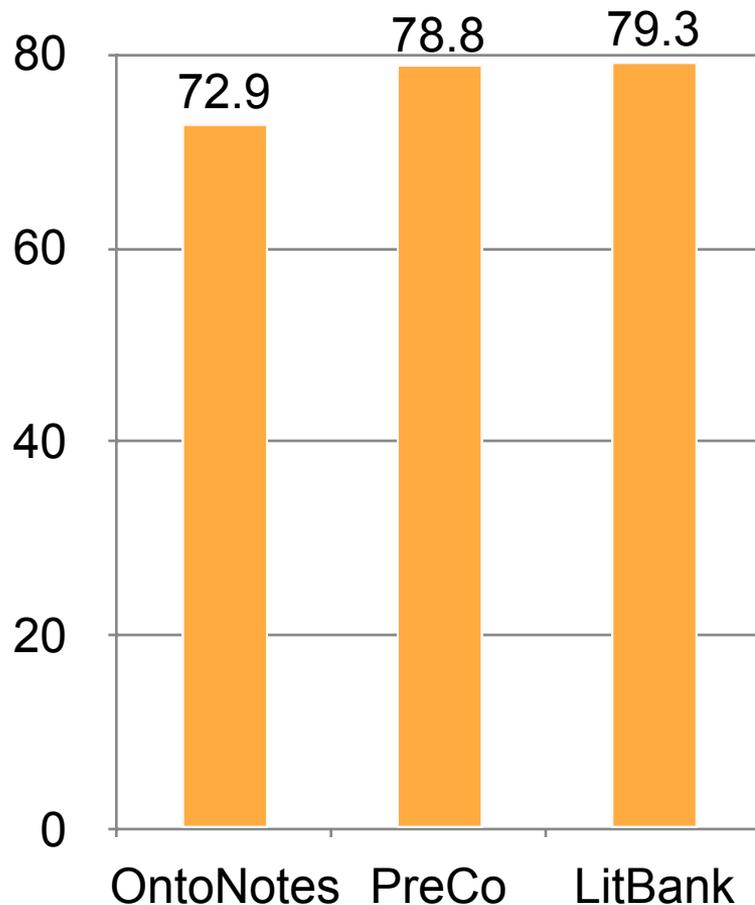
Spread between first and last mention of an entity, measured in tokens, weighted by number of mentions for each entity

- Long-range entities tend to be mentioned in bursts of attention.
- Character mentioned +100 times over 1500 tokens with lowest entropy (Basil Hallward, *Picture of Dorian Gray*) and highest entropy (narrator, *Gulliver's Travels*).

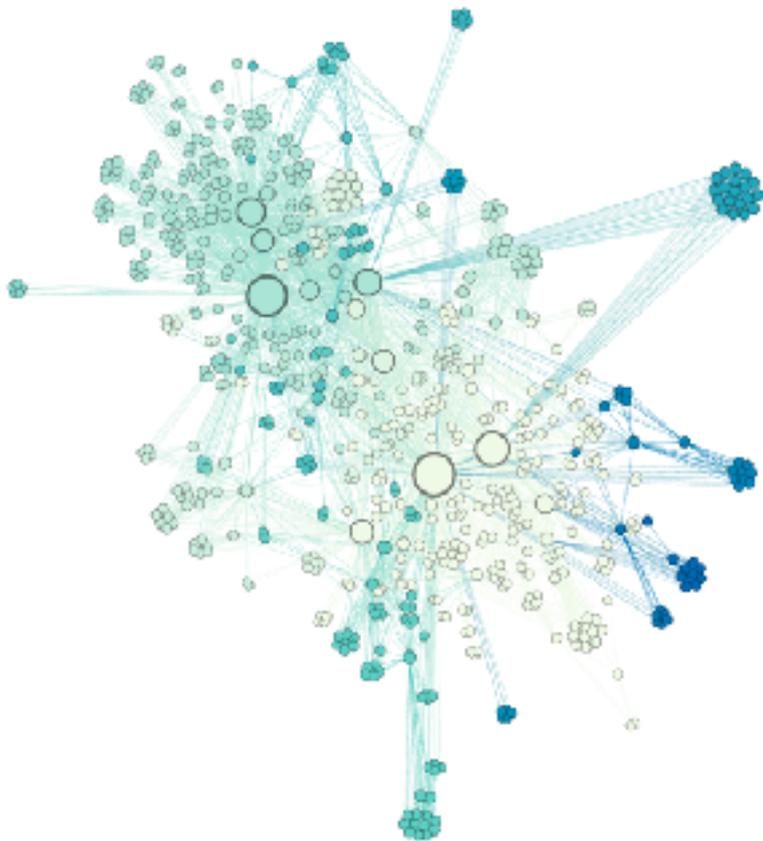


# Evaluation

- Train a BERT-based neural model on different training sources and evaluated its performance on held-out LitBank data:
  - OntoNotes (mainly news; 1.3M tokens)
  - PreCo (12.2M tokens)
  - LitBank (168K tokens)



How do we design NLP to drive *insight* into literary texts?

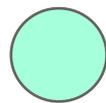
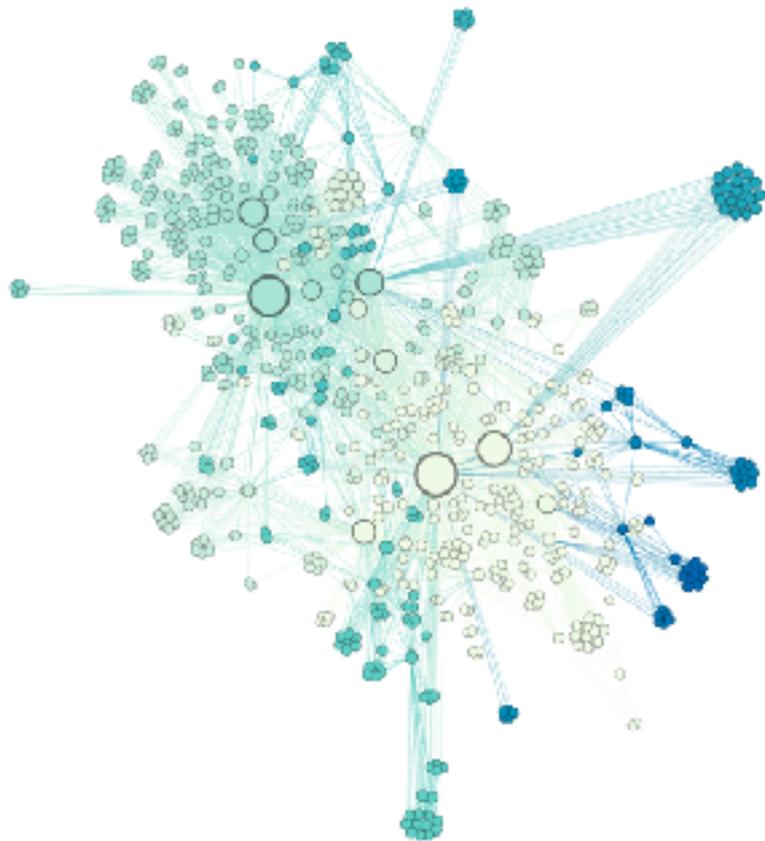


How does information propagate through **implicit** social networks?

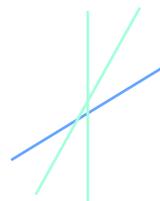


- Information diffusion in blogs (Gruhl et al., 2004; Leskovec et al., 2007)
- Spread of rumor and misinformation (Kwon et al., 2013; Friggeri et al., 2014; Del Vicario et al., 2016; Vosoughi et al., 2018)
- Textual reuse across different legislative bills (Wilkerson et al., 2015)

# *Great Expectations*

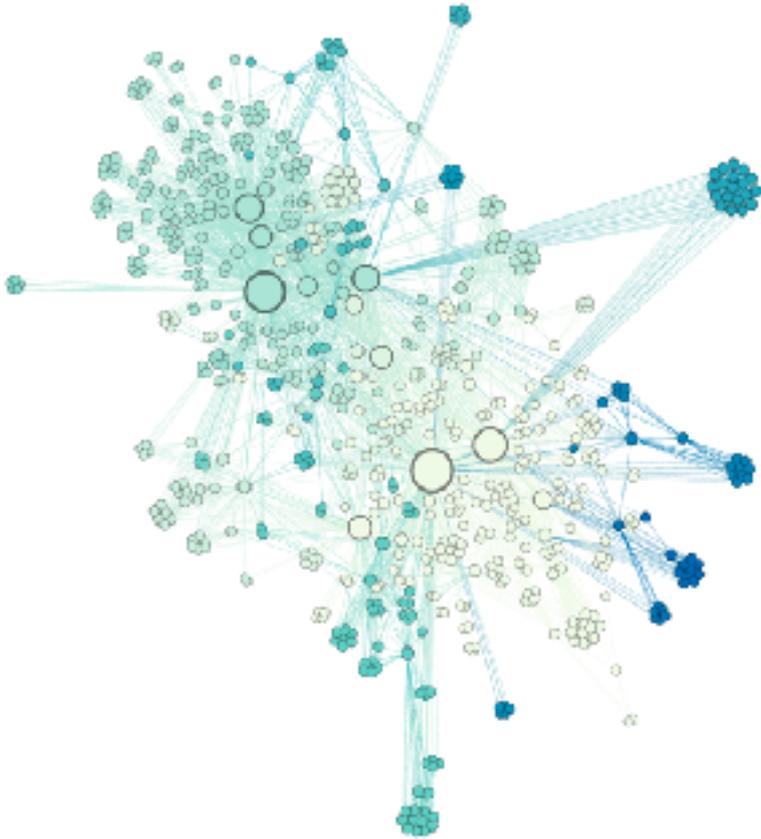


= Character



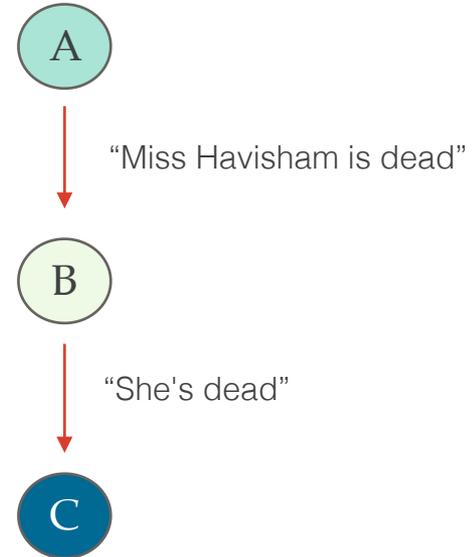
= Conversational  
Interaction

## *Great Expectations*



**Information** = Quoted Speech  
e.g.,  
“Miss Havisham is dead”

**Propagation** = Repetition of information  
across a character triad



# Research question

What are the structural properties of **information-propagating** nodes in fiction?

- “Gossip” among close friends, family, etc.: nodes that circulate information among densely-connected strong ties.
- **Information bridges**: nodes that pass information between otherwise disconnected communities.

- “Gossip” among close friends, family, etc.: nodes that circulate information among densely-connected strong ties.



Macbeth

- Information bridges: nodes that pass information between otherwise disconnected communities.



Macbeth

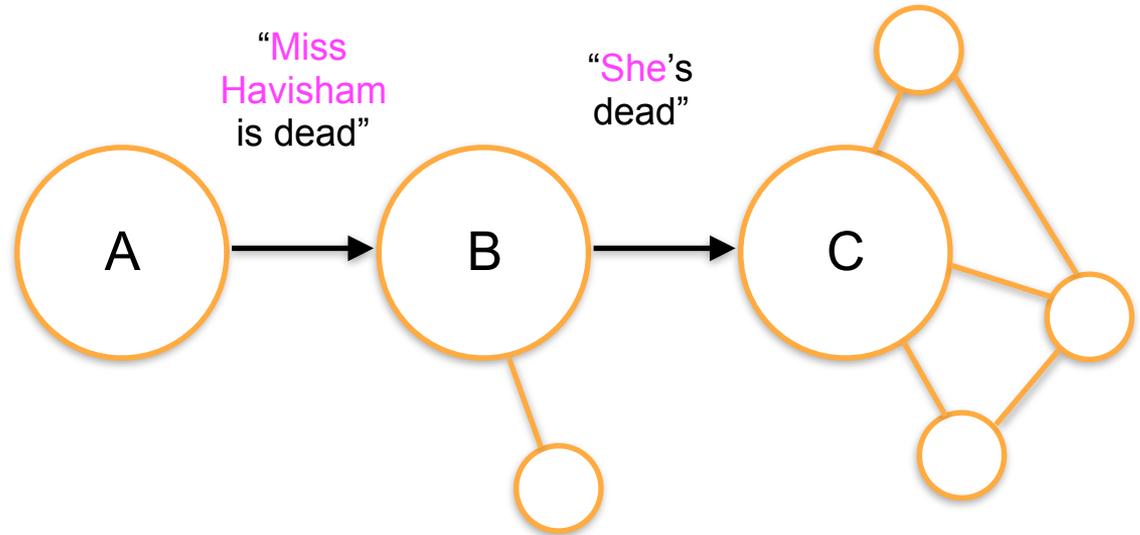
# NLP Pipeline

What are the structural properties of **information-propagating** nodes in fiction?



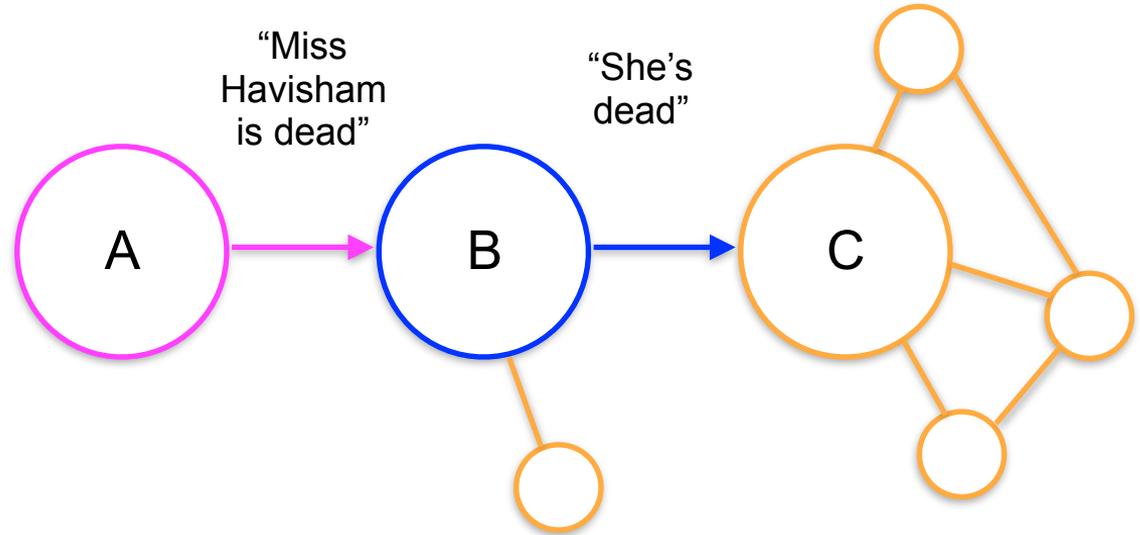
# Coreference resolution

- Identify unique characters from mentions of names, pronouns to construct network
- Identify when two mentions in quotations refer to the same individual
- $B^3 = 70.3$



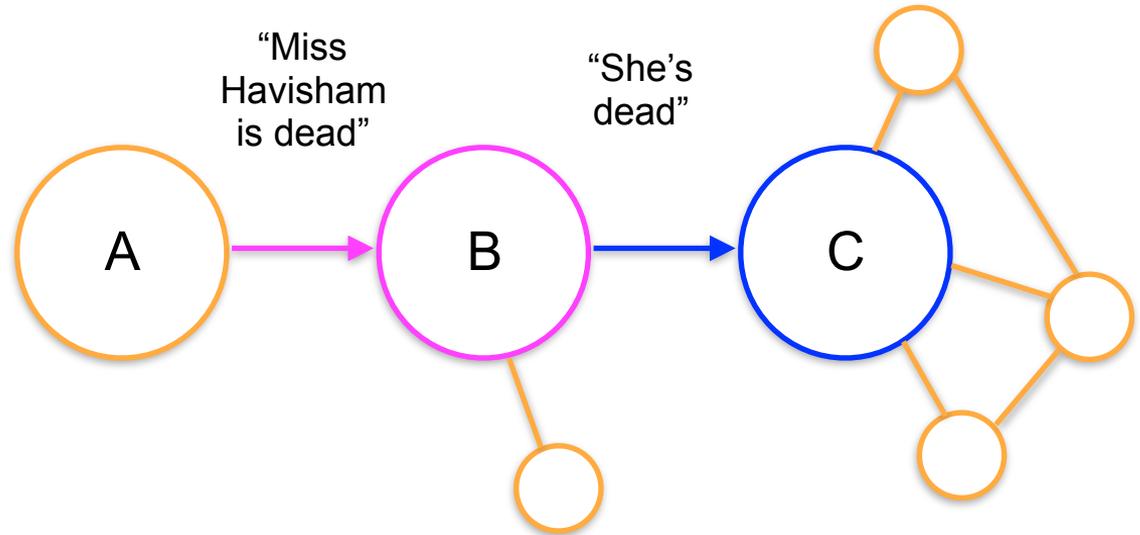
# Speaker attribution

- Link quotations to their speakers.
- $B^3 = 68.0$



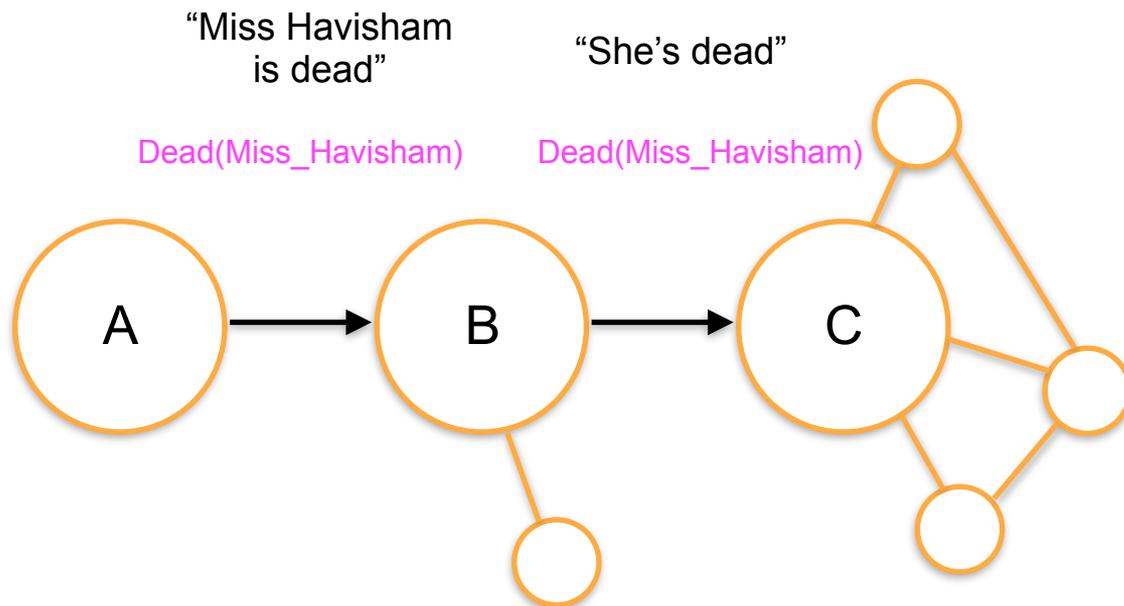
# Listener attribution

- Identify which characters were present when a given quotation was spoken.
- Identify blocks of dialogue; listeners = all characters mentioned in narrative



# Information extraction

- Extract an atomic unit of information in order to track its propagation.
- Resolve coreference and select propositional tuples of the form:  
[subject, verb, object]



# Research question

What are the structural properties of **information-propagating** nodes in fiction?

- “Gossip” among close friends, family, etc.: nodes that circulate information among densely-connected strong ties.
- Information bridges: nodes that pass information between otherwise disconnected communities.

# Data

- 15K books in Project Gutenberg (in the public domain in the US).
- Examine four high-precision topics:
- Each topic has  $> 100$  repeated tuple instances
- 4k of 15k books contain one of these repeated tuples

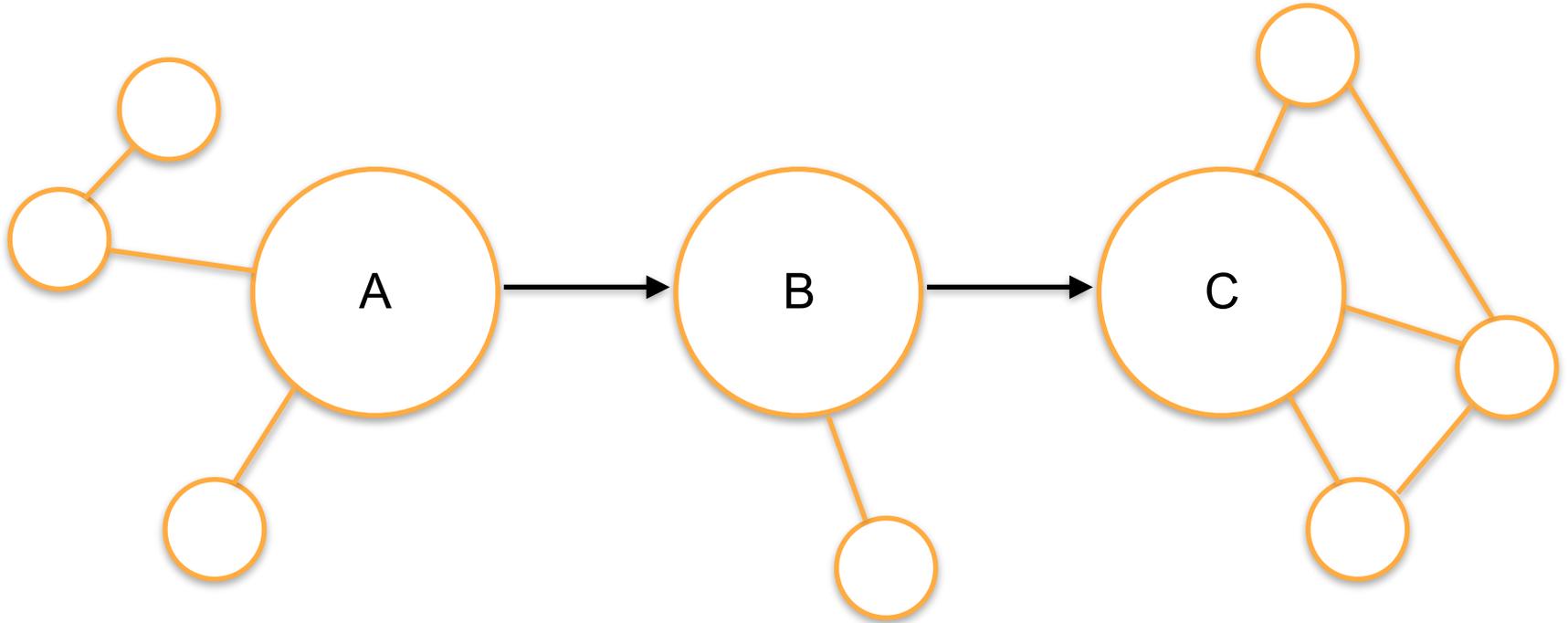
Love

Marry

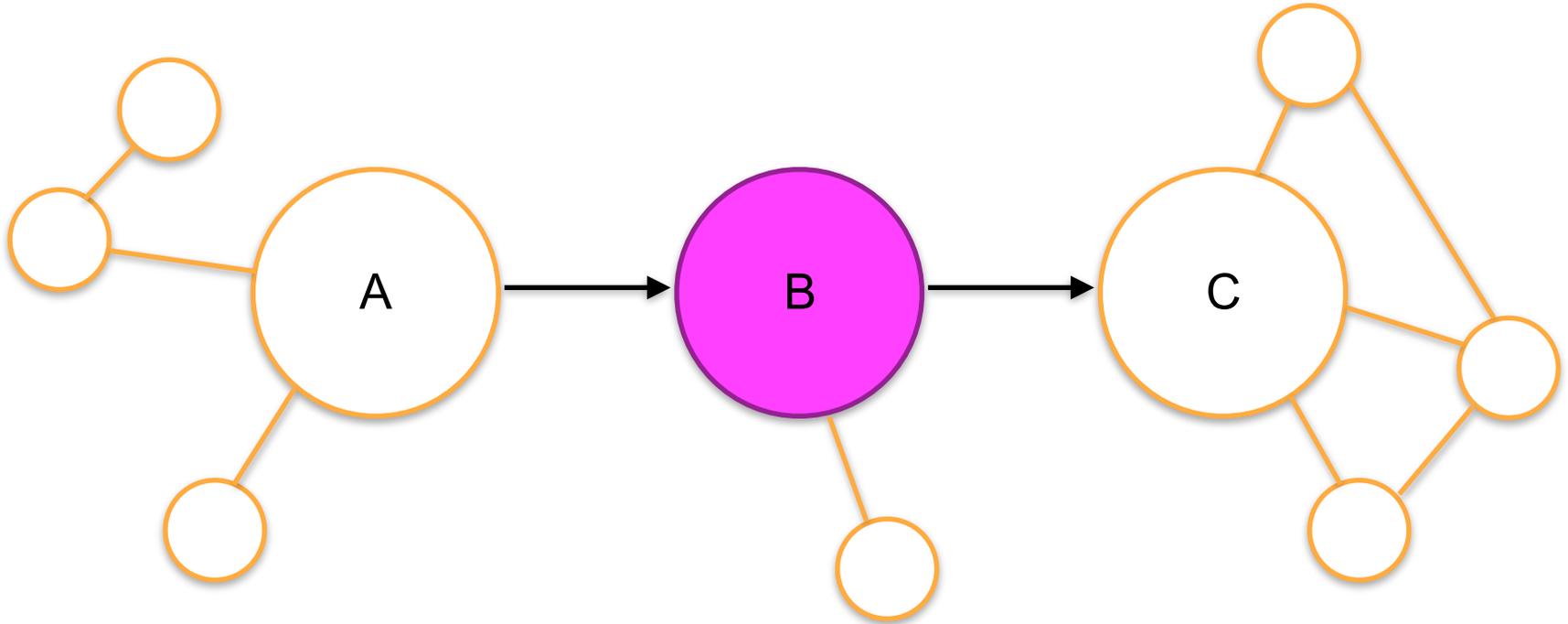
Kill

Die

# Propagation



# Propagation



What distinguishes successful propagation from unsuccessful propagation?

# Propagation



“The wicked witch is dead”

# Propagation



“You’ll never guess what happened to the wicked witch... she’s dead”

# Propagation



“You’ll never guess what happened to the wicked witch... she’s dead”

# Propagation

**A Node**



**B Nodes**



# Propagation

**A Node**



**B Nodes**



**C Node**



# Propagation

**A Node**



**B Nodes**



**C Node**



# Propagation

A Node



B Nodes



C Node



# Network Measures

## B Nodes



- Betweenness Centrality
  - Effective Size
  - Efficiency
- 
- Degree Centrality
  - Average Neighbor Degree
  - Number of Triangles

Information bridges

# Modeling

## B Nodes



x 3200



x 3200

- Use logistic regression to identify which node measures are meaningful
- Contextualize results to determine the network dynamics that are most integral to information propagation in literary fiction

# Feature Coefficients

Graph Measure	Model Coefficient	P-value < .01
Effective Size	11.6	*
Efficiency	3.8	*
Betweenness Centrality	2.3	*
Degree Centrality	-0.4	
Triangles	-6.2	*
Average Neighbor Degree	-9.5	*

# Feature Coefficients

Graph Measure	Model Coefficient	P-value < .01
Effective Size	11.6	*
Efficiency	3.8	*
Betweenness Centrality	2.3	*
Degree Centrality	-0.4	
Triangles	-6.2	*
Average Neighbor Degree	-9.5	*

“The significance of weak ties, then, would be that those which are local bridges create more, and shorter, paths .... The contention here is that removal of the average weak tie would do more ‘damage’ to transmission probabilities than would that of the average strong one.”

– Mark Granovetter

“The Strength of Weak Ties” (1973)



<https://github.com/dbamman/litbank>

README.md

## LitBank

---

LitBank is an annotated dataset of 100 works of fiction to support tasks in natural language processing and the computational humanities, described in more detail in:

David Bamman, Sejal Popat and Sheng Shen, "[An Annotated Dataset of Literary Entities](#)", NAACL 2019.



LitBank is licensed under a [Creative Commons Attribution 4.0 International License](#).

# Thanks!

David Bamman

[dbamman@berkeley.edu](mailto:dbamman@berkeley.edu)

- David Bamman, Olivia Lewke and Anya Mansoor, “An Annotated Dataset of Coreference in English Literature” (LREC 2020)
- David Bamman, Sejal Popat and Sheng Shen, “An Annotated Dataset of Literary Entities” (NAACL 2019)
- Matt Sims and David Bamman, “Information Propagation in Literary Social Networks” (in progress)

<https://github.com/dbamman/litbank>

