# Information Ordering: Entity Ordering & TSP

LING 573 — Systems & Applications Ryan Georgi — May 3rd, 2018







# Some Notes on Efficiency

- You won't be graded on efficiency
- But some thoughts re: efficiency on large files.







### DOM = "Document Object Model"

- Whole document gets parsed into tree structure
- Makes search easy
- Easy Read+Write
- SAX = "Simple API for XML"
  - Think of it as finite-state transducer
  - "Fires" on events (open/close tag, certain text)
  - Not easy to modify XML, but extremely efficient



# DOM vs SAX Parsing





# More Generally

### Object Model

- Easy to manipulate
- Can be memory heavy

### Stream Model

- Only for reading
- Triggers on events
- Very common for scalable microservices
  - e.g. Apache Kafka









- If you have large data streams
  - ... and only require a subset of information:
  - Best practice: treat as "stream" to be filtered



# Why this Matters







<DOC id="NYT ENG 19990901.0002" type="story" > <HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT></DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) /DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>





### open\_tag: DOC

### action:

*if* element.*id in* story\_list: begin\_text\_emit()

### </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> </TEXT></DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <DATELINE> (BC-EDIT-NY-MARCH-NYT) </DATELINE> <TEXT> free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <TEXT> </TEXT> </DOC>

<HEADLINE>



## SAX Example

<DOC id="NYT\_ENG\_19990901.0002" type="story" >

EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE

<P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P>

<HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE>

<P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to

<HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE>

<P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P>

PROFESSIONAL MASTER'S IN

**COMPUTATIONAL LINGUISTICS** 





### open\_tag: HEADLINE

<DOC id="NYT ENG 19990901.0002" type="story" > <HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT></DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) /DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>







### text: "EDITORIAL...

<DOC id="NYT ENG 19990901.0002" type="story" > <HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT></DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) </DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>





### close\_tag: HEADLINE

<DOC id="NYT ENG 19990901.0002" type="story" > <HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) </DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>







### open\_tag: DATELINE text:"(BC-EDIT... close\_tag: DATELINE

<DOC id="NYT ENG 19990901.0002" type="story" > <HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) </DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) </DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>







### close\_tag: DOC

### action:

end\_text\_emit()

<HEADLINE> EDITORIAL OBSERVER: NOW IT'S THE PRESIDENT'S TURN TO GAZE </HEADLINE> <DATELINE> (BC-EDIT-CLINTONS-NYT) /DATELINE> <TEXT> <P>> First Lady Hillary Rodham Clinton's campaign for Senate this week displayed what may be an unexpected new hazard - her husband. </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0003" type="story" > <HEADLINE> EDITORIAL: THE MARCH AND THE CONSTITUTION </HEADLINE> <DATELINE> (BC-EDIT-NY-MARCH-NYT) </DATELINE> <TEXT><P> The New York Times said in an editorial for Wednesday, Sept. 1: </P> <P> The U.S. District Court in New York ruled properly Tuesday that the city's attempt to prevent the so-called Million Youth March violated the constitutional right to free speech. ... </P> </TEXT> </DOC> <DOC id="NYT ENG 19990901.0004" type="story" > <HEADLINE> CLINTONS APPEAR SET TO BUY HOUSE IN CHAPPAQUA, N.Y. </HEADLINE> <DATELINE> WHITE PLAINS, N.Y. (BC-NY-CLINTONS-HOUSE-NYT) </DATELINE> <TEXT> <P> President Clinton and Hillary Rodham Clinton appeared Tuesday to be in the ... </P> </TEXT> PROFESSIONAL MASTER'S IN **COMPUTATIONAL LINGUISTICS** </DOC>



## SAX Example

<DOC id="NYT ENG 19990901.0002" type="story" >





Entity-Based Ordering





# The Entity Grid

- Rows = sentences
- Columns = discourse entities
- Values = grammatical role of mention in sentence
  - (S)ubject
  - (O)bject
  - X (other)

WASHINGTON

- (no mention)
- Multiple mentions Take highest grammatical ranking (S > O > X)

### Government **Jepartment** Microsoft Competit Competit Markets Products Brands Brands Case Netscape Software Tactics Evidenc [rial - - 0 - - X S 0 - - - - 22 - - s o - - - s o o - - - 33

6



Suit



	Department	Trial	Microsoft	Evidence	Competitors	Markets	Products	Brands	Case	Netscape	Software	Tactics	Government	Suit	Earnings	
1	S	0	S	X	0	_	_		_	—	_	_	_		_	1
2	—	_	0	_	—	X	S	0	_	—	_	_	_	_	—	2
3	_	_	S	0	_		_	_	S	0	0	_	_	_	_	3
4	—	_	S	—	—		_	_	_	_	_	S	_	_	_	4
5	_	_	_	_	_		_	_	_	_	_	_	S	0	_	5
6	_	X	S	_	_	_	_	_	_	_	_	_	_	_	0	6

- 2 competitive enough to unseat [established brands]o.
- 3 merging [browser software]o.
- 4 [Microsoft]s claims [its tactics]s are commonplace and good economically
- 5 [collusion] x is [a violation of the Sherman Act] o.
- [Microsoft]s continues to show [increased earnings]o despite [the trial]x. 6

WASHINGTON

[The Justice Department]s is conducting an [anti-trust trial]o against [Microsoft Corp.]x with [evidence]x that [the company]s is increasingly attempting to crush [competitors]o

[Microsoft]o is accused of trying to forcefully buy into [markets]x where [its own products]s are not

[The case]s revolves around [evidence]o of [Microsoft]s aggressively pressuring [Netscape]o into

[The government]s may file [a civil suit]o ruling that [conspiracy]s to curb [competition]o through



- Intuitions
  - Some columns dense: focus of text (e
    - Likely to take certain roles, e.g. S, O
  - Others **sparse**: likely other roles(X)
  - Local transitions reflect structure, topic
- local entity transitions:  $\{S,O,X,-\}^n$ 
  - Continuous column subsequences ("role n-grams"?)
  - Compute probability of sequence over grid:
    - # of occurrences of that type/# of occurrences of that length





.g. Microsoft)		Department	Trial	Microsoft	Evidence	Competitors	Markets	Products	Brands	Case	Netscape	Software	Tactics	Governmen	Suit
	1	S	0	S	X	0	—	-				_	_	_	_
	2	_	-	0	-	-	X	S	0		_	_	_	_	_
	3	_	-	S	0	-	_	-	_	S	0	0	-	_	_
c shifts	4	_	-	S	-	-	_	-	_		_	_	S	_	_
	5	_	-	—	-	-	—	-	_			_	_	S	0
	6	_	X	S	-	-	—	-	_		_	_	_	_	_





### Vector Representation 0 0 0 .03 0 0 0 $d_1$ $d_2$ ' 0 0 0 .02 0 .07 0 .02 0 0 .03 0 0 0 $d_{\mathcal{Z}}$

• Document vector:

• Length = # of transition types

• Values = Probabilities of each transition type (MLE counts over document)

• 
$$p(\mathbf{s}, -|d_1) = 0.3$$



	$\mathbf{S}$	0			$\mathbf{V}$	$\bigcirc$			
0	X	X	X	X	I		1	I	
.02	.07	0	0	.12	.02	.02	.05	.25	
.02	0	0	.06	.04	0	0	0	.36	
.06	0	0	0	.05	.03	.07	.07	.29	





# Entity-Based Ordering: Training

- I. Identify entity classes
  - a. initialize grid
- 2. Identify entity roles (Parse)
  - a. Fill out grid
- 3. Learn support vector for ranking (<u>loachims, 2002a</u>) (<u>Software</u>)

  - b. In this case:
    - Optimally coherent (original data)

c. Train vector to prefer original ordering to rearranged ordering

VASHINGTON

a. Support vector = vector that maximizes distance between two groups of data points

Non-optimally coherent (synthetically rearranged versions of original data)



# Entity-Based Ordering: Testing

- I. Identify entity classes
- 2. Identify entity roles (parse)
- 3. Using learned support vector:
  - a. Compare all new sentences matrices to support vector
  - b. Retrieve score based on distance (inner product) to support vector
  - c. Return list of ranked sentences based on this score





## Dependencies & Comparisons: **Tools Needed**

- **Coreference** system for linking entities
  - Authors used Ng & Cardie (2002)
  - Neural Coref System [link]
- **Parser** for grammatical roles:
  - Authors used CKY (Collins) to determine NP position (Subject vs. Object)
  - Might consider Stanford Parser
    - Source for Java
    - **NLTK API for Stanford Parser**







## Dependencies & Comparisons: **Tools Needed**

### • Support Vector Machine Ranking Implementation:

- For learning/applying pairwise rankings:
- (Joachims, 2002a) (Software)







• Shows accuracy of pairwise rankings between sentences

### • Settings:

- Salient:  $\geq$  2 occurrences of entity
- Transition length: 2

### • Features:

- syntax- = No S, O, X, just +/-
- coref— = don't perform coref resolution



## Results

Model	Accur
Coreference + Syntax + Salience +	80.0
Coreference + Syntax + Salience –	75.0
Coreference + Syntax – Salience +	78.8
Coreference – Syntax + Salience +	83.8
Coreference + Syntax – Salience –	71.3
Coreference – Syntax + Salience –	78.8
Coreference – Syntax – Salience +	77.5
Coreference – Syntax – Salience –	73.8

**52.5**\*\* Latent Semantic Analysis













# Discussion

### **Best Results**:

- Richer syntax and salience models
  - **not** coreference (though not significant d
- Why?
  - Synthetic nature of training data
  - unreliable coref
- Worst results:
  - Significantly worse with both simple syntax and no salience
  - Still not horrible: 71% vs. 84%
    - Much better than LSA model: 52.5%

WASHINGTON

	Model	Accura
	Coreference + Syntax + Salience +	80.0
rop)	Coreference + Syntax + Salience –	75.0
	Coreference + Syntax – Salience +	78.8
	Coreference – Syntax + Salience +	83.8
	Coreference + Syntax - Salience -	71.3*
	Coreference – Syntax + Salience –	78.8
	Coreference – Syntax – Salience +	77.5
	Coreference – Syntax – Salience –	<b>73.8</b> <sup>*</sup>
and		

Latent Semantic Analysis

**52.5**\*\*









# Comparisons

- Latent Semantic Analysis <u>Landauer et al, (1997)</u>
- "Catching the Drift" <u>Barzilay & Lee (2004)</u>







# Latent Semantic Analysis

- Distributional semantics
  - Same concept as word embeddings, different approach
- Create *term* × *document* matrix over large news corpus
  - Perform SVD to create 100-dimensional dense matrix
- Score summary as:
  - Sentence represented as mean of its word vectors
  - Average of cosine similarity scores of adjacent sentences
    - Local "concept" similarity score









### • **Intuition** — stories are:

- Composed of topics + subtopics
- Unfold in systematic, sequential manner
- Can represent ordering as sequence modeling over topics

### • Approach:

• HMM over topics



### "Catching the Drift" Barzilay and Lee, 2004





- Learn topics in unsupervised way from data
  - Assign sentences to topics
- Learn sequences from document structure
  - Given clusters, learn sequence model over them
  - No explicit topic labeling, no hand-labeling of sequence









- How to induce a set of "topics" from document set?
  - Within a given docset
- Unsupervised approach: Clustering
  - Similarity measure?
    - Cosine similarity over sentences' word bigrams
- Assume some irrelevant/off-topic sentences
  - Merge clusters with members below a threshold into "et cetera" cluster
- Result: *m* clusters



## **Topic Induction**





# Induced Topics

• Example earthquake topic

WASHINGTON

The Athens seismological institute said the temblor's epicenter was located 380 kilometers (238 miles) south of the capital.

Seismologists in Pakistan's Northwest Frontier Province said the temblor's epicenter was about 250 kilometers (155 miles) north of the provincial capital Peshawar.

The temblor was centered 60 kilometers (35 miles) north-west of the provincial capital of Kumming, about 2,200 kilometers (1,300 miles) southwest of Beijing, a bureau seismologist said.





### \*Side Note:

*p*("temblor"|*MyLanguageModel*) =



WASHINGTON





•

# Sequence Modeling

 $D(c_i, c_j)$ 

- Hidden Markov Model
  - States = Topics
    - State  $s_m$ : special "insertion" state (for digressions/other topics)
- Transition Probabilities:
  - Evidence for ordering?

 $p(s_j | s$ 

Smoothing Factors



### Document ordering — sentence from topic a appears before sentence from topic b

Count of documents in which sentence from cluster i immediately precedes one from j

Count of documents with sentence from cluster i



# Sequence Modeling

- Emission Probabilities:
  - Standard topic state:
    - Probability of observation given state (topic)
      - Probability of sentence under topic-specific bigram LM
      - Bigram probabilities
- Etcetera State:
  - Forced complementary to other states



e (topic) ecific bigram LM

$$p_{s_i}(w' | w) = \frac{f_{c_i}(ww') + \delta_1}{f_{c_i}(w) + |V|}$$

$$w' | w) = \frac{1 - \max_{i:i < m} p_{s_i}(w' | w)}{\sum_{u \in V} (1 - \max_{i:i < m} p_{s_i}(u | w))}$$



# Sequence Modeling

- Viterbi Re-estimation
  - Intuition: Refine clusters, etc. based on sequence info
  - Iterate:
    - Run Viterbi decoding over original documents
    - Assign each sentence to cluster most likely to generate it Use new clustering to recompute transition/emission
  - Until stable (or max iterations)







# Sentence Ordering Comparison

- Restricted Domain Text:
  - Separate collections of earthquake, aviation accidents
  - LSA predictions: which order has higher score
  - Topic/Context model: Highest probability under HMM

### Model

Coreference+Syntax+Salier Coreference+Syntax+Salier Coreference-Syntax-Salier Coreference-Syntax+Salier Coreference-Syntax-Salier Coreference-Syntax-Salier Coreference-Syntax-Salier

HMM-based Content Mode Latent Semantic Analysis



	Earthquakes	Accidents
ence+	87.2	90.4
nce–	88.3	90.1
nce+	86.6	88.4**
nce+	83.0**	89.9
nce–	86.1	89.2
nce–	82.3**	88.6*
nce+	83.0**	86.5**
nce–	81.4**	86.0**
els	88.0	75.8**
	81.0**	87.3**



# Example Text from Earthquakes

been evacuated from the area following Monday's quake.



• A strong earthquake hit the China-Burma border early Wednesday morning, but there were no reports of deaths, according to China's Central Seismology Bureau. The 7.3 quake hit Menglian county at 5:46am. The same area was struck by a 6.2 temblor early Monday morning, the bureau said. The bureau's Xu Wei said some buildings sustained damage and there were some injuries, but he had no further details. Communication with the remote region is difficult, and satellite phones sent from the neighboring province of Sichuan have not yet reached the site. However, he said the likelihood of deaths was low because residents should have



# **Example Text from Accidents**

weather was marginal VFR and deteriorating rapidly. Witnesses near Pine haze/fog.

WASHINGTON

• When the pilot failed to arrive for his brother's college graduation, concerned family members reported that he and his airplane were missing. A search was initiated, and the Civil Air Patrol located the airplane on top of Pine Mountain. According to the pilot's flight log, the intended destination was Pensacola, FL, with intermediate stops for fuel at Thomson, GA, and Greenville, AL. Airport personnel at Thomson confirmed that the airplane landed about 1630 on 11/6/97. They reported that the pilot purchased 26.5 gallons of 100LL fuel and departed about 1700. Witnesses at the Thomson Airport stated that when he took off, the Mountain stated that the visibility at the time of the accident was about 1/4 mile in



# Summary Coherence Scoring Comparison

- Domain Independence
  - Too little data per domain to estimate topic-content model
    - Train: 144 pairwise summary rankings
    - Test: 80 pairwise summary rankings
  - Entity grid model (best): 83.8%
  - LSA model: 52.5%
- Likely issue:
  - Bad auto summaries highly repetitive  $\rightarrow$  High inter-sentence similarity









Keeping it CLASSY ... 2006





# Ordering as Optimization

- Given a set of sentences to order
  - Define a local pairwise coherence score between sentences
  - Compute a total order optimizing local distances
- Can we do this efficiently?
  - Optimal ordering equivalent to the TSP = Traveling Salesperson Problem
    - Given a list of cities and distances between cities, find the shortest route that visits each city exactly once and returns to the origin city.
    - NP-complete (NP-hard)







# Ordering as TSP

- Can we do this practically?
  - DUC summaries are 250 words, so 6-10 sentences
  - I0 sentences have how many possible orders?
    O(n!)
  - Not impossible
- Alternatively
  - Use an approximation method
  - Take the best of a sample







### **CLASSY 2006** Conroy et al (2006)

- Formulates ordering as TSP
- Requires pairwise sentence distance measure
  - Term-based similarity: # of overlapping terms
  - **Document similarity** 
    - Multiply by a weight if in the same document (there, 1.6)
  - Normalize to between 0 and 1 (sqrt of product of self similarity)
    - Make distance: subtract from 1









# Practicalities of Ordering

- Brute force: O(n!)
  - "there are only 3,628,800 ways to order 10 sentences plus a lead sentence, so exhaustive search is feasible." (Conroy)
- Still...

WASHINGTON

- Used sample set to pick best
  - Candidates:
    - Random
    - Single-swap changes from good candidates
- 50K enough to consistently generate minimum cost order



## Results

Submission	Mean	95% CI Lower	95% CI Upper
С	0.47856	0.46005	0.49867
В	0.47284	0.45479	0.49179
D	0.46873	0.45129	0.48687
G	0.45909	0.44700	0.47187
J	0.45892	0.43983	0.47905
А	0.45816	0.44963	0.46821
Ι	0.45347	0.43629	0.46972
Н	0.44834	0.43036	0.46385
E	0.44316	0.42297	0.46286
F	0.43578	0.41054	0.45968
$15\mathrm{E}$	0.42282	0.41675	0.42845
24	0.41108	0.40488	0.41706
12	0.40488	0.39919	0.41051
23	0.40440	0.39820	0.40967
10	0.40369	0.39852	0.40874
15	0.40279	0.39649	0.40839
33	0.40206	0.39684	0.40754
8	0.40010	0.39370	0.40597
28	0.39922	0.39364	0.40460

 Table 1: Average F score of ROUGE 1 Scores

WASHINGTON

Î I

Submission	Mean	95% CI Lower	95% CI Upper
С	0.13260	0.11596	0.15197
D	0.12380	0.10751	0.14003
В	0.11788	0.10501	0.13351
G	0.11324	0.10195	0.12366
F	0.10893	0.09310	0.12780
Η	0.10777	0.09833	0.11746
J	0.10717	0.09293	0.12460
Ι	0.10634	0.09632	0.11628
E	0.10365	0.08935	0.11926
А	0.10361	0.09260	0.11617
24	0.09558	0.09144	0.09977
15E	0.09539	0.09145	0.09922
15	0.09097	0.08671	0.09478
12	0.08987	0.08583	0.09385
8	0.08954	0.08540	0.09338

 Table 2: Average F score of ROUGE 2 Scores



## Information Ordering: Conclusions

- Many cues to ordering:
  - Temporal, coherence, cohesion
  - Chronology, topic structure, entity transitions, similarity
- Strategies:
  - Heuristic, machine-learned, supervised, unsupervised
  - Incremental buildup vs. generate & rank
- lssues;

WASHINGTON

Domain independence, semantic similarity, reference 





## Digression: Information Ordering, Eggs, and Holy Wars

- "Big Endian" vs "Little Endian"
  - The information order of the most significant digit
  - Ongoing problem in standardizing computers







## Digression: Information Ordering, Eggs, and Holy Wars

- Where did it come from?
- Gulliver's Travels!





