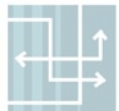


Identifying Sentence-Level Semantic Content Units with Topic Models

Leonhard Hennig, Thomas Strecker, Sascha Narr,
Ernesto William De Luca, Sahin Albayrak
DAI-Labor, TU Berlin



CC IRML
Information Retrieval
& Machine Learning



Content modeling for multi-document summarization

Gold-standard content units

A topic modeling approach to content unit discovery

Evaluation

- Word distribution similarity
- Sentence distribution similarity

Conclusion & Outlook

Summarize a set of documents related to a single event or topic

Challenge: Identification of similar information

- frequency indicates importance
- avoid redundancy



Typically solved by sentence extraction

- Compute sentence relevance based on features of interest
- Rank and select subset of important and non-redundant sentences

Previous Work:

Sentence clustering

Word frequency

Graph-based ranking

Maximum Marginal Relevance

Problems:

- Mostly lexical overlap, Vector Space Model
- Content unit granularity: sentence

But:

- Content units can be as 'small' as clauses (e.g. facts)
- Content units can be combined differently into sentences

NEW YORK, July 17 --

The U.S. Coast Guard and the Air National Guard are conducting a massive search off the coast of Long Island, N.Y. for a small plane carrying John F. Kennedy Jr., son of the 35th U.S. President, U.S. media reported Saturday.

The search began Saturday morning in an area covering some 1,000 square miles, presumably the flight path of Kennedy's plane, searchers said.

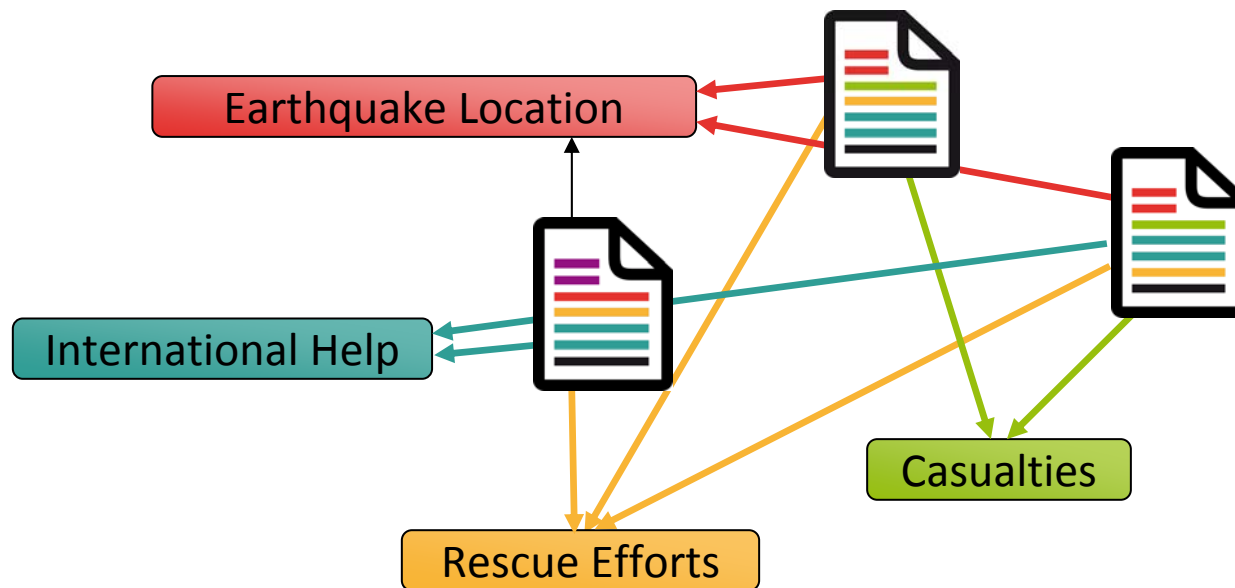
WASHINGTON, July 17 --

A small plane carrying John F. Kennedy Jr., son of the former U.S. president, was reported missing early Saturday, and a search was under way off the coast of New York's Long Island, official sources said.

The U.S. Coast Guard confirmed it was searching for the plane with help from the Air National Guard.

The search was being conducted in water off eastern tip of Long Island, along the presumed flight path of Kennedy's plane.

- Documents from same domain often exhibit similar sub-topic structuring and similar word patterns (Barzilay 04)



- Similar content units in human reference summaries are often expressed using similar word patterns, but with variations in word choice (e.g. synonymy) and word order (Nenkova 04, Harnly 05)

Can we automatically discover these similar content units?

Challenge 1: Semantic similarity

- Synonymy and polysemy
- Variations in word order and word choice

Challenge 2: (Sub-) sentential content units

- Granularity: sentence or clause
- Combined differently into sentences

Content modeling for multi-document summarization

Gold-standard content units

A topic modeling approach to content unit discovery

Evaluation

- Word distribution similarity
- Sentence distribution similarity

Conclusion & Outlook

Copies

- Sentences are verbatim copies of each other

Similar

- Sentences express the same content

- (a) The supreme court struck down as unconstitutional a law giving the president a line-item veto which lets him cancel specific items in tax and spending measures.
- (b) The U.S. supreme court Thursday struck down as unconstitutional the line-item veto law that lets the U.S. president strike out specific items in tax and spending measures.

Clause

- Sentence clauses are repeated verbatim or with similar word usage as clauses of other sentences

- (a) Germany, Azerbaijan, Greece, France, the Netherlands, Kazakhstan, Ukraine and Russia have been participating in the fight against the blaze that threatened to engulf the entire field of 30 storage tanks containing 1 million tons of crude oil.
- (b) However, he said the strong fire had destroyed seven storage tanks and damaged two other ones in the refinery which held 30 storage tanks containing 1 million tons of crude oil.
- (c) Germany, Azerbaijan, Greece, France, the Netherlands, Kazakistan, Ukraine and Russia participated in the fight against the blaze.

Unique

- Sentences express information not repeated in any other sentence

11 closely related document pairs selected from DUC 2007 multi-document summarization data set

3 annotators

$$\hat{\mathbb{H}}^{(a)} =$$

Content Units	Sentences					
	s ₁	s ₂	...	s _{m-1}	s _m	
z' ₁	1	1		0	0	
z' ₂	0	1		1	0	
...						
z' _{k-1}	0	0		1	0	
z' _k	0	0		0	1	

- Inter-annotator agreement on number of content units: **0.97**
- Inter-annotator agreement on content units: **0.69**
 (i.e. agreement on "identical" content units)

Distribution of types

53 % unique

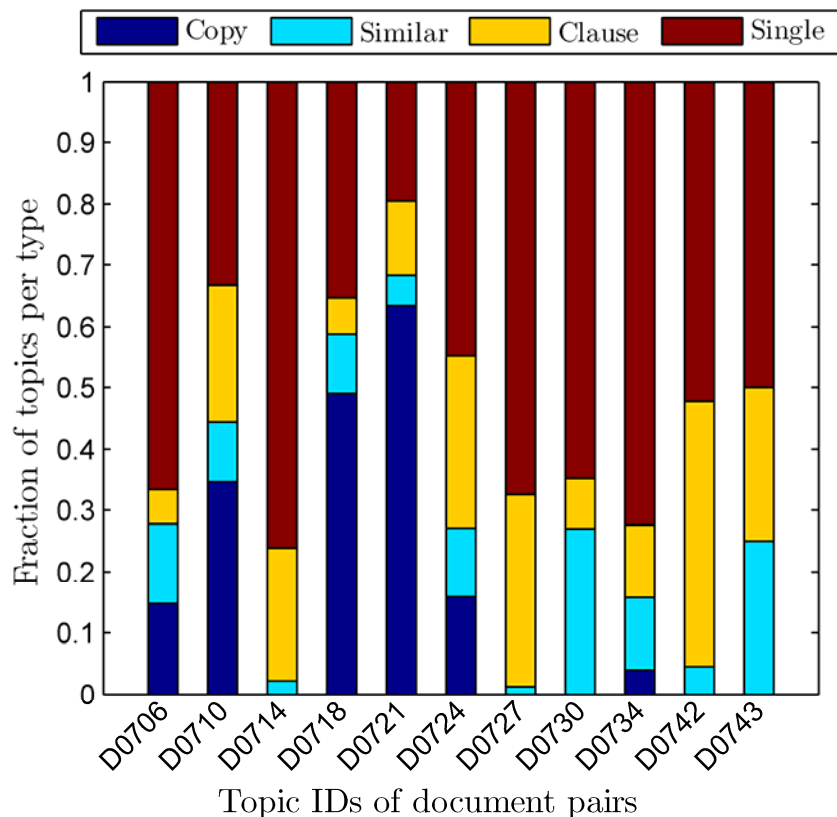
20% clause

11% similar

ML estimate for content unit word likelihood

$\hat{\Phi} = P(w | \hat{z}^{(a)}), \text{ where}$

$$p(w_i | \hat{z}_k^{(a)}) = \frac{w_{is}}{\sum_{s \in \hat{z}_k^{(a)}} \sum_i w_{is}}$$



Content modeling for multi-document summarization

Gold-standard content units

A topic modeling approach to content unit discovery

Evaluation

- Word distribution similarity
- Sentence distribution similarity

Conclusion & Outlook

Utilize Latent Dirichlet Allocation (LDA, Blei et al. 2003)

- Exploit word-document co-occurrence patterns
- Model each document as a distribution over topics, each topic as a distribution over words

For each document pair (d_1, d_2) :

- Create Term-**Sentence** Matrix \mathbf{A}
- \mathbf{A}_{ij} = frequency of term i in sentence j
- $M = |s \in d_1 \cup s \in d_2|$

$\mathbf{A} =$

Terms	Sentences				
	s_1	s_2	...	s_{m-1}	s_m
w_1	1	2		0	0
w_2	1	0		1	0
...					
w_{n-1}	0	0		0	3
w_n	0	1		0	1

Gibbs sampling, 2000 iterations

Compute $p(w|z)$ and $p(z|s)$ from single sample

$\Phi =$

		z_1	z_2	...	z_{k-1}	z_k
Terms	w_1	0.02	0.03		0.0	0.0
	w_2	0.02	0.0		0.018	0.0
	...					
	w_{n-1}	0.0	0.0		0.0	0.04
	w_n	0.0	0.015		0.0	0.017

Topics

$\Theta =$

		s_1	s_2	...	s_{m-1}	s_m
Topics	z_1	0.0	0.6		0.0	0.51
	z_2	0.98	0.38		0.0	0.0
	...					
	z_{k-1}	0.0	0.0		0.98	0.0
	z_k	0.0	0.0		0.0	0.48

Sentences

Content modeling for multi-document summarization

Gold-standard content units

A topic modeling approach to content unit discovery

Evaluation

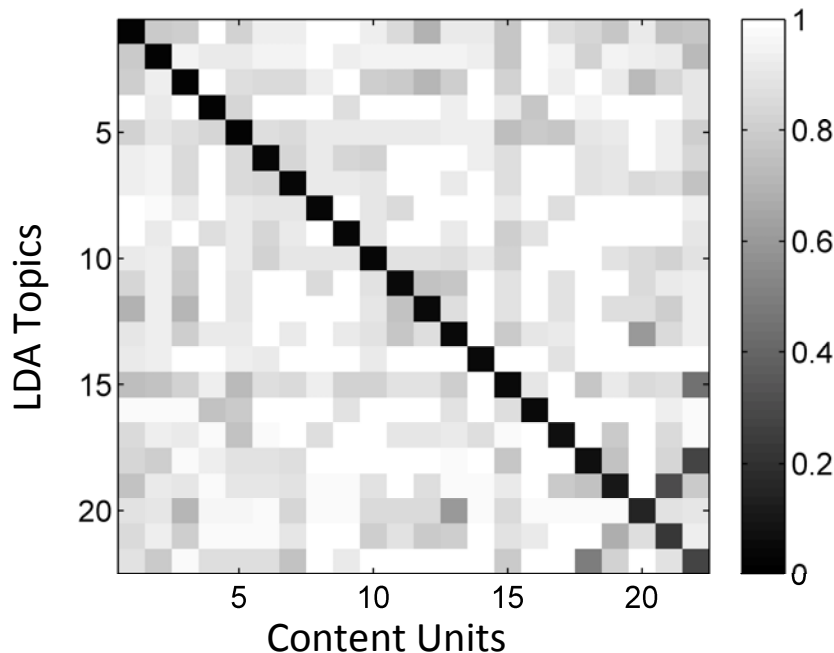
- Word distribution similarity
- Sentence distribution similarity

Conclusion & Outlook

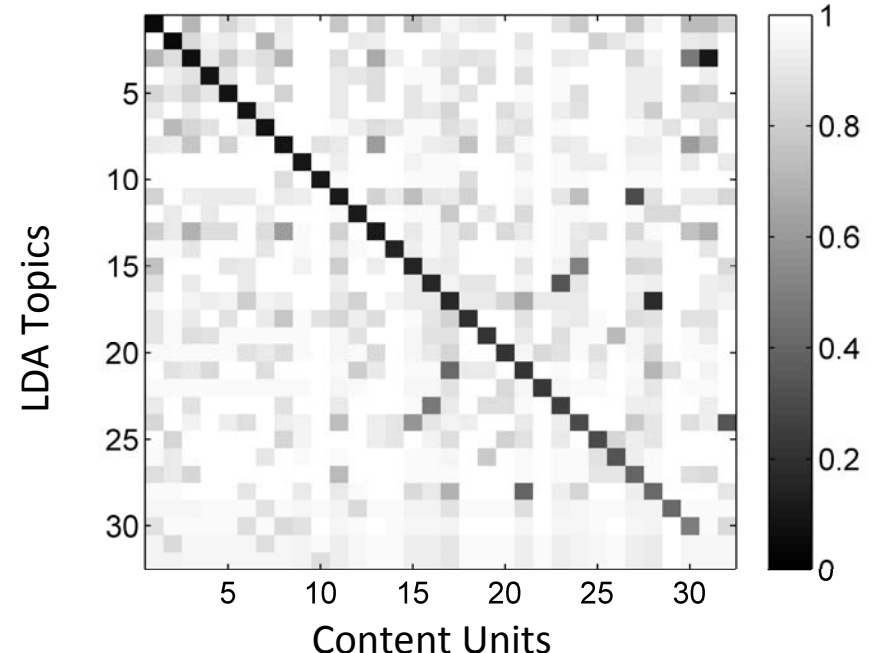
Pair-wise comparison of word distributions $\hat{\Phi}$, Φ of annotated content units and latent topics using Jensen-Shannon divergence

- Diagonal shows matches of most similar topic-content unit pairs
- Lower (darker) is better
- Reordered so that best matches are in upper left

Document pair from D0710

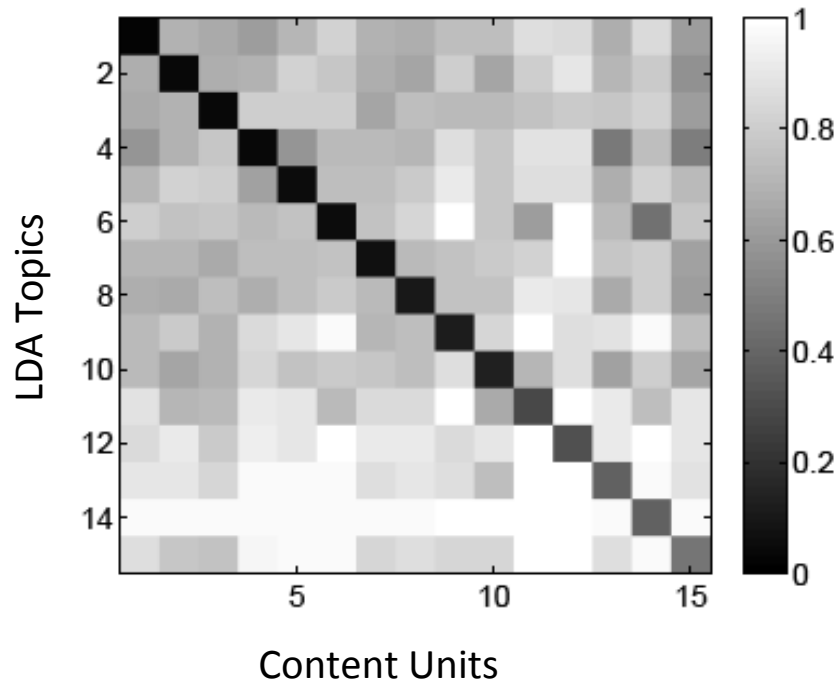


Document pair from D0734

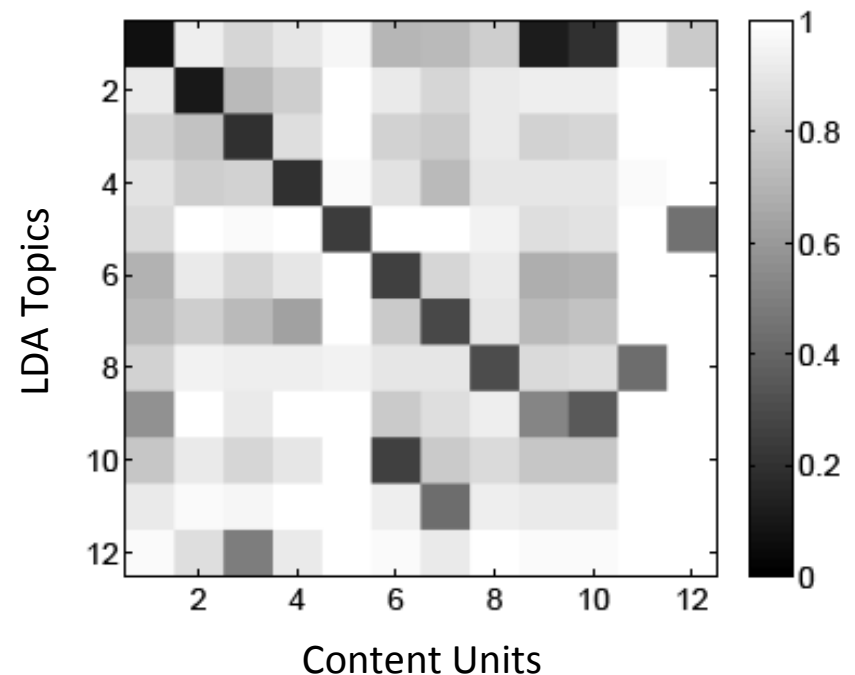


LDA does not always result in topics matching well with manually annotated content units

Document pair from D0730



Document pair from D0742



Content modeling for multi-document summarization

Gold-standard content units

A topic modeling approach to content unit discovery

Evaluation

- Word distribution similarity
- Sentence distribution similarity

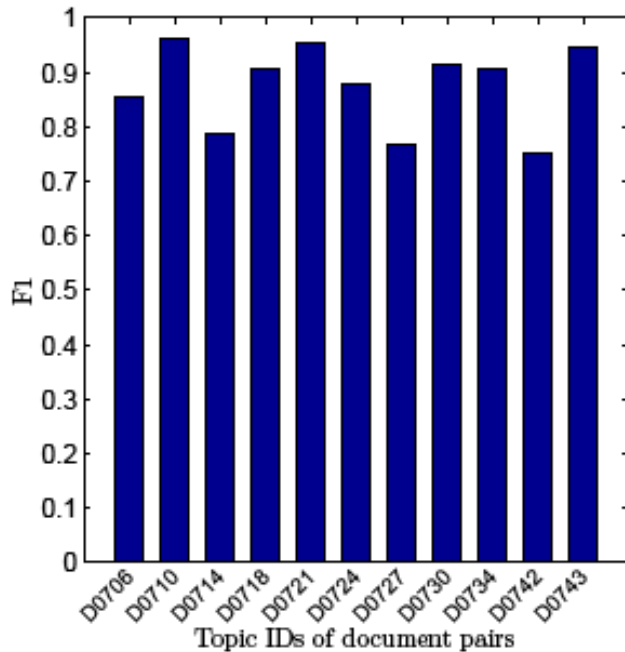
Conclusion & Outlook

Compare sentence-topic assignments

$\Theta =$	Topics	<table style="width: 100%; border-collapse: collapse;"> <thead> <tr style="background-color: #4a86e8; color: white;"> <th style="width: 10%;"></th> <th style="width: 15%;">s₁</th> <th style="width: 15%;">s₂</th> <th style="width: 15%;">...</th> <th style="width: 15%;">s_{m-1}</th> <th style="width: 15%;">s_m</th> </tr> </thead> <tbody> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z₁</td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.6</td> <td></td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.51</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z₂</td> <td style="text-align: center;">0.98</td> <td style="text-align: center;">0.38</td> <td></td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.0</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">...</td> <td></td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z_{k-1}</td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.0</td> <td></td> <td style="text-align: center;">0.98</td> <td style="text-align: center;">0.0</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z_k</td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.0</td> <td></td> <td style="text-align: center;">0.0</td> <td style="text-align: center;">0.48</td> </tr> </tbody> </table>		s ₁	s ₂	...	s _{m-1}	s _m	z ₁	0.0	0.6		0.0	0.51	z ₂	0.98	0.38		0.0	0.0	...						z _{k-1}	0.0	0.0		0.98	0.0	z _k	0.0	0.0		0.0	0.48	Sentences
	s ₁	s ₂	...	s _{m-1}	s _m																																		
z ₁	0.0	0.6		0.0	0.51																																		
z ₂	0.98	0.38		0.0	0.0																																		
...																																							
z _{k-1}	0.0	0.0		0.98	0.0																																		
z _k	0.0	0.0		0.0	0.48																																		
$\hat{\Theta}^{(a)} =$	Content Units	<table style="width: 100%; border-collapse: collapse;"> <thead> <tr style="background-color: #4a86e8; color: white;"> <th style="width: 10%;"></th> <th style="width: 15%;">s₁</th> <th style="width: 15%;">s₂</th> <th style="width: 15%;">...</th> <th style="width: 15%;">s_{m-1}</th> <th style="width: 15%;">s_m</th> </tr> </thead> <tbody> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z'₁</td> <td style="text-align: center;">0</td> <td style="text-align: center;">1</td> <td></td> <td style="text-align: center;">0</td> <td style="text-align: center;">1</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z'₂</td> <td style="text-align: center;">1</td> <td style="text-align: center;">0</td> <td></td> <td style="text-align: center;">0</td> <td style="text-align: center;">0</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">...</td> <td></td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z'_{n-1}</td> <td style="text-align: center;">0</td> <td style="text-align: center;">0</td> <td></td> <td style="text-align: center;">1</td> <td style="text-align: center;">0</td> </tr> <tr> <td style="background-color: #4a86e8; color: white; text-align: center;">z'_n</td> <td style="text-align: center;">0</td> <td style="text-align: center;">0</td> <td></td> <td style="text-align: center;">1</td> <td style="text-align: center;">1</td> </tr> </tbody> </table>		s ₁	s ₂	...	s _{m-1}	s _m	z' ₁	0	1		0	1	z' ₂	1	0		0	0	...						z' _{n-1}	0	0		1	0	z' _n	0	0		1	1	Sentences
	s ₁	s ₂	...	s _{m-1}	s _m																																		
z' ₁	0	1		0	1																																		
z' ₂	1	0		0	0																																		
...																																							
z' _{n-1}	0	0		1	0																																		
z' _n	0	0		1	1																																		

Binarize Θ to compute per-topic precision and recall:

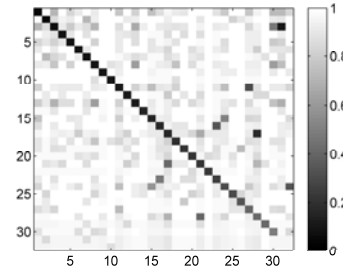
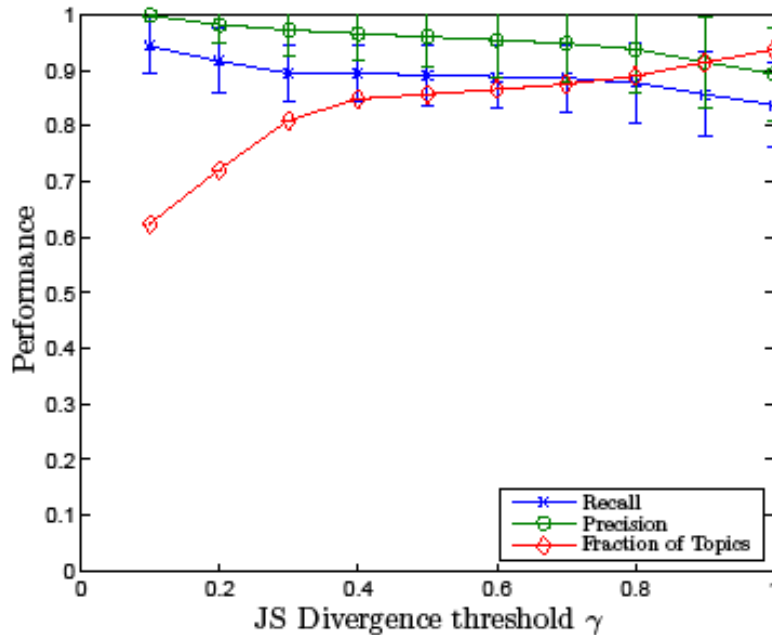
$$\Theta'_{ij} = \begin{cases} 1, & \Theta_{ij} \geq \varepsilon \\ 0, & \Theta_{ij} < \varepsilon \end{cases}$$



**F1, averaged across annotators
and using all topic-content unit
matches**

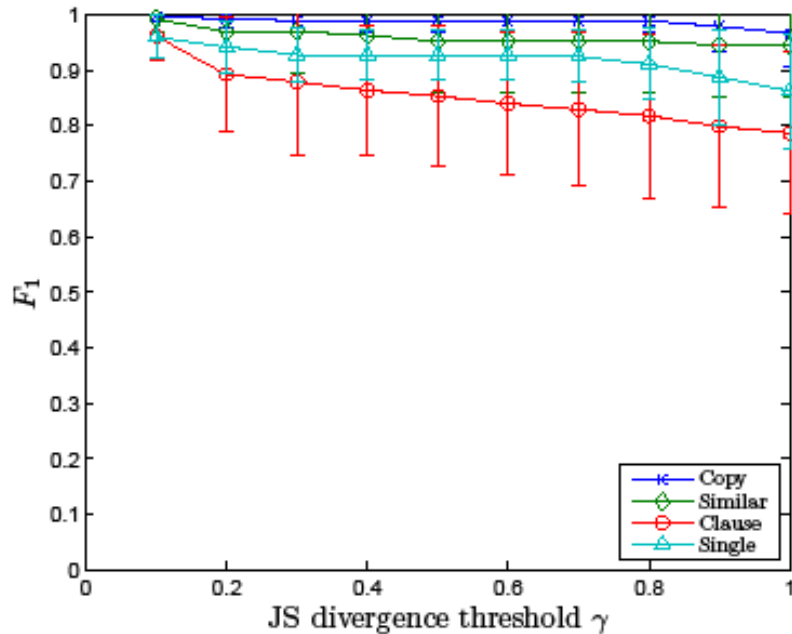
Mean F1 = 0.88

- Suggests that LDA identifies the same "associations" of topics and sentences as human annotators



Precision and recall, averaged over all document pairs and annotators

- Many matches have low JS divergence (approx. 80% ≤ 0.3)
- Precision and recall decrease as matches of lower quality are added



F1, averaged over all document pairs and annotators, per content unit type

Copies and **Similar** are easiest to identify

Clause is hardest due to word "noise"

Overall F1 quite high

LDA learns topics similar to manually annotated content units

- unsupervised
- domain- and language-independent
- benefits of probabilistic models

Initial results seem promising

- phrasal / sentence meaning similarity remains open (e.g. influence of word order, negation, ...)

Representation of sentences in terms of content units is a step towards bridging the gap between word and sentence importance modeling in multi-document summarization?

Thank you.

Questions?