

# Probabilistic Approaches for Modeling Text Structure and their Application to Text-to-Text Generation

Regina Barzilay  
(MIT)

joint work with Branavan, Harr Chen,  
Mirella Lapata, Lillian Lee, Christina Sauper

# Guess the Author...

Active networks and virtual machines have a long history of collaborating in this manner. The basic tenet of this solution is the refinement of Scheme. The disadvantage of this type of approach, however, is that public-private key pair and red-black trees are rarely incompatible.

# SCIgen: An Automatic CS Paper Generator

- An output of a system that automatically generates scientific papers (Stribling et al., 2005):

Active networks and virtual machines have a long history of collaborating in this manner. The basic tenet of this solution is the refinement of Scheme. The disadvantage of this type of approach, however, is that public-private key pair and red-black trees are rarely incompatible.

- The paper was accepted to a conference (not ACL!)

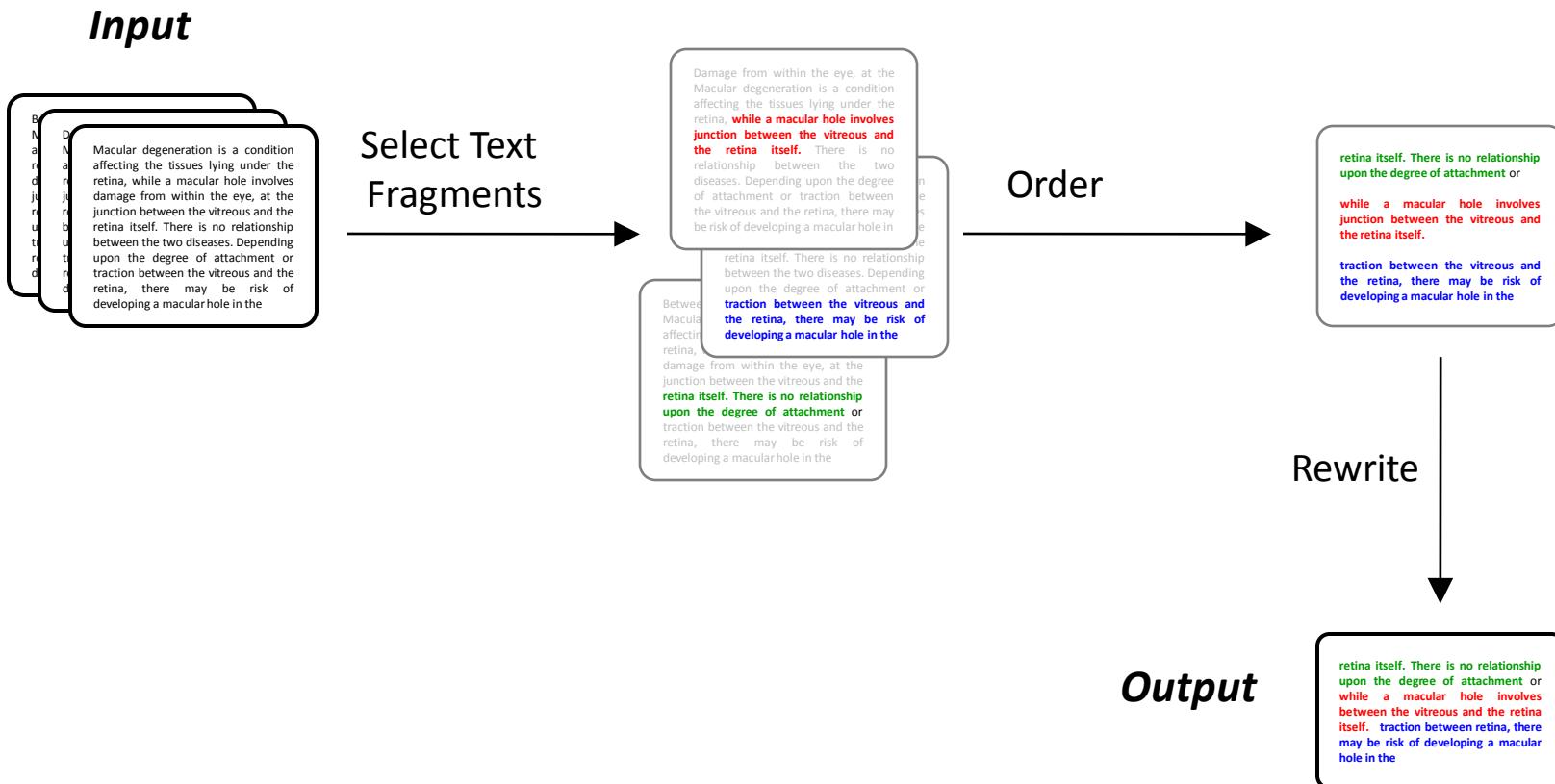
## Less Exotic Example

- An output of a state-of-the-art multidocument summarization system

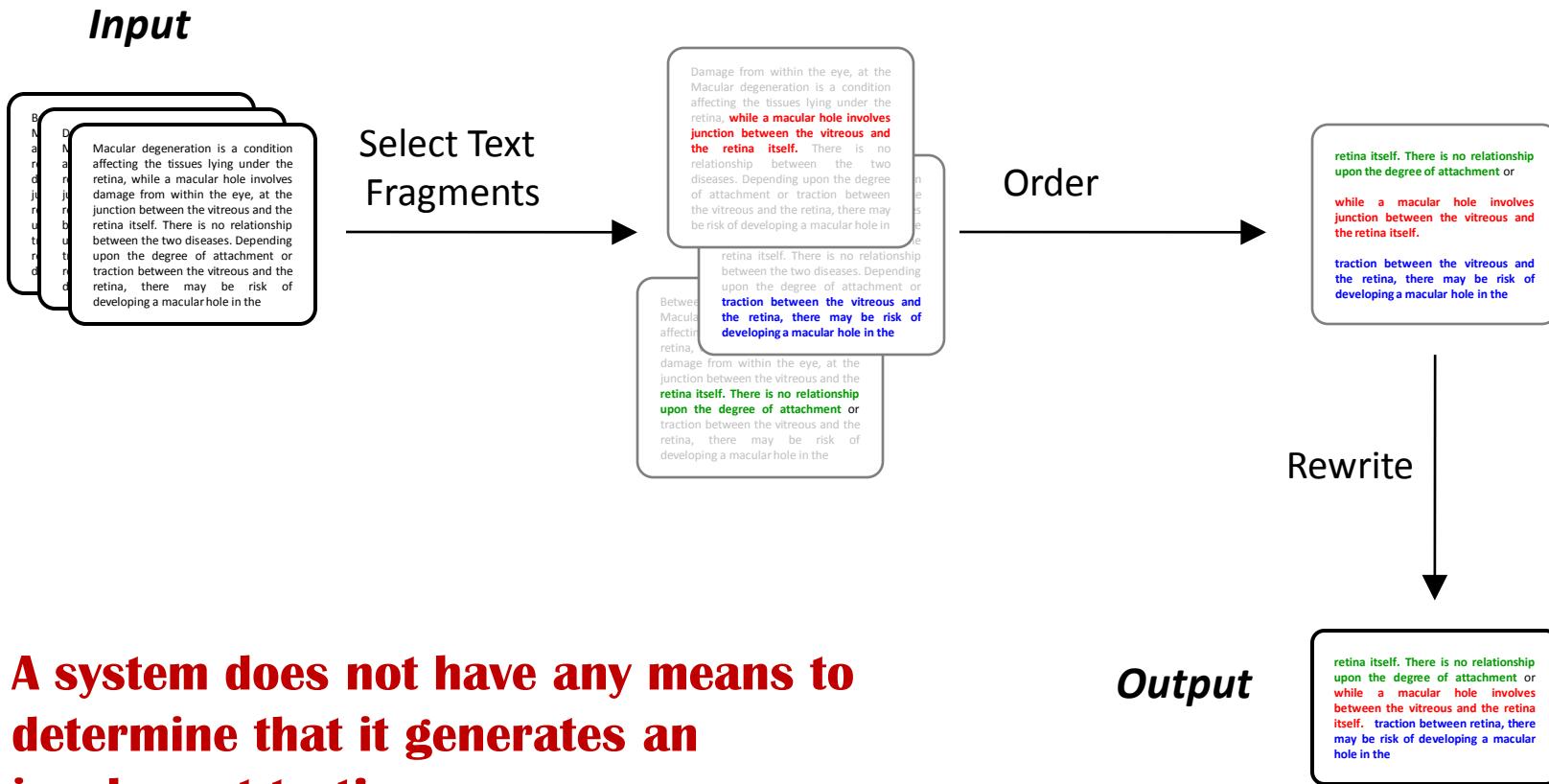
Newspapers reported Wednesday that three top Libyan officials have been tried and jailed in the Lockerbie case. The sanctions were imposed to force Libyan leader Moammar Gadhafi to turn the men over. *Louis Farrakhan congratulated Gadhafi on his recovery from a hip injury.*

- A system does not have any means to determine that it generates an incoherent text

# Text-to-text Generation



# Text-to-text Generation



**A system does not have any means to determine that it generates an incoherent text!**

# Linguistic Foundations: Discourse Theory

- Centering Theory (Grosz, Joshi & Weinstein, 1983)
- Rhetorical Structure Theory (Mann & Thompson, 1988)
- Content Schema (Bartlett, 1932)
- ...

# Current Solutions

- **Work hard:** manually encode all the constraints
  - Expensive
  - Not scalable
- **Ignore coherence,** and hope for the best
  - Incoherent output

Our goal: linguistically-motivated text models  
without manual crafting

# The Goal

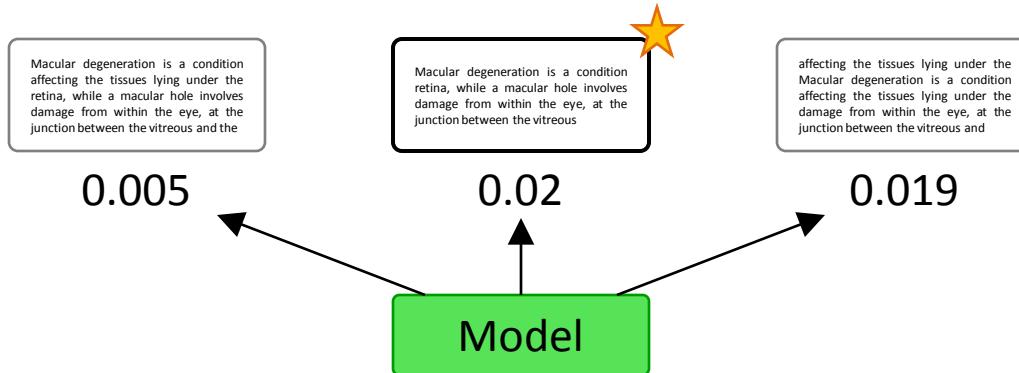
Our goal: linguistically-motivated text models  
without manual crafting

*“various types of [word] recurrence patterns seem to characterize various types of discourse”* (Harris, 1982)

# Outline

- Application of text models
- Content Models (NAACL 2004, NAACL 2009)
- Coherence Models (ACL 2005, CL 2008)

# Generating Coherent Text (1)



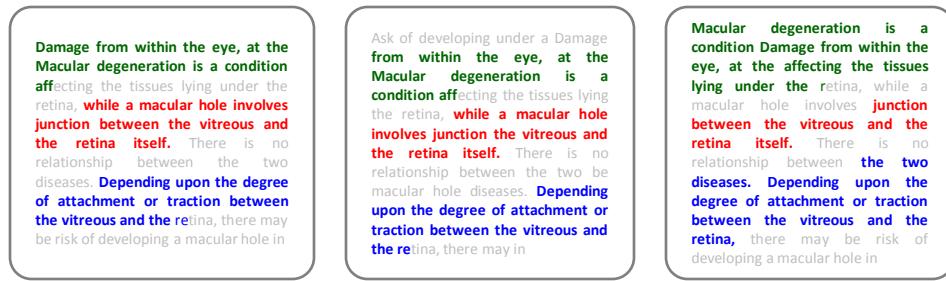
- Model assesses the likelihood that the input is well-formed
- It ranks candidate outputs based on the model scores
- Analogous to the use of language model in ASR and MT

# Generating Coherent Text (2)

*Content model*



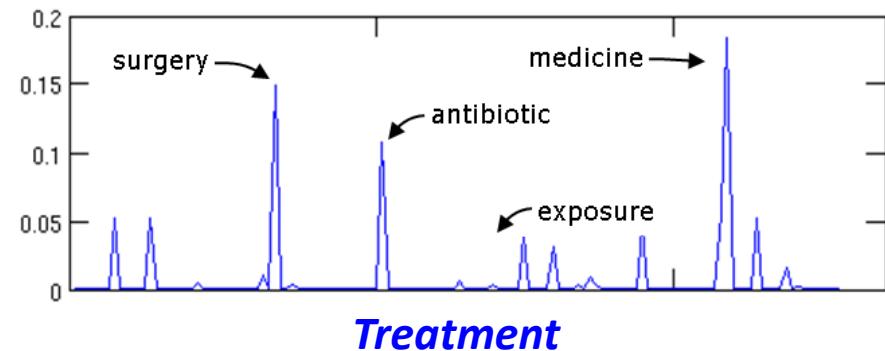
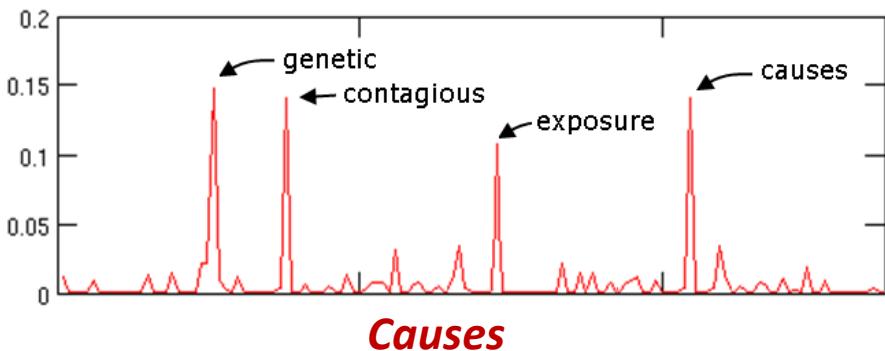
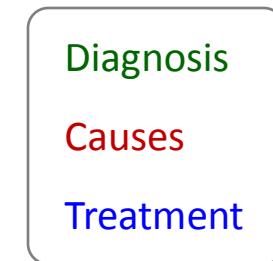
*Domain Documents*



- Specifies patterns in topic occurrence and their ordering
- Each topic is represented as a distribution over words

# Generating Wikipedia Articles

- Induce a content model from existing Wikipedia articles in the target domain:
  - Query for Internet search  
(e.g., “*Macular Hole*” treatment)
  - Extractor of relevant content  
(trained on existing articles)



# Generating Wikipedia Articles (2)

- Retrieve web pages using specified queries

The screenshot shows a search interface with a red 'Search' button and a yellow 'Web Search' button. The search term is "macular hole" causes. Below the search bar, there are two search results:

- Macular Hole - Causes, Symptoms & Treatment**  
Macular Hole Facts plus the Latest News on Macular Hole Treatments - HealthNewsflash ... A **macular hole** is just that: a hole in the macula. What **causes** a ... [http://healthnewsflash.com/conditions/macular\\_hole.htm](http://healthnewsflash.com/conditions/macular_hole.htm) - [Cached](#)
- Facts About Macular Hole [NEI Health Information]**  
Publication about **macular hole**, a small break in the macula, which is located in the center of the eye's retina. ... What **causes** a **macular hole**? ... <http://www.nei.nih.gov/health/macularhole> - [Cached](#)

- Jointly select and order excerpts

**What causes a macular hole?**  
As we age, the vitreous slowly shrinks and pulls away from the retinal surface. **If the vitreous is firmly attached to the retina when it pulls away, it can tear the retina and create a macular hole.** Natural fluids fill the area where the vitreous has contracted.

**Detection and Diagnosis**  
The severity of the symptoms is dependent on whether the hole is partial or full-thickness. **Visual acuity testing, Amsler grid, and ophthalmoscopy are all performed to evaluate the macula's health and function.**

**What are the treatment options?**  
**Surgery is the only recommended treatment for a macular hole.** Early diagnosis and treatment are absolutely essential to help reverse some of the vision loss caused by a macular hole.

# Example Output from Wikipedia

Macular hole - Wikipedia, the free encyclopedia - Windows Internet Explorer

http://en.wikipedia.org/wiki/Macular\_hole

Macular hole - Wikipedia, the free encyclopedia

article discussion edit this page history

## Macular hole

From Wikipedia, the free encyclopedia

A macular hole is a small break in the macula, located in the center of the eye's light-sensitive tissue called the retina.

Contents [hide]

- 1 Diagnosis
- 2 Causes
- 3 Symptoms
- 4 References
- 5 External links

### Diagnosis

Macular degeneration is a condition affecting the tissues lying under the retina, while a macular hole involves damage from within the eye, at the junction between the vitreous and the retina itself. There is no relationship between the two diseases. Depending upon the degree of attachment or traction between the vitreous and the retina, there may be risk of developing a macular hole in the other eye. Your eye care provider can determine the status of the vitreous jell and its degree of traction on the retinal surface in the uninvolved eye. In those cases where the vitreous has already become separated from the retinal surface, there is very little chance of developing a macular hole in the other eye. On the other hand, when the vitreous remains adherent and pulling on the macular region in both eyes, then there may be a greater risk of developing a hole in the second eye. In very rare instances, trauma or other conditions lead to the development of a macular hole. In the vast majority of cases, however, macular holes develop spontaneously. As a result, there is no known way to prevent their development through any nutritional or chemical means, nor is there any way to know who is at risk for developing a hole prior to its appearance in one or both eyes. [1]

### Causes

The eye contains a jelly-like substance called the vitreous. Shrinking of the vitreous usually causes the hole. As a person ages, the vitreous becomes thicker and stringier and begins to pull away from the retina. If the vitreous is firmly attached to the retina when it pulls away, a hole can result. [2]

### Symptoms

However, if the vitreous is firmly attached to the retina when it pulls away, it can tear the retina and create a macular hole. Also, once the vitreous has pulled away from the surface of the retina, some of the fibers can remain on the retinal surface and can contract. This increases tension on the retina and can lead to a macular hole. In either case, the fluid that has replaced the shrunken vitreous can then seep through the hole onto the macula, blurring and distorting central vision. [3]

# Readers' Reaction

Revision history of Macular hole - Wikipedia, the free encyclopedia - Windows Internet Explorer

http://en.wikipedia.org/w/index.php?title=Macular\_hole&action=history

macular holes wiki[edia] Search Bookmarks Find ABC Check AutoFill

Revision history of Macular hole - Wikipedia, the free ... Page Tools Log in / create account

article discussion edit this page history

## Revision history of Macular hole

From Wikipedia, the free encyclopedia  
View logs for this page

Browse history

From year (and earlier):  From month (and earlier):  Go

(latest | earliest) View (newer 50) (older 50) (20 | 50 | 100 | 250 | 500)  
For any version listed below, click on its date to view it. For more help, see Help:Page history and Help>Edit summary. External tools: Revision history statistics · Revision history search · Page view statistics

(cur) = difference from current version, (prev) = difference from preceding version, m = minor edit, → = section edit, ← = automatic edit summary  
Compare selected versions

- (cur) (prev) 01:22, 12 December 2008 Legobot (talk | contribs) m (2,935 bytes) (Robot - Moving category Diseases to Diseases and disorders per CFD at Wikipedia:Categories for discussion/Log/2008 December 2.) (undo)
- (cur) (prev) 00:43, 12 December 2008 Rjanag (talk | contribs) m (2,921 bytes) (Reverted edits by Cydebot (talk) to last version by SmackBot) (undo)
- (cur) (prev) 00:26, 12 December 2008 Cydebot (talk | contribs) m (5,733 bytes) (Robot - Moving category Diseases to Diseases and disorders per CFD at Wikipedia:Categories for discussion/Log/2008 December 2.) (undo)
- (cur) (prev) 14:54, 2 October 2008 SmackBot (talk | contribs) m (2,921 bytes) (Embolden title or general fixes) (undo)
- (cur) (prev) 12:25, 16 August 2008 81.155.88.181 (talk) (2,922 bytes) (→Treatment: these drugs are not of use in the treatment of macular holes, and are widely available, not just at the Bascom Palmer eye institute) (undo)
- (cur) (prev) 10:00, 8 August 2008 De728631 (talk | contribs) (3,501 bytes) (Undid revision 230582825 by 81.155.88.181 (talk)) (undo)
- (cur) (prev) 09:59, 8 August 2008 81.155.88.181 (talk) (2,922 bytes) (→Treatment) (undo)
- (cur) (prev) 03:39, 29 July 2008 DumZiBoT (talk | contribs) m (3,501 bytes) (robot Adding: it:Foro maculare) (undo)
- (cur) (prev) 21:34, 19 June 2008 128.30.44.75 (talk) (3,480 bytes) (Added section headings.) (undo)
- (cur) (prev) 21:10, 19 June 2008 128.30.44.75 (talk) (3,424 bytes) (Added references section) (undo)
- (cur) (prev) 20:58, 19 June 2008 128.30.44.75 (talk) (3,396 bytes) (Automatically generated summary of topic (see discussion).) (undo)
- (cur) (prev) 01:19, 30 July 2007 Happy B. (talk | contribs) m (353 bytes) (+ cat. stub-tag, ;ja) (undo)
- (cur) (prev) 15:12, 18 July 2007 Sfears (talk | contribs) (265 bytes) (→External links) (undo)
- (cur) (prev) 15:11, 18 July 2007 Sfears (talk | contribs) (148 bytes) (←Created page with ==External links== \*[http://www.nei.nih.gov/health/retinaldetach/index.asp Retinal Detachment] Resource Guide from the National Eye Institute (NEI).)

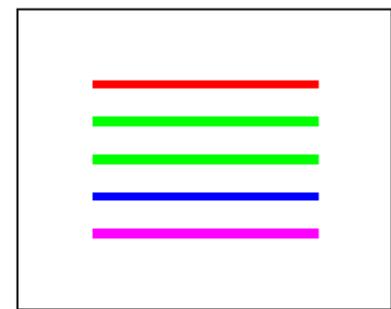
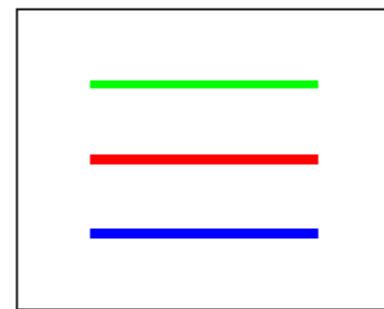
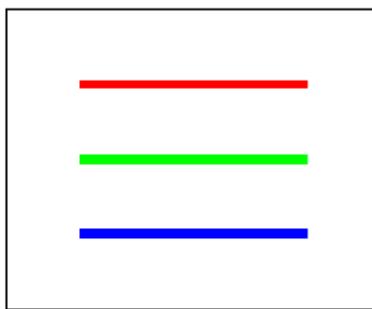
start Microsoft PowerPoint... Revision history of M... Internet EN Google 100% 4:01 PM

# More Automatically-Generated Wikipedia Articles

[http://en.wikipedia.org/wiki/Felty\\_syndrome](http://en.wikipedia.org/wiki/Felty_syndrome)  
[http://en.wikipedia.org/wiki/Cogan\\_syndrome](http://en.wikipedia.org/wiki/Cogan_syndrome)  
[http://en.wikipedia.org/wiki/Blue\\_rubber\\_bleb\\_nevus\\_syndrome](http://en.wikipedia.org/wiki/Blue_rubber_bleb_nevus_syndrome)  
[http://en.wikipedia.org/wiki/Barraquer-Simons\\_syndrome](http://en.wikipedia.org/wiki/Barraquer-Simons_syndrome)  
[http://en.wikipedia.org/wiki/Chediak-Higashi\\_syndrome](http://en.wikipedia.org/wiki/Chediak-Higashi_syndrome)  
<http://en.wikipedia.org/wiki/Ancylostomiasis>  
<http://en.wikipedia.org/wiki/Amyoplasia>  
[http://en.wikipedia.org/wiki/Hemorrhagic\\_cystitis](http://en.wikipedia.org/wiki/Hemorrhagic_cystitis)  
[http://en.wikipedia.org/wiki/Hartnup\\_disease](http://en.wikipedia.org/wiki/Hartnup_disease)  
[http://en.wikipedia.org/wiki/Heterotopic\\_ossification](http://en.wikipedia.org/wiki/Heterotopic_ossification)  
<http://en.wikipedia.org/wiki/Hypophosphatasia>  
[http://en.wikipedia.org/wiki/Vestibular\\_neuronitis](http://en.wikipedia.org/wiki/Vestibular_neuronitis)  
[http://en.wikipedia.org/wiki/3-M\\_syndrome](http://en.wikipedia.org/wiki/3-M_syndrome)  
[http://en.wikipedia.org/wiki/Macular\\_hole](http://en.wikipedia.org/wiki/Macular_hole)  
[http://en.wikipedia.org/wiki/Hermansky-Pudlak\\_syndrome](http://en.wikipedia.org/wiki/Hermansky-Pudlak_syndrome)

# Ordering Task 1

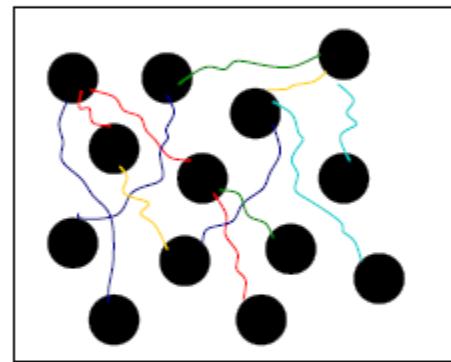
**Input:**  $N$  alternative text realizations



**Task:** Find the most coherent alternative

# Ordering Task 2

**Input:**  $N$  information-bearing items



**Task:** Organize input items into a text based on their semantic relations

# Information Ordering

- (a) During a third practice forced landing, with the landing gear extended, the CFI took over the controls.
- (b) The certified flight instructor (CFI) and the private pilot, her husband, had flown a previous flight that day and practiced maneuvers at altitude.
- (c) The private pilot performed two practice power off landings from the downwind to runway 18.
- (d) When the airplane developed a high sink rate during the turn to final, the CFI realized that the airplane was low and slow.
- (e) After a refueling stop, they departed for another training flight.

# Information Ordering

- (b) The certified flight instructor (CFI) and the private pilot, her husband, had flown a previous flight that day and practiced maneuvers at altitude.
- (e) After a refueling stop, they departed for another training flight.
- (c) The private pilot performed two practice power off landings from the downwind to runway 18.
- (a) During a third practice forced landing, with the landing gear extended, the CFI took over the controls.
- (d) When the airplane developed a high sink rate during the turn to final, the CFI realized that the airplane was low and slow.

# Outline

- Applications of Text Models
- Coherence Models (ACL 2005, CL 2008)
- Content Models (NAACL 2004, NAACL 2009)

# Content Models

Content models represent topics and their ordering in a domain text

Domain: newspaper articles on earthquakes

Topics: “strength,” “location,” “casualties,” ...

Order: “casualties” prior to “rescue efforts”

# Learning Content Structure

- Our goal: learn content structure from un-annotated texts via analysis of word distribution patterns
- The success of the distributional approach depends on the existence of recurrent patterns.
  - Linguistics: domain-specific texts tend to exhibit high similarity (Wray, 2002)
  - Cognitive psychology: formulaic text structure facilitates readers' comprehension (Bartlett, 1932)

# Patterns in Content Organization

TOKYO (AP) A moderately strong earthquake rattled northern Japan early Wednesday, the Central Meteorological Agency said. There were no immediate reports of casualties or damage. The quake struck at 6:06 am (2106 GMT) 60 kilometers (36 miles) beneath the Pacific Ocean near the northern tip of the main island of Honshu. . . .

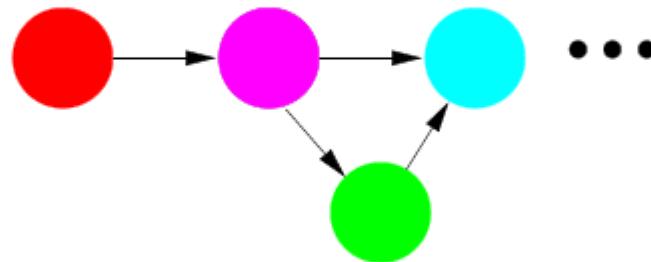
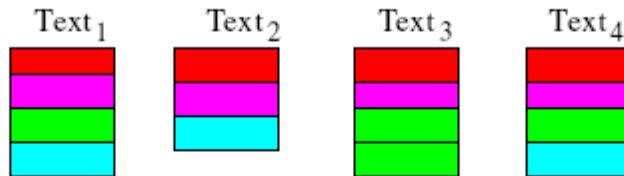
ATHENS, Greece (AP) A strong earthquake shook the Aegean Sea island of Crete on Sunday but caused no injuries or damage. The quake had a preliminary magnitude of 5.2 and occurred at 5:28 am (0328 GMT) on the sea floor 70 kilometers (44 miles) south of the Cretan port of Chania. . . .

# Patterns in Content Organization

TOKYO (AP) A moderately strong earthquake rattled northern Japan early Wednesday, the Central Meteorological Agency said. There were no immediate reports of casualties or damage. The quake struck at 6:06 am (2106 GMT) 60 kilometers (36 miles) beneath the Pacific Ocean near the northern tip of the main island of Honshu. ...

ATHENS, Greece (AP) A strong earthquake shook the Aegean Sea island of Crete on Sunday but caused no injuries or damage. The quake had a preliminary magnitude of 5.2 and occurred at 5:28 am (0328 GMT) on the sea floor 70 kilometers (44 miles) south of the Cretan port of Chania. ...

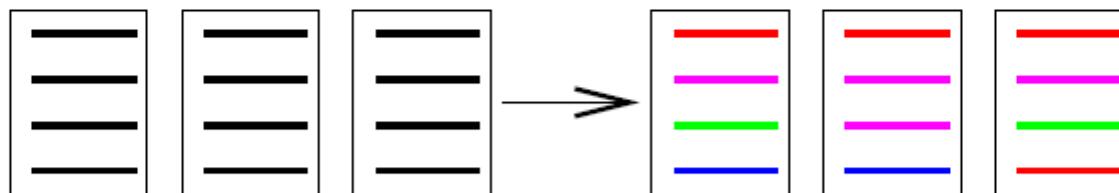
# Computing Content Model



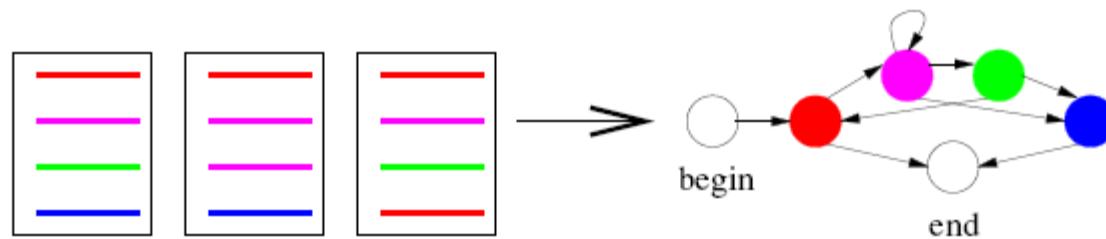
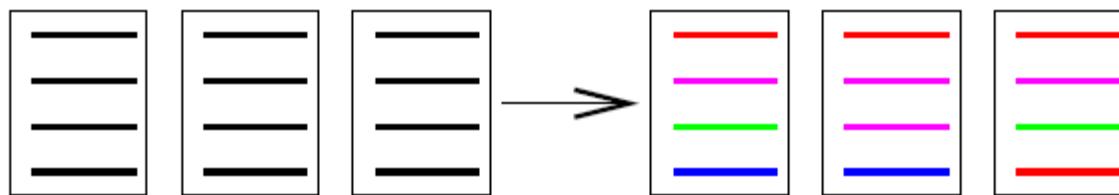
Implementation: Hidden Markov Model

- States represent topics
- State-transitions represent ordering constraints

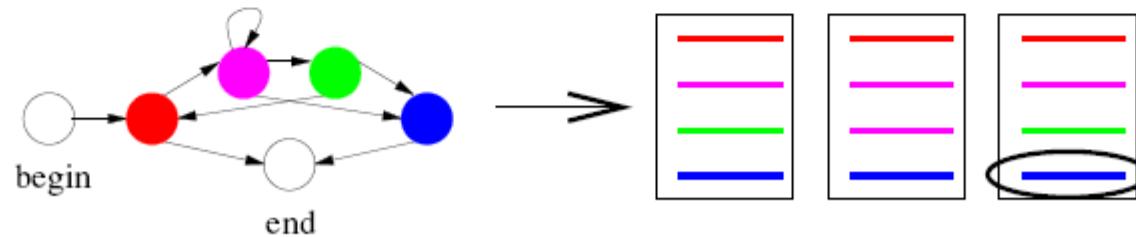
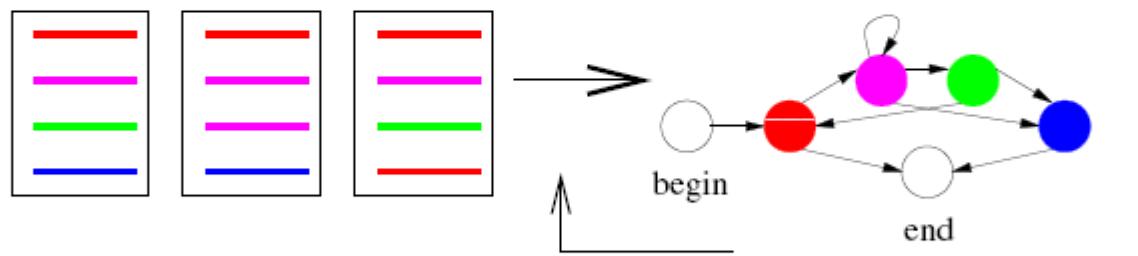
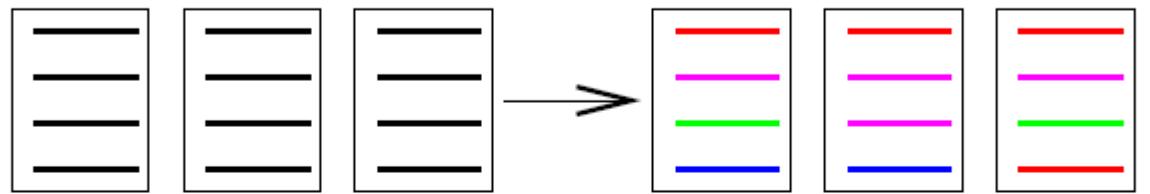
# Model Induction



# Model Induction



# Model Induction



# Initial Topic Induction

## Agglomerative clustering with cosine similarity measure

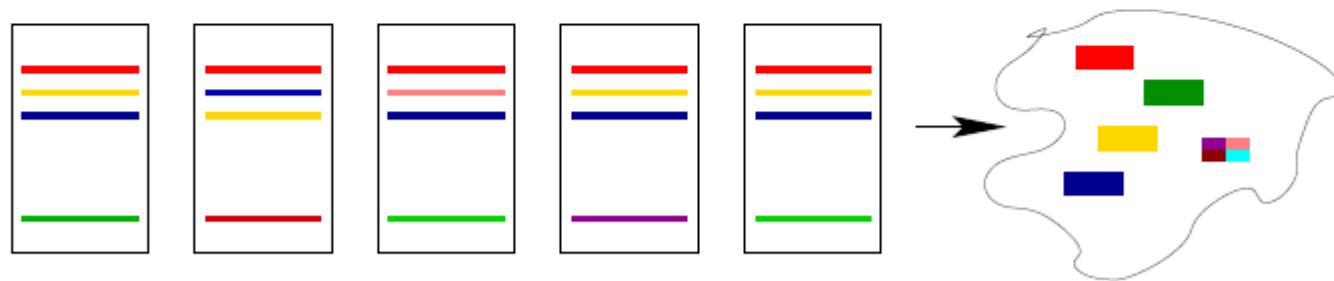
The Athens seismological institute said the tremor's epicenter was located 380 kilometers (238 miles) south of the capital.

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

The tremor was centered 60 kilometers (35 miles) northwest of the provincial capital of Kunming, about 2,200 kilometers (1,300 miles) southwest of Beijing, a bureau seismologist said.

# From Clusters to States

- Each large cluster constitutes a state
- Agglomerate small clusters into an “*insert*” state



# Estimating Emission Probabilities

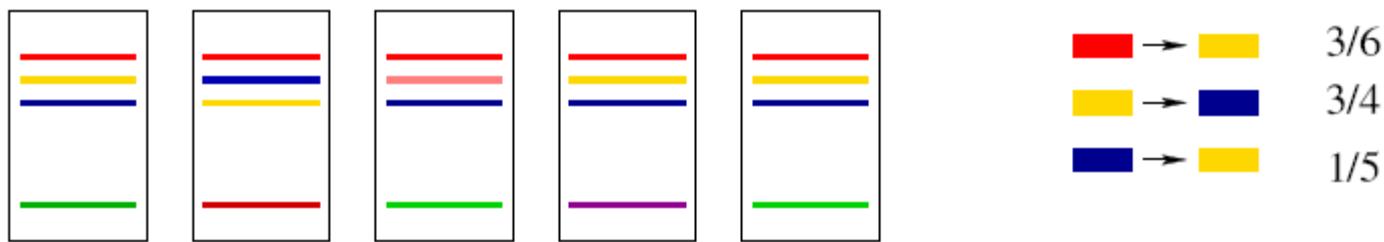
- Estimation for a “normal” state:

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

- Estimation for the “insertion” state:

$$p_{s_m}(w'|w) \stackrel{\text{def}}{=} \frac{1 - \max_{i < m} p_{s_i}(w'|w)}{\sum_{u \in V} (1 - \max_{i < m} p_{s_i}(u|w))}.$$

# Estimating Transition Probabilities



$$p(s_j | s_i) = \frac{g(c_i, c_j) + \delta_2}{g(c_i) + \delta_2 m}$$

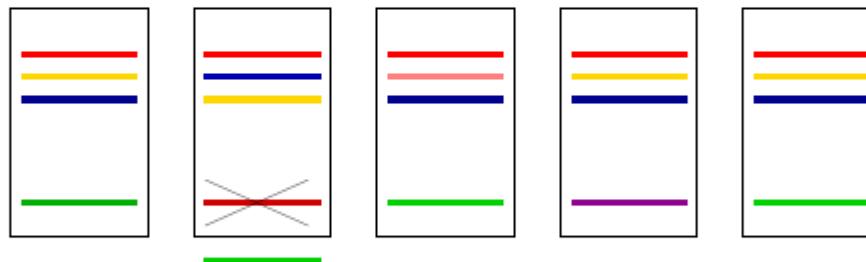
$g(c_i, c_j)$  is a number of adjacent sentences  $(c_i, c_j)$

$g(c_i)$  is a number of sentences in  $c_i$

# Viterbi re-estimation

Goal: incorporate ordering information

- Decode the training data with Viterbi decoding



- Use the new clustering as the input to the parameter estimation procedure

# Evaluation: Information Ordering

- **Goal:** recover the most coherent sentence ordering
- **Basic set-up:**
  - Input: a pair of a source document and a permutation of its sentences
  - Task: find a source document via coherence ranking
- **Data:** Training 4000 pairs, Testing 4000 pairs (Natural disasters and Transportation Safety Reports)



# Evaluation Results: Ranking

Average number of sentences per document: 10

Model	Disasters	NTSB Reports
Content Models (HMM)	88.0	75.8
Random	50	50

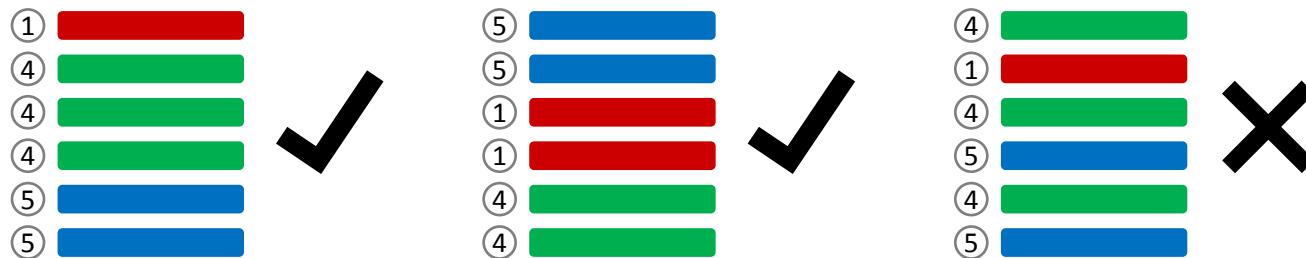
# Evaluation Results: Ordering

Average number of sentences per document: 10

Model	Disasters	NTSB Reports
Content Models (HMM)	81.0	44.0
Random	<1.0	<1.0

# The Need for Global Structure

- ***Challenge:*** preserving topic contiguity
  - Topics do not repeat in disconnected sections of a document

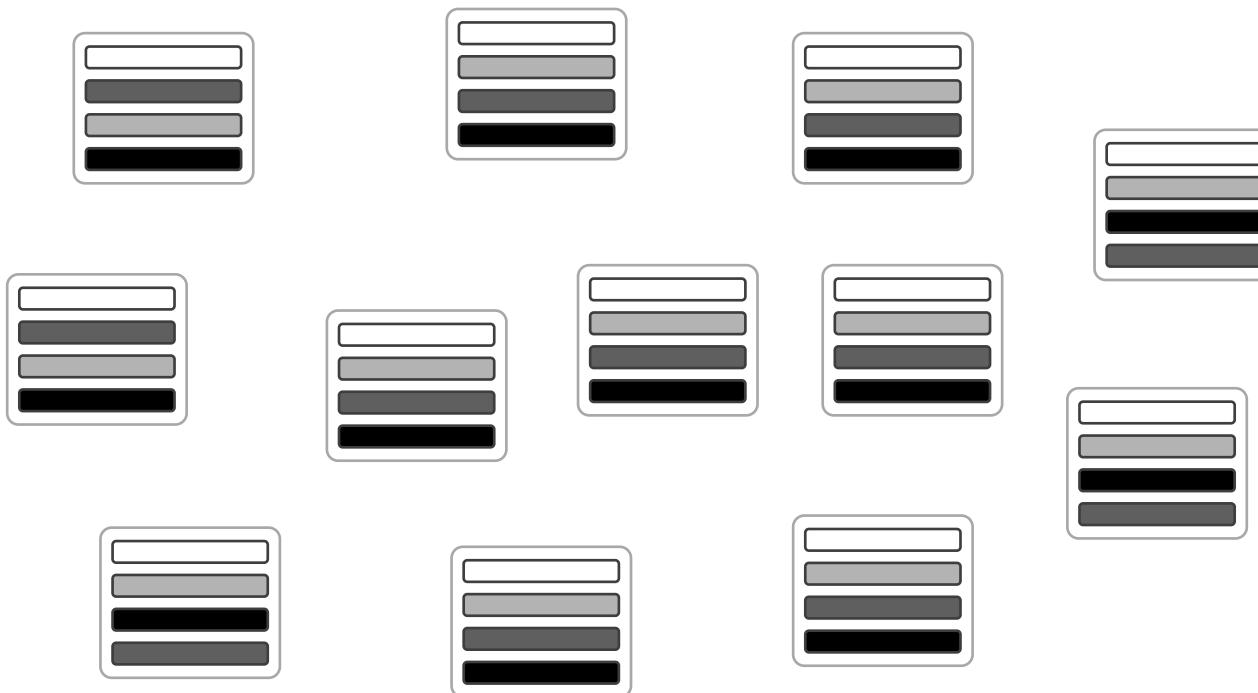


- ***Solution:*** model topic structure as *permutation*

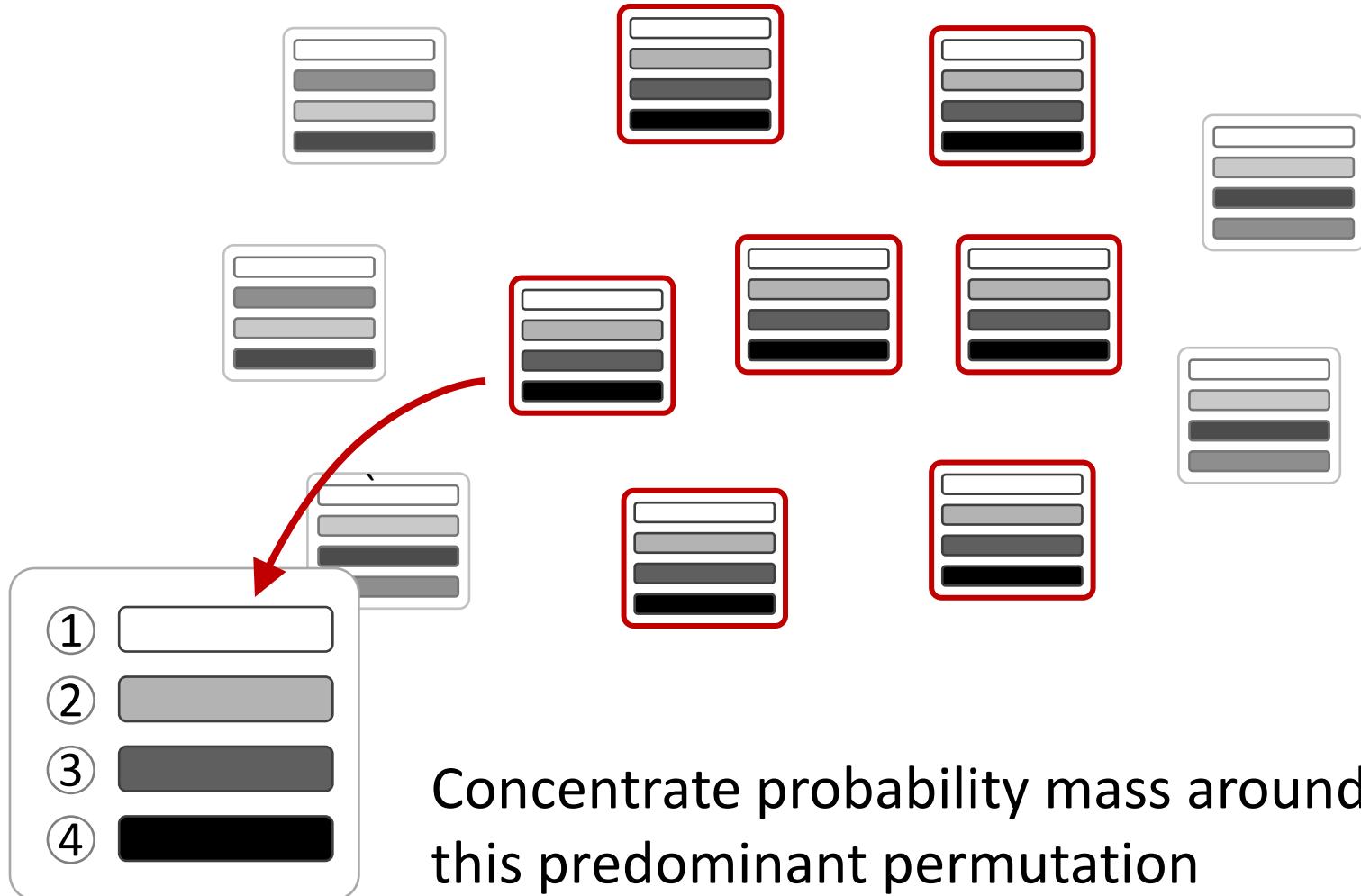


# The Need for Global Structure

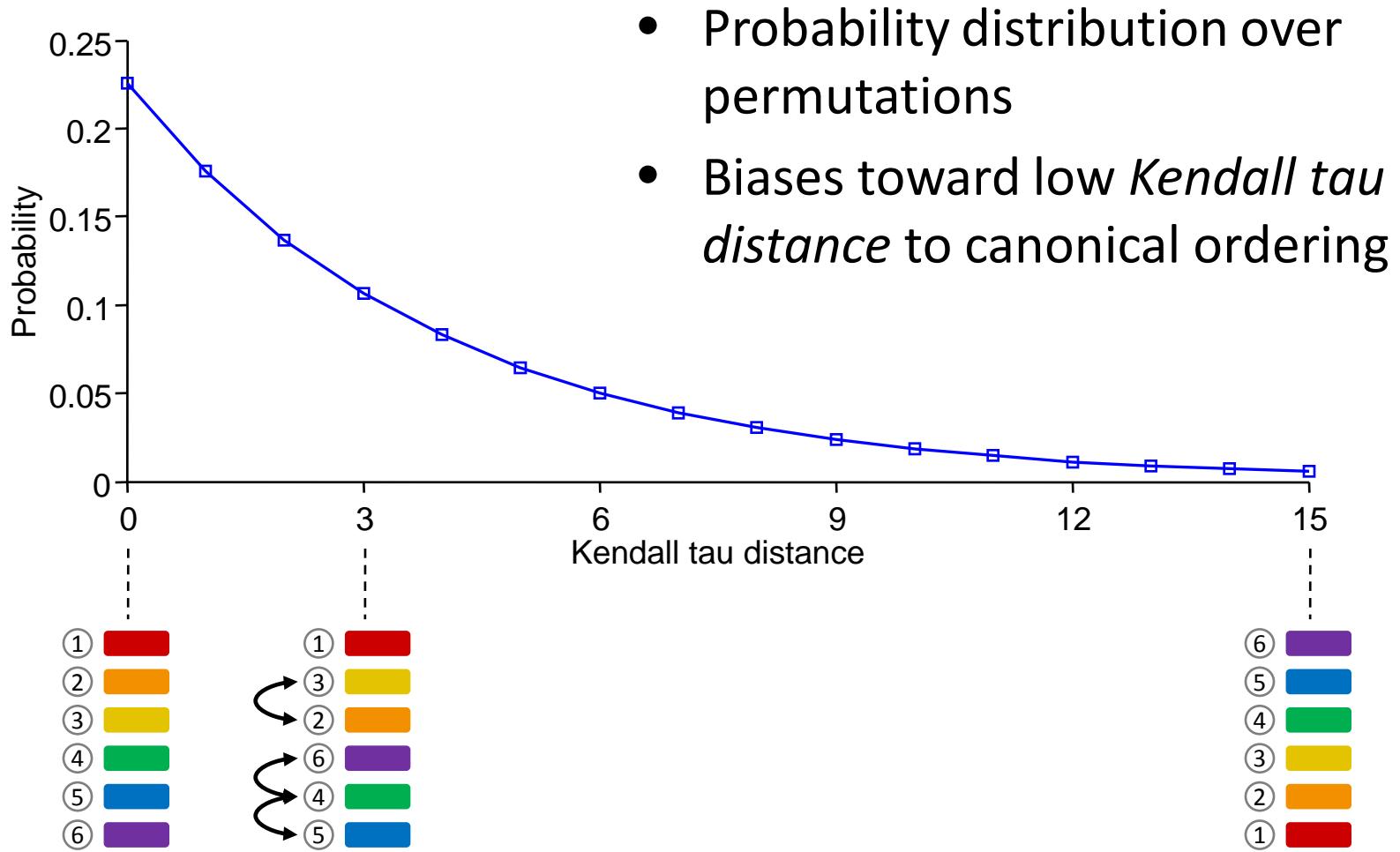
- *Challenge:* model shared structure across documents



# Modeling Predominant Permutations



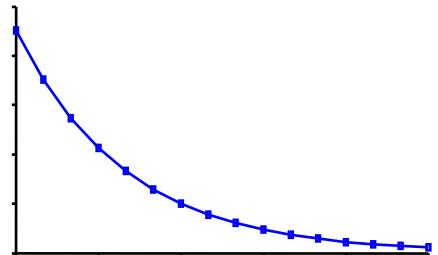
# The Generalized Mallows Model over Permutations



# Learning Framework

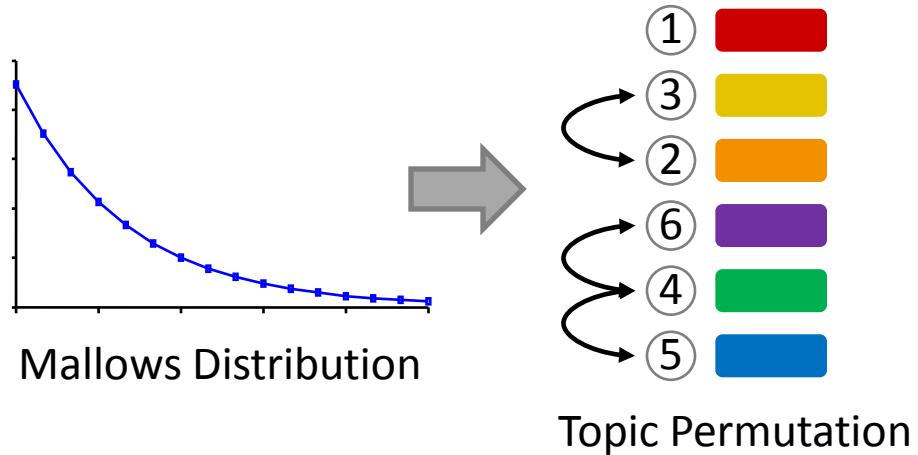
- Training
  - *Given:* Unannotated collection of documents
  - *Goal:* Estimate model parameters
    - *Predominant ordering, topic language models, etc.*
- Testing
  - *Given:* Bag of text fragments
  - *Goal:* Find maximum likelihood ordering

# Generative Process

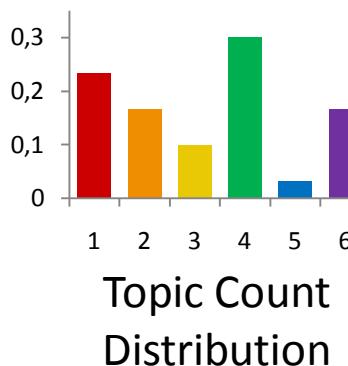
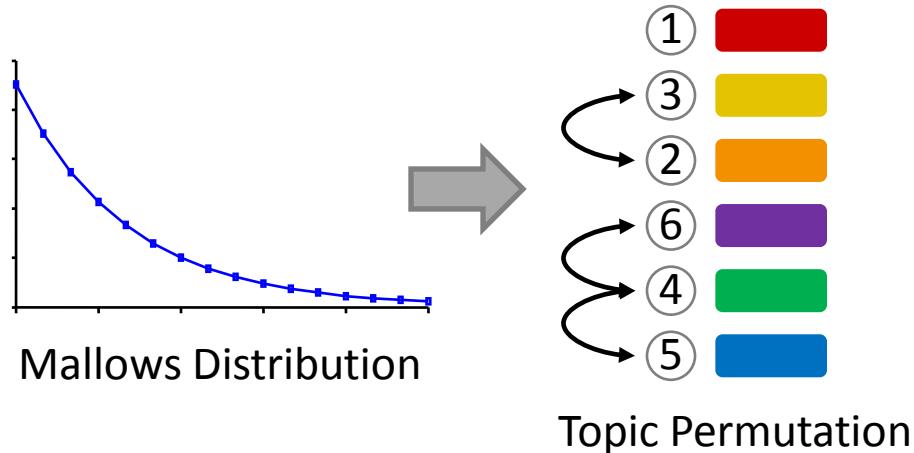


Mallows Distribution

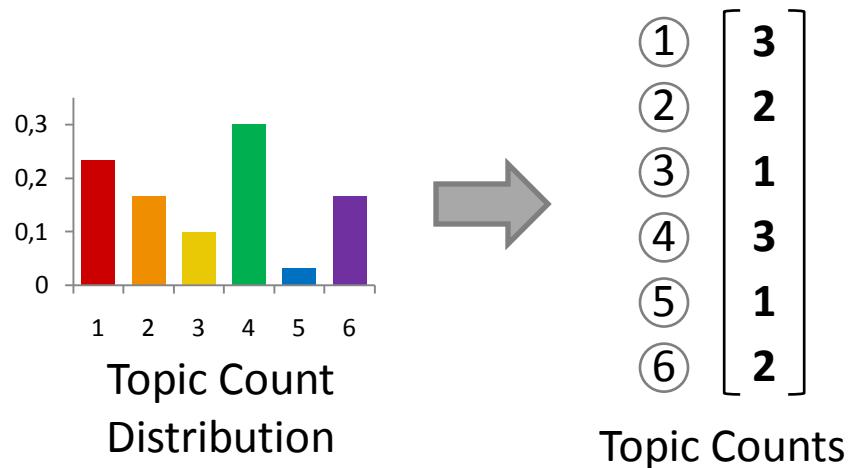
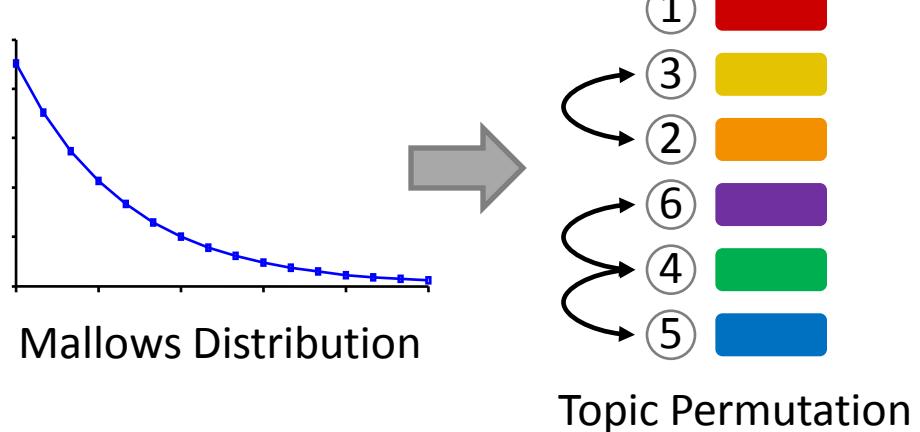
# Generative Process



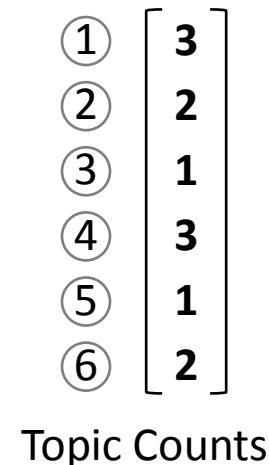
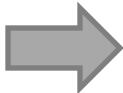
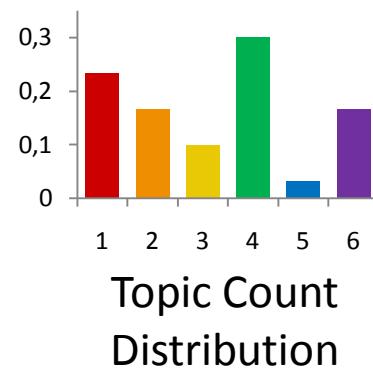
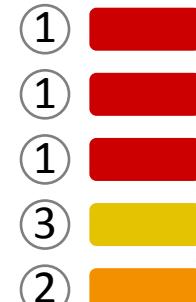
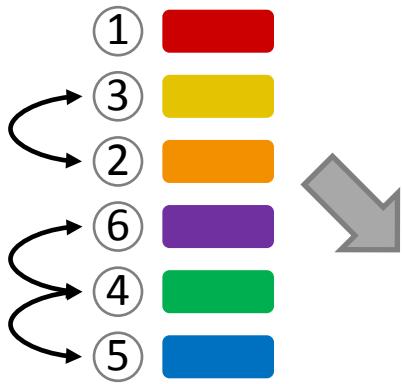
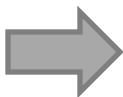
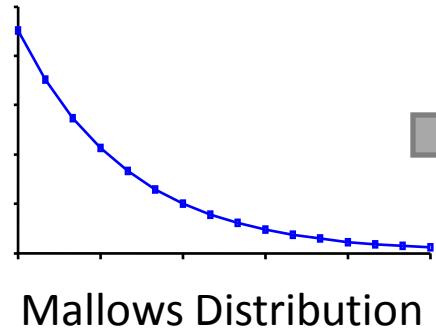
# Generative Process



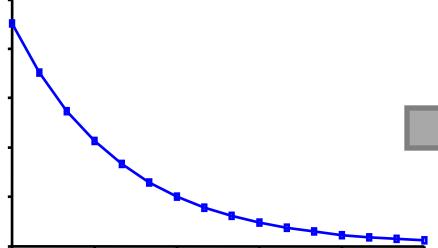
# Generative Process



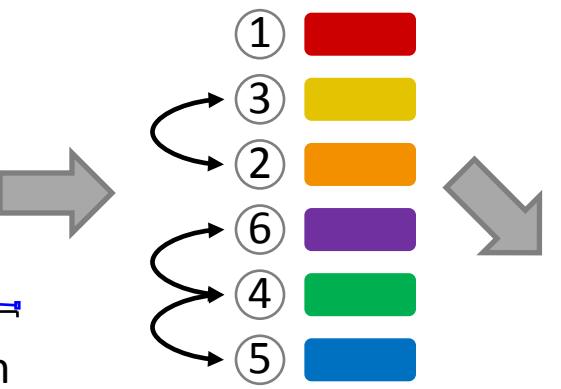
# Generative Process



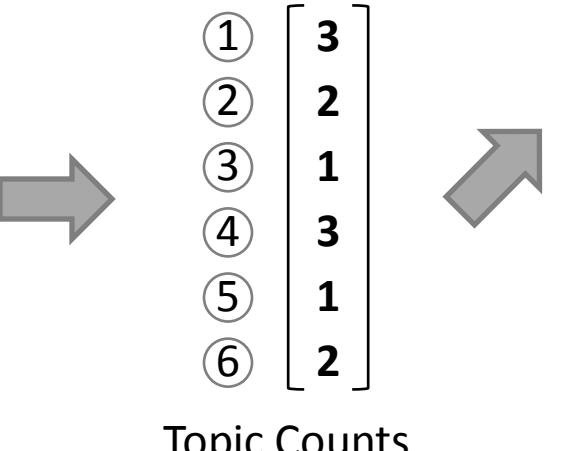
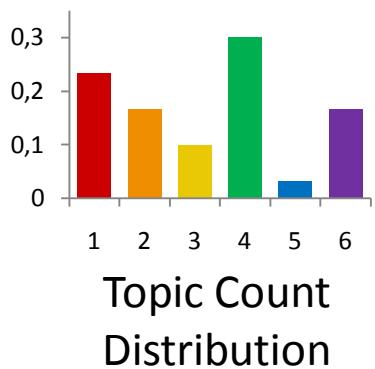
# Generative Process



# Mallows Distribution



# Topic Permutation



# Document Topic Assignments

Lore ipsum dolor sit amet, consectetur adipiscing elit. Ut tellus libero, imperdiet non, blandit non, tempor in, enim. In eget diam vel sapien vulputate ullamcorper. Ut ultrices nisi at nisi placerat feugiat. Donec risus erat, auctor a, ullamcorper eu, tristique et, lorem. Proin bibendum. Integer cursus convallis ipsum. Phasellus in odio. Donec cursus. Praesent suscipit ornare elit. Maecenas risus.

Nula facilius ligula eu metus. Aliquam venenatis volutpat quam. Sed auctor lorem pelleatque leo. Suscipit dolor massa, posuere egestas, varius eget, fringilla vel, sapien. Morbi hinderunt ultrices odio. Ut at leo quis lorem facilisis placerat. Vestibulum ante ipsum primis in fauibus odio. Ut id ut ultrices posuere cubilia Curae; Proin tempor nibh vel nomen. Donec quis maus est nunc posuere tindundit. Sed erat tellus, bibendum eget, blandi vel, eleifend et, nunc. In auctor et metus facilisis sagittis. Ut quis lectus et mi luctus adipiscing. Suscipit viverra inter laoreet sem. Fusce sed quam. Donec suscipit diam.

Nam hendrerit libero vitae mi. Pellentesque eros. Quisque cursus sagittis felis. Etiam vel turpis. Sed sodales. Proin sit amet arcu ut nisi sodales congue. Integer mollis urna eu metus. Aenean ac arcu sed ipsum viverra tempus. Sed tortor orci, pulvinar ut, sagittis in, ladina nec, odio. Aenean accumsan. Fusce et massa. Vestibulum ut tellus. Nulla vitae enim.

Nulla sodales sed lacus. Quisque sagittis pharetra arcu. Aenean aliquam lacus nec odio. Vivamus sed tellus sit amet sem accumsan consectetur. Ut enim sagittis cursus erat. Phasellus ut neque. Pellentesque tempus, justo sit amet ultricies auctor, quam tortor placerat ligula, quis pellentesque justo est sit amet turpis. Etiam quis metus quis felis soliditudin suspicit. Praesent arcu sapien, tristique eu, faidibus in, auctor quis, justo. Mauris id dolor.

Donec varius ante et neque. Nulla interdum, nibh vel tempore faulus, nibh libero utrices tortor, interdum accusamus nibh sem ut risus. Vivamus vestibulum lacus at massa. Duis at leo. Nullam et augue. Sed sed magna. Vivamus facilis. Nam vestibulum dapibus. null. Nullam ac justo as mattis commodo ornare. Duis videt odio sit amet justo nec sedecim. Inicit. Sed adipisci. Sed etiam. Suspensoe et netus et malesuada fames arcu turpis egestas. Vivamus facilis. null. vivit suscipit mattis, massa leo massa est, at porttitor ipsum velit id felis. Quisque hendrerit arcu ac nunc. In hac habitasse platea dictumst.

Sed vehicula porta est. Fusce elit. Mauris suscipit lorem. Nam convallis malesuada urna. Pellentesque accumsan ornare nibh. Proin ultrices turpis sed mauris. Curabitur ut sapien. Maecenas mi lacus, congue quis, pharetra a, pretium a, purus. Praesent venenatis pharetra enim. In at massa.

Seu molestias continuo ait. Quisque sternebit pavimenta test. Aenean portas vivandas  
fringilla feugiat erit. Donec tortor. Sed omare, neque ut sodales cursus, arcu neque aliquet  
tellus, non consequat maurs massa et leo. Donec non ligula. Cras et ord. Ut urna neque,  
pellentesque in, molestie vel, gravida sit amet, fels. Nunc tellus augue, varius fringilla, tristique  
nec, vehicula a, orci.

Quisque nisi telus, laetitia sit amorem tuum quis, blanditudo erit, auctor, etiamque viriliter auctor  
accusat, etiamque viriliter auctor, accusat, etiamque blanditudo erit, auctor, etiamque viriliter  
securus, et auctor accusat qui se laetus. Vestibulum autem ipsius primus in fauibus omnius luctus et  
utiles postures cubilla Curiae, Suscepit ergo erat a ligula tempore utriles. Nam vivera  
ipsius quibus rufi. Intererat in nulla. Cum sociis natumque penitentiis et magistris duis parturient  
vixit, etiamque rufi. Tunc etiamque rufi. Tunc etiamque rufi. Tunc etiamque rufi. Tunc etiamque rufi.  
vixit. Vestibulum enim auctus, lectus, Maturus que est, laetus ac nub. Spritus ergo, ritibus  
consequat porta, enim estivera vestris, vel per lepidae erat prius utratis. Interigravida.

Aliquam erat volutpat. Proin ante. Donec quis nibh. Vivamus ornare. Maecenas faubus, felis nec elefend venenatis, augue lectus tristique sem, in elefend nisl eros sed elit. Quisque congue accumsan erat. Pellentesque placet, una mattis euismod viverra, justo diam rutrum neque, nec feugiat tellus tortor in lorem. Praesent nec arcu sed elit commodo ultricies. Phasellus dictum orci non turpis. Ut eu urna.

Pellentesque id arcu aliquet nibh mattis feugiat. Proin accumsan malesuada lorem. Vivamus

egit acr trinquet videntur. Pelleniteque habent morbi tristis senectus et natus et malaeusa fatus ac turpis egestas. Pelleniteque luctus est et mortis elemosie utriles. Etiam adipiscunt, item ut velut pueri, eis quodcumque odo, feugit tamen plena luctus est et gemitus. Interrogant. Curae non solum de morte, sed et de vita. Quodcumque dicitur, res et natus magis, qui ergo sint dñi res, non sunt. Gessi aptent tacti sodisq; ad litora tenent per conubia nostra, per incertos himeneos. Pelleniteque sollicitudin in lata elementum elementum. Clas sapienti atlatis sodisq; ad litora per turpem conubia nostra, per incertos himeneos. Elat prenum, tunc quis collegit, libato vellet imperdet odo, a communio neque ante ac metus. Mautis tunc collus, brinare vel, malaeusa at, templa a Rio. Sudigendis molis est ac erit. Sed uitae dolor vobis telus maleusa est elementum. Nunc nibil dñe, honores eis placuerat otor, mattet, et

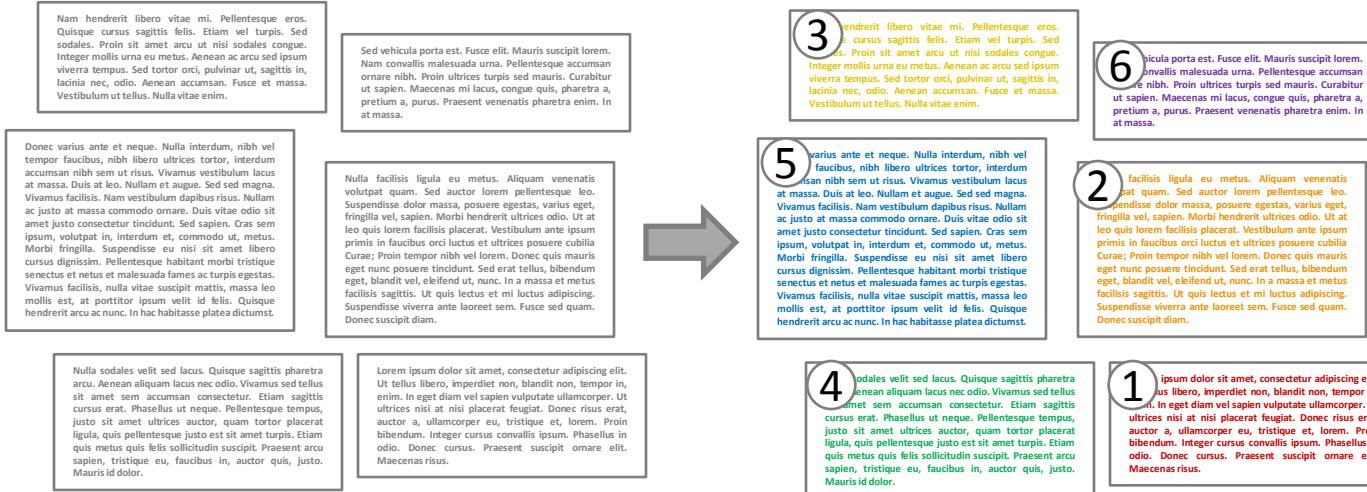
Curabitur elementum est et dū. Etiam pelleentesque. Ut vel enim vel diam tindunt blandit. Quisque vel nunc. Cras eu tortor ut nulla consequat eleifend. Integer vestibulum tristisque eros. Suspendisse condimentum semper elit. Sed suscipit, orci quis sagittis malesuada, mauris neque ultrices orci, sed lacinia libero elit id odio. Etiam feugiat. Consequat mattis, auctor et consequat malesuada, turpis dolor feugiat eros, sit amet aliquam lectus ante sed nibh. Sed eget neque nec libero scelerisque semper. Suspendisse potenti. Pellentesque dignissim pelleentesque magna. Fusce vulputate, nunc venenatis sapien sed erat. Fusce tellus mauris, porta non, lacinia nec, scelerisque sed, donec nulla facilis. Ut accumsan posuerit sapien.

*Curabitur nisl odio, sagittis in, euismod nisl, lobortis velit, augue. Fusce nibh mauris, rhoncus dictum, faucibus, dignissim non, neque. Etiam feugiat mauris eget turpis. Pellentesque tellus velit, egaszt sagittis, egaszt a, varius et, neque. Aliquam non lacus. Aliquam in elit. Fusce sagittis tellus quis nisi. Praesent vel enim ut nibh lacuna congue. Pellentesque et sapien ac nisi elevent varius. Nulla molestie. Aenean gravida urna. Aenean suscipit venenatis massa. Morbi velit augue. Curabitur non eros vitae lorem portat soleisque. Mauris ultrices suscipit egestet. Curabitur porta, tellus posuere metuere ultrices, mautis metus laculus ac, eget pelleentesque ligula erat ac dui. Morbi laculis posuere libero.*

# Training Details

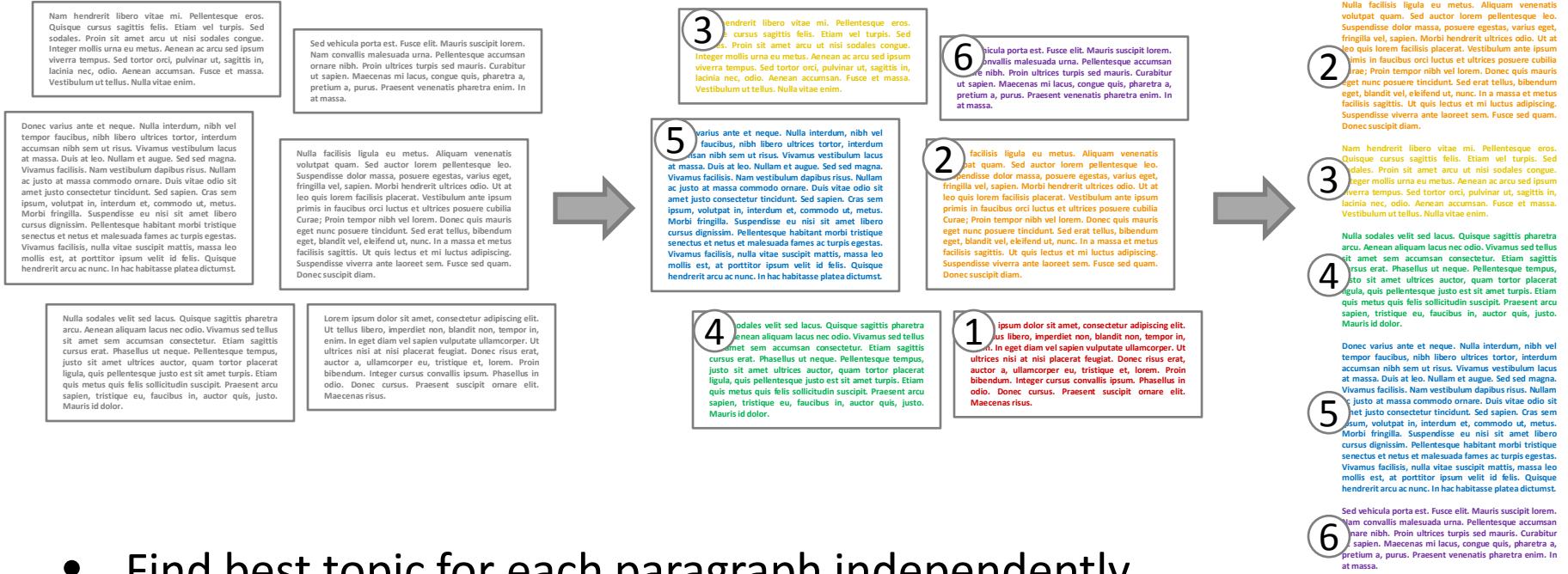
- Hidden variables:
  - *Predominant ordering, Mallows parameters*
  - *Topic counts, language models*
- Estimated from unannotated documents using Bayesian inference techniques
  - *Collapsed Gibbs sampling*
- Mallows distribution admits computationally tractable inference procedures

# Predicting Document Ordering



- Find best topic for each paragraph independently
  - Computationally simple, because marginalization over all orderings is not necessary!

# Predicting Document Ordering



# Evaluation Results: Ordering

Average number of sentences per document: 10

Model	Disasters	NTSB Reports
Content Models (HMM)	81.0	44.0
Content Models (GMM)	83.0	57.0
Random	<1.0	<1.0

See NAACL'09 for evaluation on segmentation and semantic labeling

# Outline

- Applications of Text Models
- Coherence Models (ACL 2005, CL 2008)
- Content Models (NAACL 2004, NAACL 2009)

# Modeling Coherence

Active networks and virtual machines have a long history of collaborating in this manner. The basic tenet of this solution is the refinement of Scheme. The disadvantage of this type of approach, however, is that public-private key pair and red-black trees are rarely incompatible.

- **Coherence** is a property of well-written texts that makes them easier to read and understand than a sequence of randomly strung sentences
- **Local coherence** captures text organization at the level of sentence-to-sentence transitions

# Our Approach

**Key Premise: the distribution of entities in locally coherent discourse exhibits certain regularities**

- Abstract a text into an entity-based representation that encodes syntactic and distributional information
- Learn properties of coherent texts, given a training set of coherent and incoherent texts

# Our Inspiration: Centering Theory

Grosz&Joshi&Weinstein,1983; Strube&Hahn,1999;

Poesio&Stevenson&Di Eugenio&Hitzeman,2004

- Constraints on the entity distribution in a coherent text
  - Focus is the most salient entity in a discourse segment
  - Transition between adjacent sentences is characterized in terms of focus switch
- Constraints on linguistic realization of focus
  - Focus is more likely to be realized as subject or object
  - Focus is more likely to be referred to with anaphoric expression

# Text Representation

- Entity Grid — a two-dimensional array that captures the distribution of discourse entities across text sentences
- Discourse Entity — a class of coreferent noun phrases

## Input Text

- 1 Former Chilean dictator Augusto Pinochet, was arrested in London on October 14th, 1998.
- 2 Pinochet, 82, was recovering from surgery.
- 3 The arrest was in response to an extradition warrant served by a Spanish judge.
- 4 Pinochet was charged with murdering thousands, including many Spaniards.
- 5 He is awaiting a hearing, his fate in the balance.
- 6 American scholars applauded the arrest.

# Input Text with Syntactic Annotation

Use Collins' parser(1997):

1. [Former Chilean dictator Augusto Pinochet]**S**, was arrested in [London]**X** on [October 14th]**X** 1998.
2. [Pinochet]**S**, 82, was recovering from [surgery]**X**.
3. [The arrest]**S** was in [response]**X** to [an extradition warrant]**X** served by [a Spanish judge]**S**.
4. [Pinochet]**S** was charged with murdering [thousands]**O**, including many [Spaniards]**O**.
5. [He]**S** is awaiting [a hearing]**O**, [his fate]**X** in [the balance]**X**.
6. [American scholars]**S** applauded the [arrest]**O**.

Notation: **S**=subjects, **O**=object, **X**=other

# Input Text with Coreference Information

Use noun-phrase coreference tool (Ng and Cardie, 2002):

1. [Former Chilean dictator Augusto Pinochet]**S**, was arrested in [London]**X** on [October 14]**X** 1998.
2. [Pinochet]**S**, 82, was recovering from [surgery]**X**.
3. [The arrest]**S** was in [response]**X** to [an extradition warrant]**X** served by [a Spanish judge]**S**.
4. [Pinochet]**S** was charged with murdering [thousands]**O**, including many [Spaniards]**O**.
5. [He]**S** is awaiting [a hearing]**O**, [his fate]**X** in [the balance]**X**.
6. [American scholars]**S** applauded the [arrest]**O**.

## Output Entity Grid

	Pinochet	London	October	Surgery	Arrest	Extradition	Warrant	Judge	Thousands	Spaniards	Hearing	Fate	Balance	Scholars	
1	<i>S</i>	X	X	-	-	-	-	-	-	-	-	-	-	-	1
2	<i>S</i>	-	-	X	-	-	-	-	-	-	-	-	-	-	2
3	<i>-</i>	-	-	-	<i>S</i>	X	X	<i>S</i>	-	-	-	-	-	-	3
4	<i>S</i>	-	-	-	-	-	-	-	O	O	-	-	-	-	4
5	<i>S</i>	-	-	-	-	-	-	-	-	-	O	X	X	-	5
6	<i>-</i>	-	-	-	O	-	-	-	-	-	-	-	-	<i>s</i>	6

# Comparing Grids

# Text Encoding as Feature Vector

	s s	o x		s o x		s	o x		s	o	x	
	s s	s s	s	o o o	o	x	x x x	x				
$d_{i1}$	0	0	0	.03	0	0	.02	.07	0	0	.12	.02
$d_{i2}$	.02	0	0	.03	0	0	.06	0	0	0	.05	.03
									.07	.07	.25	.29

Each grid rendering  $x_{ij}$  of a document  $d_i$  is represented by a feature vector:

$$\Phi(x_{ij}) = (p_1(x_{ij}), p_2(x_{ij}), \dots, p_m(x_{ij}))$$

where  $m$  is the number of all predefined entity transitions, and  $p_t(x_{ij})$  the probability of transition  $t$  in the grid  $x_{ij}$

# Learning a Ranking Function

- **Training Set**

Ordered pairs  $(x_{ij}, x_{ik})$ , where  $x_{ij}$  and  $x_{ik}$  are renderings of the same document  $d_i$ , and  $x_{ij}$  exhibits a higher degree of coherence than  $x_{ik}$

- **Training Procedure**

- Goal: Find a parameter vector  $\vec{w}$  that yields a “ranking score” function  $\vec{w} \cdot \Phi(x_{ij})$  satisfying:

$$\vec{w} \cdot (\Phi(x_{ij}) - \Phi(x_{ik})) > 0$$

$\forall (x_{ij}, x_{ik})$  in training set

- Method: Constraint optimization problem solved using the search technique described in Joachims (2002)

# Evaluation: Information Ordering

- **Goal:** recover the most coherent sentence ordering
- **Basic set-up:**
  - Input: a pair of a source document and a permutation of its sentences
  - Task: find a source document via coherence ranking
- **Data:** Training 4000 pairs, Testing 4000 pairs (Natural disasters and Transportation Safety Reports)



# Evaluation: Summarization

- **Goal:** select the most coherent summary among several alternatives
- **Basic set-up:**
  - Input: a pair of system summaries
  - Task: predict the ranking provided by human
- **Data:** 96 summary pairs for training, 32 pairs for testing (from DUC 2003)

# Evaluation Results

Model	Disasters	NTSB reports	Summary Ranking
Grid	87.2	90.4	80.0
Content Model (HMM)	88.0	75.8	N/A
Random	50	50	50

# Follow-Ups

- Elsner, Austerweil, and Charniak (NAACL 2007): combining content and coherence models is effective
  - **Brown Coherence Toolkit**: software for coherence models and test applications
- Pitler and Nenkova (EMNLP 2008): grid-based features help to predict readability of human-authored texts

# Conclusions

- Practical Benefits: The proposed models can be easily incorporated in text-to-text generation
- Main Findings: The proposed document models can be successfully applied to modeling text structure
- Future Research: Design statistical models that match representational power of traditional discourse theories

# References

- Harr Chen, S.R.K. Branavan, Regina Barzilay, David R. Karger. “*Latent Topic Models for Document Structure Induction*”, NAACL-HLT, 2009.
- Regina Barzilay, Mirella Lapata “*Modeling Local Coherence: An Entity-based Approach*”, Computational Linguistics, 2008.
- Regina Barzilay, Lillian Lee “*Catching the Drift: Probabilistic Content Models, with Applications to Generation and Summarization*”, NAACL-HLT, 2004