

# Robust Pseudo Feedback & HMM Passage Extraction

UIUC at TREC 2006 Genomics Track

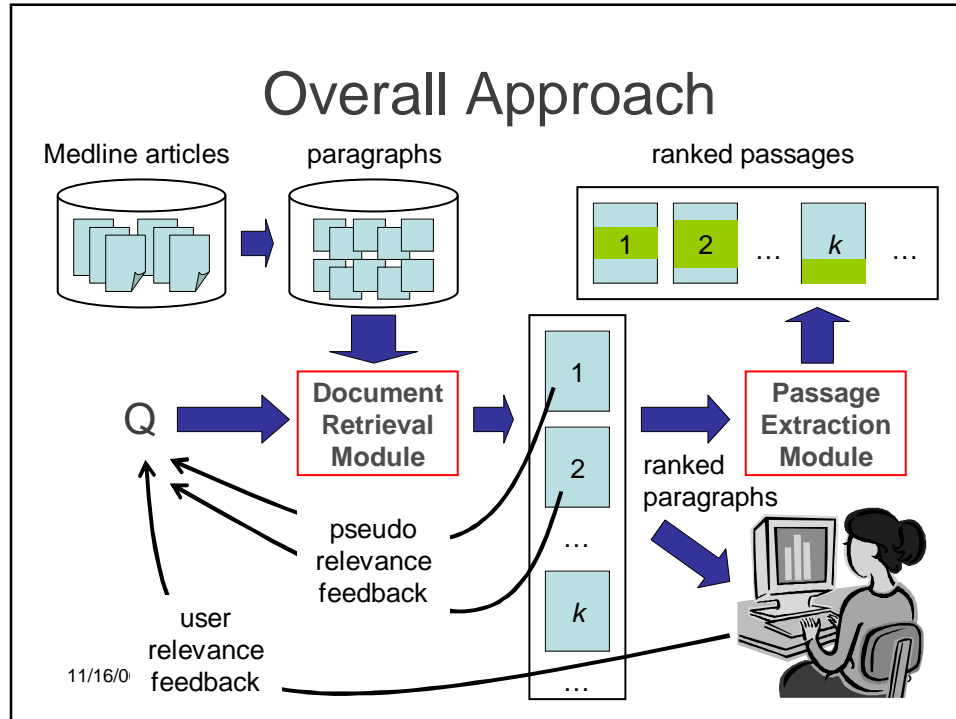
Jing Jiang, Xin He, ChengXiang Zhai  
University of Illinois at Urbana-Champaign

## Goal of Participation

- To test the effectiveness of some recent language modeling methods for genomics retrieval
  - Robust pseudo feedback [Tao & Zhai 06]
  - HMM passage extraction [Jiang & Zhai 06]
- Task at 2006 genomics track
  - Document-level retrieval
  - Passage-level retrieval
  - Aspect-level retrieval

11/16/06

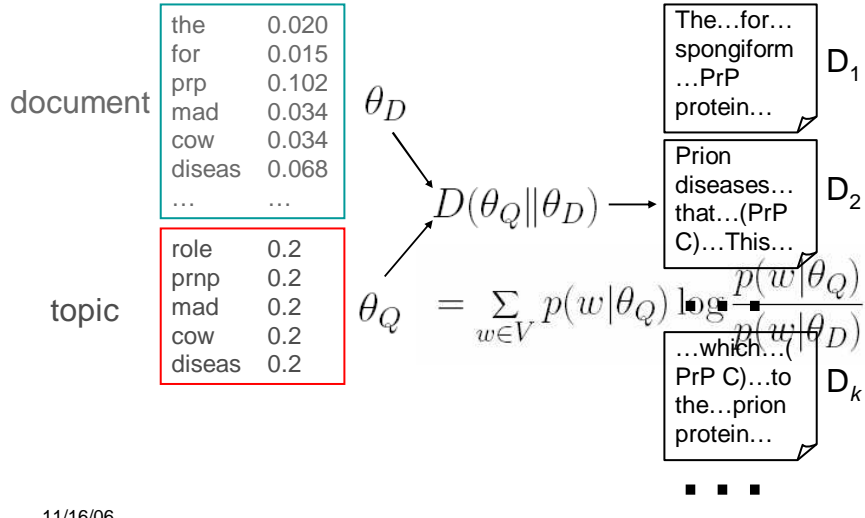
2



- ## Goal of Participation
- To test the effectiveness of some recent language modeling methods for genomics retrieval
    - Robust pseudo feedback [Tao & Zhai 06]
    - HMM passage extraction [Jiang & Zhai 06]
- 11/16/06 4

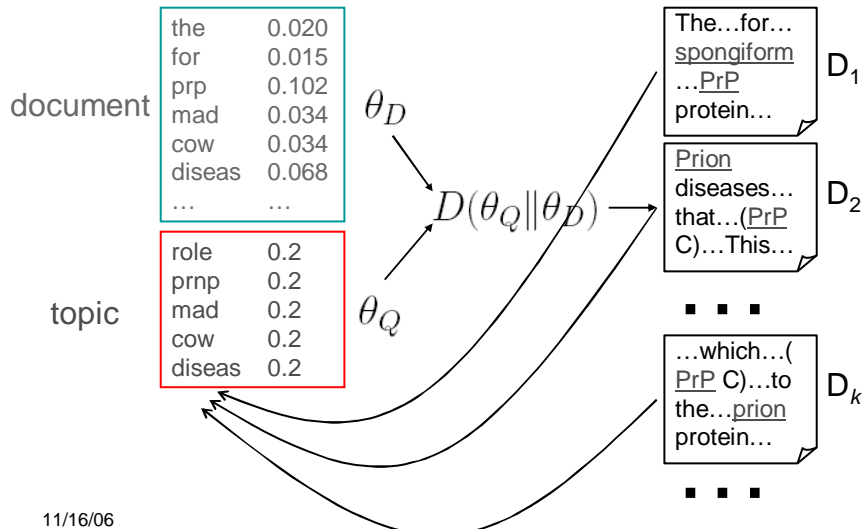
# KL-Divergence Retrieval Model

[Lafferty & Zhai 01]



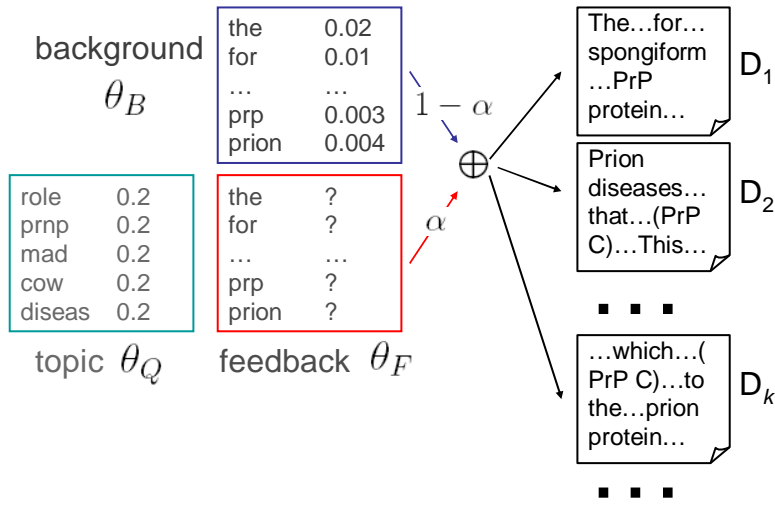
# KL-Divergence Retrieval Model

[Lafferty & Zhai 01]



# Model-Based Feedback

[Zhai & Lafferty 01]

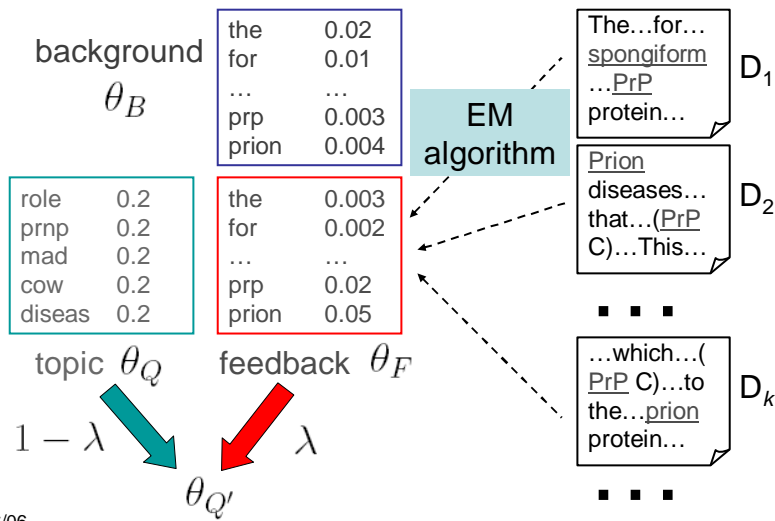


11/16/06

7

# Model-Based Feedback

[Zhai & Lafferty 01]

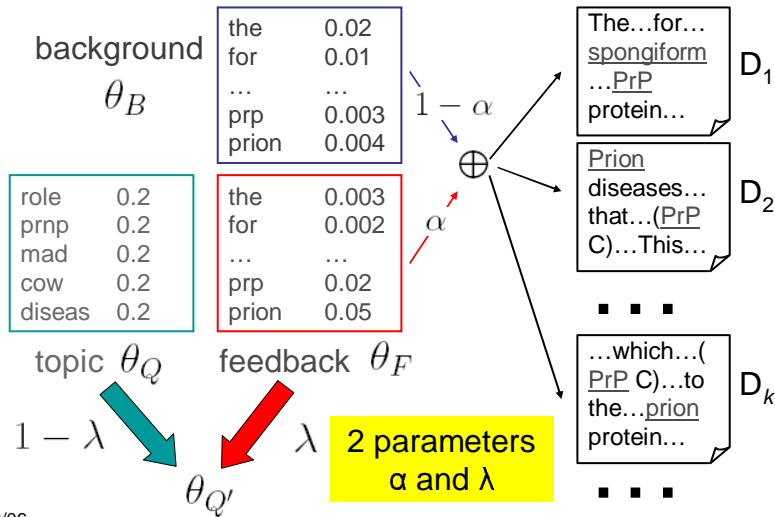


11/16/06

8

# Model-Based Feedback

[Zhai & Lafferty 01]

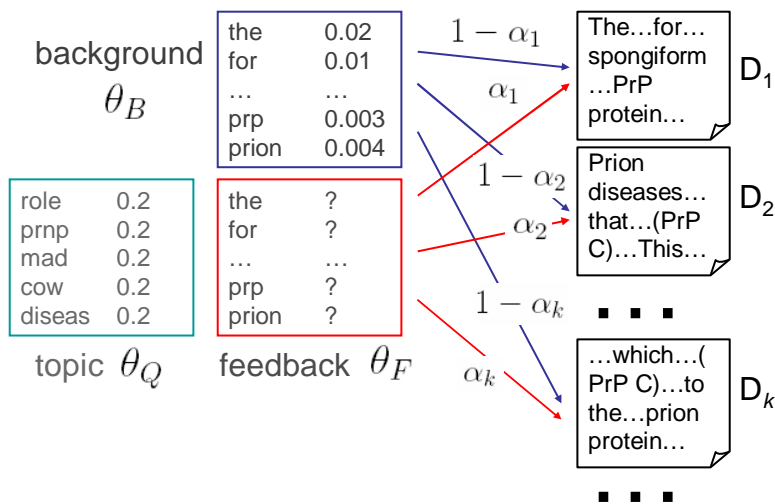


11/16/06

9

# Regularized Estimation

[Tao & Zhai 06]

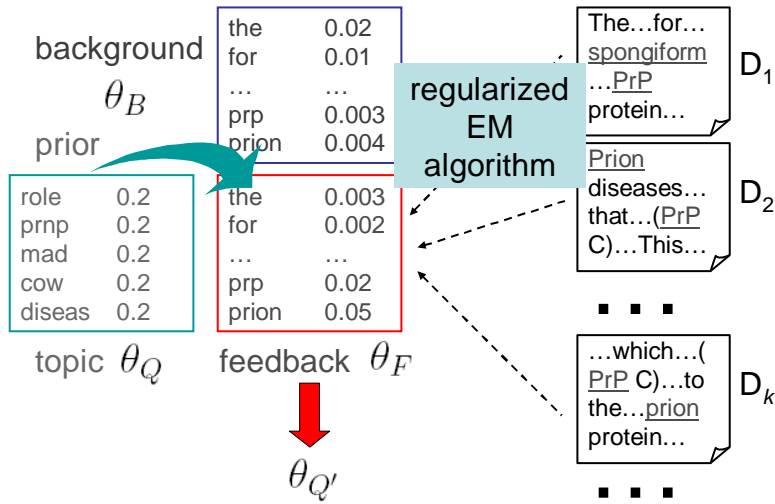


11/16/06

10

# Regularized Estimation

[Tao & Zhai 06]

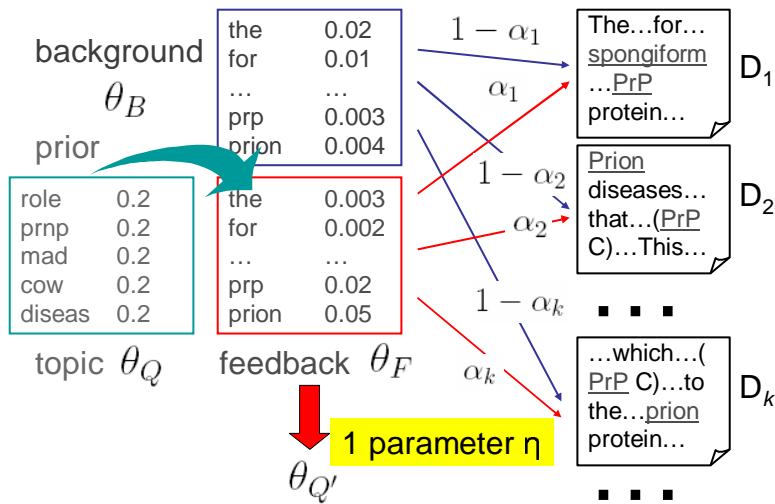


11/16/06

11

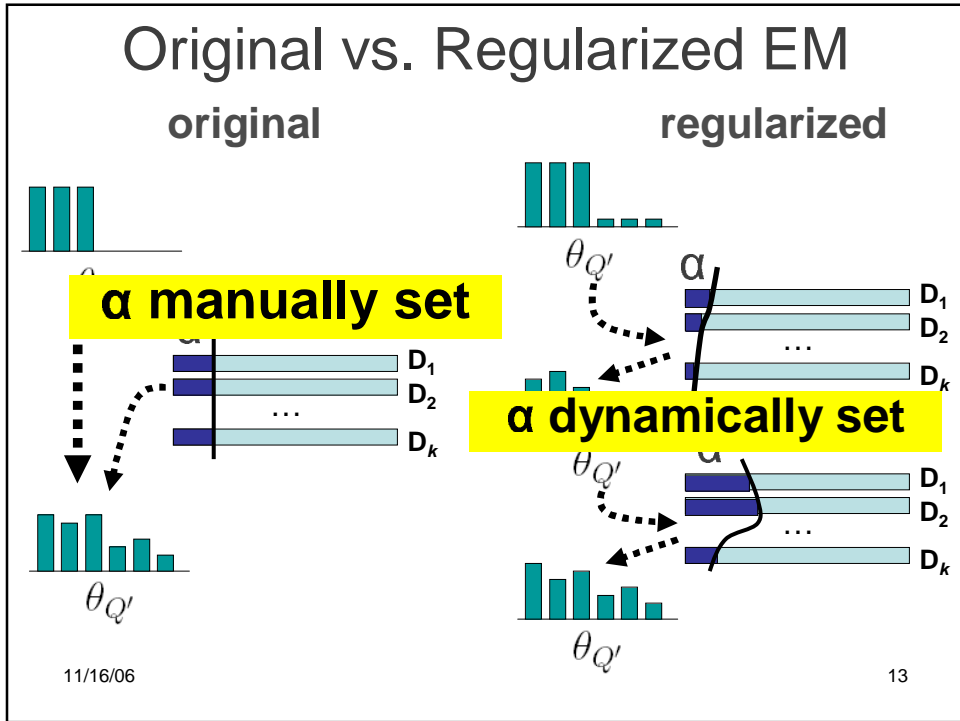
# Regularized Estimation

[Tao & Zhai 06]



11/16/06

12



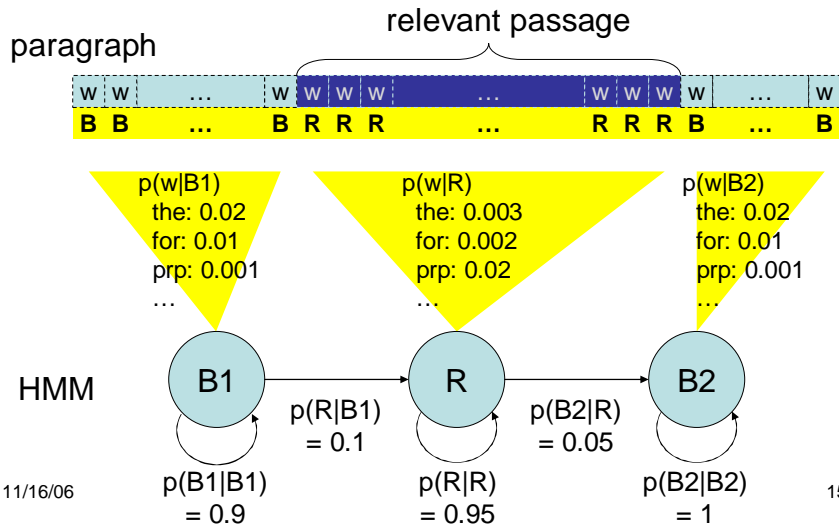
## Goal of Participation

- To test the effectiveness of some recent language modeling methods for genomics retrieval
  - Robust pseudo feedback [Tao & Zhai 06]
  - HMM passage extraction [Jiang & Zhai 06]

11/16/06 14

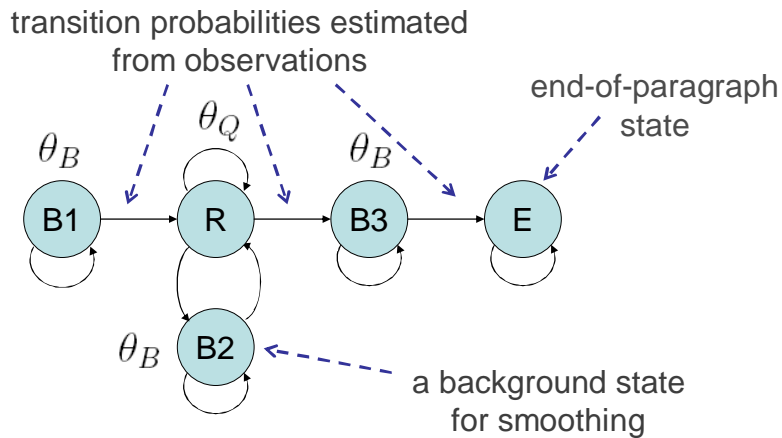
# HMM Passage Extraction

[Jiang & Zhai 06]



# HMM Passage Extraction

[Jiang & Zhai 06]



11/16/06

16

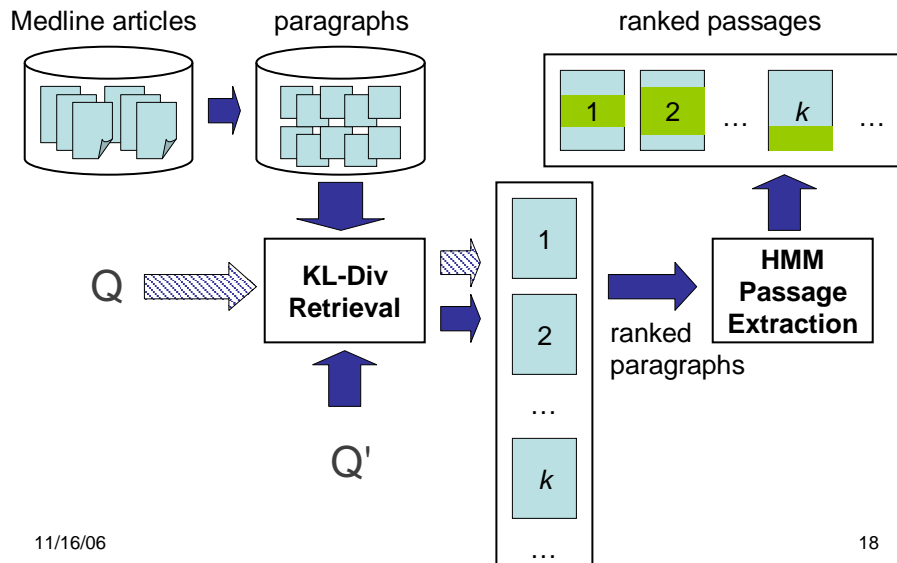
# Experiment Design

- Pre-processing
  - HTML parsing
  - paragraph boundaries
  - Tokenization
- User relevance feedback

11/16/06

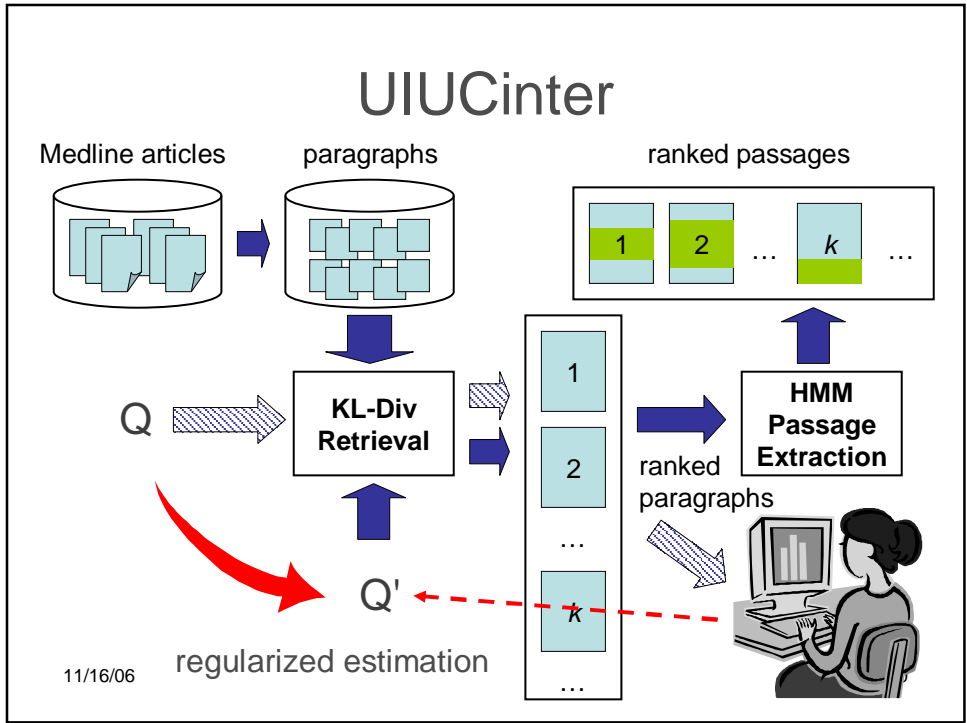
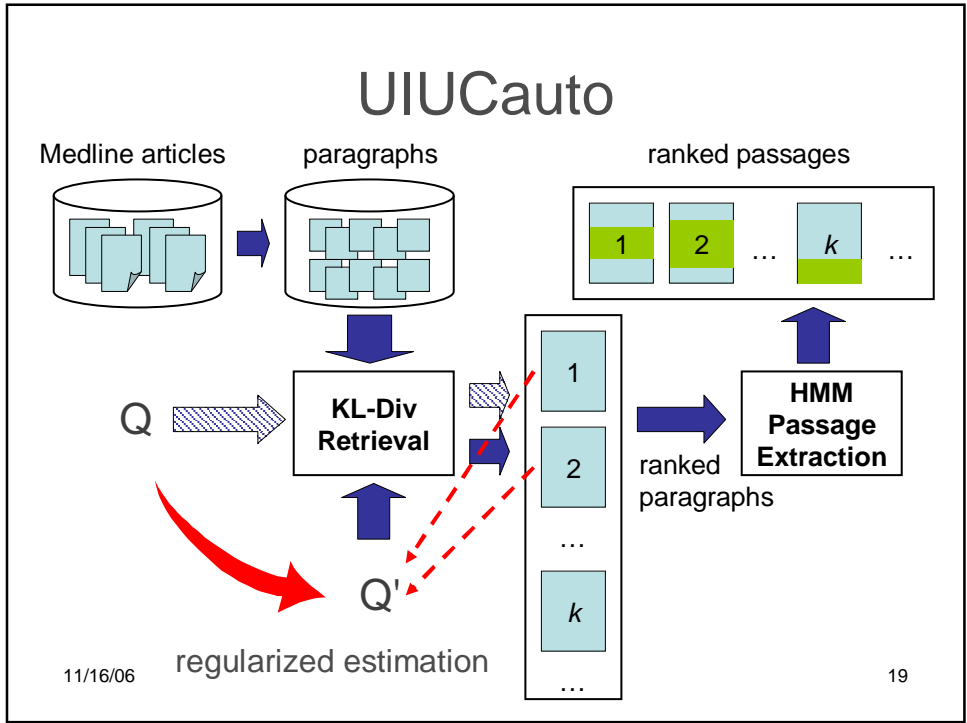
17

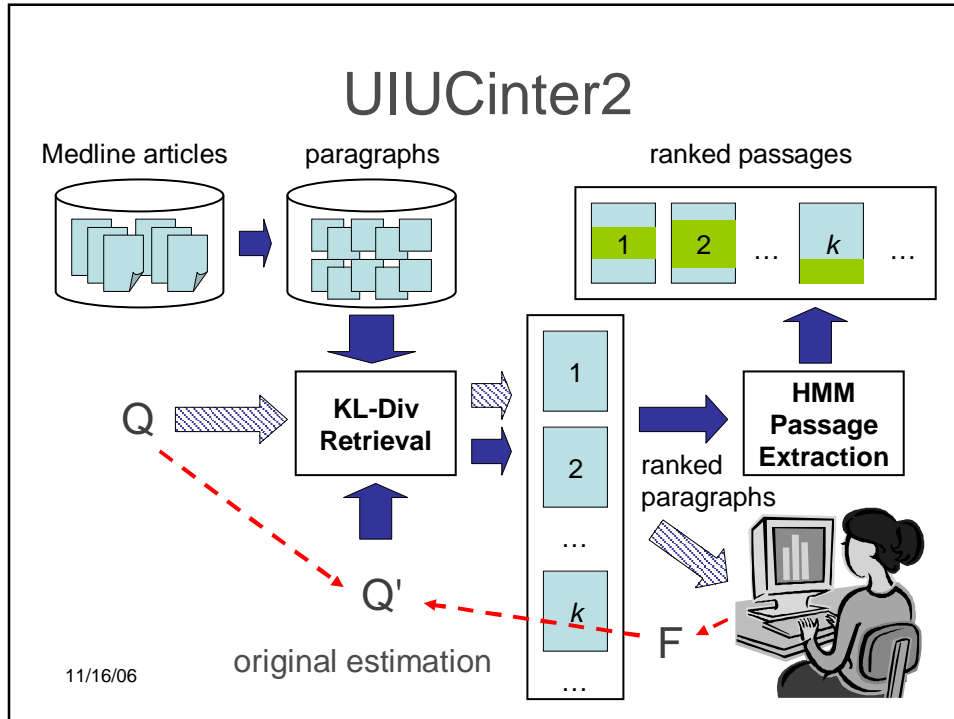
## Official Runs



11/16/06

18





## Pseudo Relevance Feedback ( $k = 10$ )

Method		Doc MAP	Rel. Impr.
<b>Baseline (no feedback)</b>		0.3484	N/A
<b>Original Estimation</b>	Def	0.3606	<b>+3.50%</b>
	Opt	0.3943	<b>+13.2%</b>
<b>Regularized Estimation</b>	Def	0.3842 (UIUCauto)	<b>+10.3%</b>
	Opt	0.3952	<b>+13.4%</b>

$\eta$  is similar to  $\lambda / (1 - \lambda)$

11/16/06 22

## Pseudo Relevance Feedback ( $k = 10$ )

Method		Doc MAP	Rel. Impr.
<b>Baseline (no feedback)</b>		0.3484	N/A
<b>Original Estimation</b>	Def	0.3606	<b>+3.50%</b>
	Opt	0.3943	<b>+13.2%</b>
<b>Regularized Estimation</b>	Def	0.3842 (UIUCauto)	<b>+10.3%</b>
	Opt	0.3952	<b>+13.4%</b>

$\eta$  is similar to  $\lambda / (1 - \lambda)$

11/16/06

23

## Pseudo Relevance Feedback ( $k = 10$ )

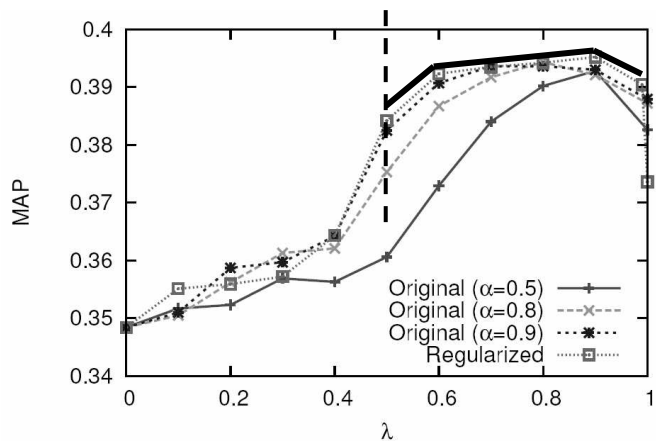
Method		Doc MAP	Rel. Impr.
<b>Baseline (no feedback)</b>		0.3484	N/A
<b>Original Estimation</b>	Def	0.3606	<b>+3.50%</b>
	Opt	0.3943	<b>+13.2%</b>
<b>Regularized Estimation</b>	Def	0.3842 (UIUCauto)	<b>+10.3%</b>
	Opt	0.3952	<b>+13.4%</b>

$\eta$  is similar to  $\lambda / (1 - \lambda)$

11/16/06

24

## Parameter Sensitivity (pseudo feedback, $k = 10$ )



11/16/06

25

## User Relevance Feedback

Method		Doc MAP		
		Pseudo Feedback	User Feedback	Rel. Impr.
<b>Original Estimation</b>	Def	0.3606	0.3986	<b>+10.5%</b>
	Opt	0.3943	0.4511	<b>+14.4%</b>
<b>Regularized Estimation</b>	Def	0.3842 (UIUCauto)	0.4261 (UIUCinter)	<b>+10.9%</b>
	Opt	0.3952	0.4515	<b>+14.2%</b>

11/16/06

26

## User Relevance Feedback

Method		Doc MAP		
		Pseudo Feedback	User Feedback	Rel. Impr.
Original Estimation	Def	0.3606	0.3986	+10.5%
	Opt	0.3943	0.4511	+14.4%
Regularized Estimation	Def	0.3842 (UIUCauto)	0.4261 (UIUCinter)	+10.9%
	Opt	0.3952	0.4515	+14.2%

11/16/06

27

## User Relevance Feedback

Method		Doc MAP		
		Pseudo Feedback	User Feedback	Rel. Impr.
Original Estimation	Def	0.3606	0.3986	+10.5%
	Opt	0.3943	0.4511	+14.4%
Regularized Estimation	Def	0.3842 (UIUCauto)	0.4261 (UIUCinter)	+10.9%
	Opt	0.3952	0.4515	+14.2%

11/16/06

28

## HMM Passage Extraction

Method		Psg MAP
UIUCauto	Paragraph	0.03753
	HMM Passage	0.04864
	Rel. Impr.	<b>+29.6%</b>
UIUCinter	Paragraph	0.04481
	HMM Passage	0.05906
	Rel. Impr.	<b>+31.8%</b>
UIUCinter2	Paragraph	0.04580
	HMM Passage	0.06038
	Rel. Impr.	<b>+31.8%</b>

11/16/06

29

## Passage Length (In Bytes)

	Max	Min	Avg	Std
<b>True Passages</b>	6928	27	<b>399.8</b>	489.4
<b>HMM Passages</b>	6955	34	<b>1525.8</b>	949.7
<b>Paragraph</b>	8670	60	<b>2105.4</b>	1136.8

**HMM passages are generally too long!**

11/16/06

30

## Example Passage

**Prion diseases, which include Creutzfeldt-Jacob disease in humans, mad cow disease in cattle, and scrapie in sheep, involve the misfolding of the benign cellular prion protein (PrP C) 1 to the infectious disease-causing scrapie isoform PrP Sc.** The prion protein (PrP C) is a copper-binding cell surface glycoprotein. The role of copper in the normal function of PrP, as well as in prion diseases, has been the subject of a number of excellent reviews. The mature cellular form of PrP consists of residues 23 to 231 and is tethered to the cell surface via a glycosylphosphatidylinositol anchor at the C terminus. There are now a number of NMR solution structures of copper-free mammalian PrPs. A crystal structure of PrP C has also been published; this structure is dimeric involving domain swapping of the monomeric form.

11/16/06

31

## Example Passage

**Prion diseases, which include Creutzfeldt-Jacob disease in humans, mad cow disease in cattle, and scrapie in sheep, involve the misfolding of the benign cellular prion protein (PrP C) 1 to the infectious disease-causing scrapie isoform PrP Sc.** The prion protein (PrP C) is a copper-binding cell surface glycoprotein. The role of copper in the normal function of PrP, as well as in prion diseases, has been the subject of a number of excellent reviews. The mature cellular form of PrP consists of residues 23 to 231 and is tethered to the cell surface via a glycosylphosphatidylinositol anchor at the C terminus. There are now a number of NMR solution structures of copper-free mammalian PrPs. A crystal structure of PrP C has also been published; this structure is dimeric involving domain swapping of the monomeric form.

11/16/06

32

## Conclusions and Future Work

- The two language modeling methods in general works well in genomics domain
  - Regularized feedback estimation can effectively eliminates parameter  $\alpha$
  - HMM passages improves over paragraphs
- User relevance feedback is effective
- Limitations and future work
  - Regularized feedback estimation still has parameter  $\eta$  to tune
    - How to eliminate  $\eta$ ?
  - The inherent coherence property of HMM passages may not suit the task well
    - Different/better HMM architecture?

11/16/06

33

## The End

- Questions?

11/16/06

34