

YAHOO!

Storylines from Streaming Text

The Infinite Topic Cluster Model

Amr Ahmed, Jake Eisenstein, Qirong Ho

Alex Smola, Choon Hui Teo, Eric Xing

Carnegie Mellon University

Yahoo! Research

Outline

- Visualizing a news stream
- Goals
 - Clusters & Content analysis
- The story-topic model
 - Recurrent Chinese Restaurant Process
 - Latent Dirichlet Allocation
- Examples

News Stream



Add-ons turn tax cut bill into 'Christmas tree'

AP - 1 hr 32 mins ago
WASHINGTON - In the

BEYOND FOSSIL FUELS

Using Waste, Swedish



As part of its citywide system, Kristianstad burns wood waste like tree prunings and scraps from flooring factories to power an underground district heating grid.

China says inflation up 5.1 per cent

Associated Press

Buzz up! 19 votes | Share



Wall Street Video: [Charting Consumer Sentiment](#) CNBC



Wall Street Video: [Bright Future](#) TheStreet.com

RELATED QUOTES

^DJI	11,410.32	+40.26
^GSPC	1,240.40	+7.40
^IXIC	2,637.54	+20.87

By CARA ANNA, Associated Press

BEIJING - China's inflation surged Saturday, despite supplies and end diesel shortages

The 5.1 percent inflation rate was driven by a 11.7 percent jump in food prices year on year.

The news comes as China's leaders meet for the top economic planning conference of the year and as financial markets watch for a widely anticipated [interest rate hike](#) to help bring rapid economic growth to a more sustainable level.

"I think this means that an interest rate hike of 25 basis points is very likely by the end of the year," said CLSA analyst Andy Rothman.

Suit to Recover Madoff's Money Calls Austrian an Accomplice

By DIANA B. HENRIQUES and PETER LATTMAN

Sonja Kohn, an Austrian banker, is accused of masterminding a 23-year conspiracy that played a central role in financing the gigantic Ponzi scheme.

Post a Comment

er

Print

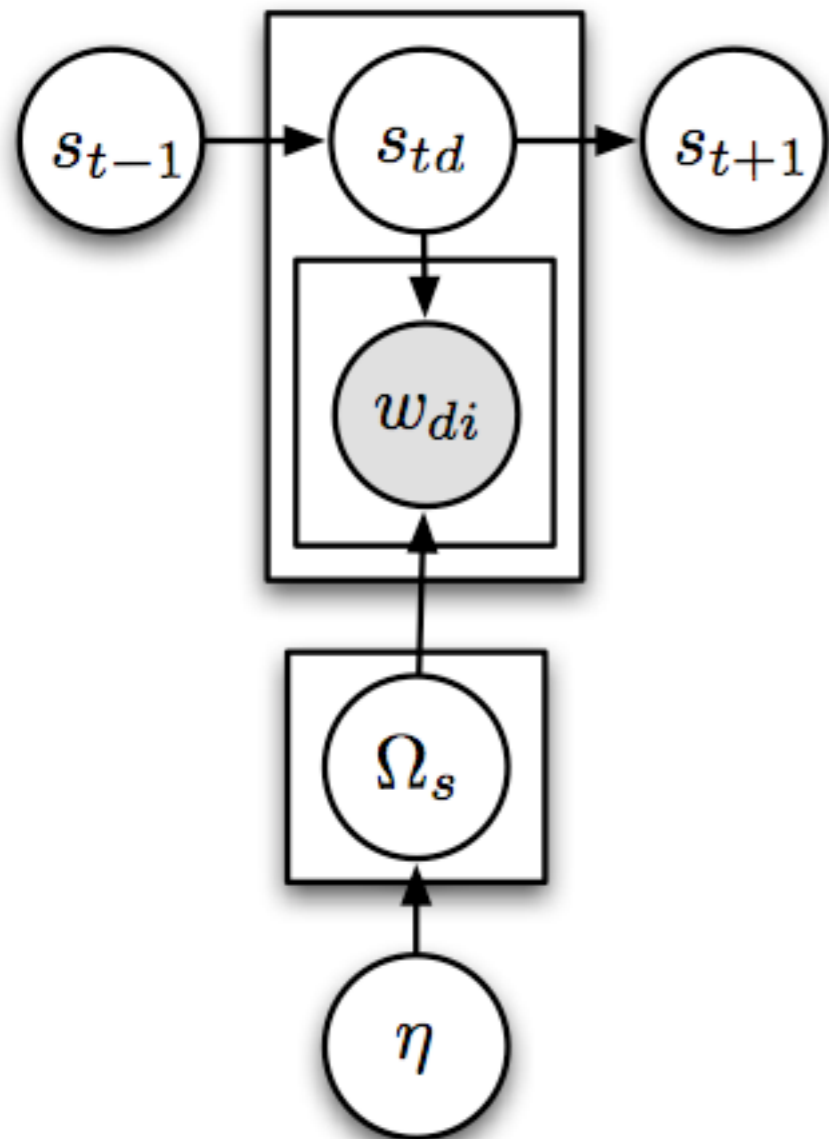
November, base food

News Stream

- **Realtime news stream**
 - **Multiple sources (Reuters, AP, CNN, ...)**
 - **Same story from multiple sources**
 - **Stories are related**
- **Goals**
 - **Aggregate articles into a storyline**
 - **Analyze the storyline (topics, entities)**

Precursors

Evolutionary Clustering / RCRP



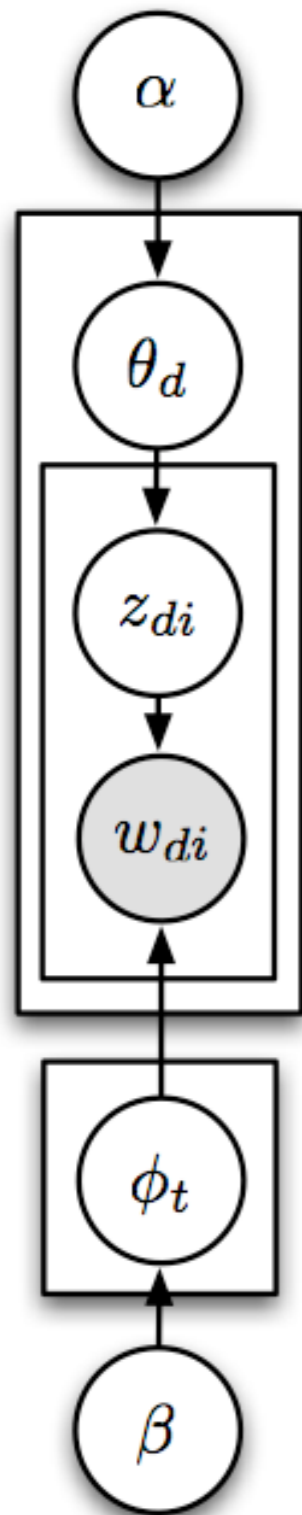
- Assume active story distribution at time t
- Draw story indicator
- Draw words from story distribution
- Down-weight story counts for next day

Ahmed & Xing, 2008

Clustering / RCRP

- Pro
 - Nonparametric model of story generation (no need to model frequency of stories)
 - No fixed number of stories
 - Efficient inference via collapsed sampler
- Con
 - We learn nothing!
 - No content analysis

Latent Dirichlet Allocation



- Generate topic distribution per article
- Draw topics per word from topic distribution
- Draw words from topic specific word distribution

Blei, Ng, Jordan, 2003

Latent Dirichlet Allocation

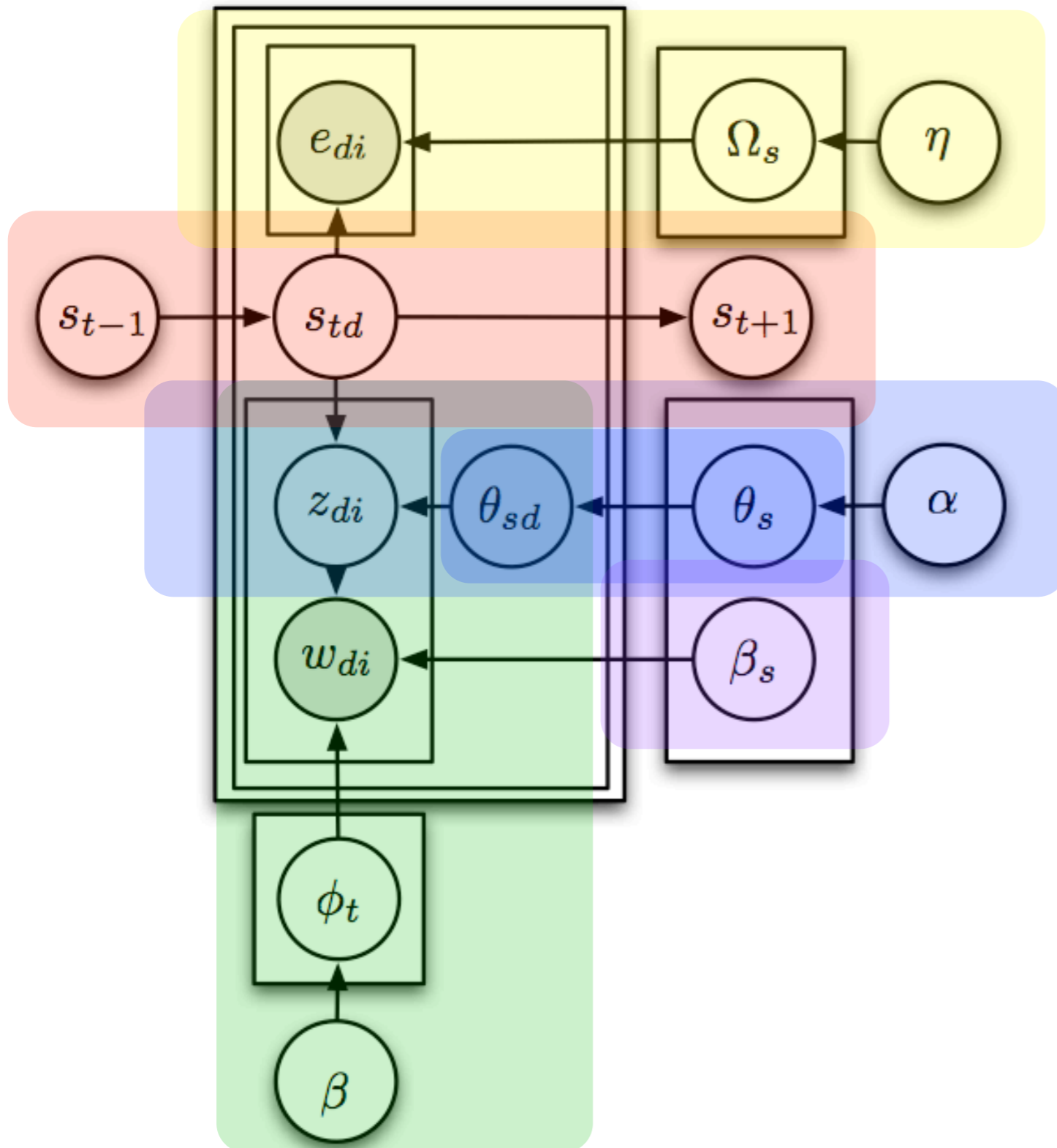
- Pro
 - Topical analysis of stories
 - Topical analysis of words (meaning, saliency)
 - More documents improve estimates
- Con
 - No clustering

More Issues



Storylines

Storylines Model



- Topic model
- Topics per cluster
- RCRP for cluster
- Hierarchical DP for article
- Separate model for named entities
- Story specific correction

Inference

- We receive articles as a stream
 - Want topics & stories now
- Variational inference infeasible
 - (RCRP, sparse to dense, vocabulary size)
- We have a 'tracking problem'
 - Sequential Monte Carlo
 - Use sampled variables of surviving particle
 - Use ideas from Cannini et al. 2009

Estimation

- Proposal distribution - draw stories s , topics z

$$p(s_{t+1}, z_{t+1} | x_{1..t+1}, s_{1..t}, z_{1..t})$$

using Gibbs Sampling for each particle

- Reweight particle via

past state

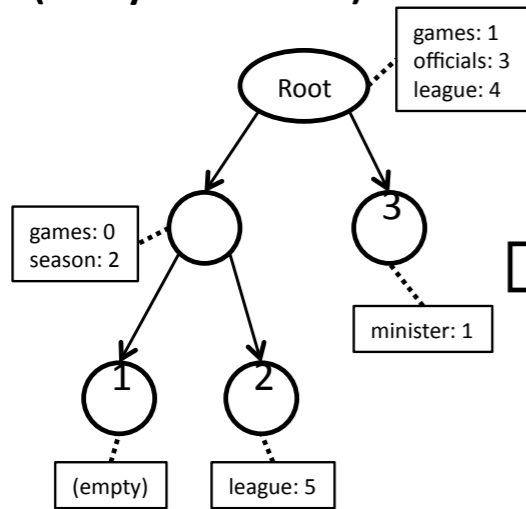
new data

$$p(x_{t+1} | x_{1..t}, s_{1..t}, z_{1..t})$$

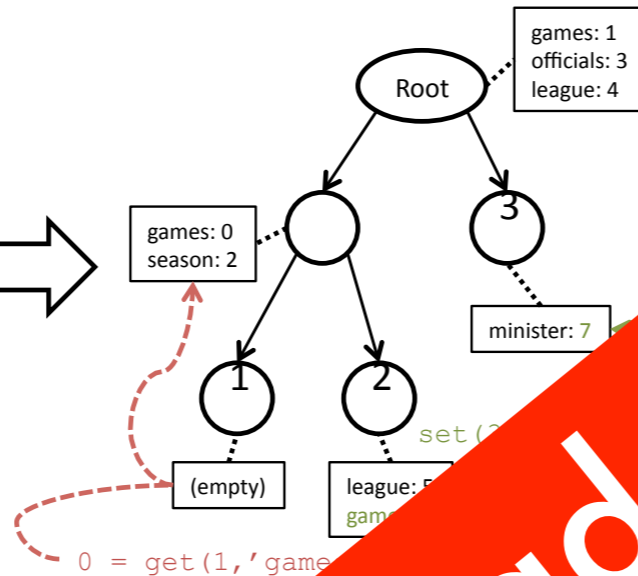
- Resample particles if l_2 norm too large
(resample some assignments for diversity, too)
- Compare to multiplicative updates algorithm
In our case predictive likelihood yields weights

Inheritance Tree

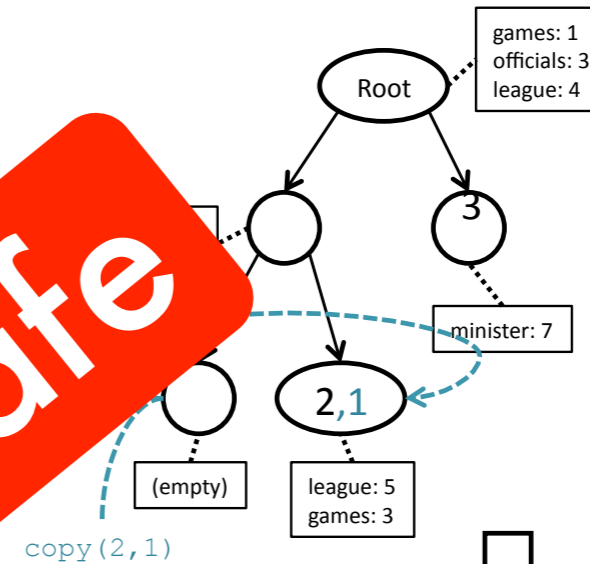
Initial tree
(ready for threads)



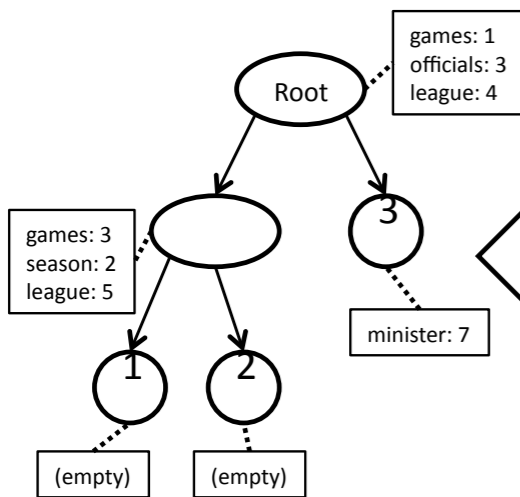
Filter threads *update* particles



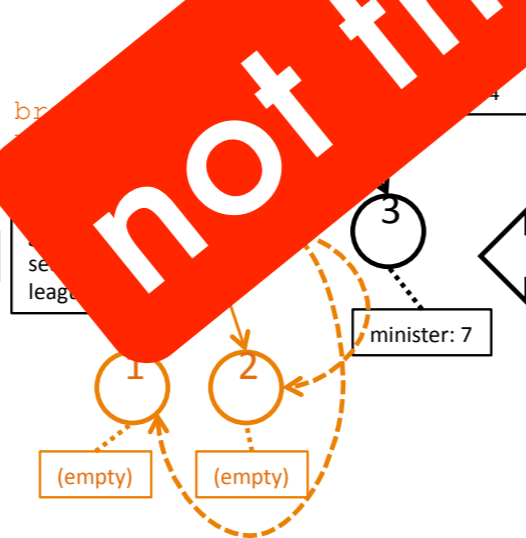
Resampling *copies* particles



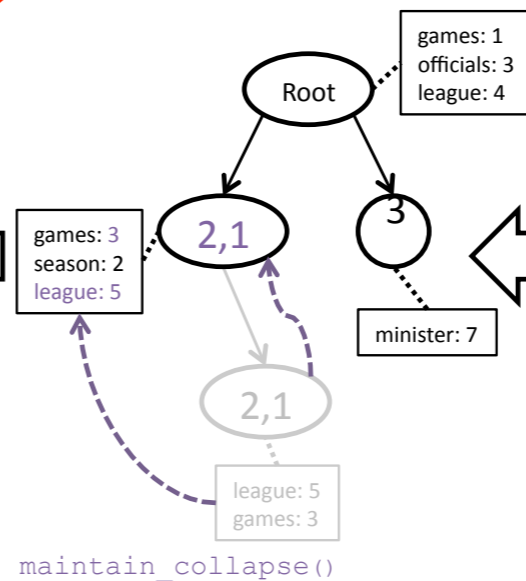
New initial tree
(ready for threads)



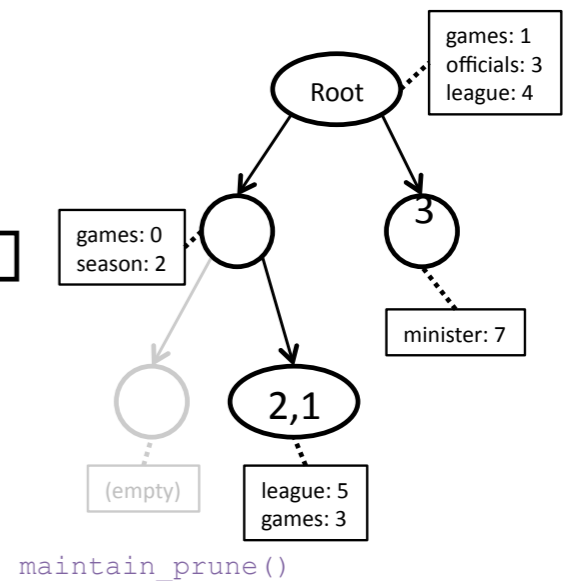
Create *new* particles



collapse long branches



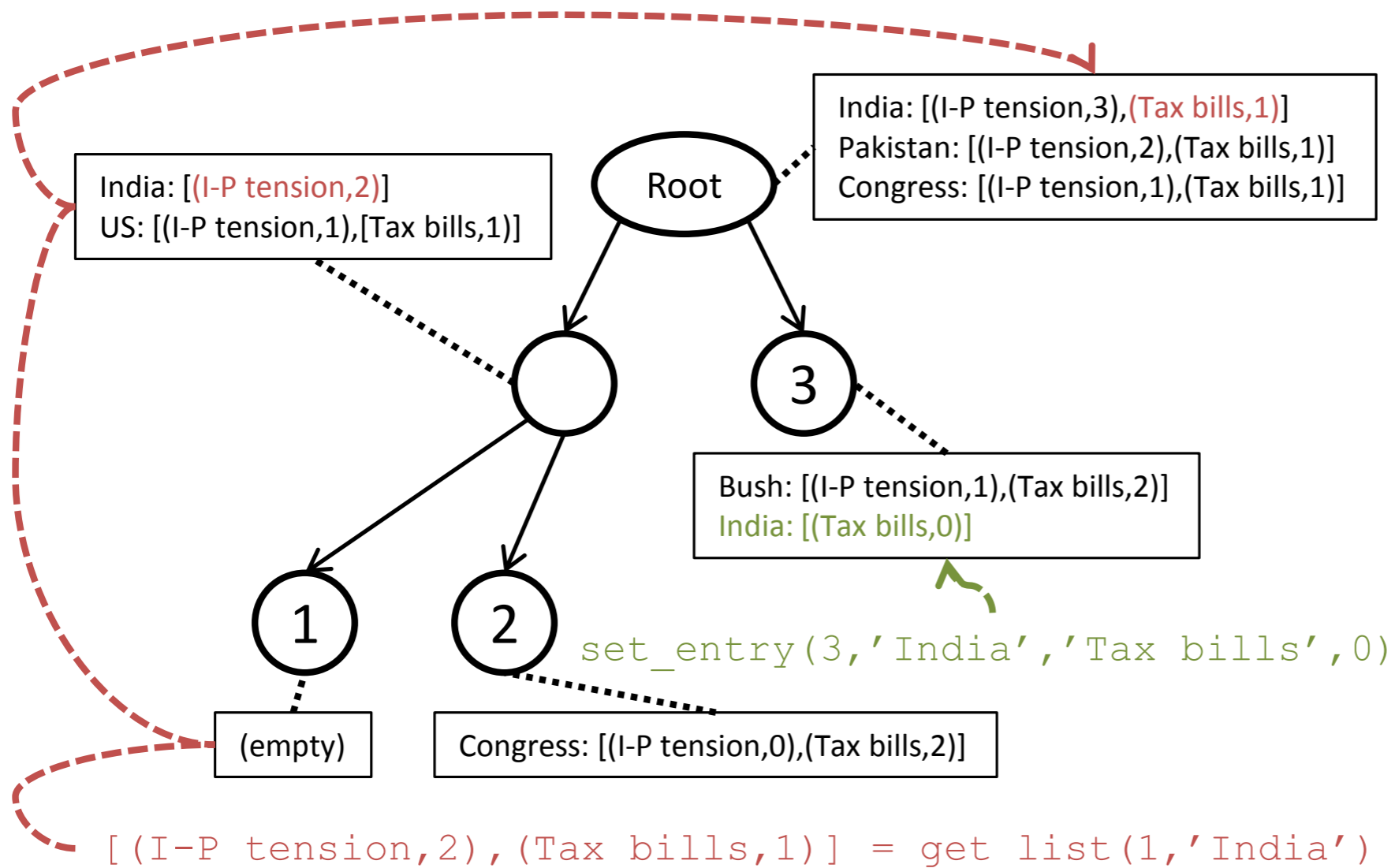
Prune unused branches



not thread safe

Extended Inheritance Tree

Extended Inheritance Tree



write only in
the leaves
(per thread)

Note: "I-P tension" is short for "India-Pakistan tension"

Results

Numbers ...

- **TDT5 (Topic Detection and Tracking)**
macro-averaged minimum detection cost: 0.714

time	entities	topics	story words
0.84	0.90	0.86	0.75

This is the best performance on TDT5!

- **Yahoo News data**
... beats all other clustering algorithms

Stories

TOPICS

Sports

games
won
team
final
season
league
held

Politics

government
minister
authorities
opposition
officials
leaders
group

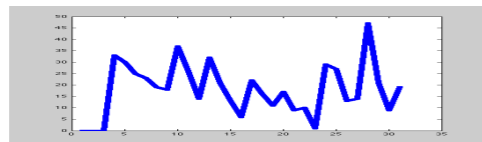
Unrest

police
attack
run
man
group
arrested
move

STORYLINES

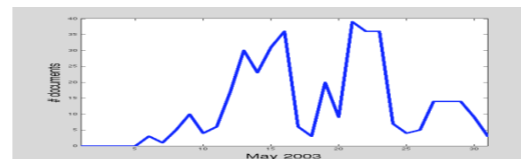
UEFA-soccer

champions	<i>Juventus</i>
goal	<i>AC Milan</i>
leg	<i>Real Madrid</i>
coach	<i>Milan</i>
striker	<i>Lazio</i>
midfield	<i>Ronaldo</i>
penalty	<i>Lyon</i>



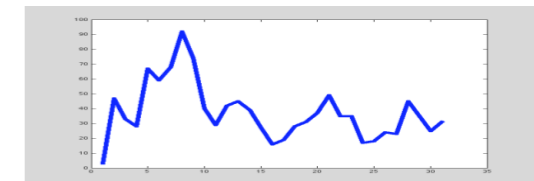
Tax bills

tax	<i>Bush</i>
billion	<i>Senate</i>
cut	<i>US</i>
plan	<i>Congress</i>
budget	<i>Fleischer</i>
economy	<i>White House</i>
lawmakers	<i>Republican</i>



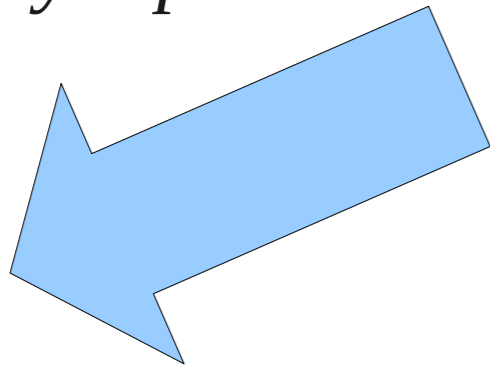
India-Pakistan tension

nuclear	<i>Pakistan</i>
border	<i>India</i>
dialogue	<i>Kashmir</i>
diplomatic	<i>New Delhi</i>
militant	<i>Islamabad</i>
insurgency	<i>Musharraf</i>
missile	<i>Vajpayee</i>



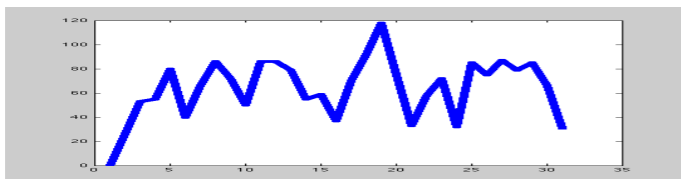
Related Stories

“Show similar stories by topic”



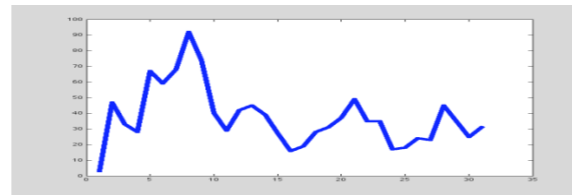
Middle-east conflict

Peace	<i>Israel</i>
Roadmap	<i>Palestinian</i>
Suicide	<i>West bank</i>
Violence	<i>Sharon</i>
Settlements	<i>Hamas</i>
bombing	<i>Arafat</i>

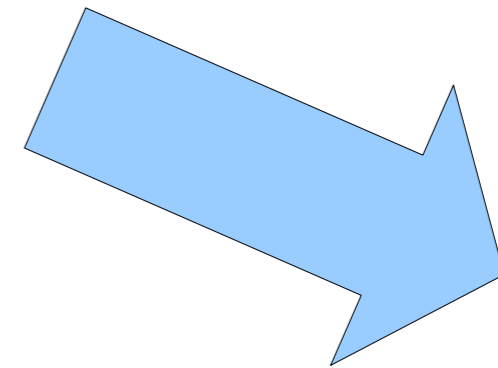


India-Pakistan tension

nuclear	<i>Pakistan</i>
border	<i>India</i>
dialogue	<i>Kashmir</i>
diplomatic	<i>New Delhi</i>
militant	<i>Islamabad</i>
insurgency	<i>Musharraf</i>
missile	<i>Vajpayee</i>

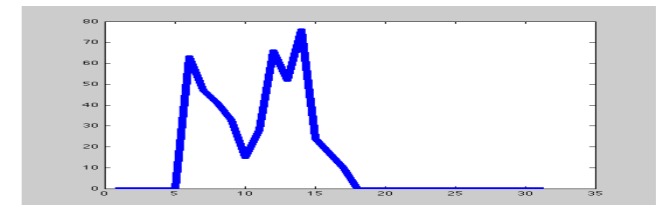


“Show similar stories, require the word nuclear”



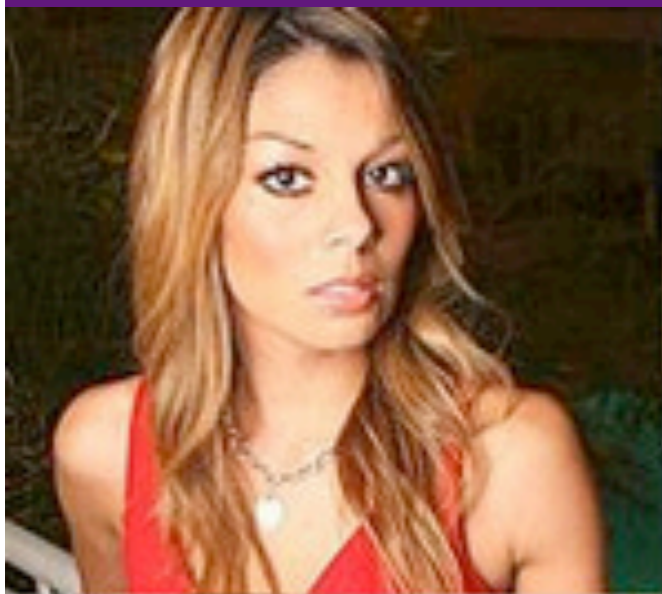
North Korea nuclear

nuclear	<i>North Korea</i>
summit	<i>South Korea</i>
warning	<i>U.S</i>
policy	<i>Bush</i>
missile	<i>Pyongyang</i>
program	



Issues

To Do



?

