

Distributed MAP Inference for Undirected Graphical Models

Sameer Singh¹ Amarnag Subramanya²
Fernando Pereira² Andrew McCallum¹

¹University of Massachusetts, Amherst MA

²Google Research, Mountain View CA

Workshop on Learning on Cores, Clusters and Clouds (LCCC)
Neural Information Processing Systems (NIPS) 2010



UMASS
AMHERST



Motivation

- Graphical models are used in a number of information extraction tasks
- Recently, models are getting larger and denser
 - Coreference Resolution [CULOTTA ET AL. NAACL 2007]
 - Relation Extraction [RIEDEL ET AL. EMNLP 2010, POON & DOMINGOS EMNLP 2009]
 - Joint Inference [FINKEL & MANNING. NAACL 2009, SINGH ET AL. ECML 2009]
- Inference is difficult, and approximations have been proposed
 - LP-Relaxations [MARTINS ET AL. EMNLP 2010]
 - Dual Decomposition [RUSH ET AL. EMNLP 2010]
 - MCMC-Based [MCCALLUM ET AL. NIPS 2009, POON ET AL. AAAI 2008]

Motivation

- Graphical models are used in a number of information extraction tasks
- Recently, models are getting larger and denser
 - Coreference Resolution [CULOTTA ET AL. NAACL 2007]
 - Relation Extraction [RIEDEL ET AL. EMNLP 2010, POON & DOMINGOS EMNLP 2009]
 - Joint Inference [FINKEL & MANNING. NAACL 2009, SINGH ET AL. ECML 2009]
- Inference is difficult, and approximations have been proposed
 - LP-Relaxations [MARTINS ET AL. EMNLP 2010]
 - Dual Decomposition [RUSH ET AL. EMNLP 2010]
 - MCMC-Based [MCCALLUM ET AL. NIPS 2009, POON ET AL. AAAI 2008]

Without parallelization, these approaches have restricted scalability

Motivation

Contributions:

- ➊ Distribute MAP Inference for a large, dense factor graph
 - 1 million variables, 250 machines
- ➋ Incorporate [sharding](#) as variables in the model

Outline

① Model and Inference

- Graphical Models

- MAP Inference

- Distributed Inference

② Cross-Document Coreference

- Coreference Problem

- Pairwise Model

- Inference and Distribution

③ Hierarchical Models

- Sub-Entities

- Super-Entities

④ Large-Scale Experiments

Factor Graphs

Represent distribution over variables Y using factors ψ .

$$p(Y = y) \propto \exp \sum_{y_c \subseteq y} \psi_c(y_c)$$

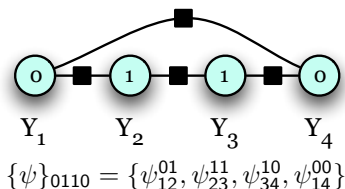
Note: Set of factors is different of every assignment $Y = y$ ($\{\psi\}_y$)

Factor Graphs

Represent distribution over variables Y using factors ψ .

$$p(Y = y) \propto \exp \sum_{y_c \subseteq Y} \psi_c(y_c)$$

Note: Set of factors is different of every assignment $Y = y$ ($\{\psi\}_y$)

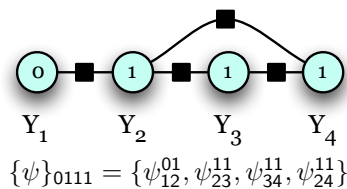
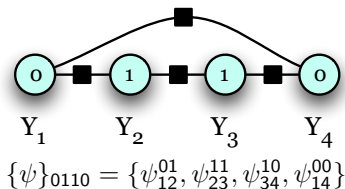


Factor Graphs

Represent distribution over variables Y using factors ψ .

$$p(Y = y) \propto \exp \sum_{y_c \subseteq Y} \psi_c(y_c)$$

Note: Set of factors is different of every assignment $Y = y$ ($\{\psi\}_y$)



MAP¹ Inference

We want to find the **best** configuration according to the model,

$$\begin{aligned}\hat{y} &= \arg \max_y p(Y = y) \\ &= \arg \max_y \exp \sum_{y_c \subseteq y} \psi_c(y_c)\end{aligned}$$

¹MAP = maximum a posteriori

MAP¹ Inference

We want to find the **best** configuration according to the model,

$$\begin{aligned}\hat{y} &= \arg \max_y p(Y = y) \\ &= \arg \max_y \exp \sum_{y_c \subseteq y} \psi_c(y_c)\end{aligned}$$

Computational bottlenecks:

- ① Space of Y is usually enormous (exponential)
- ② Even evaluating $\sum_{y_c \subseteq y} \psi_c(y_c)$ for each y may be polynomial

¹MAP = maximum a posteriori

MCMC for MAP Inference

Initial Configuration $y = y_0$

for (num_samples):

① **Propose** a change to y to get configuration y'
(Usually a *small* change)

② Acceptance probability: $\alpha(y, y') = \min \left(1, \left(\frac{p(y')}{p(y)} \right)^{1/t} \right)$
(Only involve computations local to the change)

③ if Toss(α): **Accept** the change, $y = y'$

return y

MCMC for MAP Inference

Initial Configuration $y = y_0$

for (num_samples):

① **Propose** a change to y to get configuration y'
(Usually a *small* change)

② Acceptance probability: $\alpha(y, y') = \min \left(1, \left(\frac{p(y')}{p(y)} \right)^{1/t} \right)$
(Only involve computations local to the change)

③ if Toss(α): **Accept** the change, $y = y'$

return y

$$\frac{p(y')}{p(y)} = \exp \left\{ \sum_{y'_c \subseteq y'} \psi_c(y'_c) - \sum_{y_c \subseteq y} \psi_c(y_c) \right\}$$

Mutually Exclusive Proposals

Let $\{\psi\}_y^{y'}$ be the set of factors used to evaluate a proposal $y \rightarrow y'$

$$\text{i.e. } \{\psi\}_y^{y'} = (\{\psi\}_y \cup \{\psi\}_{y'}) - (\{\psi\}_y \cap \{\psi\}_{y'})$$

Consider two proposals $y \rightarrow y_a$ and $y \rightarrow y_b$ such that,

$$\{\psi\}_y^{y_a} \cap \{\psi\}_y^{y_b} = \{\}$$

Completely different set of factors are required to evaluate these proposals.

Mutually Exclusive Proposals

Let $\{\psi\}_y^{y'}$ be the set of factors used to evaluate a proposal $y \rightarrow y'$

$$\text{i.e. } \{\psi\}_y^{y'} = (\{\psi\}_y \cup \{\psi\}_{y'}) - (\{\psi\}_y \cap \{\psi\}_{y'})$$

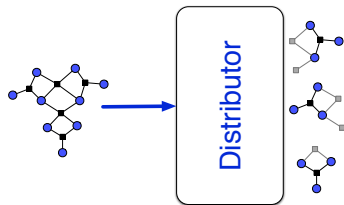
Consider two proposals $y \rightarrow y_a$ and $y \rightarrow y_b$ such that,

$$\{\psi\}_y^{y_a} \cap \{\psi\}_y^{y_b} = \{\}$$

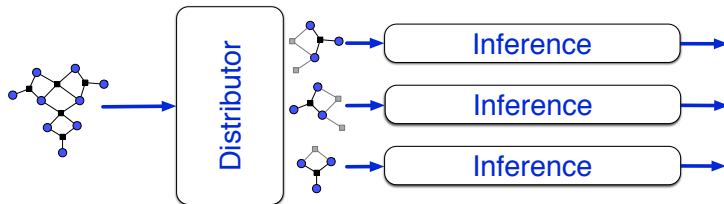
Completely different set of factors are required to evaluate these proposals.

These two proposals can be evaluated (and accepted) in parallel.

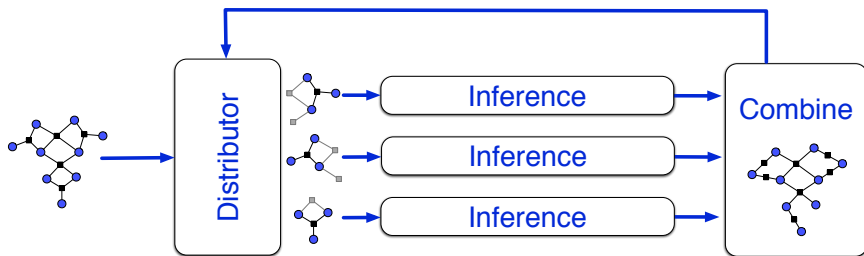
Distributed Inference



Distributed Inference



Distributed Inference



Outline

1 Model and Inference

- Graphical Models

- MAP Inference

- Distributed Inference

2 Cross-Document Coreference

- Coreference Problem

- Pairwise Model

- Inference and Distribution

3 Hierarchical Models

- Sub-Entities

- Super-Entities

4 Large-Scale Experiments

Coreference Problem

... The Physiological Basis of Politics," by **Kevin B. Smith**, Douglas Oxley, Matthew Hibbing...

...during the late 60's and early 70's, **Kevin Smith** worked with several local...

...the term hip-hop is attributed to **Lovebug Starski**. What does it actually mean...

The filmmaker **Kevin Smith** returns to the role of Silent Bob...

Nothing could be more irrelevant to **Kevin Smith's** audacious "Dogma" than ticking off...

Firefighter **Kevin Smith** spent almost 20 years preparing for Sept. 11. When he...

Like Back in 2008, the Lions drafted **Kevin Smith**, even though Smith was badly...

...shorthanded backfield in the wake of **Kevin Smith's** knee injury, and the addition of Haynesworth...

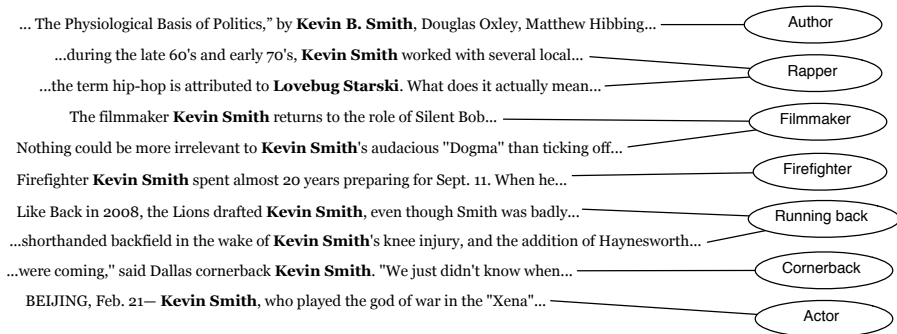
...were coming," said Dallas cornerback **Kevin Smith**. "We just didn't know when...

BELJING, Feb. 21— **Kevin Smith**, who played the god of war in the "Xena" ...

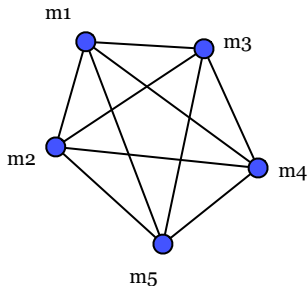
Coreference Problem



Coreference Problem



Input Features



Define similarity between mentions, $\phi : \mathcal{M}^2 \rightarrow \mathcal{R}$

- $\phi(m_i, m_j) > 0$: m_i, m_j are similar
- $\phi(m_i, m_j) < 0$: m_i, m_j are dissimilar

We use cosine similarity of the context bag of words:

$$\phi(m_i, m_j) = \text{cosSim}(\{c\}_i, \{c\}_j) - b$$

Graphical Model

The random variables in our model are entities (E) and mentions (M)

Graphical Model

The random variables in our model are entities (E) and mentions (M)
 For any assignment to these entities ($E = e$), we define the model score:

$$p(E = e) \propto \exp \left\{ \sum_{m_i \sim m_j} \psi_a(m_i, m_j) + \sum_{m_i \approx m_j} \psi_r(m_i, m_j) \right\}$$

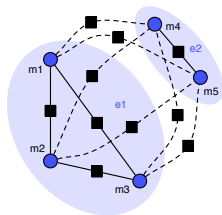
where $\psi_a(m_i, m_j) = w_a \phi(m_i, m_j)$, and
 $\psi_r(m_i, m_j) = -w_r \phi(m_i, m_j)$

Graphical Model

The random variables in our model are entities (E) and mentions (M)
For any assignment to these entities ($E = e$), we define the model score:

$$p(E = e) \propto \exp \left\{ \sum_{m_i \sim m_j} \psi_a(m_i, m_j) + \sum_{m_i \approx m_j} \psi_r(m_i, m_j) \right\}$$

where $\psi_a(m_i, m_j) = w_a \phi(m_i, m_j)$, and
 $\psi_r(m_i, m_j) = -w_r \phi(m_i, m_j)$



For the following configuration,

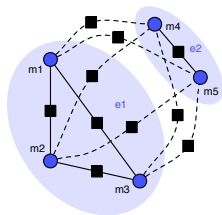
$$p(e_1, e_2) \propto \exp \left\{ \begin{aligned} &w_a (\phi_{12} + \phi_{13} + \phi_{23} + \phi_{45}) \\ &- w_r (\phi_{15} + \phi_{25} + \phi_{35} \\ &\quad + \phi_{14} + \phi_{24} + \phi_{34}) \end{aligned} \right\}$$

Graphical Model

The random variables in our model are entities (E) and mentions (M)
 For any assignment to these entities ($E = e$), we define the model score:

$$p(E = e) \propto \exp \left\{ \sum_{m_i \sim m_j} \psi_a(m_i, m_j) + \sum_{m_i \approx m_j} \psi_r(m_i, m_j) \right\}$$

where $\psi_a(m_i, m_j) = w_a \phi(m_i, m_j)$, and
 $\psi_r(m_i, m_j) = -w_r \phi(m_i, m_j)$

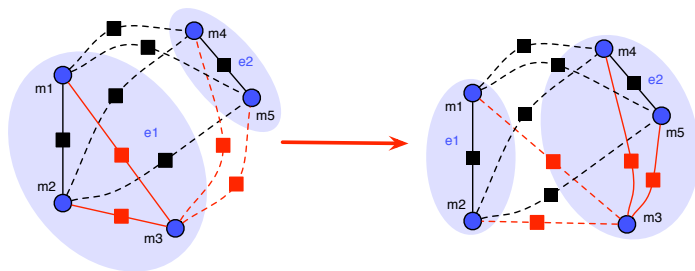


For the following configuration,

$$p(e_1, e_2) \propto \exp \left\{ \begin{aligned} &w_a (\phi_{12} + \phi_{13} + \phi_{23} + \phi_{45}) \\ &- w_r (\phi_{15} + \phi_{25} + \phi_{35} \\ &\quad + \phi_{14} + \phi_{24} + \phi_{34}) \end{aligned} \right\}$$

- ① Space of E is Bell Number(n) in number of mentions
- ② Evaluating model score for each $E = e$ is $O(n^2)$

MCMC for MAP Inference

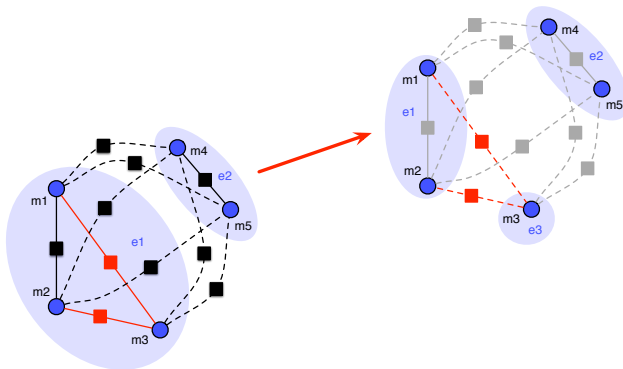


$$p(e) \propto \exp\{w_a(\phi_{12} + \phi_{13} + \phi_{23} + \phi_{45}) - w_r(\phi_{15} + \phi_{25} + \phi_{35} + \phi_{14} + \phi_{24} + \phi_{34})\}$$

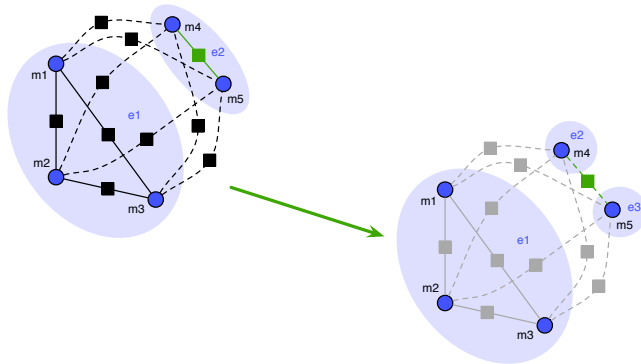
$$p(\acute{e}) \propto \exp\{w_a(\phi_{12} + \phi_{34} + \phi_{35} + \phi_{45}) - w_r(\phi_{15} + \phi_{25} + \phi_{13} + \phi_{14} + \phi_{24} + \phi_{23})\}$$

$$\log \frac{p(\acute{e})}{p(e)} = w_a(\phi_{34} + \phi_{35} - \phi_{13} - \phi_{23}) - w_r(\phi_{13} + \phi_{23} - \phi_{34} - \phi_{35})$$

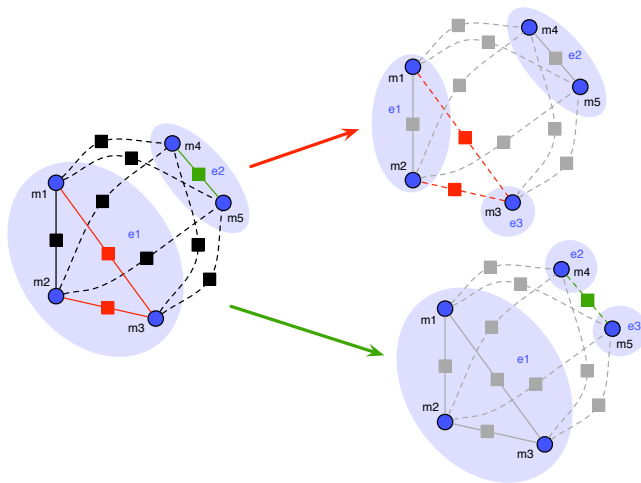
Mutually Exclusive Proposals



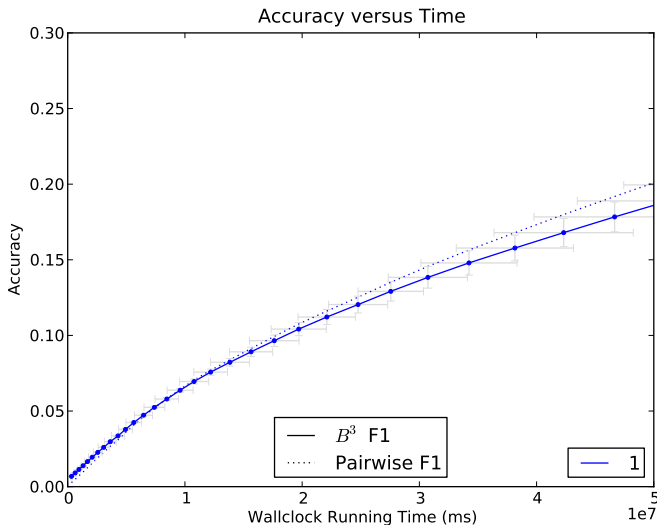
Mutually Exclusive Proposals



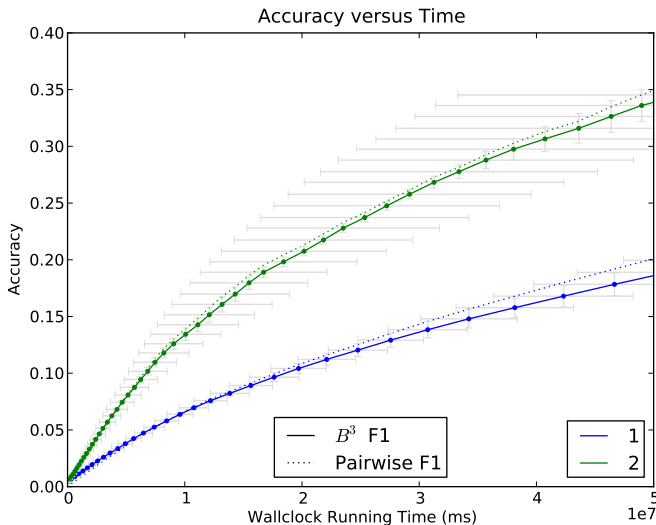
Mutually Exclusive Proposals



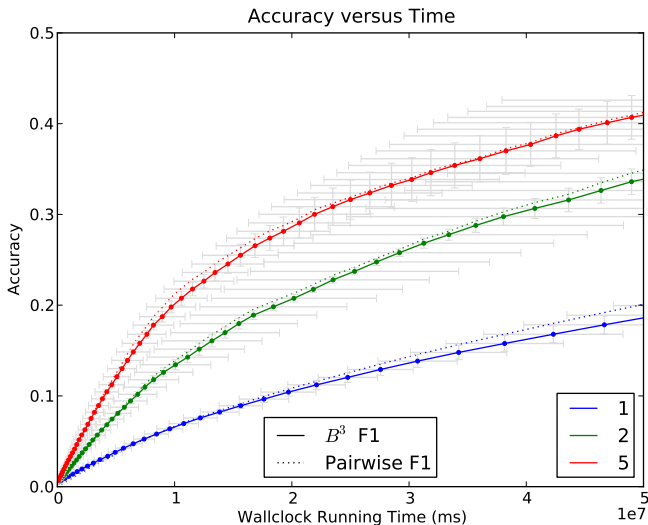
Results



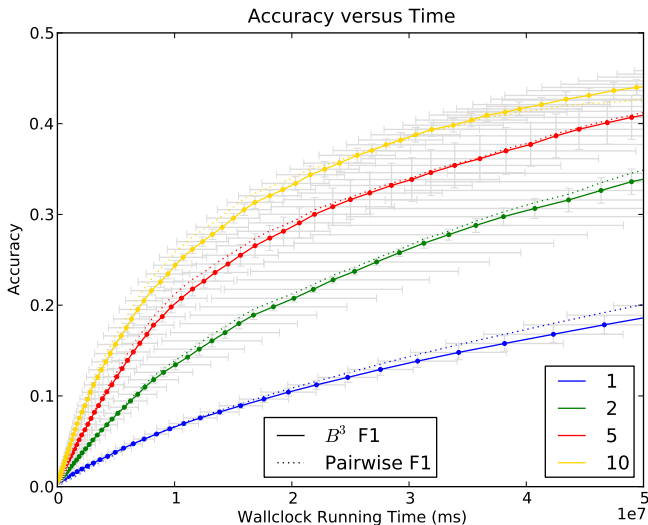
Results



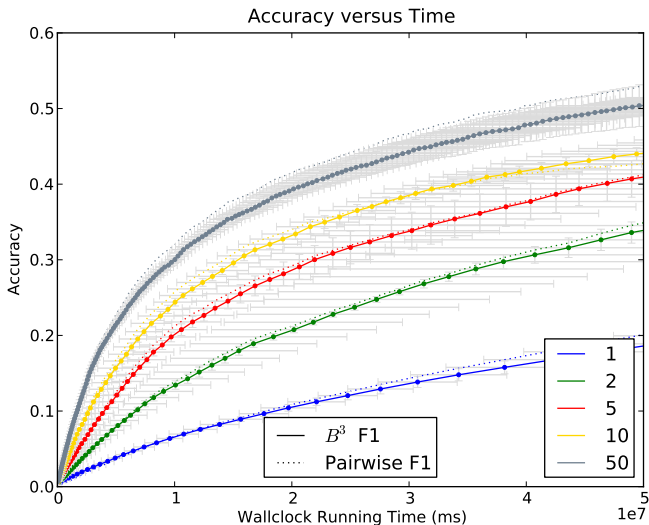
Results



Results



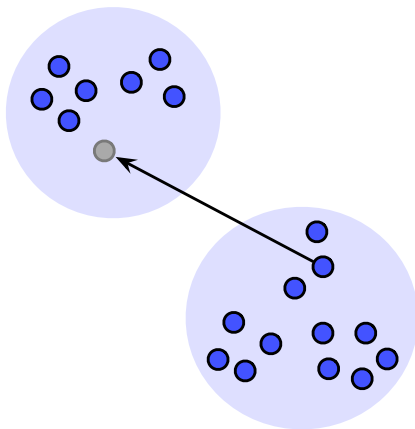
Results



Outline

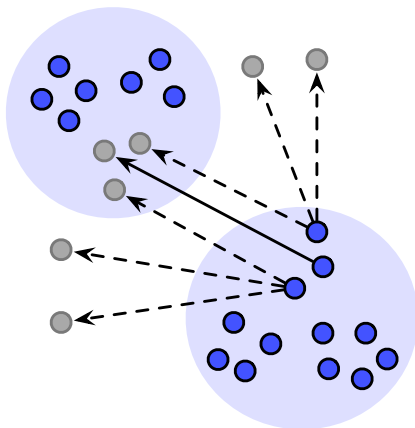
- ① Model and Inference
 - Graphical Models
 - MAP Inference
 - Distributed Inference
- ② Cross-Document Coreference
 - Coreference Problem
 - Pairwise Model
 - Inference and Distribution
- ③ Hierarchical Models
 - Sub-Entities
 - Super-Entities
- ④ Large-Scale Experiments

Sub-Entities



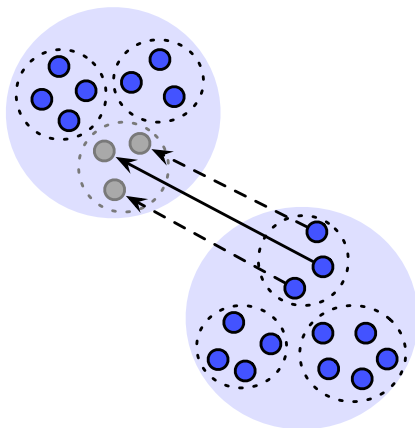
- Consider an **accepted** move for a mention

Sub-Entities



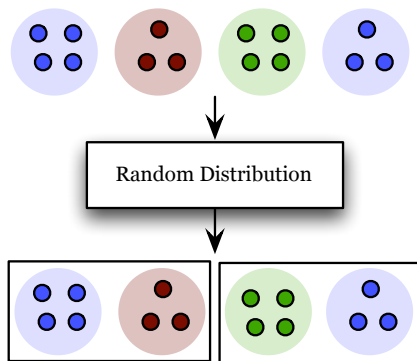
- Ideally, *similar* mentions should also move to the same entity
- Default proposal function does not utilize this
- *Good* proposals become more rare with larger datasets

Sub-Entities



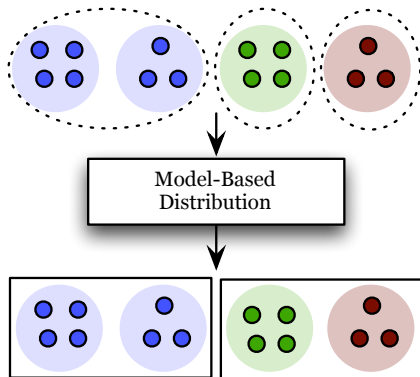
- Include **Sub-Entity** variables
- Model score is used to sample sub-entity variables
- Propose moves of mentions in a sub-entity simultaneously

Super-Entities



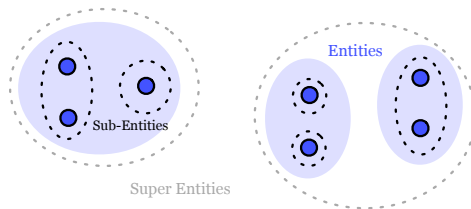
- Random distribution may not assign *similar* entities to the same machine
- Probability that similar entities will be assigned to the same machine is small

Super-Entities



- Augment model with **Super-Entities** variables
- Entities in the same super-entity are assigned the same machine
- Model score is used to sample super-entity variables

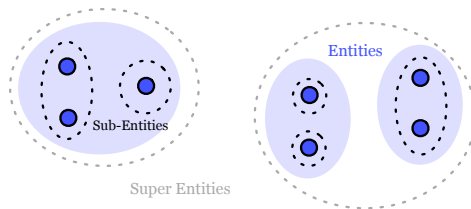
Hierarchical Representation



• Factors

- Affinity factors between **mentions** **sub-entities** in the same **sub-entities** **entities** **super-entities**
- Repulsion factors are similarly symmetric across levels

Hierarchical Representation

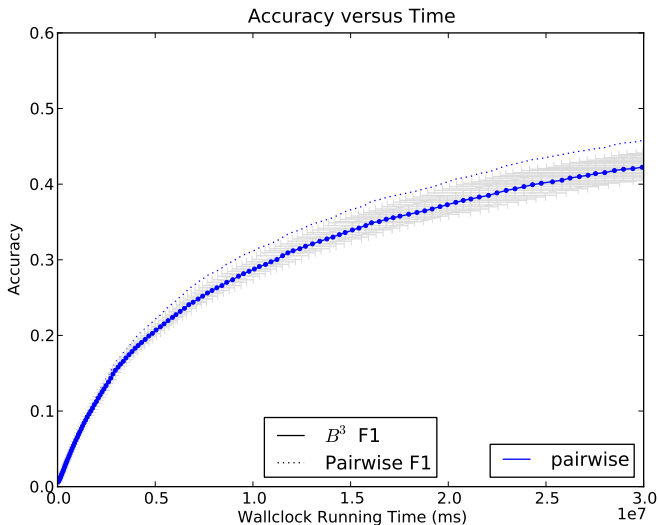


- **Factors**

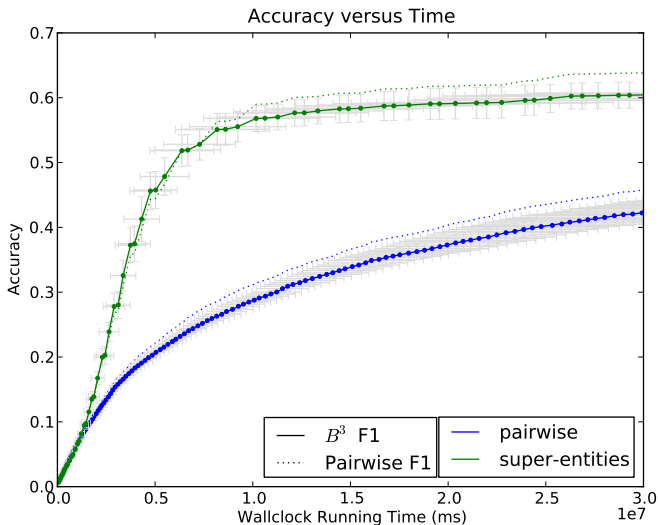
- Affinity factors between **mentions** **sub-entities** in the same **entities** **sub-entities** **entities** **super-entities**
- Repulsion factors are similarly symmetric across levels

- **Sampling:** Fix variables of two levels, sample the remaining level

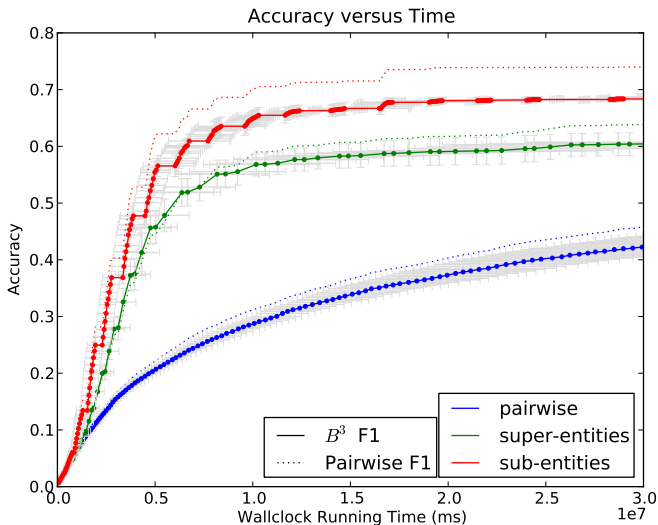
Evaluation



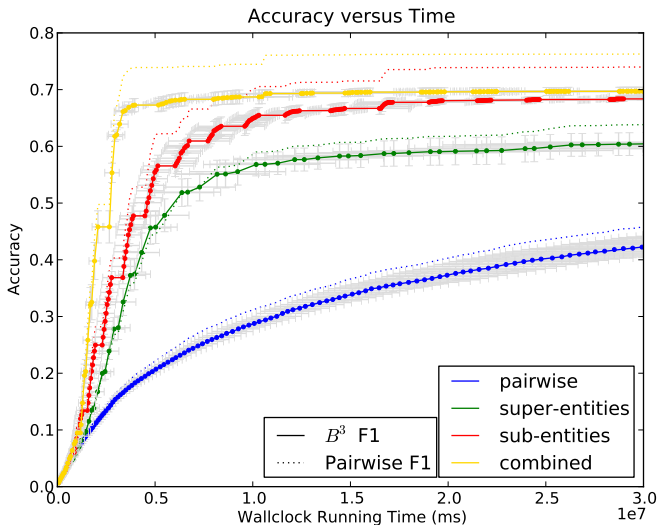
Evaluation



Evaluation



Evaluation



Outline

① Model and Inference

Graphical Models

MAP Inference

Distributed Inference

② Cross-Document Coreference

Coreference Problem

Pairwise Model

Inference and Distribution

③ Hierarchical Models

Sub-Entities

Super-Entities

④ Large-Scale Experiments

Preliminary Large-Scale Experiments

Data

- *New York Times Annotated Corpus* [SANDHOU LDC 2008]
20 years of articles (1987-2007)
- prune rare names (<1000): ~ 1 million **person name** mentions

Preliminary Large-Scale Experiments

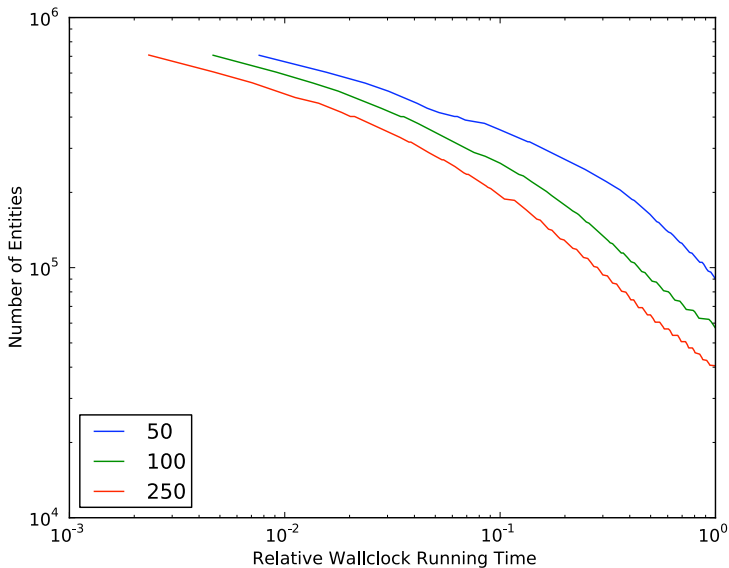
Data

- *New York Times Annotated Corpus* [SANDHUS LDC 2008]
20 years of articles (1987-2007)
- prune rare names (<1000): ~ 1 million **person name** mentions

Evaluation

- Automated labels are too noisy for evaluation
- Instead, we estimate the **speed of inference**
 - trust the model to accept good proposals
 - observe the number of predicted entities

Speed of Inference



Related Work

- GraphLab [\[LOW ET AL. UAI 2010\]](#)
 - how do we represent dynamic graphs
 - how do we represent hierarchical models

Related Work

- GraphLab [LOW ET AL. UAI 2010]
 - how do we represent dynamic graphs
 - how do we represent hierarchical models
- Graph Splashing [GONZALEZ ET AL. UAI 2009]
 - graph structure changes with every configuration
 - BP messages are enormous for exponential-domain variables

Related Work

- GraphLab [LOW ET AL. UAI 2010]
 - how do we represent dynamic graphs
 - how do we represent hierarchical models
- Graph Splashing [GONZALEZ ET AL. UAI 2009]
 - graph structure changes with every configuration
 - BP messages are enormous for exponential-domain variables
- Topic Models [SMOLA & NARAYANMURTHY. VLDB 2010, ASUNCION ET AL. NIPS 2009]
 - restrictions since they are calculating probabilities
 - we allow non-random distribution and customized proposals

Conclusions

- 1 propose **distributed inference** for graphical models
- 2 enable distributed **cross-document coreference**
- 3 improve sharding with latent **hierarchical** variables
- 4 demonstrate utility on **large** datasets

Conclusions

- 1 propose [distributed inference](#) for graphical models
- 2 enable distributed [cross-document coreference](#)
- 3 improve sharding with latent [hierarchical](#) variables
- 4 demonstrate utility on [large](#) datasets

Future Work:

- more [scalability](#) experiments
- study [mixing](#) and [convergence](#) properties
- add more expressive [factors](#)
- [supervision](#): labeled data, noisy evidences

Thanks!

Sameer Singh

sameer@cs.umass.edu

Amarnag Subramanya

asubram@google.com

Fernando Pereira

pereira@google.com

Andrew McCallum

mccallum@cs.umass.edu



UMASS
AMHERST

