

A Discriminative Latent Variable Model for Online Clustering

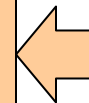
Rajhans Samdani, **Kai-Wei Chang**, Dan Roth
Department of Computer Science
University of Illinois at Urbana-Champaign

Motivating Example: Coreference

- Coreference resolution: cluster denotative noun phrases (*mentions*) in a document based on underlying entities

[Bill Clinton], recently elected as the [President of the USA], has been invited by the [Russian President], [Vladimir Putin], to visit [Russia]. [President Clinton] said that [he] looks forward to strengthening ties between [USA] and [Russia].

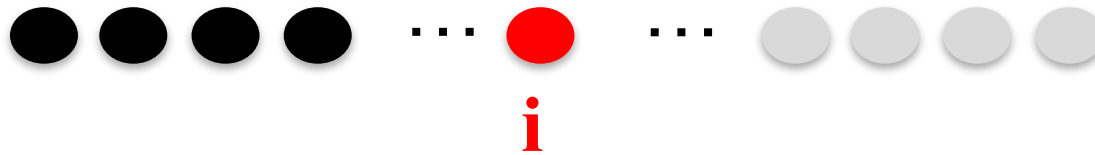
- The task: learning a clustering function from training data
 - Used expressive features between mention pairs (e.g. string similarity).
 - Learn a similarly metric between mentions.
 - Cluster mentions based on the metric.
- The mention arrives in a left-to-right order



Learning this metric using a joint distribution over clustering

Online Clustering

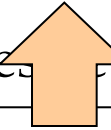
- Online clustering: items arrive in a given order



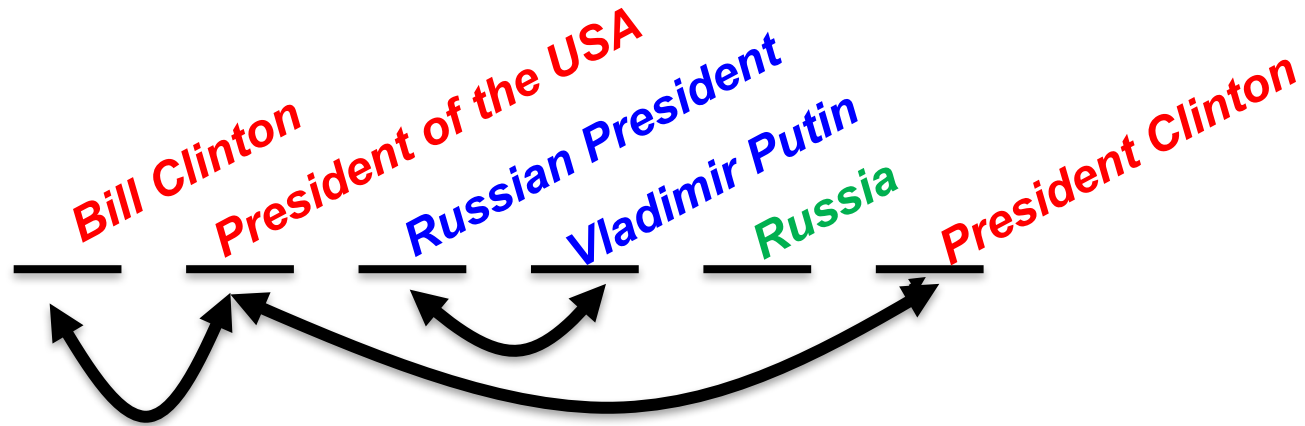
- **Motivating property:** cluster item **i** with no access to future items on the right, only the previous items to the left
- This setting is **general** and is natural in many tasks.
 - E.g., cluster posts in a forum, cluster network attack
- An online clustering algorithm is likely to be more **efficient** than a batch algorithm under such setting.

Greedy Best-Left-Link Clustering

[Bill Clinton], recently elected as the [President of the USA], has been invited by the [Russian President], [Vladimir Putin], to visit [Russia]. [President Clinton] said that [he] looks forward to strengthening ties between [USA] and [Russia].



- Best-Left-Linking decoding: (Bengtson and Roth '08).



- A Naïve way to learn the model:
 - decouple (i) learning a similarity metric between pairs; (ii) hard clustering of mentions using this metric.

Our Contribution

- A novel discriminative latent variable model, **Latent Left-Linking Model (L^3M)**, for jointly learning metric and clustering, that outperforms existing models
- Training the pair-wise similarity metric for clustering using a latent variable structured prediction
- Relaxing the single best-link: consider a distribution over links
- Efficient learning algorithm that decomposes over individual items in the training stream

Outline

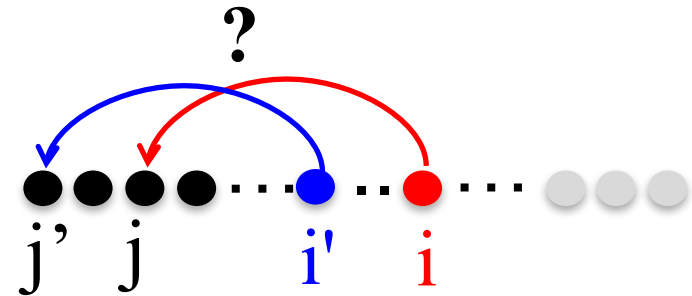
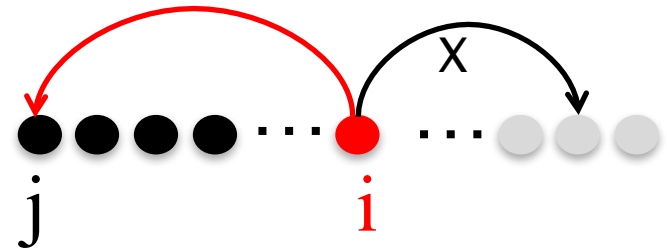
- Motivation, examples and problem description
- Latent Left-Linking Model (L^3M)
 - Likelihood computation
 - Inference
 - Role of temperature
 - Alternate latent variable perspective
- Learning
 - Discriminative structured prediction learning view
 - Stochastic gradient based decomposed learning
- Empirical study



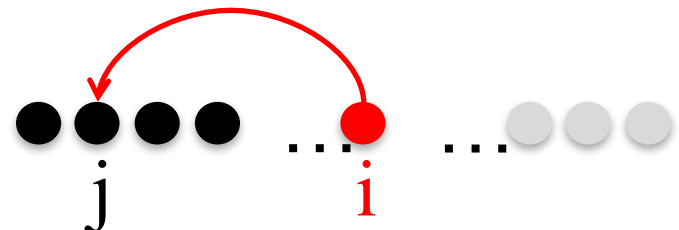
Latent Left-Linking Model (L³M)

Modeling Axioms

- Each item can link only to some item on its left (creating a *left-link*)
- Event **i** linking to **j** is ? Of **i'** linking to **j'**
- Probability of **i** linking to **j**
 $\Pr[\mathbf{j} \tilde{\mathbf{A}} \mathbf{i}] / \exp(\mathbf{w} \cdot \tilde{\mathbf{A}}(\mathbf{i}, \mathbf{j}) / \circ)$
 - $\circ \in [0,1]$ Is a temperature-like user-tuned parameter



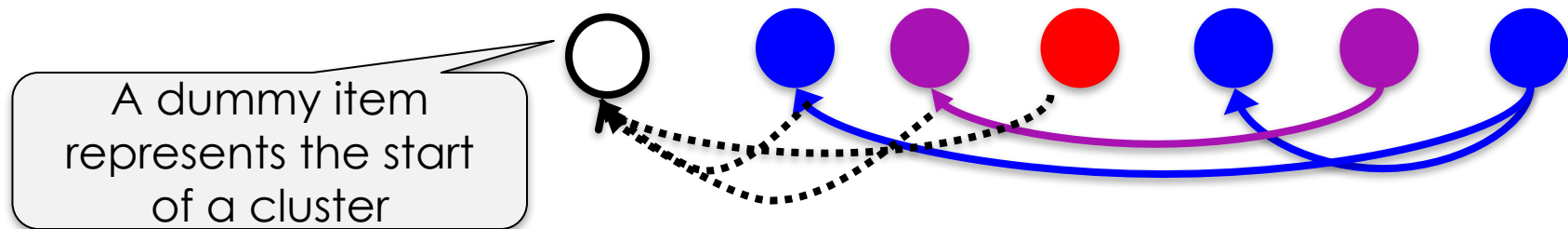
$$\exp(\mathbf{w} \cdot \tilde{\mathbf{A}}(\mathbf{i}, \mathbf{j}) / \circ)$$



L³M: Likelihood of Clustering

- **C** is a clustering of data stream **d**

□ **C**(**i**, **j**) = 1 if **i** and **j** co-clustered else 0



- Prob. of **C** : multiply prob. of items connecting as per **C**

$$\Pr[\mathbf{C}; \mathbf{w}] = \prod_i \Pr[\mathbf{i}, \mathbf{C}; \mathbf{w}] = \prod_i (\sum_{\mathbf{j} < \mathbf{i}} \Pr[\mathbf{j} \tilde{\mathbf{A}} \mathbf{i}] \mathbf{C}(\mathbf{i}, \mathbf{j}))$$

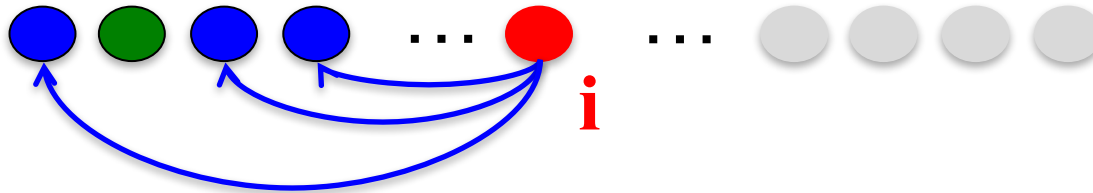
$$/ \prod_i (\sum_{\mathbf{j} < \mathbf{i}} \exp(\mathbf{w} \notin \mathbf{A}(\mathbf{i}, \mathbf{j}) / \circ) \mathbf{C}(\mathbf{i}, \mathbf{j}))$$

- Partition/normalization function efficient to compute

$$\mathbf{Z}_d(\mathbf{w}) = \prod_i (\sum_{\mathbf{j} < \mathbf{i}} \exp(\mathbf{w} \notin \mathbf{A}(\mathbf{i}, \mathbf{j}) / \circ))$$

L³M: Greedy Inference/Clustering

- Sequential arrival of items:



Prob. of **i** connecting to previously formed cluster **c**

= sum of probs. of **i** connecting to items in **c**

$$\Pr[\mathbf{c}^- \mathbf{i}] = \sum_{j \in \mathbf{c}} \Pr[\mathbf{j} \tilde{\mathbf{A}} \mathbf{i}; \mathbf{w}] / \sum_{j \in \mathbf{c}} \exp(\mathbf{w} \cdot \mathbf{A}(\mathbf{i}, \mathbf{j}) / \sigma)$$

- Greedy clustering:

- Compute $\mathbf{c}^* = \operatorname{argmax}_{\mathbf{c}} \Pr[\mathbf{c}^- \mathbf{i}]$
- Connect **i** to \mathbf{c}^* if $\Pr[\mathbf{c}^*^- \mathbf{i}] > \mathbf{t}$ (**threshold**) otherwise **i** starts a new cluster
- May not yield the most likely clustering

Inference: role of temperature \circ

- Prob. of i connecting to previous item j

$$\Pr[j \tilde{A} i] / \exp(\mathbf{w} \cdot \tilde{A}(i, j) / \circ)$$

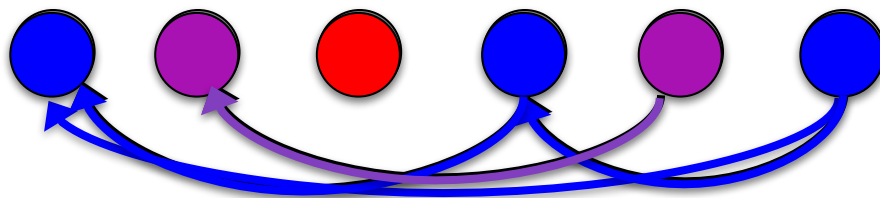
- \circ tunes the importance of high-scoring links
 - As \circ decreases from 1 to 0, high-scoring links become more important
 - For $\circ = 0$, $\Pr[j \tilde{A} i]$ is a Kronecker delta function centered on the argmax link (assuming no ties)

$$\Pr[\mathbf{c} \tilde{A} i] / \sum_j \exp(\mathbf{w} \cdot \tilde{A}(i, j) / \circ)$$

- For $\circ = 0$, clustering considers only the “best-left-link” and **greedy clustering is exact**

Latent Variables: Left-Linking Forests

- Left-linking forest, \mathbf{f} : the parent (arrow directions reversed) of each item on its left

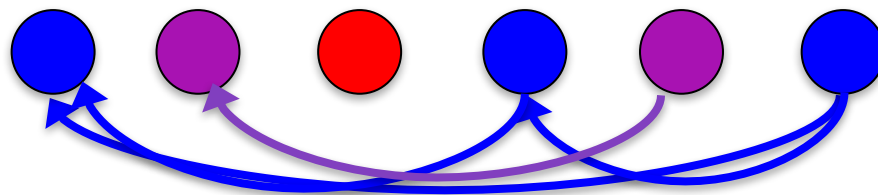


- Probability of forest \mathbf{f} based on sum of edge weights in \mathbf{f}

$$\Pr[\mathbf{f}; \mathbf{w}] / \exp(\sum_{(i,j) \in \mathbf{f}} w_{ij} \phi(\mathbf{A}(i,j)))$$

- L^3M : same as expressing the probability of \mathbf{C} as the sum of probabilities of all consistent (**latent**) Left-linking forests

$$\Pr[\mathbf{C}; \mathbf{w}] = \sum_{\mathbf{f} \in \mathcal{F}(\mathbf{C})} \Pr[\mathbf{f}; \mathbf{w}]$$



Outline

- Motivation, examples and problem description
- Latent Left-Linking Model
 - Inference
 - Role of temperature
 - Likelihood computation
 - Alternate latent variable perspective
- Learning
 - Discriminative structured prediction learning view
 - Stochastic gradient based decomposed learning
- Empirical study



L³M: Likelihood-based Learning

- Learn \mathbf{w} from annotated clustering \mathbf{C}_d for data $d \in D$

- L³M: Learn \mathbf{w} via regularized neg. log-likelihood

$$LL(\mathbf{w}) = -\|\mathbf{w}\|^2 \quad \text{Regularization}$$

$$+ \sum_d \log Z_d(\mathbf{w}) \quad \text{Partition Function}$$

$$- \sum_d \sum_i \log \left(\sum_{j < i} \exp(\mathbf{w} \cdot \mathbf{A}(i, j)) / \sum_{j < i} \exp(\mathbf{w} \cdot \mathbf{A}(i, j)) \right) \mathbf{C}_d(i, j)$$

Un-normalized Probability


- Relation to other latent variable models:

- Learn by marginalizing underlying latent left-linking forests
- $\alpha=1$: Hidden Variable CRFs (Quattoni et al, 07)
- $\alpha=0$: Latent Structural SVMs (Yu and Joachims, 09)

Training Algorithms: Discussion

- The objective function $LL(\mathbf{w})$ is non-convex
- Can use **Concave-Convex Procedure (CCCP)** (Yuille and Rangarajan 03; Yu and Joachims, 09)
 - **Pros:** guaranteed to converge to a local minima (Sriperumbudur et al, 09)
 - **Cons:** requires entire data stream to compute single gradient update
- Online updates based on **Stochastic (sub-)gradient descent (SGD)**
 - Sub-gradient can be decomposed to a **per-item basis**
 - **Cons:** no theoretical guarantees for SGD with non-convex functions
 - **Pros:** can learn in an **online fashion**; Converge much **faster** than CCCP
 - Great empirical performance

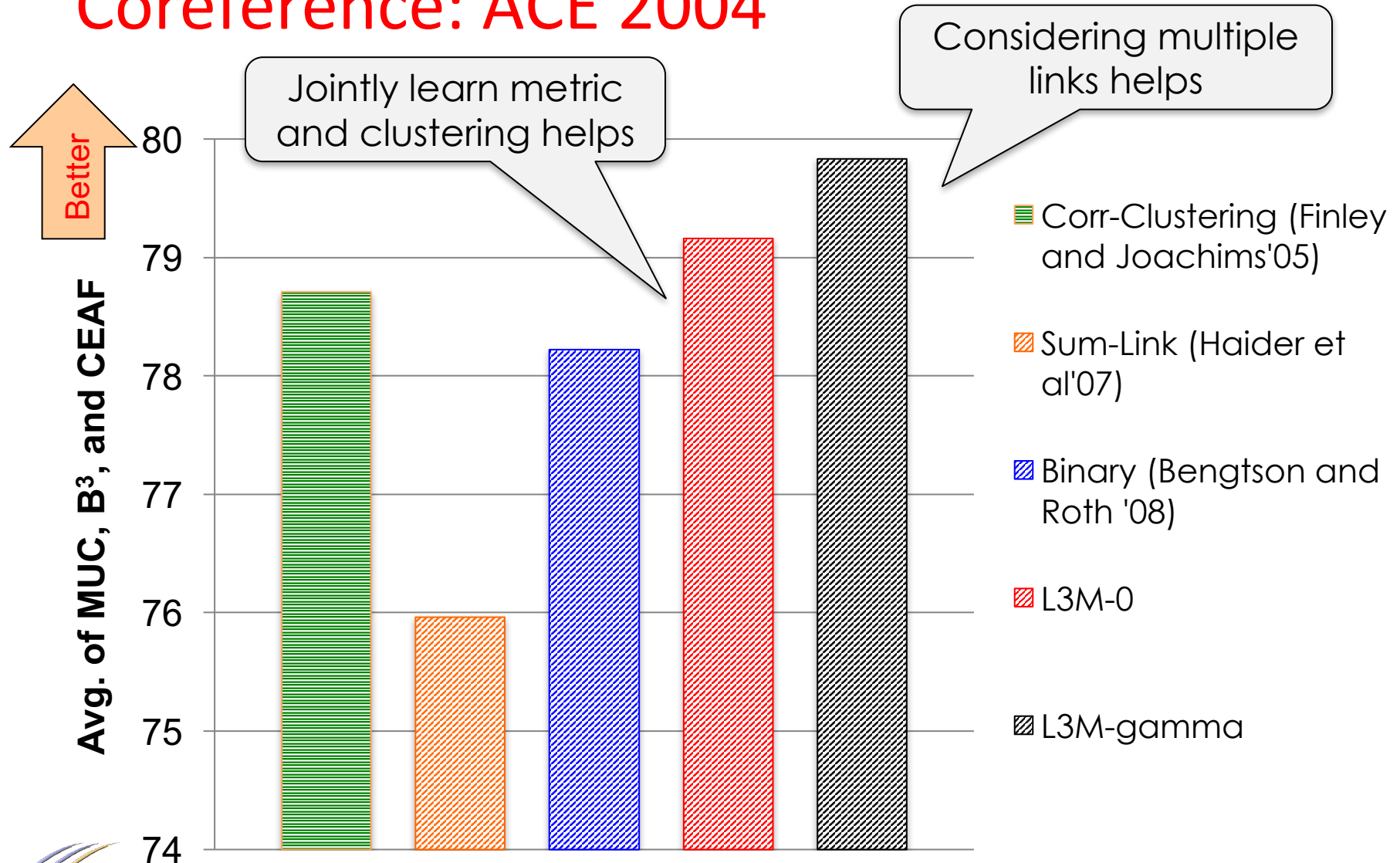
Outline

- Motivation, examples and problem description
- Latent Left-Linking Model
 - Inference
 - Role of temperature
 - Likelihood computation
 - Alternate latent variable perspective
- Learning
 - Discriminative structured prediction learning view
 - Stochastic gradient based decomposed learning
- Empirical study 

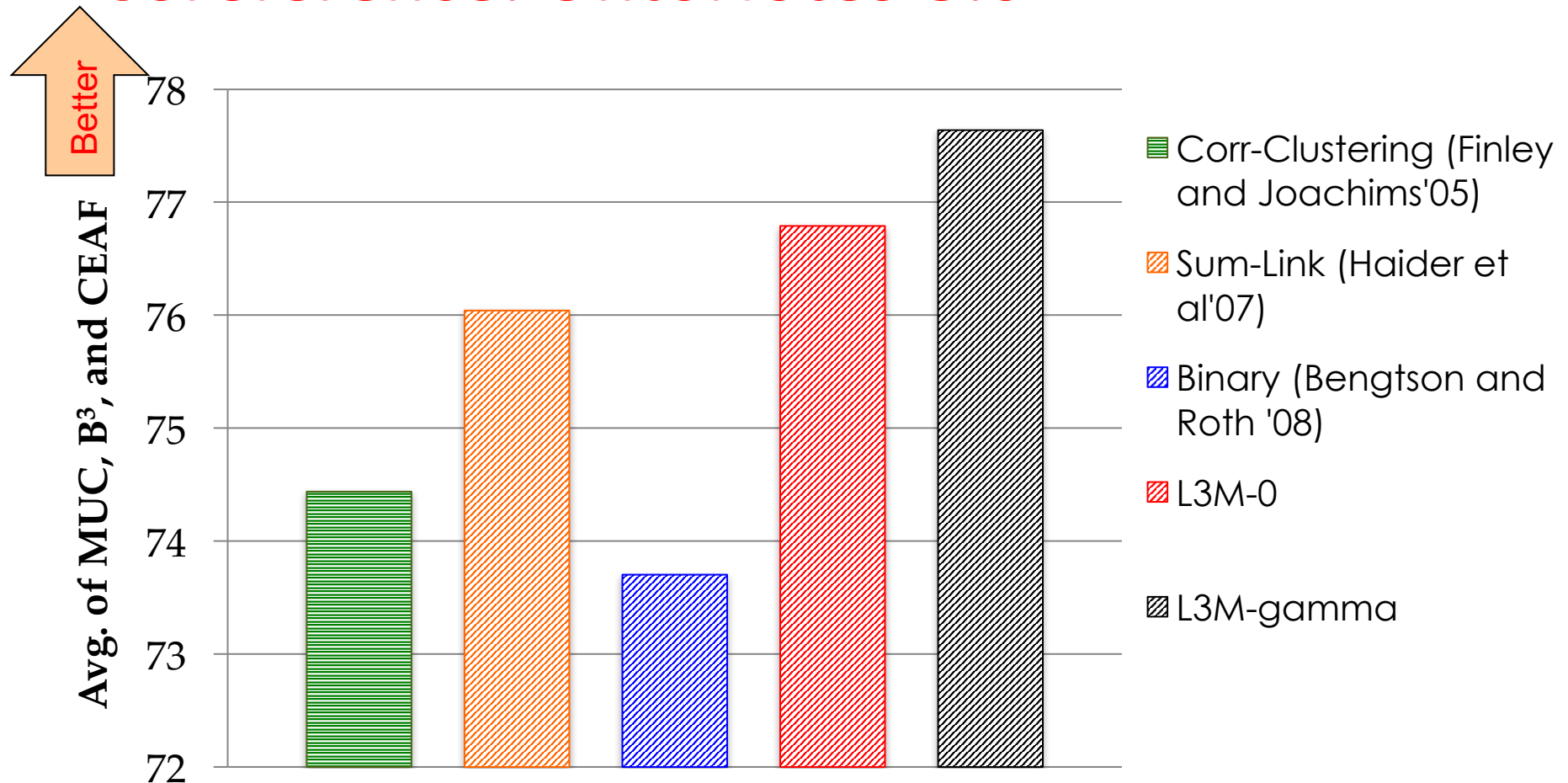
Experiment: Coreference Resolution

- Cluster denotative noun phrases called *mentions*
- Mentions follow a left-to-right order
- Features: mention distance, substring match, gender match, etc.
- Experiments on **ACE 2004** and **OntoNotes-5.0**.
- Report average of three popular coreference clustering evaluation metrics: MUC, B^3 , and CEAF

Coreference: ACE 2004



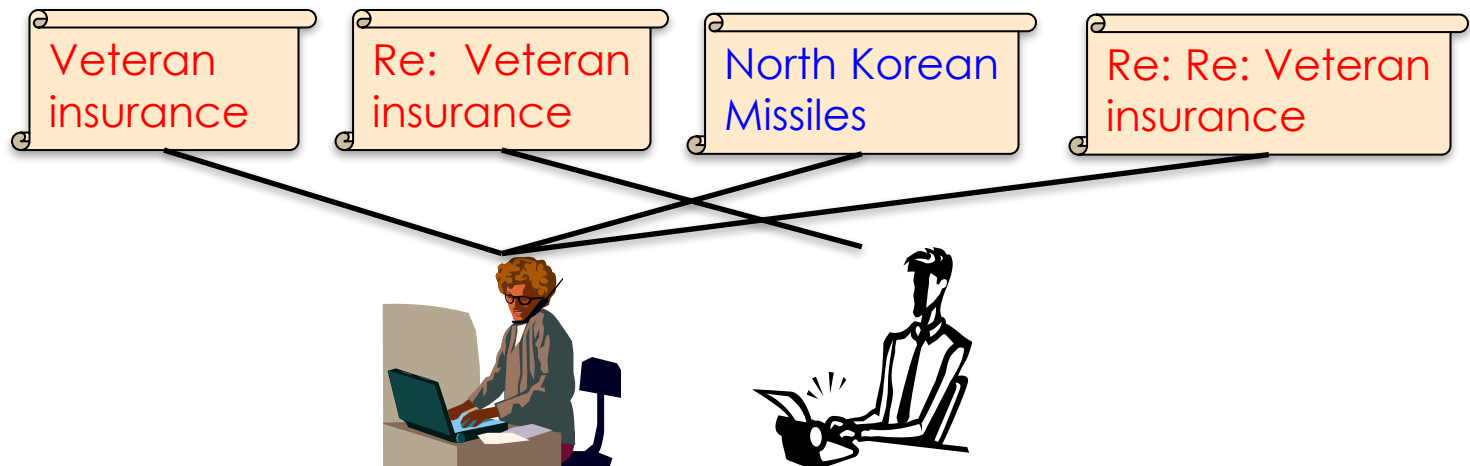
Coreference: OntoNotes-5.0



By incorporating with domain knowledge constraints, **L³M** achieves the **state of the art** performance on OntoNotes-5.0 (Chang et al. 13)

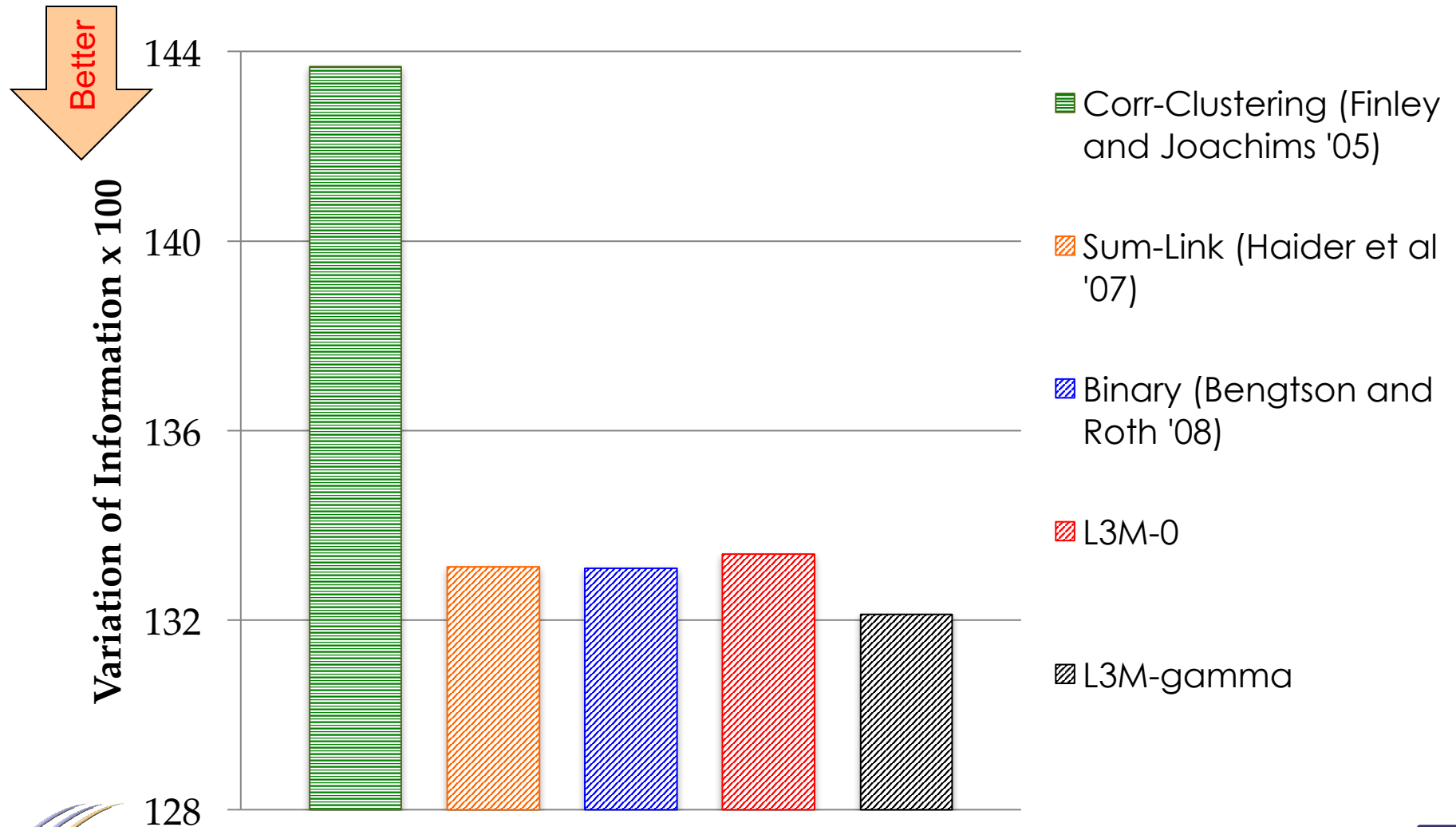
Experiments: Document Clustering

- Cluster the posts in a forum based on **authors** or **topics**.
- Dataset: discussions from www.militaryforum.com
- The posts in the forum arrive in **a time order**:

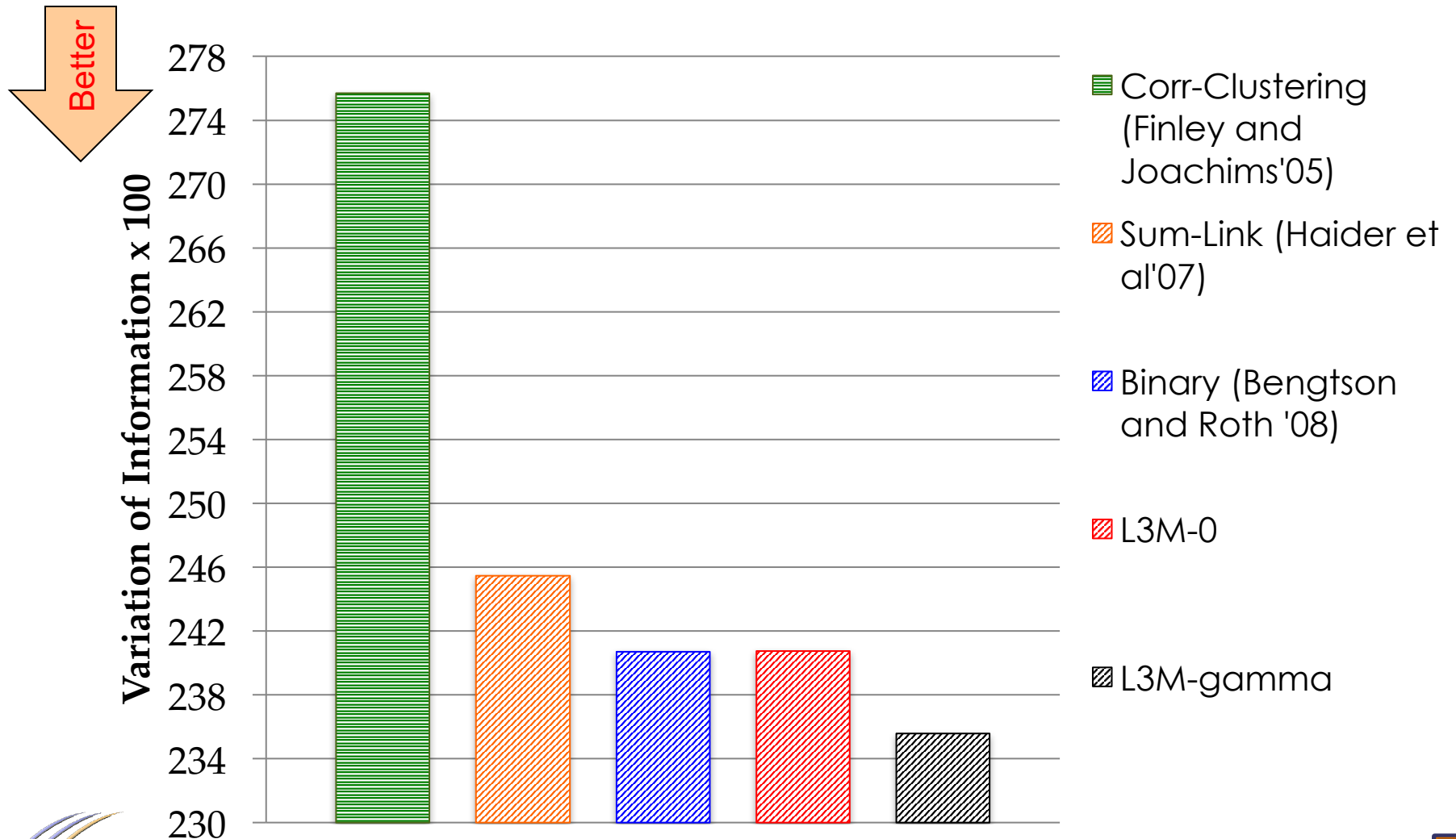


- Features: **common words**, **tf-idf similarity**, **time between arrival**
- Evaluate with Variation-of-Information (Meila, 07)

Author Based Clustering



Topic Based Clustering



Conclusions

- Latent Left-Linking Model
 - Principled probabilistic modeling for online clustering tasks
 - Marginalizes underlying latent link structures
 - Tuning \circ helps – considering multiple links helps
 - Efficient greedy inference
- SGD-based learning
 - Decompose learning into smaller gradient updates over individual items
 - Rapid convergence and high accuracy
- Solid empirical performance on problems with a natural streaming order