# Optimizing Sampling-based Entity Resolution over Streaming Documents

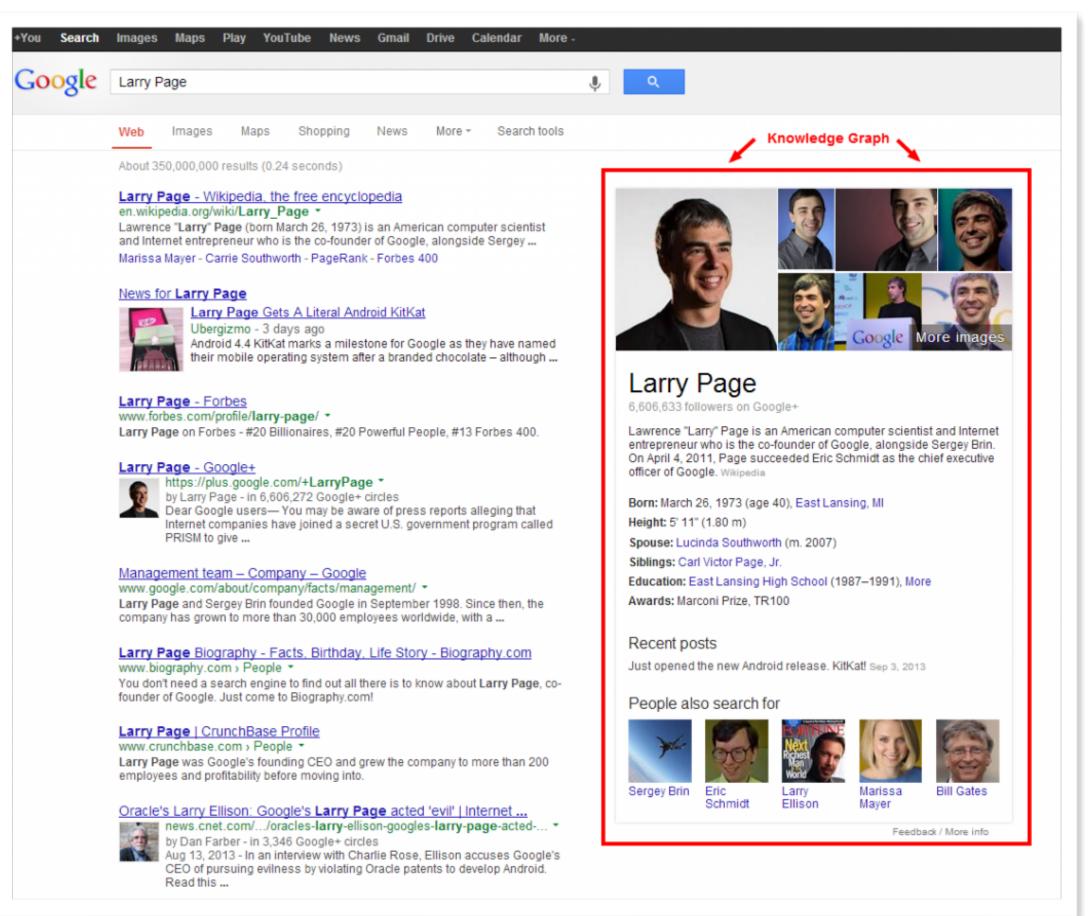**Christan Grant** **and Daisy Zhe Wang**
University of Florida

**SIAM BSA Workshop 2015**

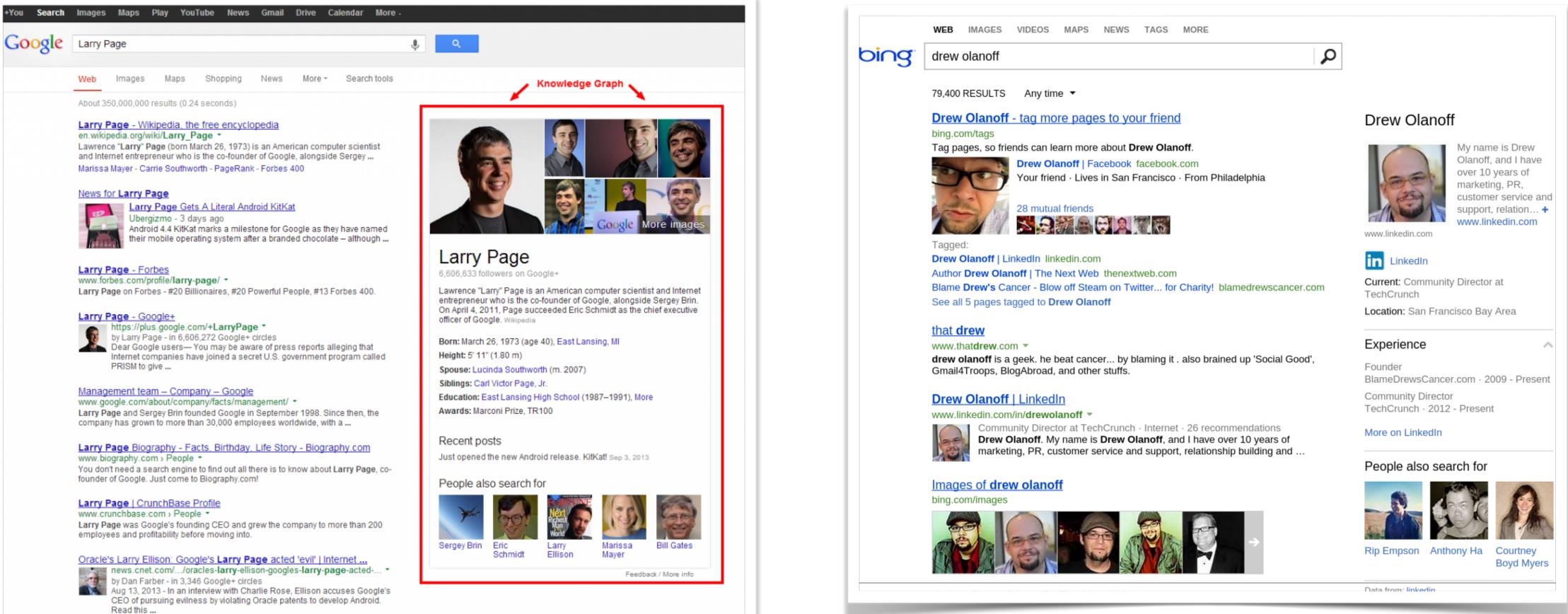Knowledge Bases are important structure for organizing and categorizing information.

# Knowledge Bases are important structure for organizing and categorizing information.

# Knowledge Bases are important structure for organizing and categorizing information.

- Many of these knowledge bases and new knowledge bases are bootstrapped using **Wikipedia**/Freebase.

- Many of these knowledge bases and new knowledge bases are bootstrapped using **Wikipedia**/Freebase.

- All Wikipedia information is based on facts from (reputable?) web sources.

# Society for Industrial and Applied Mathematics

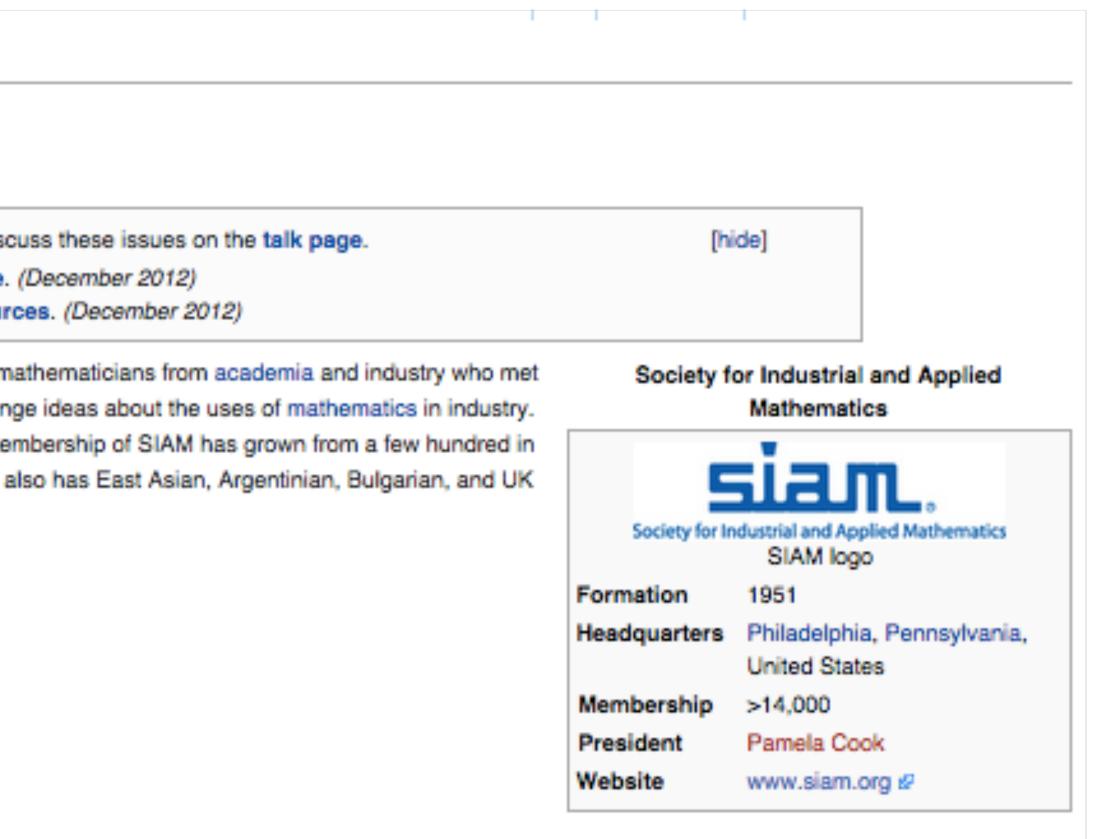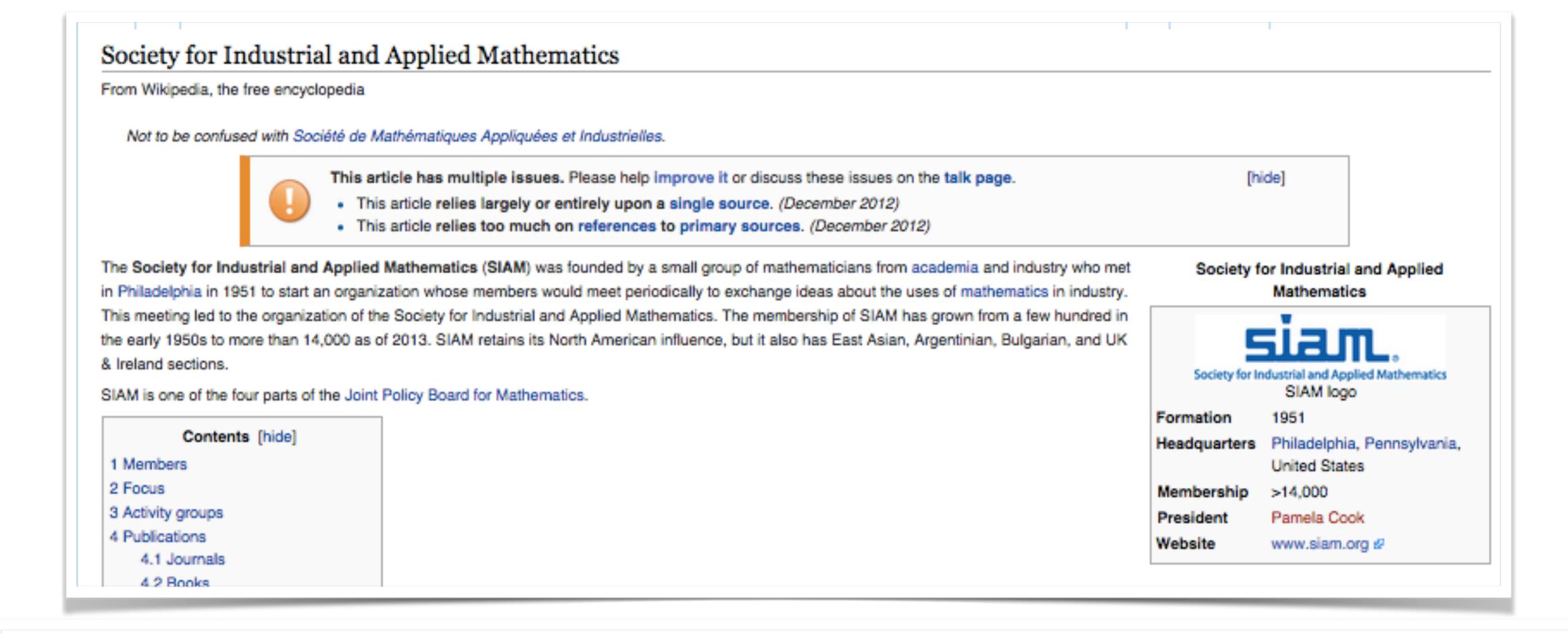*Not to be confused with Société de Mathématiques Appliquées et Industrielles.*

> ⚠️ **This article has multiple issues.** Please help **improve it** or discuss these issues on the **talk page.**  [hide]
> - This article **relies largely or entirely upon a single source.** *(December 2012)*
> - This article **relies too much on references to primary sources.** *(December 2012)*

The **Society for Industrial and Applied Mathematics (SIAM)** was founded by a small group of mathematicians from academia and industry who met in Philadelphia in 1951 to start an organization whose members would meet periodically to exchange ideas about the uses of mathematics in industry. This meeting led to the organization of the Society for Industrial and Applied Mathematics. The membership of SIAM has grown from a few hundred in the early 1950s to more than 14,000 as of 2013. SIAM retains its North American influence, but it also has East Asian, Argentinian, Bulgarian, and UK & Ireland sections.

SIAM is one of the four parts of the Joint Policy Board for Mathematics.

**Contents** [hide]

**Society for Industrial and Applied Mathematics**



SIAM logo

| | |
|---|---|
| Formation | 1951 |
| Headquarters | Philadelphia, Pennsylvania, United States |
| Membership | >14,000 |
| President | Pamela Cook |
| Website | www.siam.org |

# Society for Industrial and Applied Mathematics

*Not to be confused with Société de Mathématiques Appliquées et Industrielles.*

> **This article has multiple issues.** Please help **improve it** or discuss these issues on the **talk page.**     [hide]
> - This article **relies largely or entirely upon a single source.** *(December 2012)*
> - This article **relies too much on references to primary sources.** *(December 2012)*

The **Society for Industrial and Applied Mathematics (SIAM)** was founded by a small group of mathematicians from academia and industry who met in Philadelphia in 1951 to start an organization whose members would meet periodically to exchange ideas about the uses of mathematics in industry. This meeting led to the organization of the Society for Industrial and Applied Mathematics. The membership of SIAM has grown from a few hundred in the early 1950s to more than 14,000 as of 2013. SIAM retains its North American influence, but it also has East Asian, Argentinian, Bulgarian, and UK & Ireland sections.

SIAM is one of the four parts of the Joint Policy Board for Mathematics.

| **Contents** [hide] |
| --- |
| 1 Members |
| 2 Focus |
| 3 Activity groups |
| 4 Publications |
|     4.1 Journals |
|     4.2 Books |

**Society for Industrial and Applied Mathematics**

siam.
Society for Industrial and Applied Mathematics
SIAM logo

| | |
| --- | --- |
| **Formation** | 1951 |
| **Headquarters** | Philadelphia, Pennsylvania, United States |
| **Membership** | >14,000 |
| **President** | Pamela Cook |
| **Website** | www.siam.org |

## SIAM Fellows    [edit]

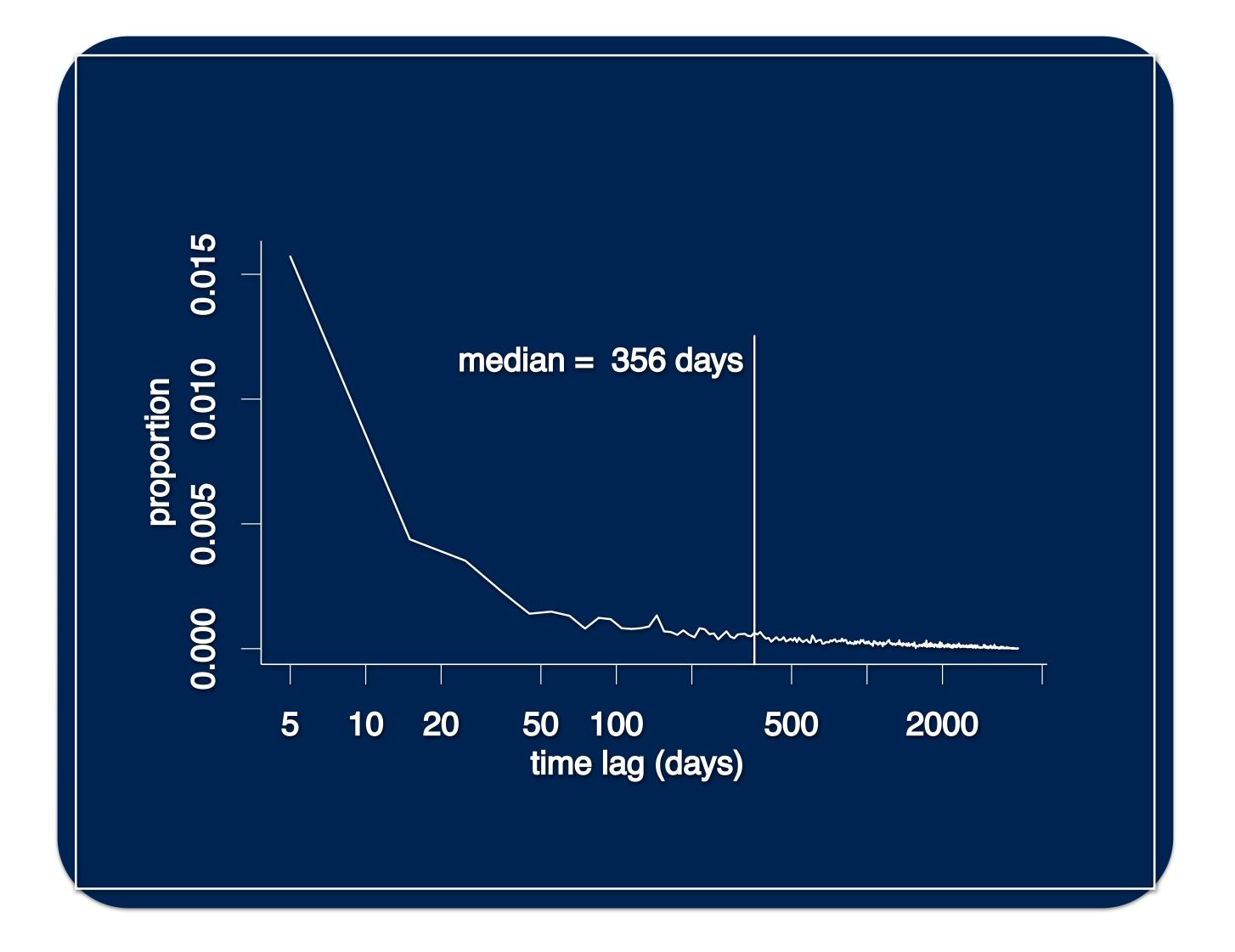- In 2009 SIAM instituted a Fellows program to recognize certain members who have made outstanding contributions to the fields SIAM serves[14]

# Society for Industrial and Applied Mathematics

*Not to be confused with Société de Mathématiques Appliquées et Industrielles.*

> **⚠** **This article has multiple issues.** Please help **improve it** or discuss these issues on the **talk page**.　　　　[hide]
> - This article **relies largely or entirely upon a single source.** *(December 2012)*
> - This article **relies too much on references to primary sources.** *(December 2012)*

The **Society for Industrial and Applied Mathematics (SIAM)** was founded by a small group of mathematicians from academia and industry who met in Philadelphia in 1951 to start an organization whose members would meet periodically to exchange ideas about the uses of mathematics in industry. This meeting led to the organization of the Society for Industrial and Applied Mathematics. The membership of SIAM has grown from a few hundred in the early 1950s to more than 14,000 as of 2013. SIAM retains its North American influence, but it also has East Asian, Argentinian, Bulgarian, and UK & Ireland sections.

SIAM is one of the four parts of the Joint Policy Board for Mathematics.

**Society for Industrial and Applied Mathematics**

SIAM logo

| | |
|---|---|
| Formation | 1951 |
| Headquarters | Philadelphia, Pennsylvania, United States |
| Membership | >14,000 |
| President | Pamela Cook |
| Website | www.siam.org |

**Contents** [hide]

1 Members
2 Focus
3 Activity groups
4 Publications
　　4.1 Journals
　　4.2 Books

## SIAM Fellows   [edit]

- In 2009 SIAM instituted a Fellows program to recognize certain members who have made outstanding contributions to the fields SIAM serves[14]

14. ^ "Fellows Program" . SIAM. Retrieved 2012-12-04.

The average time between an event and its appearance on Wikipedia is **356** days.

J. Frank et al. 2012

# Knowledge Base Acceleration

# Knowledge Base Acceleration
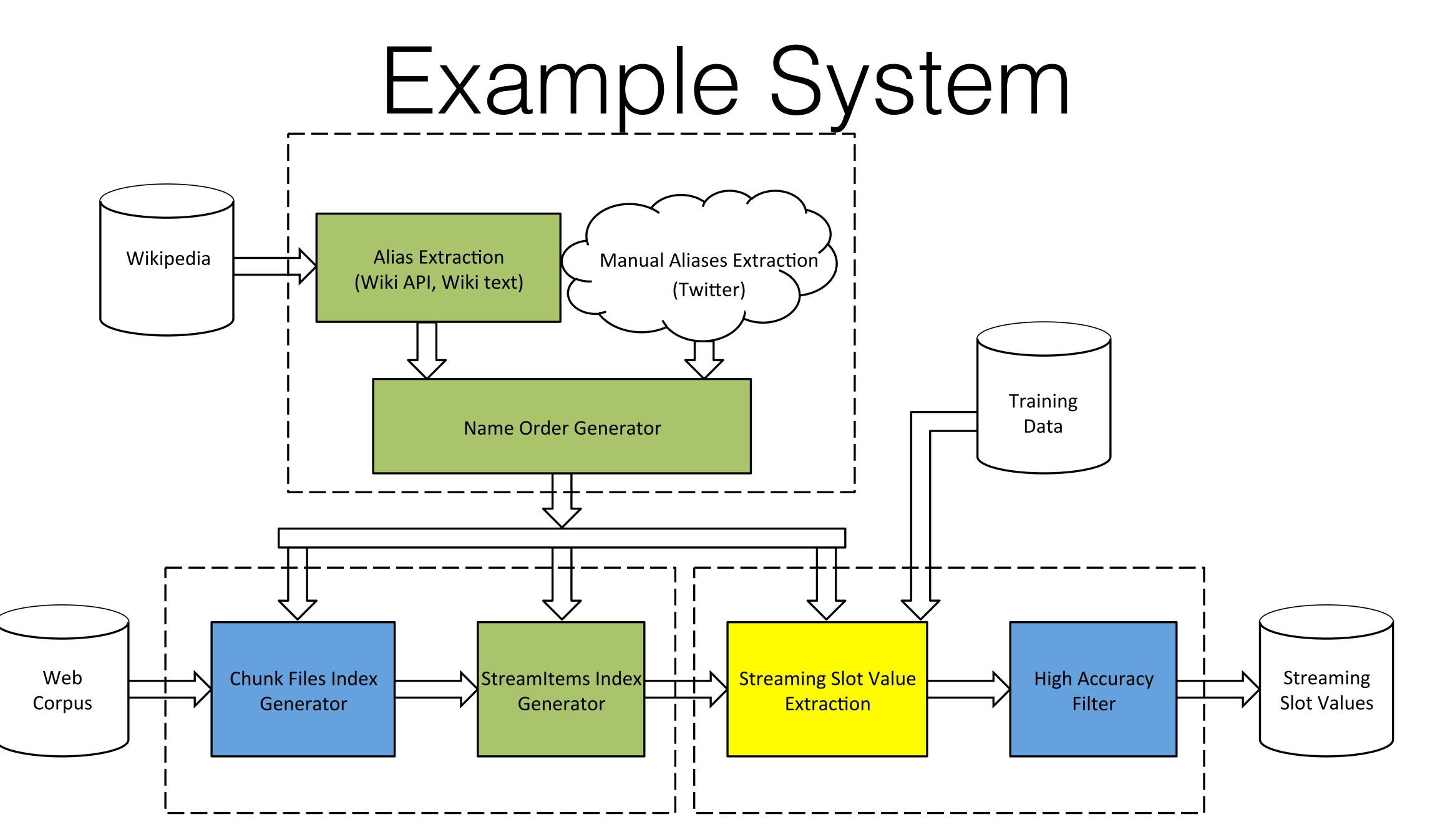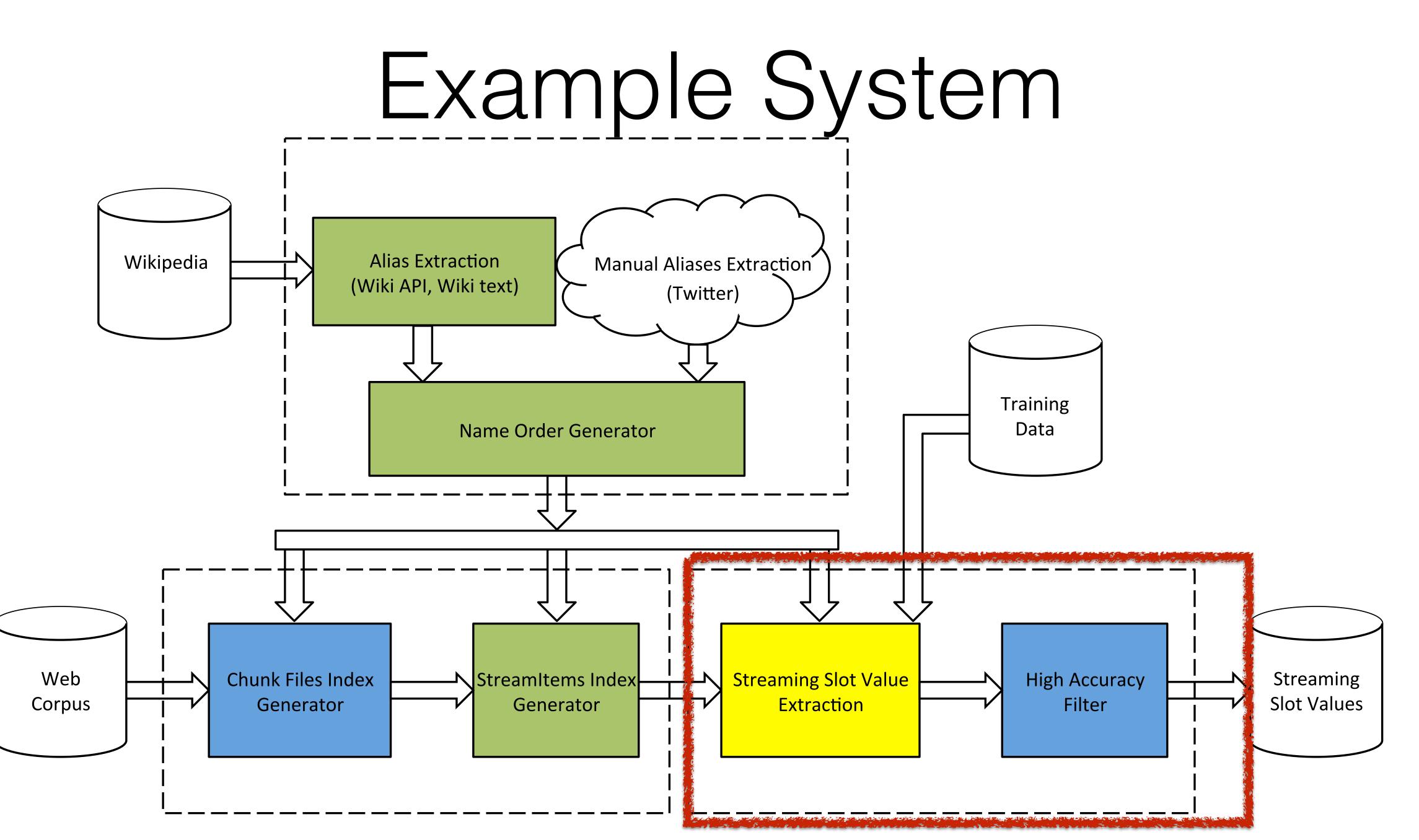
NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

Challenges:

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

Challenges:

    1) A large amount of documents

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

Challenges:

    1) A large amount of documents

    2) Ambiguous text

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

Challenges:

   1) A large amount of documents

   2) Ambiguous text

   3) Ambiguous Entities

# Knowledge Base Acceleration

NIST TREC created a track that reads in streaming documents and a set of entities and suggests citations for wikipedia entities.

Challenges:

1) A large amount of documents

2) Ambiguous text

3) Ambiguous Entities

4) Finding *relevant* facts

# Example System

# Example System

# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.
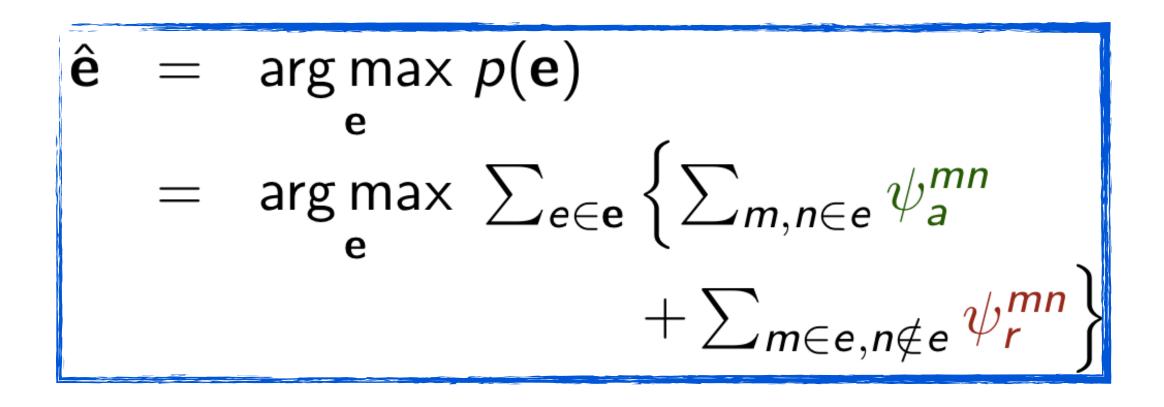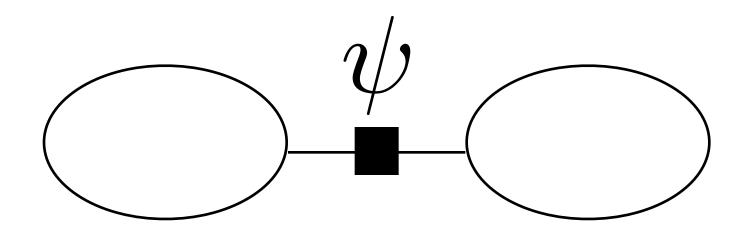
# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.

  - Difficult because of ambiguity

# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.

  - Difficult because of ambiguity

  Same Name, Different Person

# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.

  - Difficult because of ambiguity

    Same Name, Different Person

    Different Name, Same Person

# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.

  - Difficult because of ambiguity



    Same Name, Different Person

    Different Name, Same Person

# Entity Resolution

- Entity resolution is the process of identifying and clustering different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object.



- Difficult because of ambiguity

Same Name, Different Person

Different Name, Same Person

# Entity Resolution Model

# Entity Resolution Model

# Entity Resolution Model

$$\hat{\mathbf{e}} = \arg\max_{\mathbf{e}} p(\mathbf{e})$$

$$= \arg\max_{\mathbf{e}} \sum_{e \in \mathbf{e}} \left\{ \sum_{m,n \in e} \psi_a^{mn} + \sum_{m \in e, n \notin e} \psi_r^{mn} \right\}$$

# Entity Resolution Model

Find the best arrangement.

$$\hat{\mathbf{e}} = \arg\max_{\mathbf{e}} p(\mathbf{e})$$

$$= \arg\max_{\mathbf{e}} \sum_{e \in \mathbf{e}} \left\{ \sum_{m,n \in e} \psi_a^{mn} + \sum_{m \in e, n \notin e} \psi_r^{mn} \right\}$$

# Entity Resolution Model

Find the best arrangement.

$$\hat{\mathbf{e}} = \arg\max_{\mathbf{e}} p(\mathbf{e})$$

$$= \arg\max_{\mathbf{e}} \sum_{e \in \mathbf{e}} \left\{ \sum_{m,n \in e} \psi_a^{mn} + \sum_{m \in e, n \notin e} \psi_r^{mn} \right\}$$

$\psi$

# Entity Resolution Model

Find the best arrangement.

$$\hat{\mathbf{e}} = \underset{\mathbf{e}}{\arg\max}\, p(\mathbf{e})$$

$$= \underset{\mathbf{e}}{\arg\max} \sum_{e \in \mathbf{e}} \left\{ \sum_{m,n \in e} \psi_a^{mn} + \sum_{m \in e, n \notin e} \psi_r^{mn} \right\}$$

# Entity Resolution Algorithm

The Baseline ER metropolis hastings takes a random mention and adds it to a random entity.

# Entity Resolution Algorithm

The Baseline ER metropolis hastings takes a random mention and adds it to a random entity.



Random Number Generator

# Entity Resolution Algorithm

# Entity Resolution Algorithm

1.Select a source mention at *random*.

# Entity Resolution Algorithm

1.Select a source mention at *random*.
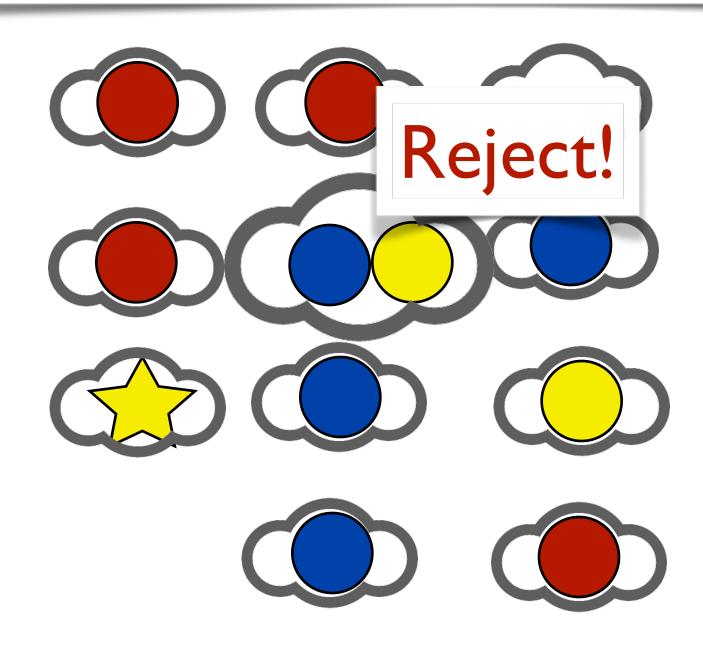
# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
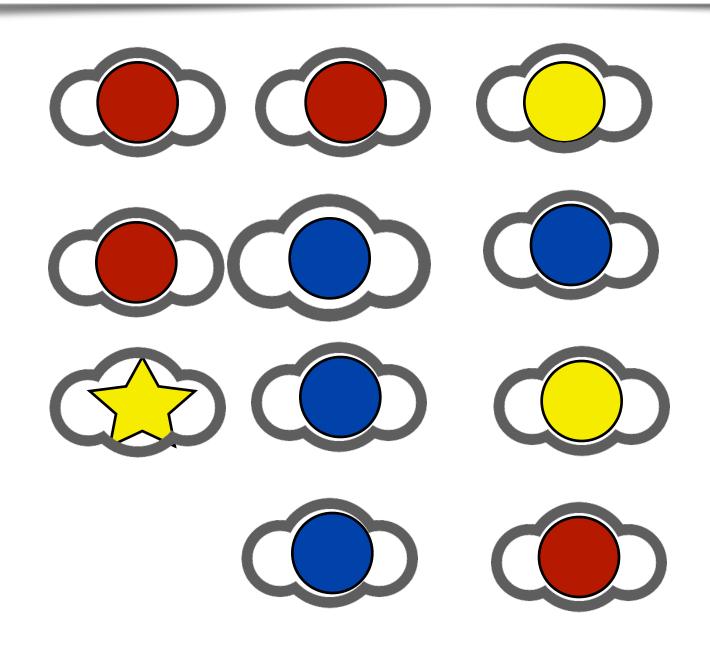3. Propose a merge.

Reject!

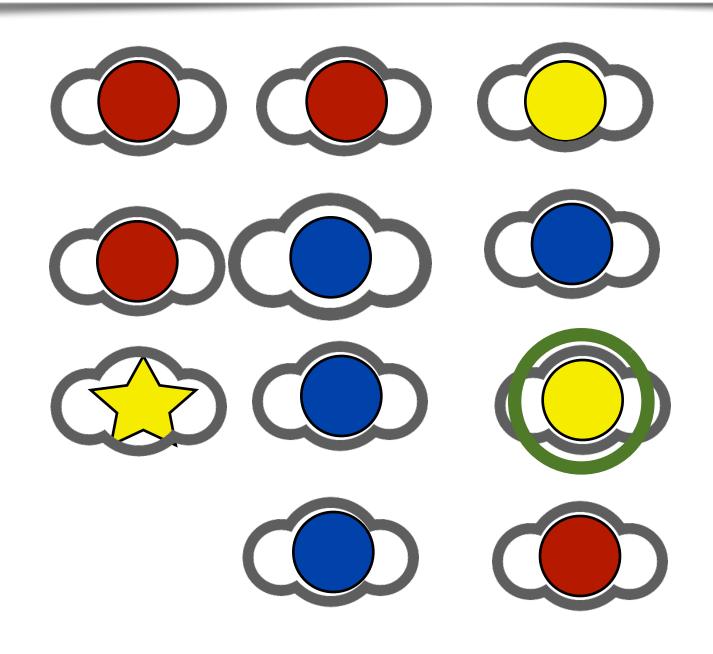$$\alpha(\mathbf{e}, \mathbf{e}') = \min\left(1, \frac{p(\mathbf{e}')}{p(\mathbf{e})}\right)$$

# Entity Resolution Algorithm
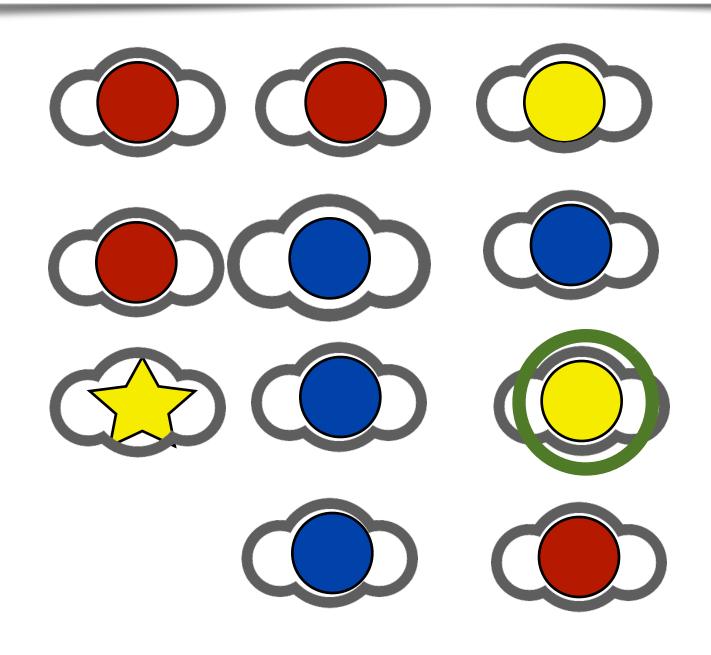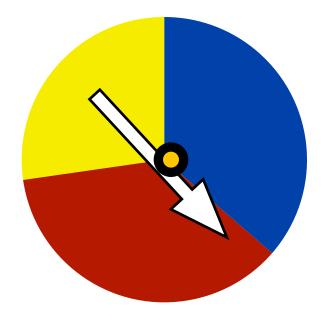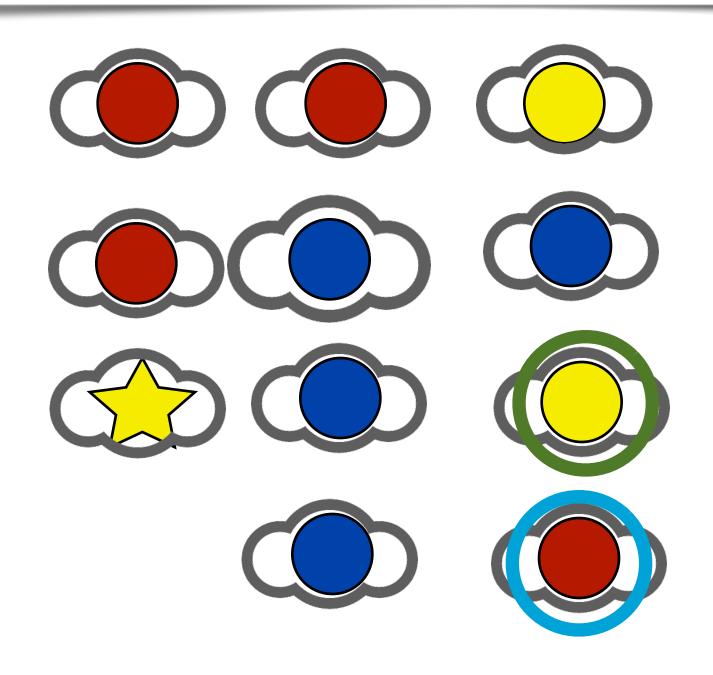
1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.



Reject!

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.

# Entity Resolution Algorithm

1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.
4. Accept when it improves the state.
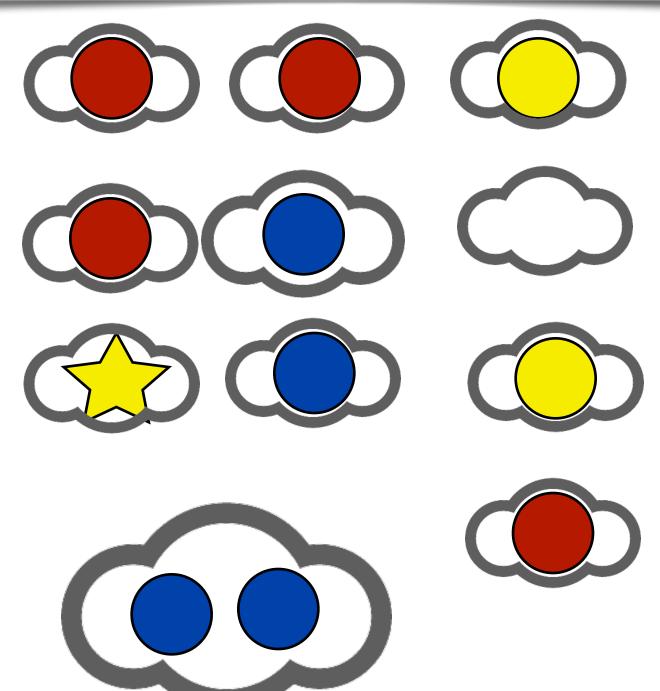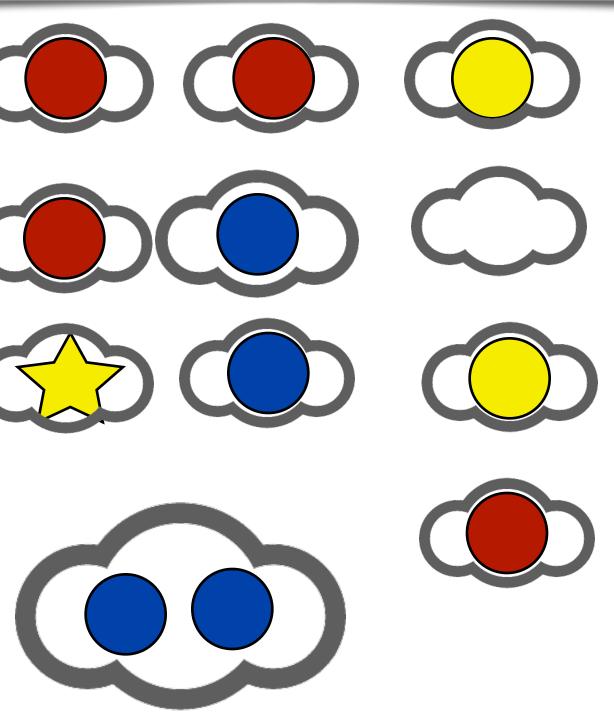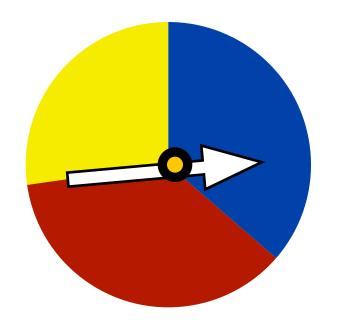
# Entity Resolution Algorithm
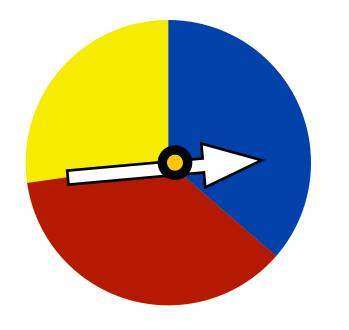
1. Select a source mention at *random*.
2. Select a destination mention at *random*.
3. Propose a merge.
4. Accept when it improves the state.

Accept!

# Entity Resolution Algorithm

Eventually **converges**. (State does not oscillate or vary)

# Entity Resolution Algorithm

Eventually **converges**. (State does not oscillate or vary)



**Markov Chain Monte Carlo Metropolis Hastings!**

# Sampling Optimizations

Distributed Computations (Singh et al. 2011)

# Sampling Optimizations

Distributed Computations (Singh et al. 2011)

Query-Driven Computation (Grant et al. 2015)

# Sampling Inefficiencies

# Sampling Inefficiencies

1. Large clusters are the slowest.

# Sampling Inefficiencies

1. Large clusters are the slowest.

   • Pairwise comparisons are expensive.

# Sampling Inefficiencies

1. Large clusters are the slowest.

   - Pairwise comparisons are expensive. $\Theta(n^2)$

# Sampling Inefficiencies

1. Large clusters are the slowest.

   Pairwise comparisons are expensive. $\Theta(n^2)$

2. Excessive computation on unambiguous entities

# Sampling Inefficiencies

1. Large clusters are the slowest.

   Pairwise comparisons are expensive. $\Theta(n^2)$

2. Excessive computation on unambiguous entities

   Entities such as *Carnegie Mellon* are relatively unambiguous.

# Sampling Inefficiencies

1. Large clusters are the slowest.

   Pairwise comparisons are expensive. $\Theta(n^2)$

2. Excessive computation on unambiguous entities

   Entities such as *Carnegie Mellon* are relatively unambiguous.

   Streaming documents exacerbates these problems.

# Optimizer for MCMC Sampling *

Database style optimizer for streaming MCMC.

# Optimizer for MCMC Sampling ⭐

Database style optimizer for streaming MCMC.

This optimizer makes two decisions:

# Optimizer for MCMC Sampling ⭐

Database style optimizer for streaming MCMC.

This optimizer makes two decisions:

1. Can I approximate the state score calculation?

# Optimizer for MCMC Sampling ⭐

Database style optimizer for streaming MCMC.

This optimizer makes two decisions:

1. Can I approximate the state score calculation?

2. Should I compress an Entity?

# Experiments

- **Wikilink Data Set** *(Singh, Subramaniya, Pereira, McCallum, 2011)*

  - Largest fully-labeled data set

  - 40 Million Mentions

  - 180 GBs of data



Figure 1: Links to Wikipedia as Entity Labels

# Large Entity Sizes



Size Histogram

Size Histogram — Large Entity Sizes. Log-log scatter plot of Number of Entities (y-axis, from 1 to 1e+07) versus Entity Sizes (x-axis, from 1 to 100000).

# Entity Compression

# Entity Compression

- Known matches can be compressed into a representative mention.

# Entity Compression

- Known matches can be compressed into a representative mention.

- Entity compression can reduce the number of mentions (**$n$**).

# Entity Compression

- Known matches can be compressed into a representative mention.

- Entity compression can reduce the number of mentions (**$n$**).

- Compression of large and *popular* entities is **costly**.

# Entity Compression

- Known matches can be compressed into a representative mention.

- Entity compression can reduce the number of mentions (***n***).
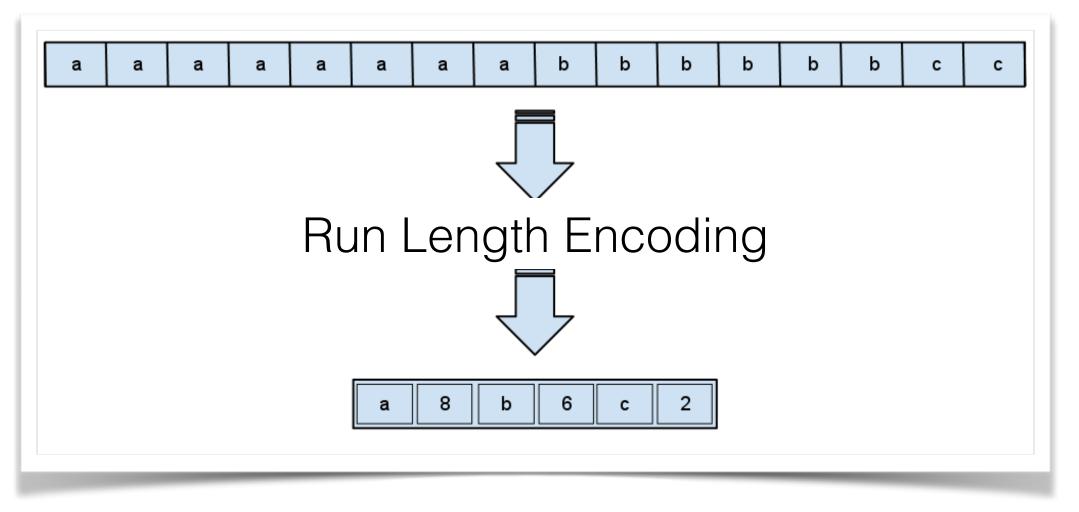
- Compression of large and *popular* entities is **costly**.

- Compression errors are permanent.

# Compression Types



- Run-Length Encoding

- Hierarchical Compression (Wick et al.)

# Early Stopping

- Can we estimate the computation of the features?

Singh et al. EMNLP'12

# Early Stopping

- Can we estimate the computation of the features?

- Given a *p* value, randomly select less values.

Singh et al. EMNLP'12

# Early Stopping



- Can we estimate the computation of the features?

- Given a **p** value, randomly select less values.

Singh et al. EMNLP'12

# Optimizer

Current work

1. Classifier for deciding when to perform *early stopping*.

2. Classifier for the decision to *compress*.

# When should it compress?

# When should it compress?

Power law says there are only a small number of very large clusters.

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

Examining the Wiki Link data set.

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

Examining the Wiki Link data set.

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

Examini      Exact String Match Initialization
set.

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

Examining set.



Exact String Match Initialization

Ground Truth

# When should it compress?

Power law says there are only a small number of very large clusters.

We can treat these in a special way.

Examining set.



Exact String Match Initialization

Ground Truth

# When should it compress?

Power law says there are only a small number of very large clusters.
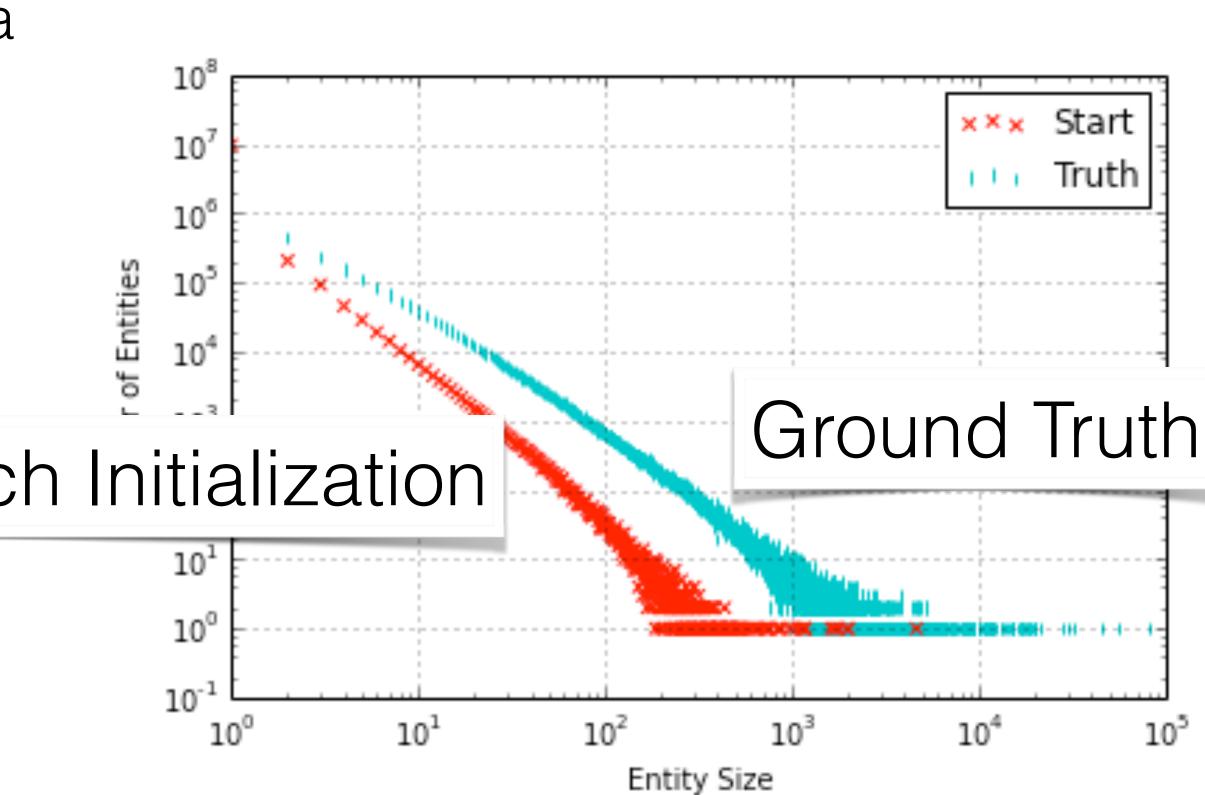
We can treat these in a special way.

Examining set.



Exact String Match Initialization

Ground Truth

# When should it compress?

We could make 100,000 insertions in the time it take to to compress a 300K mention cluster.



Compression time across entity sizes

- Cardinality 0.8
- Cardinality 0.6
- Cardinality 0.4
- Cardinality 0.2
- Insertion time

# When should it compress?

We could make 100,000 insertions in the time it take to to compress a 300K mention cluster.

Compression must be worth it.

# When should it compress?

We could make 100,000 insertions in the time it take to to compress a 300K mention cluster.

Compression must be worth it.



Compression time across entity sizes

# When should we approximate?

# When should we approximate?

- Early stopping only makes sense for clusters of medium size.
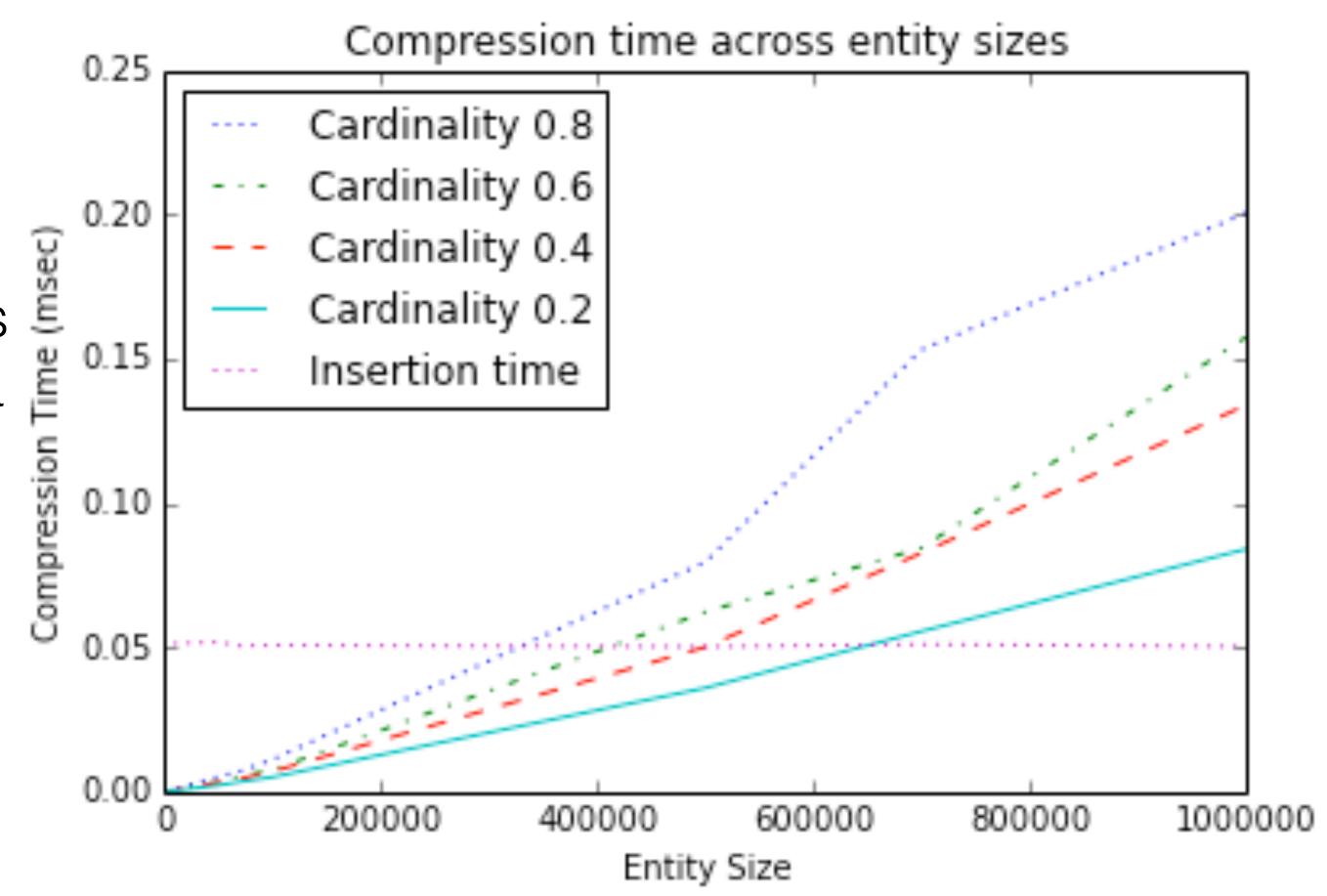
- It is better to do full comparison for small and large cluster sizes.

# When should we approximate?

- Early stopping only makes sense for clusters of medium size.

- It is better to do full comparison for small and large cluster sizes.



#dimensions=5; algorithms=0; querynodes=1; conf=0.8; iterations=5; clocks_per_sec=1000000

Study of clustering and early stopping methods

Legend: early stopping, baseline

Y-axis: Clock Ticks ($10^1$ to $10^{11}$)

X-axis: (Source Cluster, Destination Cluster)
(10, 10), (10, 100), (100, 10), (100, 100), (1000, 10), (10, 1000), (1000, 100), (100, 1000), (1000, 1000), (10, 10000), (10000, 10), (10000, 100), (100, 10000), (10000, 1000), (1000, 10000), (10000, 10000), (100000, 10), (10, 100000), (100000, 100), (100, 100000), (100000, 1000), (1000, 100000), (10000, 100000), (100000, 100000), (1000000, 1000000)
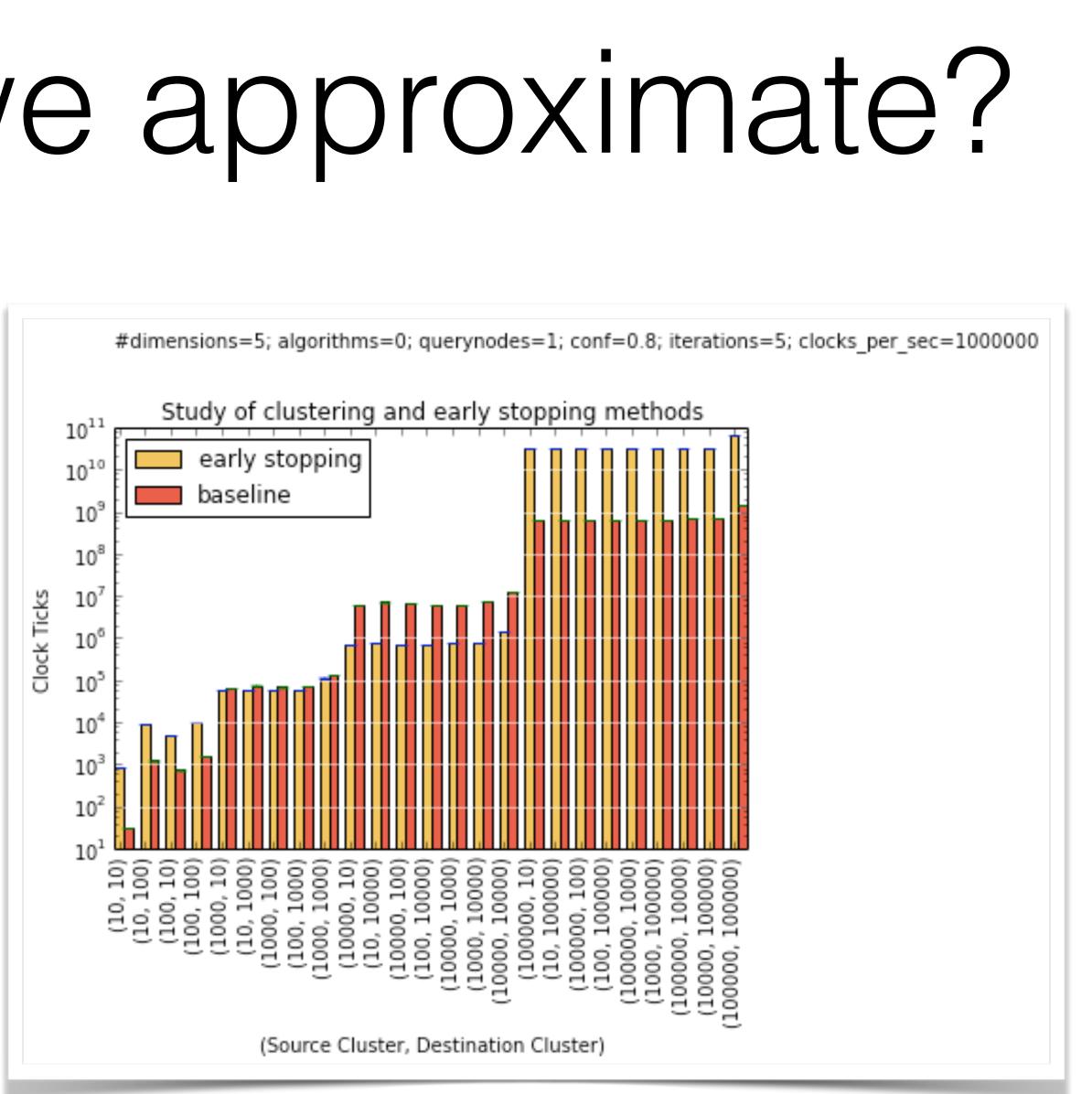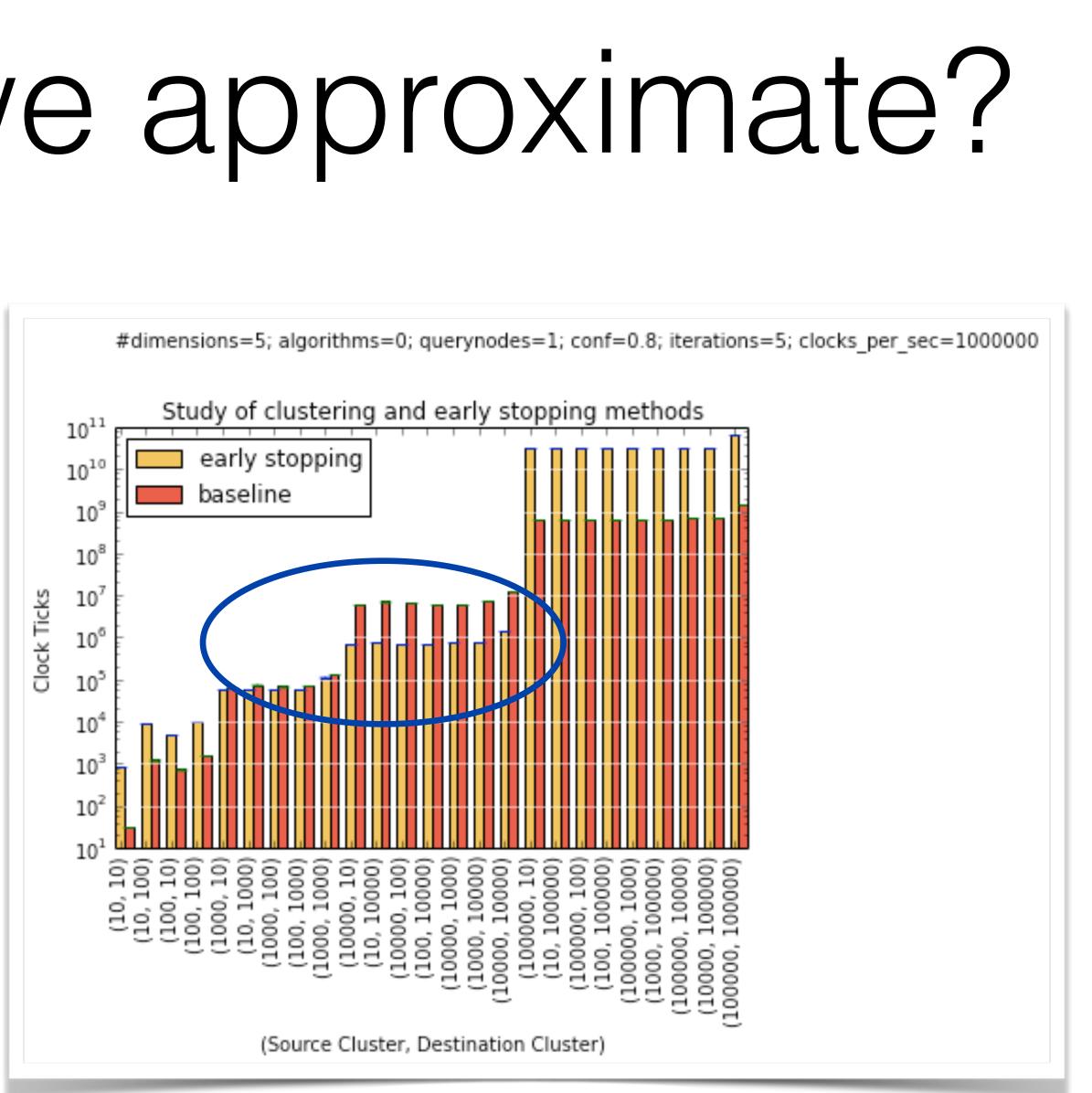
# When should we approximate?

- Early stopping only makes sense for clusters of medium size.

- It is better to do full comparison for small and large cluster sizes.

# Optimizer for Query-Driven Sampling

Optimizer needs to know:
- Current Cardinality of Items in each entity.
- Memory/CPU configuration for estimating baseline time

```
while samples--  > 0:
    m ~ Mentions
    e ~ Entities
    state' = move(state, m, e)
    o = Optimize(state, state', m, e)
    if (!score(state',state, o)):
        state = state'
    doCompress(state, m, e, o)
```

# Summary

- We motivated the need and discussed the open space for optimization of MCMC sampling methods.

- We plan to use the newly released labeled TREC stream corpus.

- Want to collaborate?!

- Lets talk if you want to do a Ph.D. at the University of Oklahoma!

# Thank you!