

Fully Distributed EM for Very Large Datasets

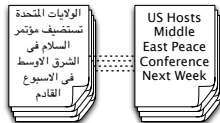
Jason Wolfe Aria Haghighi Dan Klein

Computer Science Division
UC Berkeley

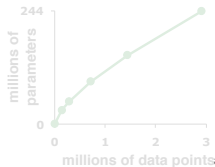


Overview

- **Task:** unsupervised learning via EM



- **Focus:** models w/ many local parameters (relevant to few datums)

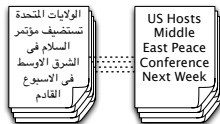


- **Approach:** fully distributed, localized EM
★ *parameter locality* → less bandwidth

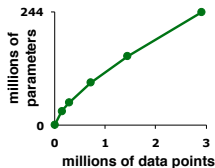


Overview

- **Task:** unsupervised learning via EM



- **Focus:** models w/ many local parameters
(relevant to few datums)

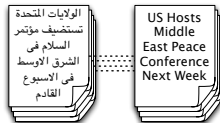


- **Approach:** fully distributed, localized EM
★ *parameter locality* → less bandwidth

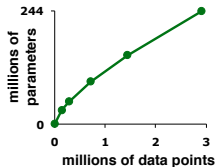


Overview

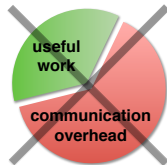
- **Task:** unsupervised learning via EM



- **Focus:** models w/ many local parameters
(relevant to few datums)



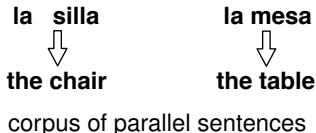
- **Approach:** fully distributed, localized EM
★ *parameter locality* → less bandwidth



- Running example: IBM Model 1 for word alignment
- Naive distributed EM
- Efficiently distributed EM

Word alignment for machine translation

- **Goal:** parallel sentences \rightarrow word-level translation model
- Parameters $\theta_{s \rightarrow t}$: probability that Spanish word s translates to English word t

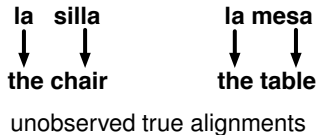
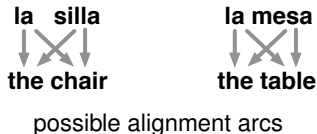
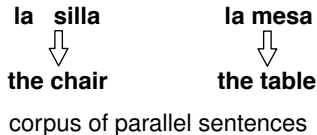


$$\theta = \left\{ \begin{array}{l} \theta_{\text{la} \rightarrow \text{the}} \\ \theta_{\text{la} \rightarrow \text{chair}} \\ \theta_{\text{la} \rightarrow \text{table}} \\ \hline \theta_{\text{silla} \rightarrow \text{the}} \\ \theta_{\text{silla} \rightarrow \text{chair}} \\ \hline \theta_{\text{mesa} \rightarrow \text{the}} \\ \theta_{\text{mesa} \rightarrow \text{table}} \end{array} \right.$$

Word alignment for machine translation

- **Goal:** parallel sentences \rightarrow word-level translation model
- Parameters $\theta_{s \rightarrow t}$: probability that Spanish word s translates to English word t

$$\theta = \left\{ \begin{array}{l} \theta_{la \rightarrow the} = 1.0 \\ \theta_{la \rightarrow chair} = 0.0 \\ \theta_{la \rightarrow table} = 0.0 \\ \hline \theta_{silla \rightarrow the} = 0.0 \\ \theta_{silla \rightarrow chair} = 1.0 \\ \hline \theta_{mesa \rightarrow the} = 0.0 \\ \theta_{mesa \rightarrow table} = 1.0 \end{array} \right.$$



IBM Model 1 for word alignment

a Steve no le gustan las ferias grandes

? ? ? ? ? ? ?

each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

IBM Model 1 for word alignment

a Steve no le gustan las ferias grandes



? ? ? ? ? ? ?

each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

IBM Model 1 for word alignment

a Steve no le gustan las ferias grandes
↓
Steve ? ? ? ? ? ?

each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

IBM Model 1 for word alignment

a Steve no le gustan las ferias grandes
↓ ↙
Steve ? ? ? ? ? ?

each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

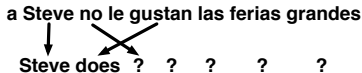
IBM Model 1 for word alignment

a Steve no le gustan las ferias grandes
Steve does ? ? ? ? ?

each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

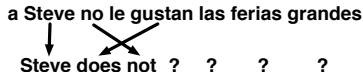
IBM Model 1 for word alignment



each target word is generated by exactly one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

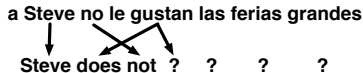
IBM Model 1 for word alignment



each target word is generated by exactly one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

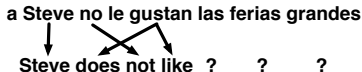
IBM Model 1 for word alignment



each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

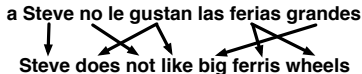
IBM Model 1 for word alignment



each target word is generated by exactly one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

IBM Model 1 for word alignment



each target word is generated by exactly
one source word chosen u.a.r

- **IBM Model 1**: a simple generative model
 - For each target position i , independently
 - choose a source index a_i u.a.r.
 - choose a target word $T_i \sim \theta_{S_{a_i} \rightarrow}$.

EM algorithm for IBM Model 1

- $\theta \leftarrow$ some initial guess

$$\theta_{\text{la} \rightarrow \text{the}} = .33, \theta_{\text{la} \rightarrow \text{chair}} = .33,$$
$$\theta_{\text{la} \rightarrow \text{table}} = .33, \theta_{\text{silla} \rightarrow \text{the}} = .5, \dots$$

EM algorithm for IBM Model 1

- $\theta \leftarrow$ some initial guess

$$\theta_{\text{la} \rightarrow \text{the}} = .33, \theta_{\text{la} \rightarrow \text{chair}} = .33, \\ \theta_{\text{la} \rightarrow \text{table}} = .33, \theta_{\text{silla} \rightarrow \text{the}} = .5, \dots$$

- Iterate:

- 1 **E-step:** estimate alignment counts η
- 1 compute posteriors $p(a_i | \theta)$

$$\frac{.33}{.33 + .5} = .4 \rightarrow \begin{array}{c} \text{la} \quad \text{silla} \\ \downarrow \quad \swarrow \\ \text{the} \quad \text{chair} \end{array} \leftarrow .6 = \frac{.5}{.33 + .5}$$

EM algorithm for IBM Model 1

- $\theta \leftarrow$ some initial guess

$$\theta_{la \rightarrow the} = .33, \theta_{la \rightarrow chair} = .33, \\ \theta_{la \rightarrow table} = .33, \theta_{silla \rightarrow the} = .5, \dots$$

- Iterate:

① **E-step:** estimate alignment counts η

① compute posteriors $p(a_i | \theta)$

$$\frac{.33}{.33 + .5} = .4 \rightarrow \begin{array}{cc} \text{la} & \text{silla} \\ \downarrow & \swarrow \\ \text{the} & \text{chair} \end{array} \leftarrow .6 = \frac{.5}{.33 + .5}$$

$$\begin{array}{cc} \text{la} & \text{silla} \\ \downarrow & \swarrow \\ \text{the} & \text{chair} \end{array} \quad \begin{array}{cc} \text{la} & \text{mesa} \\ \downarrow & \swarrow \\ \text{the} & \text{table} \end{array}$$

EM algorithm for IBM Model 1

- $\theta \leftarrow$ some initial guess

$$\theta_{la \rightarrow the} = .33, \theta_{la \rightarrow chair} = .33, \\ \theta_{la \rightarrow table} = .33, \theta_{silla \rightarrow the} = .5, \dots$$

- Iterate:

- 1 **E-step:** estimate alignment counts η

- 1 compute posteriors $p(a_i | \theta)$
- 2 aggregate into **expected counts** $\eta_{s \rightarrow t}$
(expected # times $s \rightarrow t$ under θ)

$$\eta_{s \rightarrow t} \leftarrow \sum_c \frac{\theta_{s \rightarrow t}}{\sum_{i'} \theta_{s_{i'} \rightarrow t}}$$

$$\frac{.33}{.33 + .5} = .4 \rightarrow \begin{array}{c} \text{la silla} \\ \downarrow \swarrow \\ \text{the chair} \end{array} \leftarrow .6 = \frac{.5}{.33 + .5}$$

$$\begin{array}{c} \text{la silla} \\ \downarrow \times \downarrow \\ \text{the chair} \end{array} \leftarrow .6 \quad \begin{array}{c} \text{la mesa} \\ \downarrow \times \downarrow \\ \text{the table} \end{array} \leftarrow .6$$

$$\eta_{la \rightarrow the} = .8, \eta_{la \rightarrow chair} = .4, \\ \eta_{la \rightarrow table} = .4, \eta_{silla \rightarrow the} = .6, \dots$$

EM algorithm for IBM Model 1

- $\theta \leftarrow$ some initial guess

$$\theta_{la \rightarrow the} = .33, \theta_{la \rightarrow chair} = .33, \\ \theta_{la \rightarrow table} = .33, \theta_{silla \rightarrow the} = .5, \dots$$

- Iterate:

- 1 **E-step:** estimate alignment counts η

- 1 compute posteriors $p(a_i | \theta)$
- 2 aggregate into **expected counts** $\eta_{s \rightarrow t}$
(expected # times $s \rightarrow t$ under θ)

$$\eta_{s \rightarrow t} \leftarrow \sum_c \frac{\theta_{s \rightarrow t}}{\sum_{i'} \theta_{s_{i'} \rightarrow t}}$$

- 2 **M-step:** normalize η to get **new ML θ**

$$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

$$\frac{.33}{.33 + .5} = .4 \rightarrow \begin{matrix} la & silla \\ \downarrow & \swarrow \\ the & chair \end{matrix} \leftarrow .6 = \frac{.5}{.33 + .5}$$

$$\begin{matrix} la & silla \\ \downarrow & \swarrow \\ the & chair \end{matrix} \quad \begin{matrix} la & mesa \\ \downarrow & \swarrow \\ the & table \end{matrix}$$

$$\eta_{la \rightarrow the} = .8, \eta_{la \rightarrow chair} = .4, \\ \eta_{la \rightarrow table} = .4, \eta_{silla \rightarrow the} = .6, \dots$$

$$\theta_{la \rightarrow the} = .5, \theta_{la \rightarrow chair} = .25, \\ \theta_{la \rightarrow table} = .25, \theta_{silla \rightarrow the} = .5, \dots$$

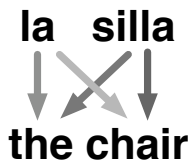
E-Step 1

la silla
↓ ↓ ↓
the chair

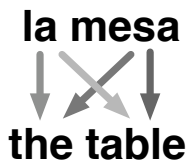
la mesa
↓ ↓ ↓
the table

E-Step 2

la silla
↓ ↓ ↓
the chair

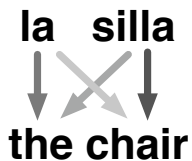


la mesa
↓ ↓ ↓
the table

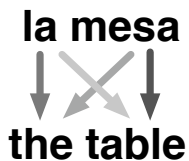


E-Step 3

la silla
↓ ↓ ↓
the chair



la mesa
↓ ↓ ↓
the table



E-Step 4

la silla
↓ ↓ ↓
the chair

la mesa
↓ ↓ ↓
the table

E-Step 5

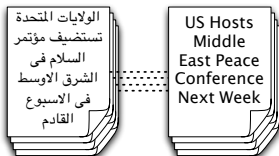
la silla
↓ ↓ ↓
the chair

la mesa
↓ ↓ ↓
the table

E-Step ∞

la silla
↓ ↓
the chair

la mesa
↓ ↓
the table



- 2.9 million sentence pairs from UN proceedings
- 243 million *unique word pairs*
(translations possible in some sentence pair)
 - 243 M parameters in θ
 - 243 M counts in η
- Even fitting all (indexed) parameters in 32-bit memory can be challenging

- Running example: IBM Model 1 for word alignment
- Naive distributed EM
- Efficiently distributed EM

Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- Remaining problems:
 - Memory at Reduce node
 - C-step (communication) bandwidth:
5.5 B numbers per iteration
(on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

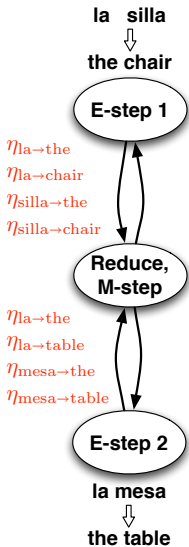
Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- Remaining problems:
- Memory at Reduce node
 - C-step (communication) bandwidth: 5.5 B numbers per iteration (on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

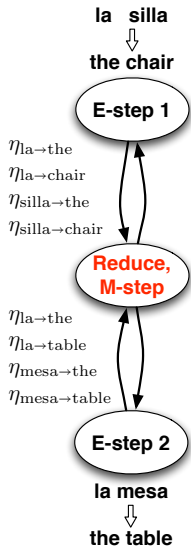
Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- Remaining problems:
 - Memory at Reduce node
 - C-step (communication) bandwidth: 5.5 B numbers per iteration (on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

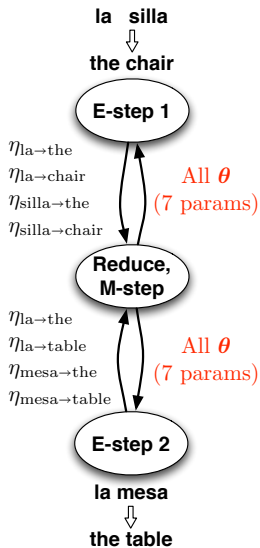
Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- Remaining problems:
 - Memory at Reduce node
 - C-step (communication) bandwidth: 5.5 B numbers per iteration (on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

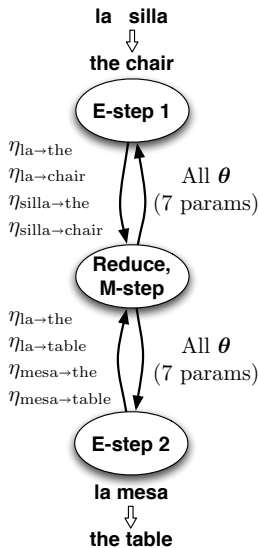
Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- Remaining problems:
 - Memory at Reduce node
 - C-step (communication) bandwidth: 5.5 B numbers per iteration (on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

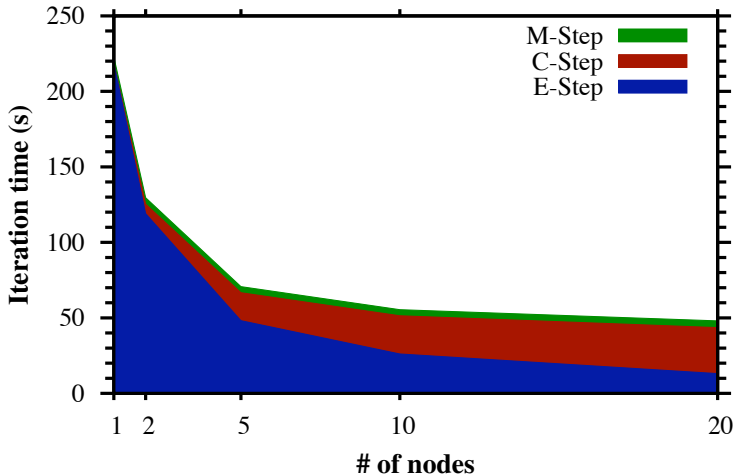
Previous approach: distributing the E-step



- 1 E-step computations distribute easily
 - partition data over k nodes
 - alignments independent given θ
 - 2 Nodes communicate partial counts to central Reduce node
 - 3 Reduce node does global M-step
 - 4 Reduce sends new parameters back
- **Remaining problems:**
 - Memory at Reduce node
 - **C-step** (communication) bandwidth: 5.5 B numbers per iteration (on full dataset with 20 nodes)

(Chu *et al.* 2006, Dyer *et al.* 2008, Newman *et al.* 2008, ...)

Iteration time vs. # of E-step nodes



Common practical solutions

- Memory and bandwidth are real problems in practice
- Workarounds
 - Use less data
 - Ignore rare words
 - Train on independent chunks
 - Swap to disk
 - Distribute over multiple machines

- Running example: IBM Model 1 for word alignment
- Naive distributed EM
- Efficiently distributed EM

Distributing the M-step locally



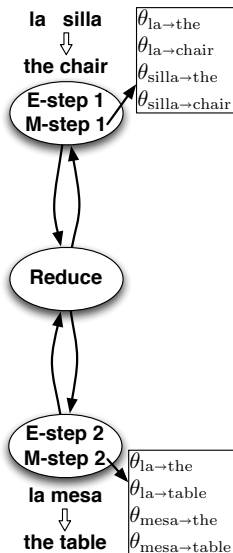
- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
- Reduce passes back **counts**
- Don't need to hear about **irrelevant** source words
- Don't need to tell (or hear) about **purely local** source words
- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s-t} + \frac{\eta_{s-t}}{\sum_{t'} \eta_{s-t'}}$$

- Bandwidth savings: 30%

Distributing the M-step locally



- Distribute M-step alongside E-step

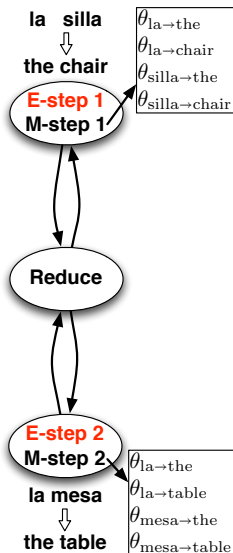
- Nodes store only needed params, compute them **locally**

- Reduce passes back counts
- Don't need to hear about irrelevant source words
- Don't need to tell (or hear) about purely local source words
- Need to hear everything about each source word: M-step denominator

$$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_r \eta_{s \rightarrow r}}$$

- Bandwidth savings: 30%

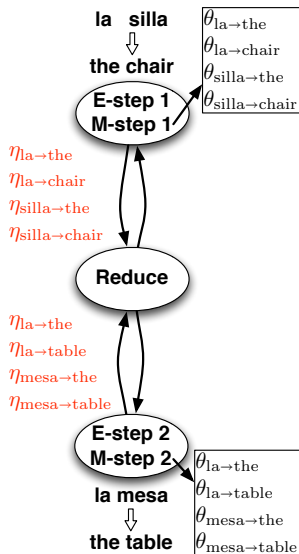
Distributing the M-step locally



- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
 - Reduce passes back **counts**
 - Don't need to hear about **irrelevant** source words
 - Don't need to tell (or hear) about **purely local** source words
 - Need to hear **everything** about each source word: M-step denominator
- $$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_r \eta_{s \rightarrow r}}$$
- Bandwidth savings: 30%

Distributing the M-step locally



- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**

- Reduce passes back **counts**

- Don't need to hear about **irrelevant** source words

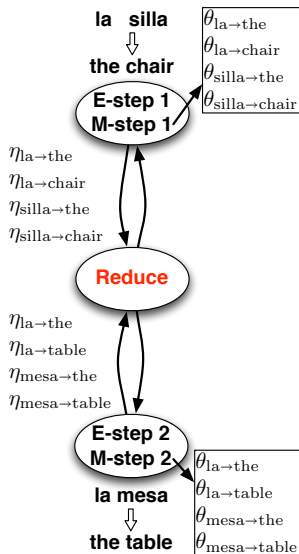
- Don't need to tell (or hear) about **purely local** source words

- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

- Bandwidth savings: **30%**

Distributing the M-step locally



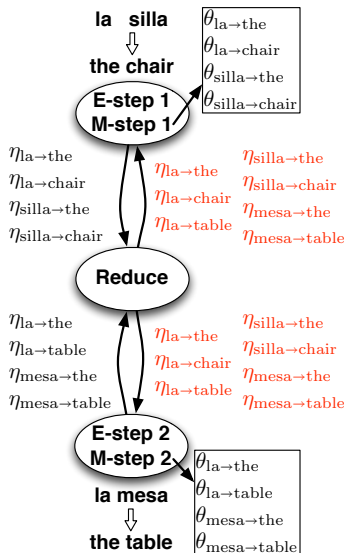
- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
- Reduce passes back **counts**
- Don't need to hear about **irrelevant source words**
- Don't need to tell (or hear) about **purely local source words**
- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

- Bandwidth savings: **30%**

Distributing the M-step locally



- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**

- Reduce passes back **counts**

- Don't need to hear about **irrelevant** source words

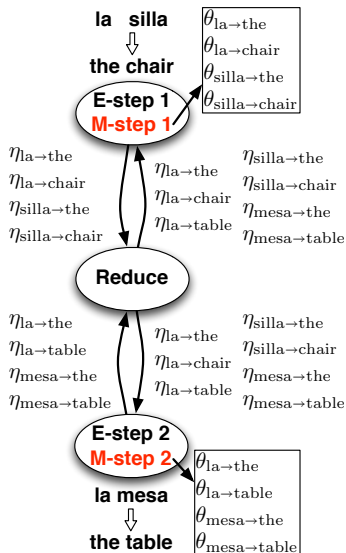
- Don't need to tell (or hear) about **purely local** source words

- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_r \eta_{s \rightarrow r}}$$

- Bandwidth savings: 30%

Distributing the M-step locally

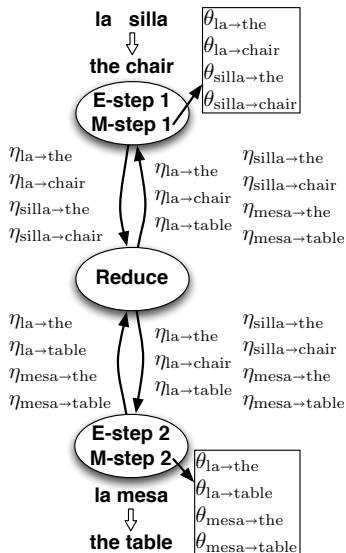


- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
- Reduce passes back **counts**
- Don't need to hear about **irrelevant** source words
- Don't need to tell (or hear) about **purely local** source words
- Need to hear **everything** about each source word: M-step denominator
- Bandwidth savings: 30%

$$\theta_{s \rightarrow t} + \frac{\eta_{s \rightarrow t}}{\sum_r \eta_{s \rightarrow r}}$$

Distributing the M-step locally



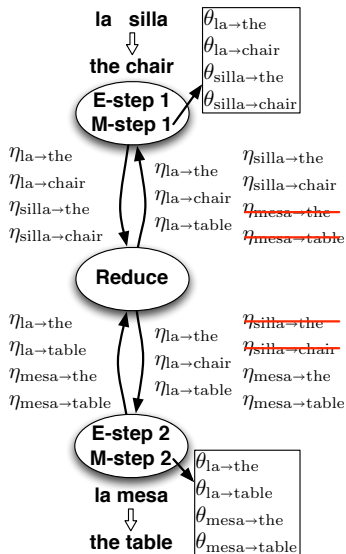
• Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
- Reduce passes back **counts**
- Don't need to hear about **irrelevant** source words
- Don't need to tell (or hear) about **purely local** source words
- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

- Bandwidth savings: **30%**

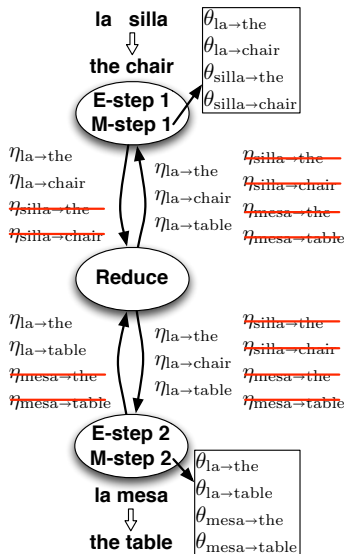
Distributing the M-step locally



- Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
 - Reduce passes back **counts**
 - Don't need to hear about **irrelevant** source words
 - Don't need to tell (or hear) about **purely local** source words
 - Need to hear **everything** about each source word: M-step denominator
- $$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$
- Bandwidth savings: **30%**

Distributing the M-step locally



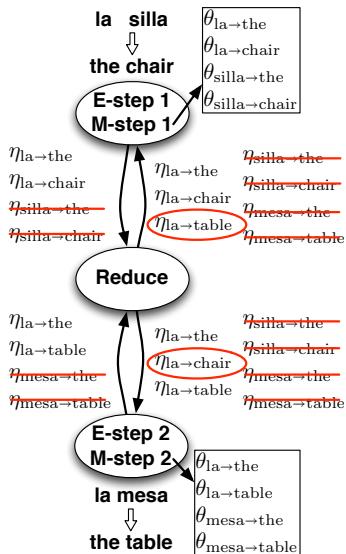
• Distribute M-step alongside E-step

- Nodes store only needed params, compute them **locally**
- Reduce passes back **counts**
- Don't need to hear about **irrelevant** source words
- Don't need to tell (or hear) about **purely local** source words
- Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

- Bandwidth savings: **30%**

Distributing the M-step locally

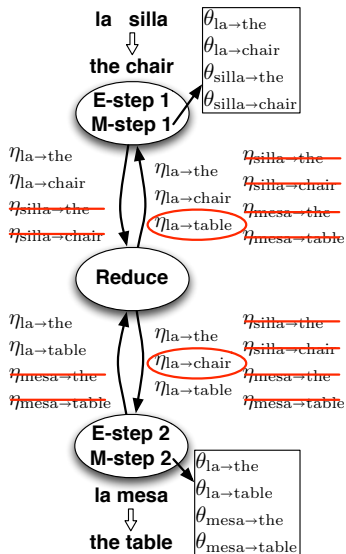


- Distribute M-step alongside E-step
 - Nodes store only needed params, compute them **locally**
 - Reduce passes back **counts**
 - Don't need to hear about **irrelevant** source words
 - Don't need to tell (or hear) about **purely local** source words
 - Need to hear **everything** about each source word: M-step denominator

$$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$

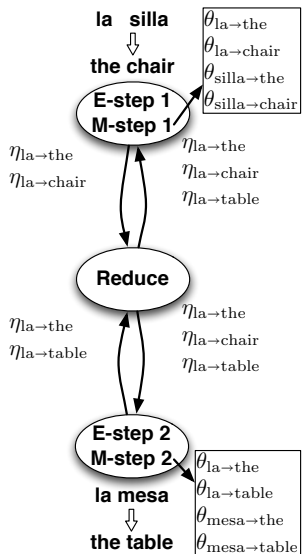
- Bandwidth savings: **30%**

Distributing the M-step locally



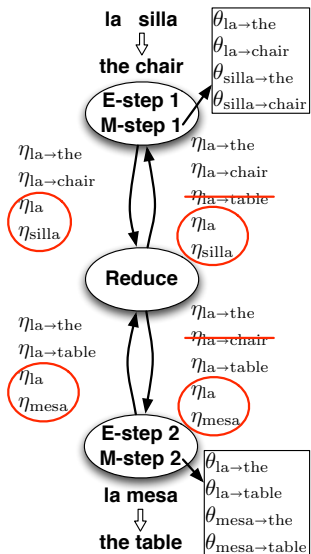
- Distribute M-step alongside E-step
 - Nodes store only needed params, compute them **locally**
 - Reduce passes back **counts**
 - Don't need to hear about **irrelevant** source words
 - Don't need to tell (or hear) about **purely local** source words
 - Need to hear **everything** about each source word: M-step denominator
- $$\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\sum_{t'} \eta_{s \rightarrow t'}}$$
- Bandwidth savings: **30%**

Augmenting η to increase locality



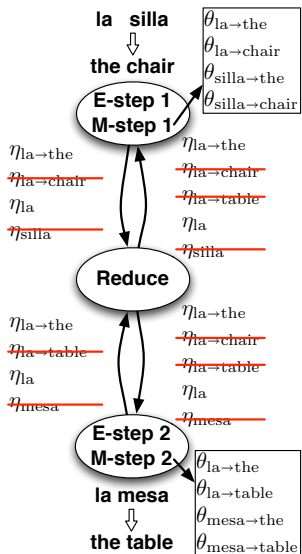
- Augment η with redundant $\eta_s = \sum_{t'} \eta_{s \rightarrow t'}$ in E-step
- M-step becomes $\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\eta_s}$
- Increases locality
- Total bandwidth savings: 84% (bigger if more nodes)
- Similar tricks for other models

Augmenting η to increase locality



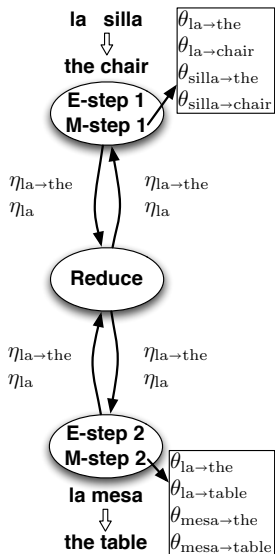
- **Augment η** with redundant $\eta_s = \sum_{t'} \eta_{s \rightarrow t'}$ in E-step
- M-step becomes $\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\eta_s}$
- Increases locality
- Total bandwidth savings: 84% (bigger if more nodes)
- Similar tricks for other models

Augmenting η to increase locality



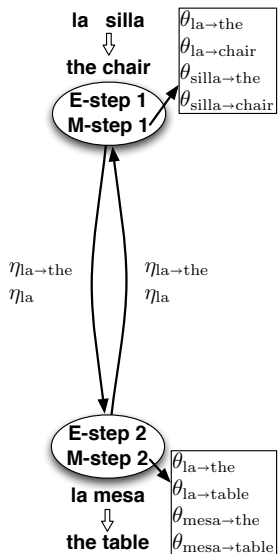
- **Augment η** with redundant $\eta_s = \sum_{t'} \eta_{s \rightarrow t'}$ in E-step
- M-step becomes $\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\eta_s}$
- Increases locality
- Total bandwidth savings: 84% (bigger if more nodes)
- Similar tricks for other models

Augmenting η to increase locality



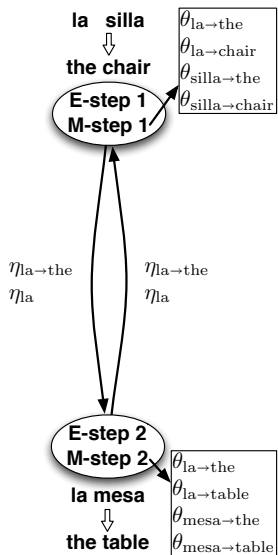
- **Augment η** with redundant $\eta_s = \sum_{t'} \eta_{s \rightarrow t'}$ in E-step
- M-step becomes $\theta_{s \rightarrow t} \leftarrow \frac{\eta_{s \rightarrow t}}{\eta_s}$
- Increases locality
- Total bandwidth savings: **84%** (bigger if more nodes)
- Similar tricks for other models

Choice of C-step topology

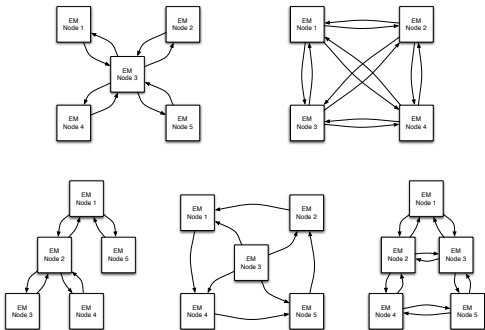


- No need for separate Reduce nodes
- By choosing connectivity, can trade off
 - bandwidth
 - latency
 - locality
 - ...

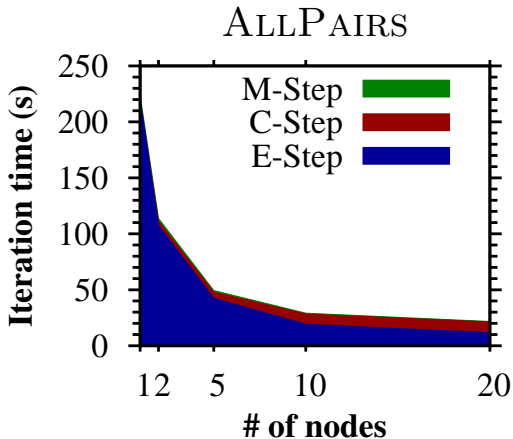
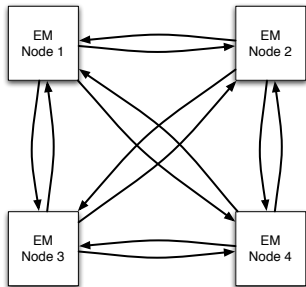
Choice of C-step topology



- No need for separate Reduce nodes
- By choosing connectivity, can trade off
 - bandwidth
 - latency
 - locality
 - ...

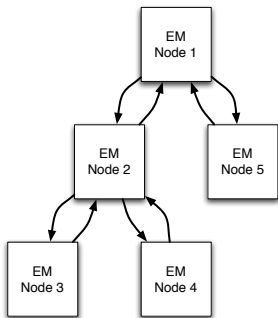


ALLPAIRS topology



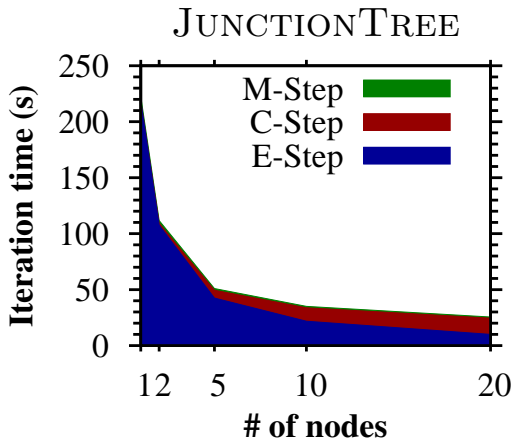
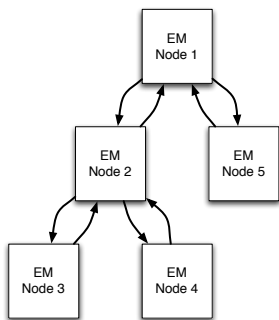
Total Bandwidth: 3.6 B counts per iteration

JUNCTIONTREE topology



- Nodes embedded in **arbitrary tree structure**
- Messages contain counts needed by nodes in both **subtrees**
- Tree can **optimize** for
 - bandwidth
 - locality
 - ...
- We use maximum spanning tree to heuristically minimize bandwidth
- Future work: multiple trees

JUNCTIONTREE topology



Total Bandwidth: 1.4 B counts per iteration

Locality in other models

- Ex: Latent Dirichlet Allocation (LDA) for **topic modeling**
 - Parameters: unigram distributions for each topic $p(w|t)$
 - Topic-word parameters local
 - Similar augmentation trick to Model 1
 - Details and results in paper
- Also applies to other EM models, beyond EM
 - **Word locality** is extremely common in NLP applications
 - Variational inference
 - Other computations that make sparse use of expectations

Conclusion

- A fully distributed, maximally localized EM algorithm
 - exploits **parameter locality** for significant **speedup**
 - is **general**; just define η for each datum
 - is **flexible** with respect to communication topology

- Many further improvements possible
 - intelligent partitioning of data
 - running E- and C-steps in parallel
 - better topologies (e.g., multiple trees)
 - exploiting approximate sparsity/locality
 - ...