# Inverted Indexing, MT, MapReduce

# Game plan for today:

Quick overview of inverted indexing
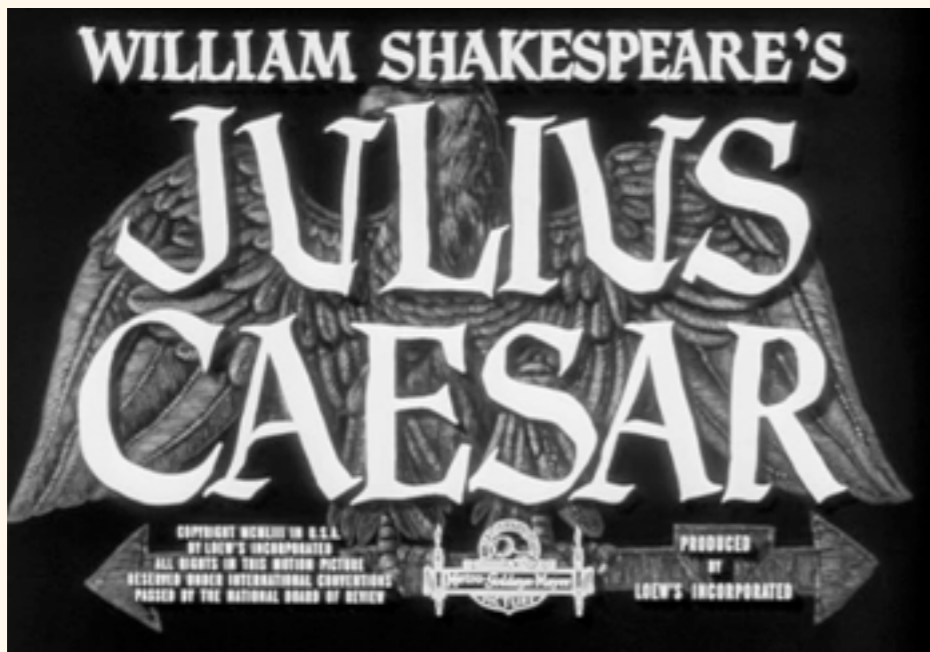
Inverted indexing & Map-Reduce

Quick MT overview

Map-Reduce & MT Model Estimation

# Information Retrieval 101:

We have a set of documents...

... and we want to be able to search them by the terms that they contain.



| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | ... |
|---|---|---|---|---|---|---|---|
| Antony | 1 | 1 | 0 | 0 | 0 | 1 | |
| Brutus | 1 | 1 | 0 | 1 | 0 | 0 | |
| Caesar | 1 | 1 | 0 | 1 | 1 | 1 | |
| Calpurnia | 0 | 1 | 0 | 0 | 0 | 0 | |
| Cleopatra | 1 | 0 | 0 | 0 | 0 | 0 | |
| mercy | 1 | 0 | 1 | 1 | 1 | 1 | |
| worser | 1 | 0 | 1 | 1 | 1 | 0 | |
| ... | | | | | | | |

To answer the query Brutus AND Caesar AND NOT Calpurnia, we take the vectors for Brutus, Caesar and Calpurnia, complement the last, and then do a bitwise AND:

110100 AND 110111 AND 101111 = 100100

# Information Retrieval 101:
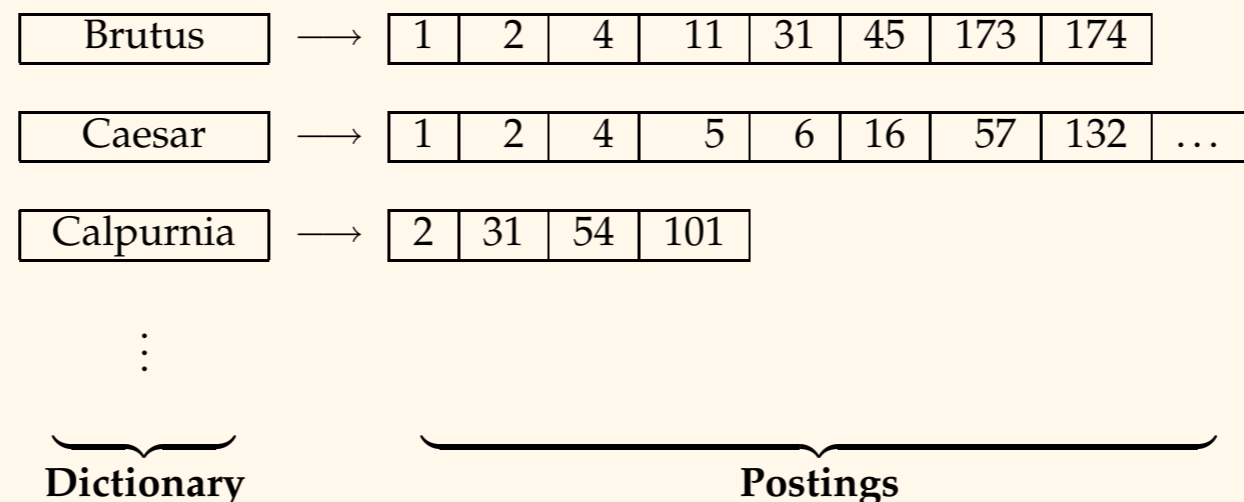
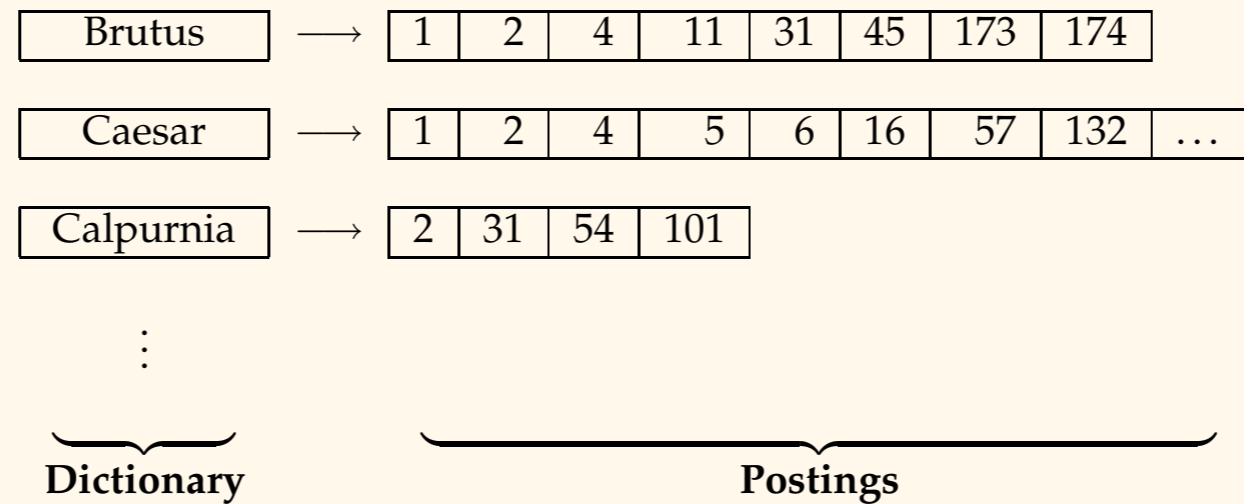| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | … |
|---|---|---|---|---|---|---|---|
| Antony | 1 | 1 | 0 | 0 | 0 | 1 | |
| Brutus | 1 | 1 | 0 | 1 | 0 | 0 | |
| Caesar | 1 | 1 | 0 | 1 | 1 | 1 | |
| Calpurnia | 0 | 1 | 0 | 0 | 0 | 0 | |
| Cleopatra | 1 | 0 | 0 | 0 | 0 | 0 | |
| mercy | 1 | 0 | 1 | 1 | 1 | 1 | |
| worser | 1 | 0 | 1 | 1 | 1 | 0 | |
| … | | | | | | | |

Obviously, this approach can only scale so far…

(as it happens, "so far" is actually not very far in this case)

One solution: build an index mapping terms to the documents in which they may be found.

| Brutus | ⟶ | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 | |
|---|---|---|---|---|---|---|---|---|---|---|

| Caesar | ⟶ | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | … |
|---|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | ⟶ | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

⋮

**Dictionary**          **Postings**

| Brutus | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 | |
|---|---|---|---|---|---|---|---|---|---|---|

| Caesar | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | … |
|---|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | → | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

⋮

**Dictionary**                    **Postings**

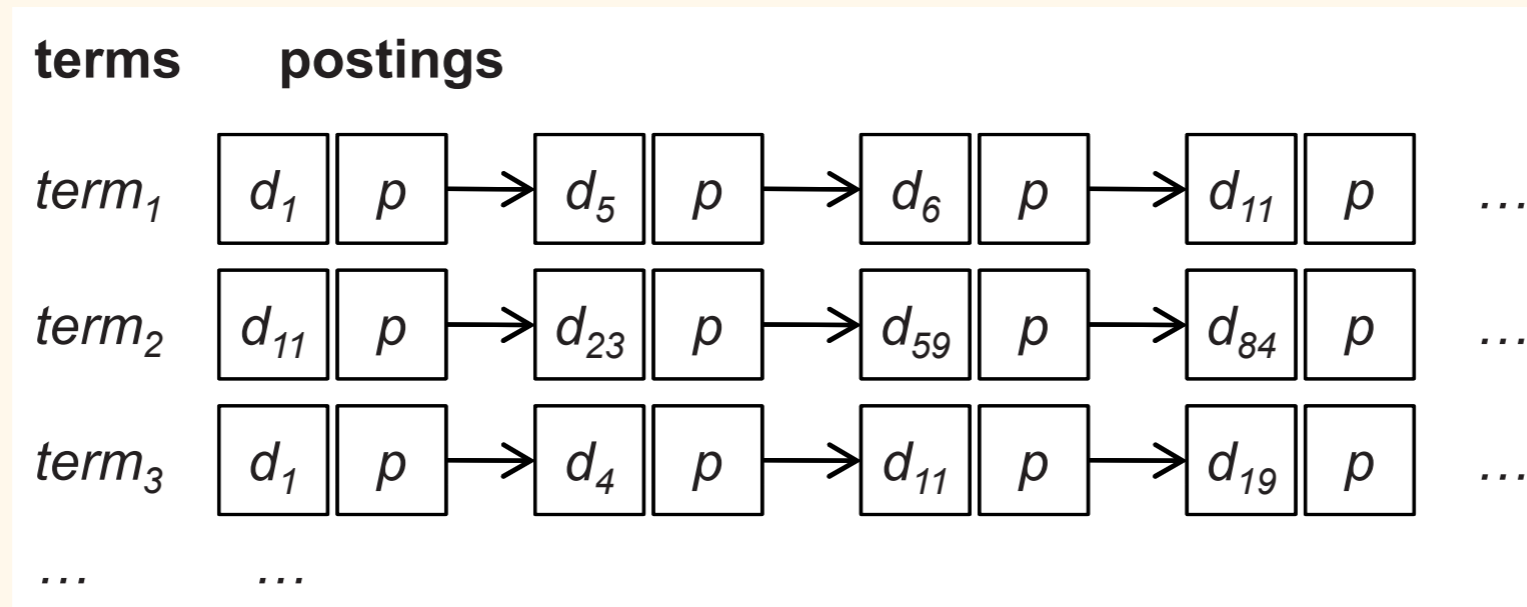Resolving a query now becomes a *set operation* on the posting lists…

… i.e., resolving the query "Brutus AND Calpurnia" would simply require intersecting their respective posting lists.

Note the order of the postings!

| terms | postings | | | | | | |
|-------|----------|--|--|--|--|--|--|
| $term_1$ | $d_1$ | $p$ → | $d_5$ | $p$ → | $d_6$ | $p$ → | $d_{11}$ | $p$ ... |
| $term_2$ | $d_{11}$ | $p$ → | $d_{23}$ | $p$ → | $d_{59}$ | $p$ → | $d_{84}$ | $p$ ... |
| $term_3$ | $d_1$ | $p$ → | $d_4$ | $p$ → | $d_{11}$ | $p$ → | $d_{19}$ | $p$ ... |
| ... | ... | | | | | | |

Postings generally include some sort of "payload" or "metadata:"

- Term frequency
- Term positions within the document
- Context surrounding the term
- PoS information
- etc.

DocIDs can be assigned randomly, or according to some scheme (documents from same domain get similar IDs, higher PageRank gets lower IDs, etc.)

# Indexing considerations:

A "web scale" corpus will involve billions of pages...

... and remember: some pages (news sites, etc.) become "stale" quickly, and must be reindexed often.

Bottleneck with inverted indexing: having to visit each document.

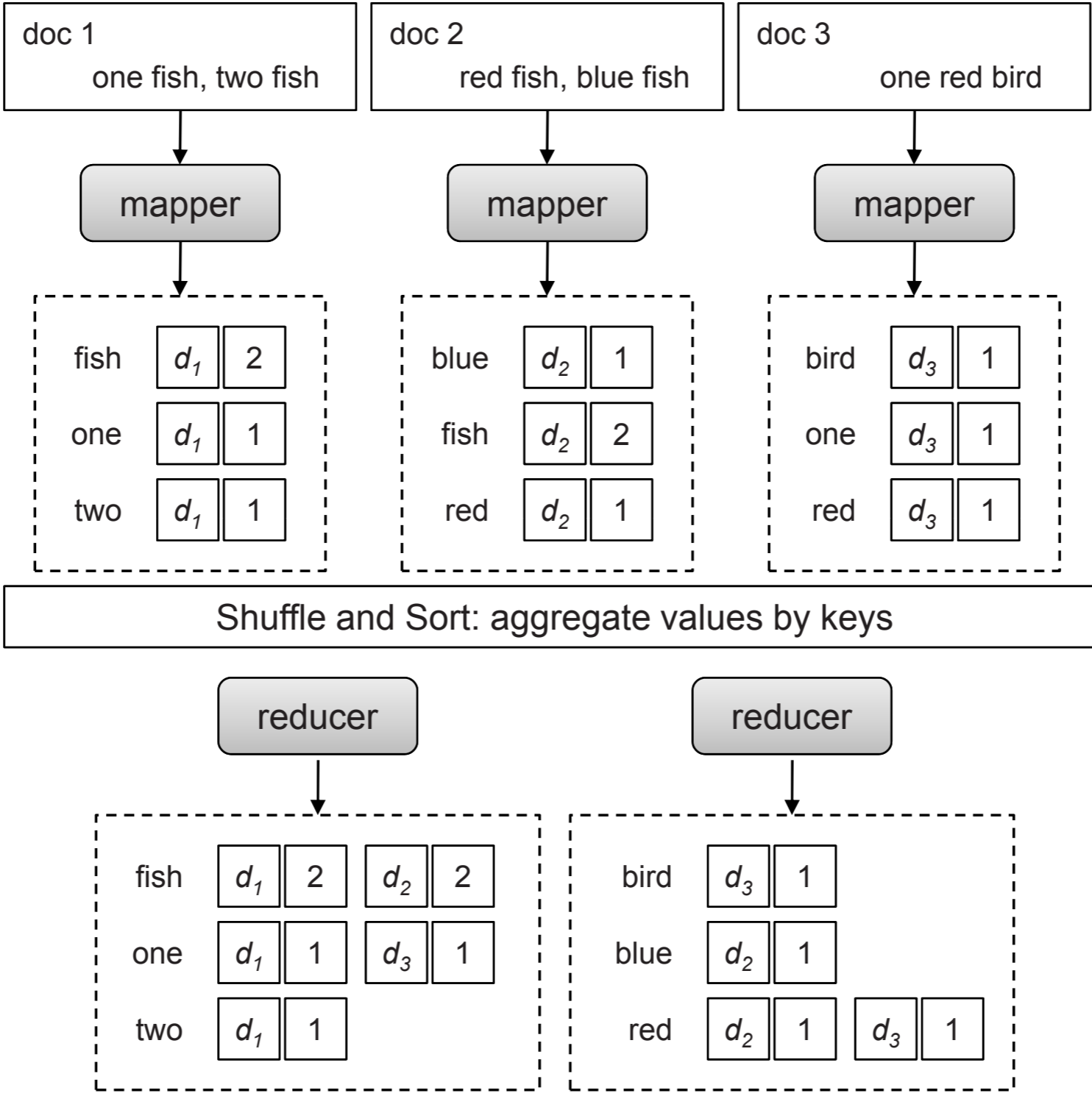Can we speed it up by distributing the work?

# Of course we can!



This is what MapReduce was *invented* for...

# Of course we can!

```
1:  class Mapper
2:      procedure Map(docid n, doc d)
3:          H ← new AssociativeArray
4:          for all term t ∈ doc d do
5:              H{t} ← H{t} + 1
6:          for all term t ∈ H do
7:              Emit(term t, posting ⟨n, H{t}⟩)

1:  class Reducer
2:      procedure Reduce(term t, postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ . . .])
3:          P ← new List
4:          for all posting ⟨a, f⟩ ∈ postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ . . .] do
5:              P.Add(⟨a, f⟩)
6:          P.Sort()
7:          Emit(term t, postings P)
```

```
1: class MAPPER
2:     procedure MAP(docid n, doc d)
3:         H ← new ASSOCIATIVEARRAY
4:         for all term t ∈ doc d do
5:             H{t} ← H{t} + 1
6:         for all term t ∈ H do
7:             EMIT(term t, posting ⟨n, H{t}⟩)

1: class REDUCER
2:     procedure REDUCE(term t, postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ . . .])
3:         P ← new LIST
4:         for all posting ⟨a, f⟩ ∈ postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ . . .] do
5:             P.ADD(⟨a, f⟩)
6:         P.SORT()
7:         EMIT(term t, postings P)
```

```
 1: class MAPPER
 2:     method MAP(docid n, doc d)
 3:         H ← new ASSOCIATIVEARRAY
 4:         for all term t ∈ doc d do
 5:             H{t} ← H{t} + 1
 6:         for all term t ∈ H do
 7:             EMIT(tuple ⟨t, n⟩, tf H{t})

 1: class REDUCER
 2:     method INITIALIZE
 3:         t_{prev} ← ∅
 4:         P ← new POSTINGSLIST
 5:     method REDUCE(tuple ⟨t, n⟩, tf [f])
 6:         if t ≠ t_{prev} ∧ t_{prev} ≠ ∅ then
 7:             EMIT(term t, postings P)
 8:             P.RESET()
 9:         P.ADD(⟨n, f⟩)
10:         t_{prev} ← t
11:     method CLOSE
12:         EMIT(term t, postings P)
```

Being clever with our keys lets the runtime do more of the work.

Remember, keys arrive at the reducer in sorted order, guaranteed.

For very large corpora, it's important to be clever about how we store our postings lists on disk...

Out of scope for today, but worth reading Lin & Dyer's summary.

# What about retrieval?

There are two main strategies for breaking up an index onto multiple machines:

1. Document sharding

2. Term sharding

| | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ | $d_7$ | $d_8$ | $d_9$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| $t_1$ | | | 2 | | | | | 3 | | |
| $t_2$ | | | 1 | | | 1 | | | 4 | partition$_a$ |
| $t_3$ | 1 | 1 | | | 2 | | | | | |
| $t_4$ | | | | 5 | | | | 2 | 2 | |
| $t_5$ | | | 1 | | | 1 | 3 | | | partition$_b$ |
| $t_6$ | 2 | | | | 1 | | | | | |
| $t_7$ | | | 2 | | 1 | | | 4 | | |
| $t_8$ | | 1 | | | | 2 | 3 | | | partition$_c$ |
| $t_9$ | | | 1 | | | 2 | | | 1 | |
| | partition$_1$ | | | partition$_2$ | | | partition$_3$ | | | |

In document sharding, each machine has a portion of the documents, and queries are processed by each machine.

# What about retrieval?

There are two main strategies for breaking up an index onto multiple machines:

1. Document sharding

2. Term sharding

| | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ | $d_7$ | $d_8$ | $d_9$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| $t_1$ | | | 2 | | | | | 3 | | |
| $t_2$ | | | 1 | | | 1 | | | 4 | $partition_a$ |
| $t_3$ | 1 | 1 | | | 2 | | | | | |
| $t_4$ | | | | 5 | | | | 2 | 2 | |
| $t_5$ | | | 1 | | | 1 | 3 | | | $partition_b$ |
| $t_6$ | 2 | | | | 1 | | | | | |
| $t_7$ | | | 2 | | 1 | | | 4 | | |
| $t_8$ | | 1 | | | | 2 | 3 | | | $partition_c$ |
| $t_9$ | | | 1 | | | 2 | | | 1 | |
| | $partition_1$ | | | $partition_2$ | | | $partition_3$ | | | |

In term sharding, each machine has a the full posting list for a subset of the terms, and only sees the part of the query that is relevant to it.

# What about retrieval?

There are tradeoffs to each:

*Document sharding* is simpler to process, but requires a lot of wasted work...

- Each shard can process independent queries

- Easy to keep around more per-document information

- Each query must involve every shard ($O(k*N)$ disk accesses, for $k$ query terms and $N$ document shards)

# What about retrieval?

There are tradeoffs to each:

*Term sharding* is arguably more efficient to process, but requires a lot more complexity.

    - *Far* less time spent doing I/O to access index (only $k$ shards are involved for a $k$-word query)...

    - Harder to have per-doc word-level information

    - Much more network bandwidth needed (data from each matching doc must be aggregated)

# What about retrieval?

There's no single right answer for all cases, but in general, document sharding seems better...

... largely due to lower search latency, which is the most relevant metric.

# Game plan for today:

Quick overview of inverted indexing

Inverted indexing & Map-Reduce

Quick MT overview

Map-Reduce & MT Model Estimation

# Quick MT overview:

$$P(f_1^m|e_1^l) = \sum_{a_1^m} P(f_1^m, a_1^m|e_1^l)$$

Foreign Sentence

$$= \sum_{a_1^m} P(a_1^m|e_1^l, f_1^m) \prod_{j=1}^{m} P(f_j|e_{a_j})$$

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# Quick MT overview:

$$P(f_1^m|e_1^l) = \sum_{a_1^m} P(f_1^m, a_1^m|e_1^l)$$

$$= \sum_{a_1^m} P(a_1^m|e_1^l, f_1^m) \prod_{j=1}^{m} P(f_j|e_{a_j})$$

Original Sentence

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# Quick MT overview:

$$P(f_1^m|e_1^l) \ = \ \sum_{a_1^m} P(f_1^m, a_1^m|e_1^l)$$

$$= \ \sum_{a_1^m} P(a_1^m|e_1^l, f_1^m) \prod_{j=1}^{m} P(f_j|e_{a_j})$$

Language model

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# Quick MT overview:

$$P(f_1^m|e_1^l) = \sum_{a_1^m} P(f_1^m, a_1^m|e_1^l)$$

$$= \sum_{a_1^m} P(a_1^m|e_1^l, f_1^m) \prod_{j=1}^{m} P(f_j|e_{a_j})$$

Translation model

$$\hat{a}_1^m = \arg\max_{a_1^m} P(a_1^m|e_1^l, f_1^m) \prod_{j=1}^{m} P(f_j|e_{a_j})$$

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# Quick MT overview:

$$P(f_1^m | e_1^l) = \sum_{a_1^m} P(f_1^m, a_1^m | e_1^l)$$

$$= \sum_{a_1^m} P(a_1^m | e_1^l, f_1^m) \prod_{j=1}^m P(f_j | e_{a_j})$$

1-best translation guess

$$\hat{a}_1^m = \arg\max_{a_1^m} P(a_1^m | e_1^l, f_1^m) \prod_{j=1}^m P(f_j | e_{a_j})$$

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# More data == better LM (and MT)



Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# More data == better LM (and MT)



Large Language Models in Machine Translation. Thorsten Brants et. al., Proc. Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP), pg. 858-867, 2007.

# How to estimate the LM?

$$P(w_1^L) = \prod_{i=1}^{L} P(w_i|w_1^{i-1}) \approx \prod_{i=1}^{L} \hat{P}(w_i|w_{i-n+1}^{i-1})$$

$$P_{MLE}(B|A) = \frac{c(A,B)}{c(A)} = \frac{c(A,B)}{\sum_{B'} c(A,B')}$$

When corpora get big...

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

Large Language Models in Machine Translation. Thorsten Brants et. al., Proc. Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP), pg. 858-867, 2007.

# MLE estimation in Map-Reduce:

| **Method 1** | |
|---|---|
| $\text{Map}_1$ | $\langle A, B \rangle \rightarrow \langle \langle A, B \rangle, 1 \rangle$ |
| $\text{Reduce}_1$ | $\langle \langle A, B \rangle, c(A, B) \rangle$ |
| $\text{Map}_2$ | $\langle \langle A, B \rangle, c(A, B) \rangle \rightarrow \langle \langle A, ^* \rangle, c(A, B) \rangle$ |
| $\text{Reduce}_2$ | $\langle \langle A, ^* \rangle, c(A) \rangle$ |
| $\text{Map}_3$ | $\langle \langle A, B \rangle, c(A, B) \rangle \rightarrow \langle A, \langle B, c(A, B) \rangle \rangle$ |
| $\text{Reduce}_3$ | $\langle A, \langle B, \frac{c(A,B)}{c(A)} \rangle \rangle$ |

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# MLE estimation in Map-Reduce:

**Method 1**

| | |
|---|---|
| $\text{Map}_1$ | $\langle A, B \rangle \to \langle \langle A, B \rangle, 1 \rangle$ |
| $\text{Reduce}_1$ | $\langle \langle A, B \rangle, c(A, B) \rangle$ |
| $\text{Map}_2$ | $\langle \langle A, B \rangle, c(A, B) \rangle \to \langle \langle A, {}^* \rangle, c(A, B) \rangle$ |
| $\text{Reduce}_2$ | $\langle \langle A, {}^* \rangle, c(A) \rangle$ |
| $\text{Map}_3$ | $\langle \langle A, B \rangle, c(A, B) \rangle \to \langle A, \langle B, c(A, B) \rangle \rangle$ |
| $\text{Reduce}_3$ | $\langle A, \langle B, \frac{c(A,B)}{c(A)} \rangle \rangle$ |

**Method 2**

| | |
|---|---|
| $\text{Map}_1$ | $\langle A, B \rangle \to \langle \langle A, B \rangle, 1 \rangle; \langle \langle A, {}^* \rangle, 1 \rangle$ |
| $\text{Reduce}_1$ | $\langle \langle A, B \rangle, \frac{c(A,B)}{c(A)} \rangle$ |

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides
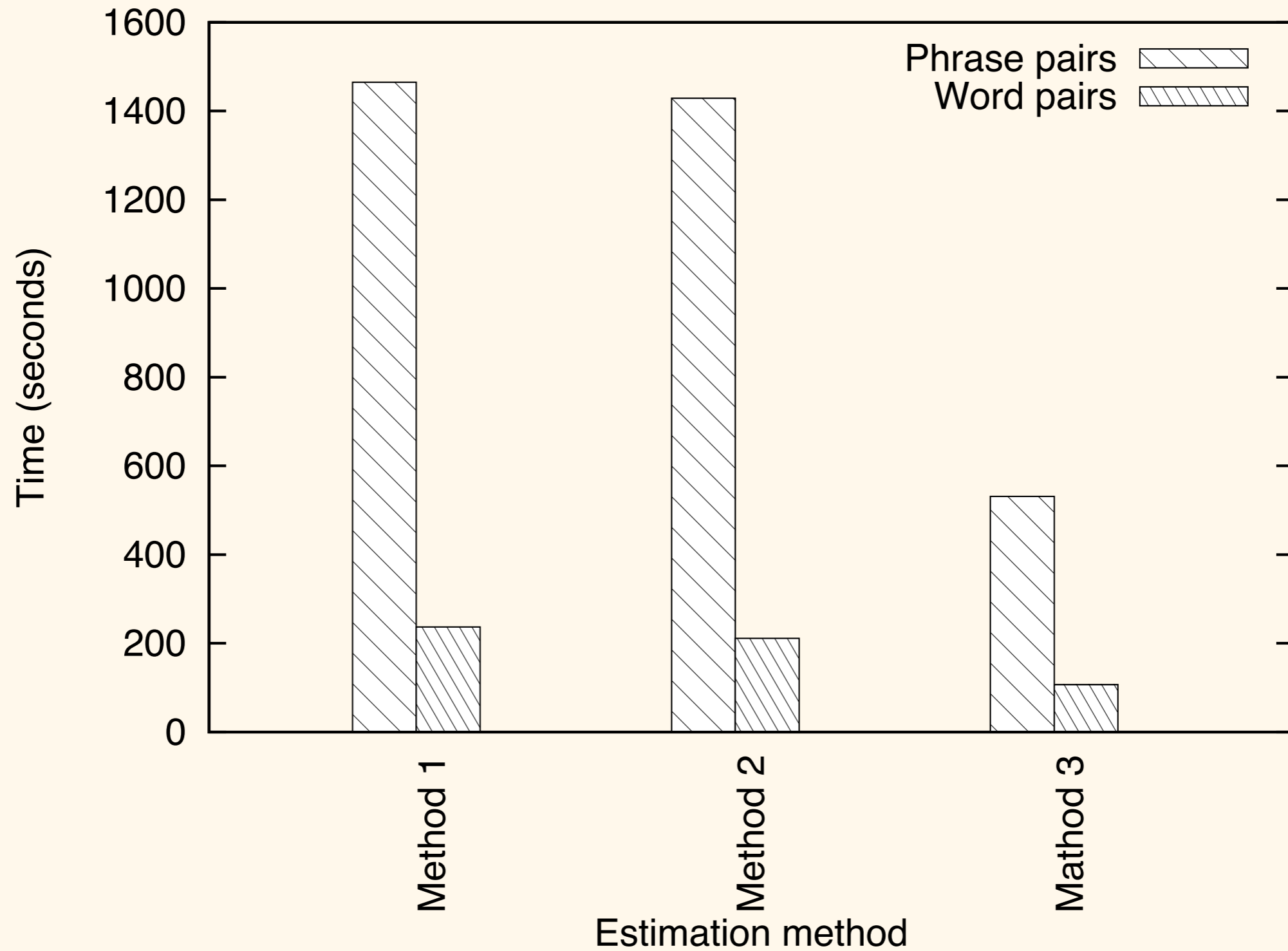
# on in Map-Reduce:

**Method 1**

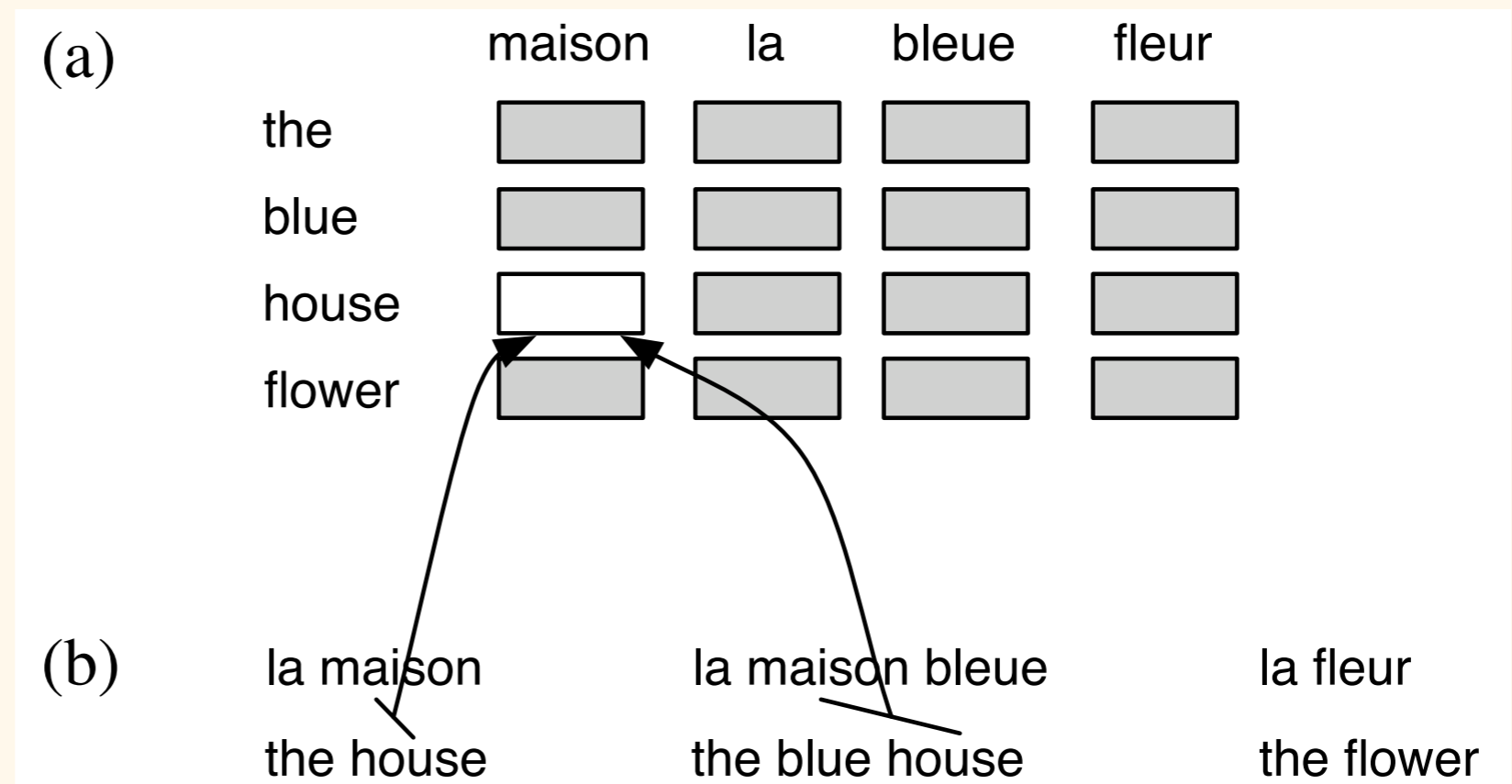| Map$_1$ | $\langle A, B \rangle \rightarrow \langle \langle A, B \rangle, 1 \rangle$ |
|---|---|
| Reduce$_1$ | $\langle \langle A, B \rangle, c(A, B) \rangle$ |
| Map$_2$ | $\langle \langle A, B \rangle, c(A, B) \rangle \rightarrow \langle \langle A, {}^* \rangle, c(A, B) \rangle$ |
| Reduce$_2$ | $\langle \langle A, {}^* \rangle, c(A) \rangle$ |
| Map$_3$ | $\langle \langle A, B \rangle, c(A, B) \rangle \rightarrow \langle A, \langle B, c(A, B) \rangle \rangle$ |
| Reduce$_3$ | $\langle A, \langle B, \frac{c(A,B)}{c(A)} \rangle \rangle$ |

**Method 2**

| Map$_1$ | $\langle A, B \rangle \rightarrow \langle \langle A, B \rangle, 1 \rangle; \langle \langle A, {}^* \rangle, 1 \rangle$ |
|---|---|
| Reduce$_1$ | $\langle \langle A, B \rangle, \frac{c(A,B)}{c(A)} \rangle$ |

**Method 3**

| Map$_1$ | $\langle A, B_i \rangle \rightarrow \langle A, \langle B_i : 1 \rangle \rangle$ |
|---|---|
| Reduce$_1$ | $\langle A, \langle B_1 : \frac{c(A,B_1)}{c(A)} \rangle, \langle B_2 : \frac{c(A,B_2)}{c(A)} \rangle \cdots \rangle$ |

map

Fast, Easy, and Cheap: Construction of Statistical Machine Translation Models with MapReduce, Christopher Dyer et. al., Proc. ACL Workshop on Statistical Machine Translation, pg. 199-207, 2008. slides

# We can use a similar method for estimating parameters for alignments:

We can use a similar method for estimating parameters for alignments:

$$P(a_1^m | e_1^l, f_1^m)$$

(a)

|        | maison | la | bleue | fleur |
|--------|--------|-----|-------|-------|
| the    |        |     |       |       |
| blue   |        |     |       |       |
| house  |        |     |       |       |
| flower |        |     |       |       |

(b)

la maison          la maison bleue          la fleur

the house          the blue house           the flower

*Figure: Time (seconds) vs Corpus size (sentences) comparing Optimal Model 1 (Giza/38), Optimal HMM (Giza/38), MapReduce Model 1 (38 M/R), and MapReduce HMM (38 M/R).*

There

is — lots
is — many
may — be
was
are — lots of
are — ε
are — many
seem
be
were — many
were — concerns

lots of — people
lots of — reasons
ε — onlookers
many — people
many — people
many — ε

*t-1*          *t*          *t+1*          *t+2*