

بسم الله الرحمن الرحيم

CMPS 653

Big Data Analytics

Session 4:

MapReduce & Unstructured Data (Text)

Some slides are adopted from Jimmy Lin's by permission



Tamer Elsayed

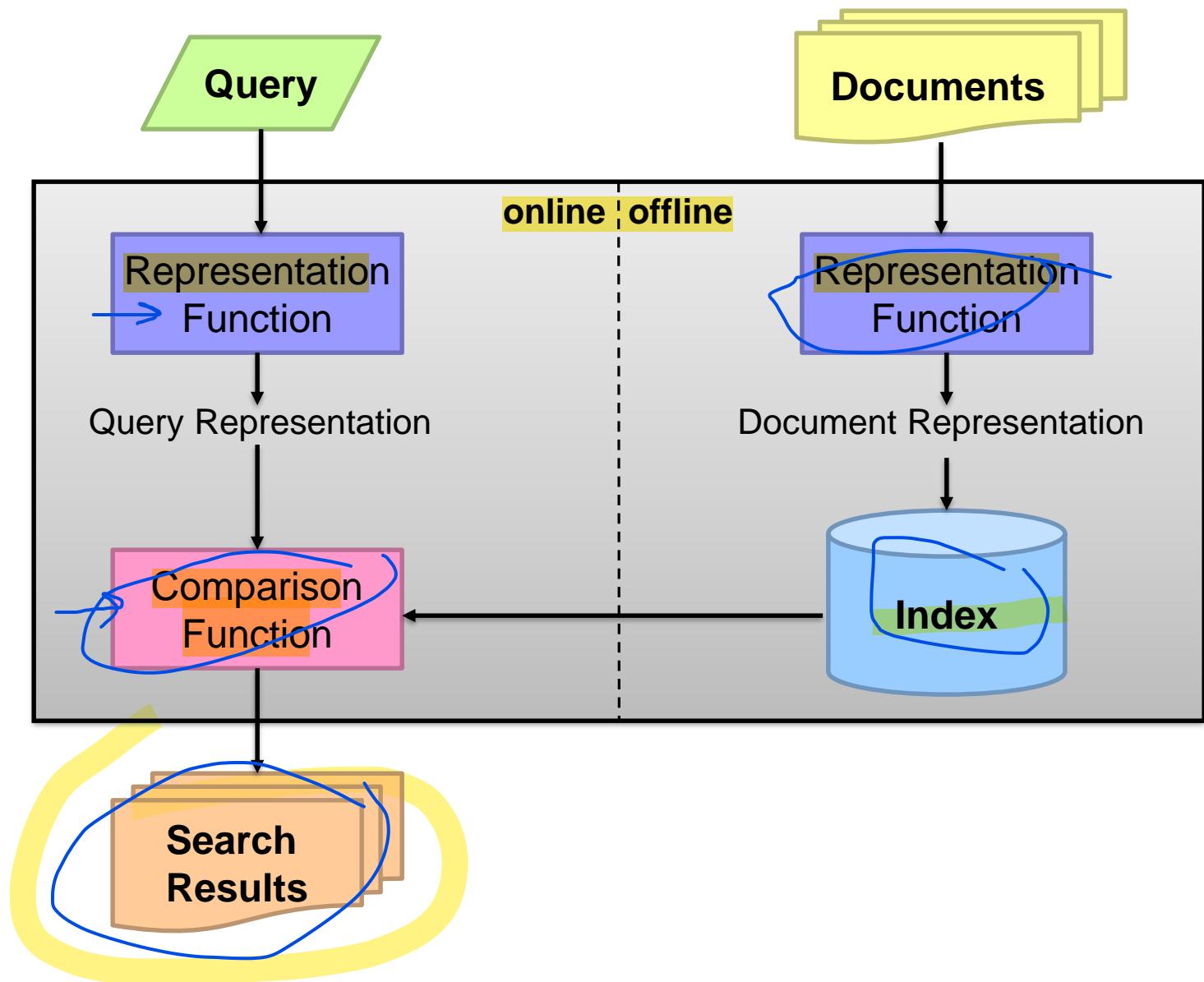


INFORMATION RETRIEVAL (IR) IN A NUTSHELL!

Information Retrieval is all about ...



Learn IR in 2 Minutes 😊





Indexing

Inverted Index: Boolean

Doc 1

one fish, two fish

Doc 2

red fish, blue fish

Doc 3

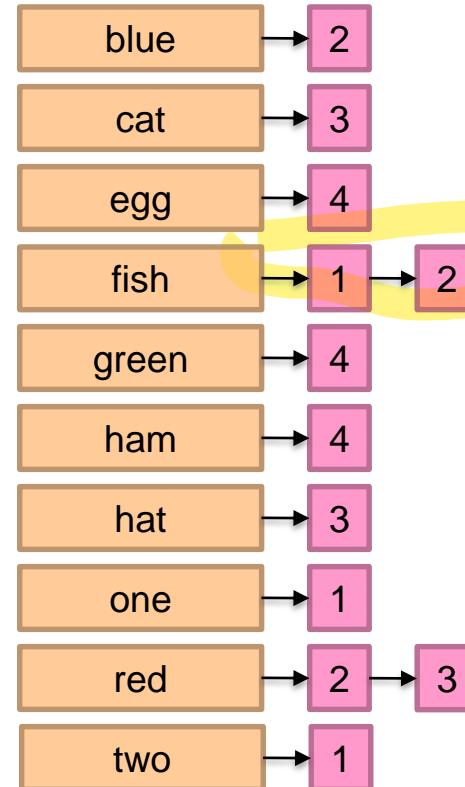
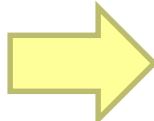
red cat in the hat

Doc 4

green eggs and ham

1 2 3 4

blue		1		
cat			1	
egg				1
fish	1	1		
green				1
ham				1
hat			1	
one	1			
red		1	1	
two	1			



Inverted Index: TF & DF

Doc 1
one fish, two fish

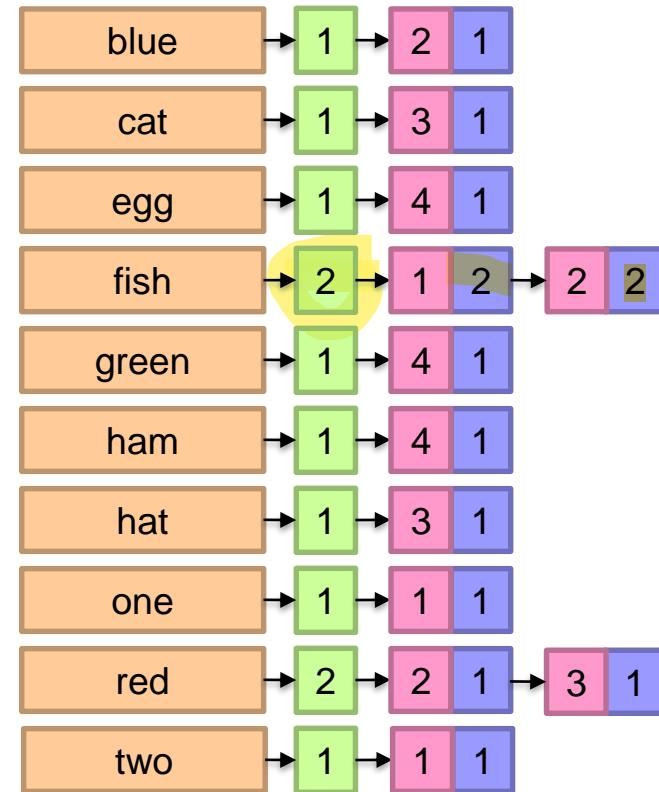
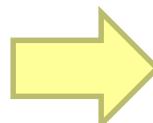
Doc 2

red fish, blue fish

Doc 3

Doc 4 green eggs and ham

	1	2	3	4	tf
blue		1			1
cat			1		1
egg				1	1
fish	2	2			2
green				1	1
ham				1	1
hat			1		1
one	1				1
red		1	1		2
two	1				1



Inverted Index: Positions

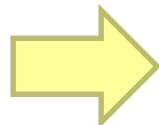
Doc 1
one fish, two fish

Doc 2
red fish, blue fish

Doc 3
red cat in the hat

Doc 4
green eggs and ham

	tf				df
	1	2	3	4	
blue		1			1
cat			1		1
egg				1	1
fish	2	2			2
green				1	1
ham				1	1
hat			1		1
one	1				1
red		1	1		2
two	1				1

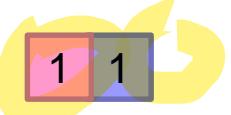


blue	→ 1 →	2	1	[3]
cat	→ 1 →	3	1	[2]
egg	→ 1 →	4	1	[2]
fish	→ 2 →	1	2	[2,4] → 2 2 [2,4]
green	→ 1 →	4	1	[1]
ham	→ 1 →	4	1	[3]
hat	→ 1 →	3	1	[3]
one	→ 1 →	1	1	[1]
red	→ 1 →	2	1	[1] → 3 1 [1]
two	→ 1 →	1	1	[3]

Indexing with MapReduce

Inverted Indexing

Doc 1
one fish, two fish

one	
two	
fish	

Doc 2
red fish, blue fish

red	
blue	
fish	

Doc 3
red cat in the hat

red	
cat	
hat	

Map

Shuffle and Sort: aggregate values by keys

Reduce

cat	
fish	
one	
red	

blue	
hat	
two	

MapReduce: Index Construction

- Map over all documents
 - Emit *term* as key, $(docno, tf)$ as value
 - Emit other information as necessary (e.g., term position)
- Sort/shuffle: group postings by term
- Reduce
 - Gather and sort the postings (e.g., by *docno* or *tf*)
 - Write postings to disk
- MapReduce does all the heavy lifting!

Inverted Indexing: Pseudo-Code

```

1: class MAPPER
2:   method MAP(docid n, doc d)
3:     H  $\leftarrow$  new ASSOCIATIVEARRAY ▷ histogram to hold term frequencies
4:     for all term t  $\in$  doc d do ▷ processes the doc, e.g., tokenization and stopword removal
5:       H{t}  $\leftarrow$  H{t} + 1
6:     for all term t  $\in$  H do
7:       EMIT(term t, posting  $\langle n, H\{t\} \rangle$ ) ▷ emits individual postings
↓

1: class REDUCER
2:   method REDUCE(term t, postings [ $\langle n_1, f_1 \rangle \dots$ ])
3:     P  $\leftarrow$  new LIST
4:     for all  $\langle n, f \rangle \in$  postings [ $\langle n_1, f_1 \rangle \dots$ ] do
5:       P.APPEND( $\langle n, f \rangle$ ) ▷ appends postings unsorted
6:     P.SORT() ▷ sorts for compression
7:     EMIT(term t, postingsList P)

```

Positional Indexes

Doc 1
one fish, two fish

one	1	1	[1]
two	1	1	[3]
fish	1	2	[2,4]

Doc 2
red fish, blue fish

red	2	1	[1]
blue	2	1	[3]
fish	2	2	[2,4]

Doc 3
red cat in the hat

red	3	1	[1]
cat	3	1	[2]
hat	3	1	[3]

Map

Shuffle and Sort: aggregate values by keys

Reduce

cat	3	1	[2]
fish	1	2	[2,4]
one	1	1	[1]
red	2	1	[1]

blue	2	1	[3]
hat	3	1	[3]
two	1	1	[3]

What's the Problem Here?!

```

1: class MAPPER
2:   method MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY           ▷ histogram to hold term frequencies
4:     for all term  $t \in$  doc  $d$  do      ▷ processes the doc, e.g., tokenization and stopword removal
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , posting  $\langle n, H\{t\} \rangle$ )           ▷ emits individual postings

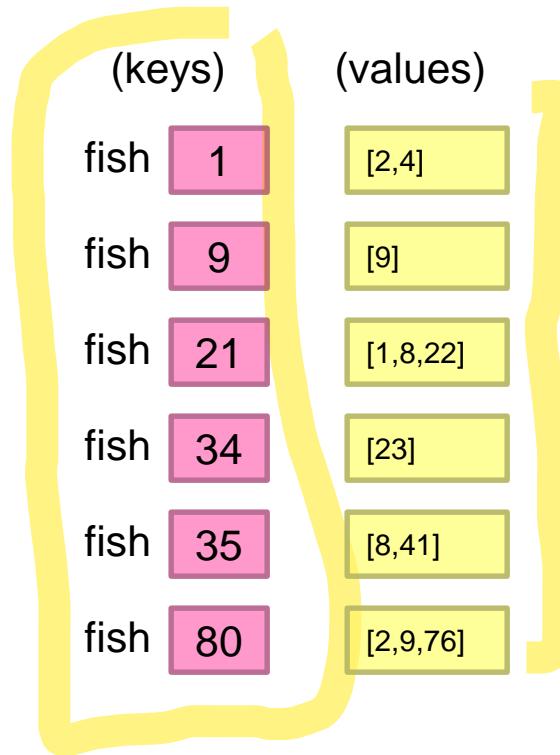
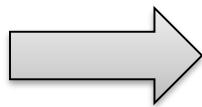
1: class REDUCER
2:   method REDUCE(term  $t$ , postings [ $\langle n_1, f_1 \rangle \dots$ ])
3:      $P \leftarrow$  new LIST
4:     for all  $\langle n, f \rangle \in$  postings [ $\langle n_1, f_1 \rangle \dots$ ] do
5:        $P.\text{APPEND}(\langle n, f \rangle)$            ▷ appends postings unsorted
6:     P.SORT()                                ▷ sorts for compression
7:     EMIT(term  $t$ , postingsList  $P$ )
  
```

Scalability Bottleneck!

- Reducers must buffer all postings assoc. with key (to sort)
- What if we run out of memory to buffer postings?

Another Try...

(key)	(values)
fish	1 2 [2,4]
	34 1 [23]
	21 3 [1,8,22]
	35 2 [8,41]
	80 3 [2,9,76]
	9 1 [9]



How is this different?

Let the framework do the sorting



Where have we seen this before?

Inverted Indexing: Pseudo-Code

```

1: class MAPPER
2:   method MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do                                 $\triangleright$  builds a histogram of term frequencies
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(tuple  $\langle t, n \rangle$ , tf  $H\{t\}$ )           $\triangleright$  emits individual postings, with a tuple as the key
1: class PARTITIONER
2:   method PARTITION(tuple  $\langle t, n \rangle$ , tf  $f$ )
3:     return HASH( $t$ ) mod  $NumOfReducers$             $\triangleright$  keys of same term are sent to same reducer

```

Inverted Indexing: Pseudo-Code

```

1: class MAPPER
2:   method MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do                                 $\triangleright$  builds a histogram of term frequencies
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(tuple  $\langle t, n \rangle$ , tf  $H\{t\}$ )            $\triangleright$  emits individual postings, with a tuple as the key

1: class PARTITIONER
2:   method PARTITION(tuple  $\langle t, n \rangle$ , tf  $f$ )
3:     return HASH( $t$ ) mod  $NumOfReducers$            $\triangleright$  keys of same term are sent to same reducer

1: class REDUCER
2:   method INITIALIZE
3:      $t_{prev} \leftarrow \emptyset$ 
4:      $P \leftarrow$  new POSTINGSLIST
5:   method REDUCE(tuple  $\langle t, n \rangle$ , tf  $[f]$ )
6:     if  $t \neq t_{prev} \wedge t_{prev} \neq \emptyset$  then
7:       EMIT(term  $t$ , postings  $P$ )                   $\triangleright$  emits postings list of term  $t_{prev}$ 
8:        $P.RESET()$ 
9:        $P.APPEND(\langle n, f \rangle)$                        $\triangleright$  appends postings in sorted order
10:       $t_{prev} \leftarrow t$ 
11:   method CLOSE
12:     EMIT(term  $t$ , postings  $P$ )            $\triangleright$  emits last postings list from this reducer

```

Chicken and Egg?

(key)	(value)
fish	1 [2,4]
fish	9 [9]
fish	21 [1,8,22]
fish	34 [23]
fish	35 [8,41]
fish	80 [2,9,76]

...



We' need df for compression
of the postings list

But we don't know the df until
we've seen all postings!

Write postings

Sound familiar?
₁₉

Getting the *df*

- In the mapper:
 - Emit “special” key-value pairs to keep track of *df*
- In the reducer:
 - Make sure “special” key-value pairs come first: process them to determine *df*
- Remember: proper partitioning!

Getting the *df*: Modified Mapper

Doc 1

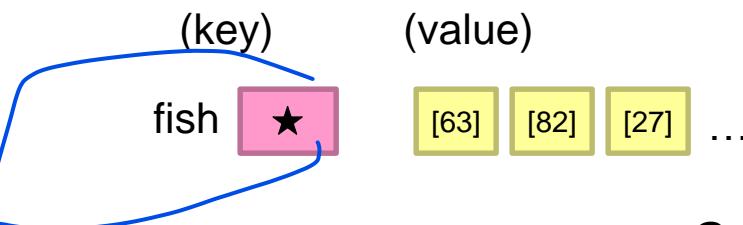
one fish, two fish

Input document...

(key) (value)

fish	1	[2,4]	Emit normal key-value pairs...
one	1	[1]	
two	1	[3]	
fish	★	[1]	Emit “special” key-value pairs to keep track of <i>df</i> ...
one	★	[1]	
two	★	[1]	

Getting the *df*: Modified Reducer



First, compute the df by summing contributions from all “special” key-value pair...

Get *df*...

fish	1	[2,4]
fish	9	[9]
fish	21	[1,8,22]
fish	34	[23]
fish	35	[8,41]
fish	80	[2,9,76]

Important: properly define sort order to make sure “special” key-value pairs come first!

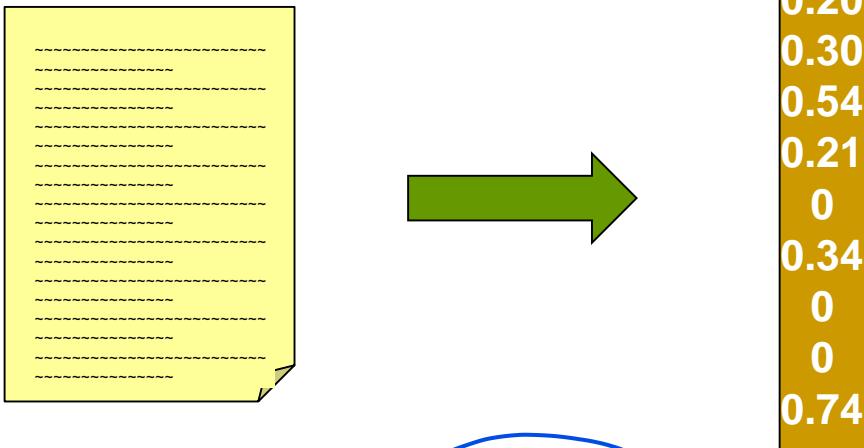
Write postings

Where have we seen this before?

Pairwise Document Similarity

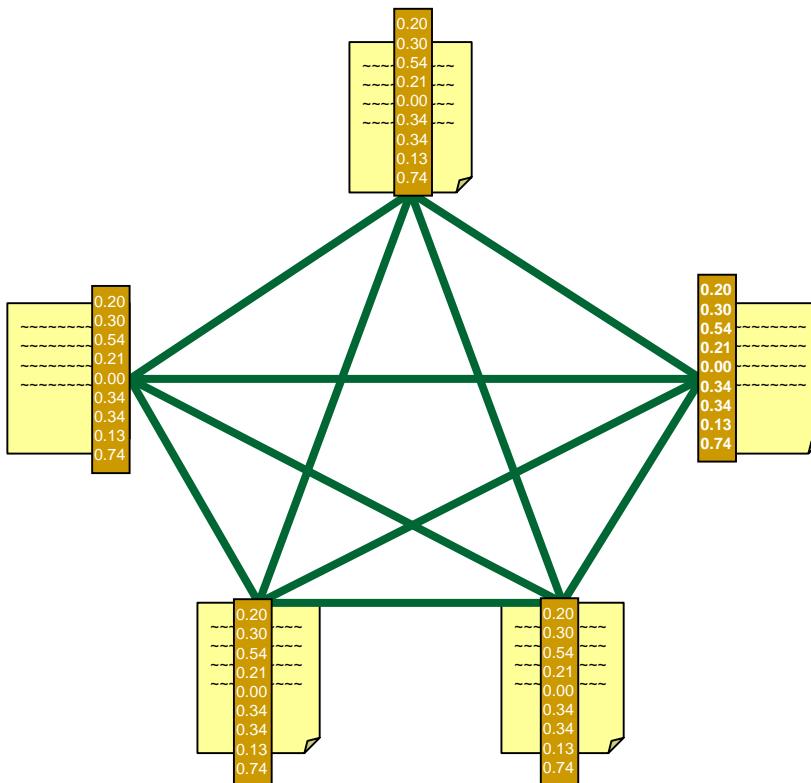


Vector Space Model (VSM)



- Represent documents as vectors
- Each term has a weight (0 if not exist)
 - term weights indicate importance
- Sparse
- [Generally: feature vector]

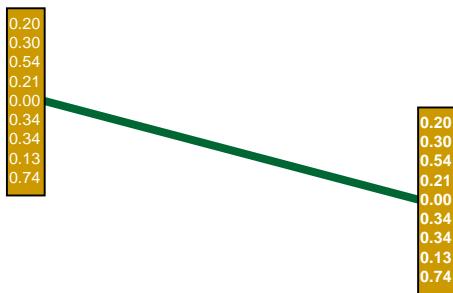
Abstract Problem



- Applications:

- Near-Duplicate Detection
- Clustering
- Friends recommendation in Facebook/Twitter
- “more-like-that” queries

Similarity in VSM



$$sim(d_i, d_j) = \sum_{t \in V} w_{t,d_i} w_{t,d_j}$$

- Dot product between vectors

- Represent a large class of similarity metrics (e.g.,
cosine similarity, BM25, etc.)

Better Solution

$$sim(d_i, d_j) = \sum_{t \in V} w_{t,d_i} w_{t,d_j}$$

Each term contributes only if appears in $d_i \cap d_j$

$$sim(d_i, d_j) = \sum_{t \in d_i \cap d_j} w_{t,d_i} w_{t,d_j}$$

$$sim(d_i, d_j) = \sum_{t \in d_i \cap d_j} term_contrib(t, d_i, d_j)$$

Can we benefit from that?

Algorithm

Algorithm 1 Compute Pairwise Similarity Matrix

```

1:  $\forall i, j : sim[i, j] \leftarrow 0$ 
2: for all  $t \in V$  do
3:    $p_t \leftarrow postings(t)$ 
4:   for all  $d_i, d_j \in p_t$  do
5:      $sim[i, j] \leftarrow sim[i, j] + w_{t, d_i} \cdot w_{t, d_j}$ 

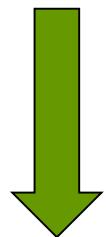
```

- Load weights for each term once
- Each term contributes $O(df_t^2)$ partial scores
- Matrix must fit in memory
 - Works for small collections
- Otherwise: disk access optimization

Large Collections!

Goal

**scalable and efficient solution
for large collections**



MapReduce ?!

Algorithm

Algorithm 1 Compute Pairwise Similarity Matrix

```

1:  $\forall i, j : sim[i, j] \leftarrow 0$  0 Indexing for each term in parallel
2: for all  $t \in V$  do
3:    $p_t \leftarrow postings(t)$  1 Partial contributions
4:   for all  $d_i, d_j \in p_t$  do
5:      $sim[i, j] \leftarrow sim[i, j] + w_{t, d_i} \cdot w_{t, d_j}$  2 Aggregation for each doc-pair in parallel

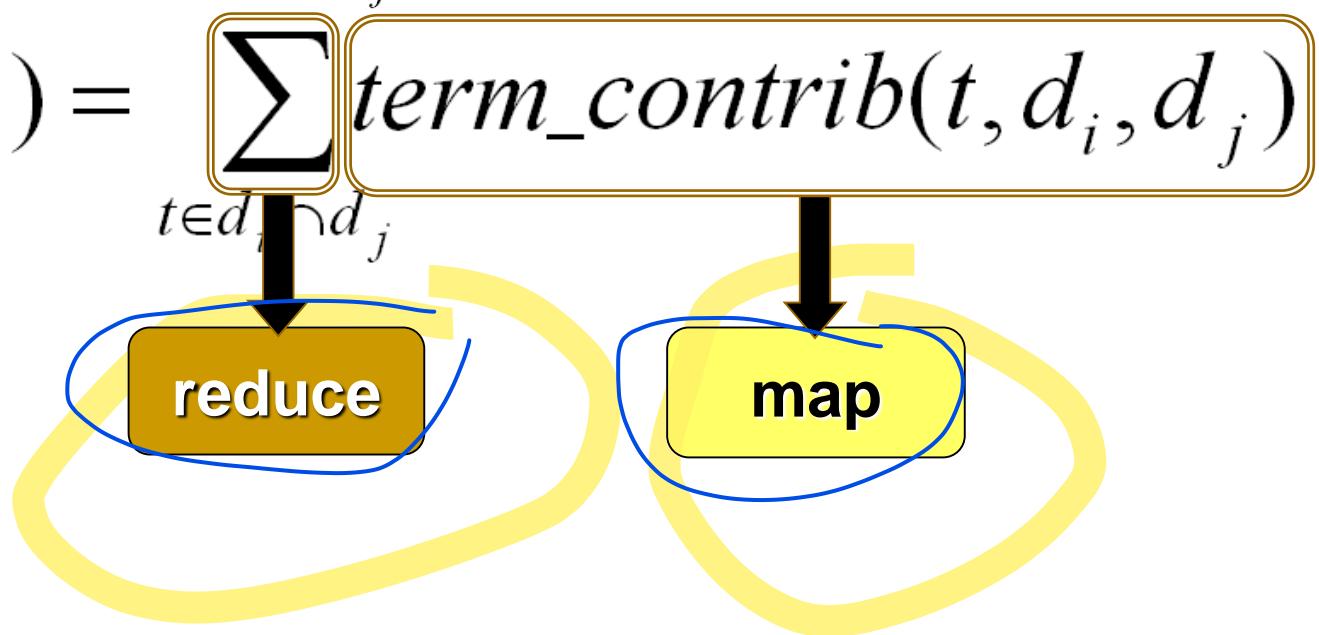
```

Decomposition

Each term contributes only if appears in $d_i \cap d_j$

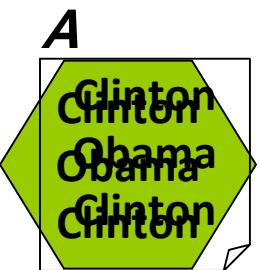
$$sim(d_i, d_j) = \sum_{t \in d_i \cap d_j} w_{t,d_i} w_{t,d_j}$$

$$sim(d_i, d_j) = \sum_{t \in d_i \cap d_j} term_contrib(t, d_i, d_j)$$



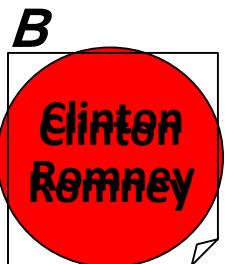
Step 1: Indexing

(a) Map



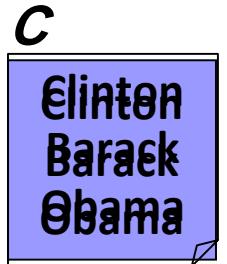
map

[Clinton ]
 [Obama ]
 [Clinton ]



map

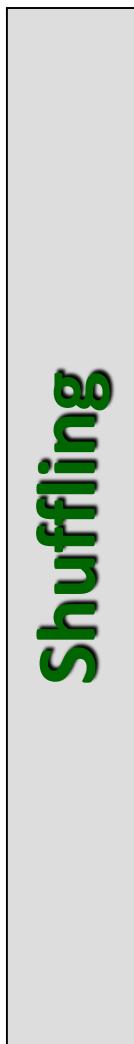
[Clinton ]
 [Romney ]



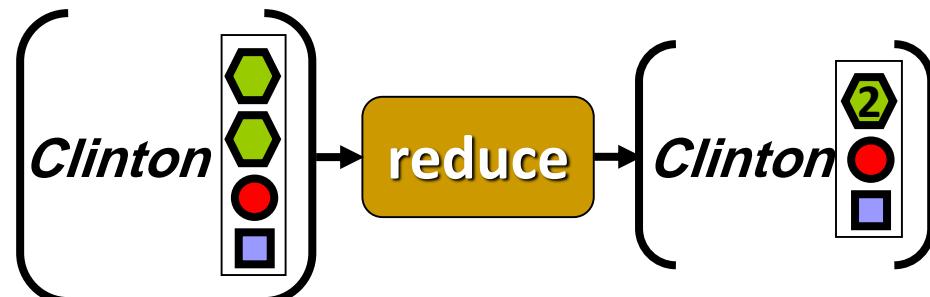
map

[Clinton ]
 [Barack ]
 [Obama ]

(b) Shuffle



(c) Reduce

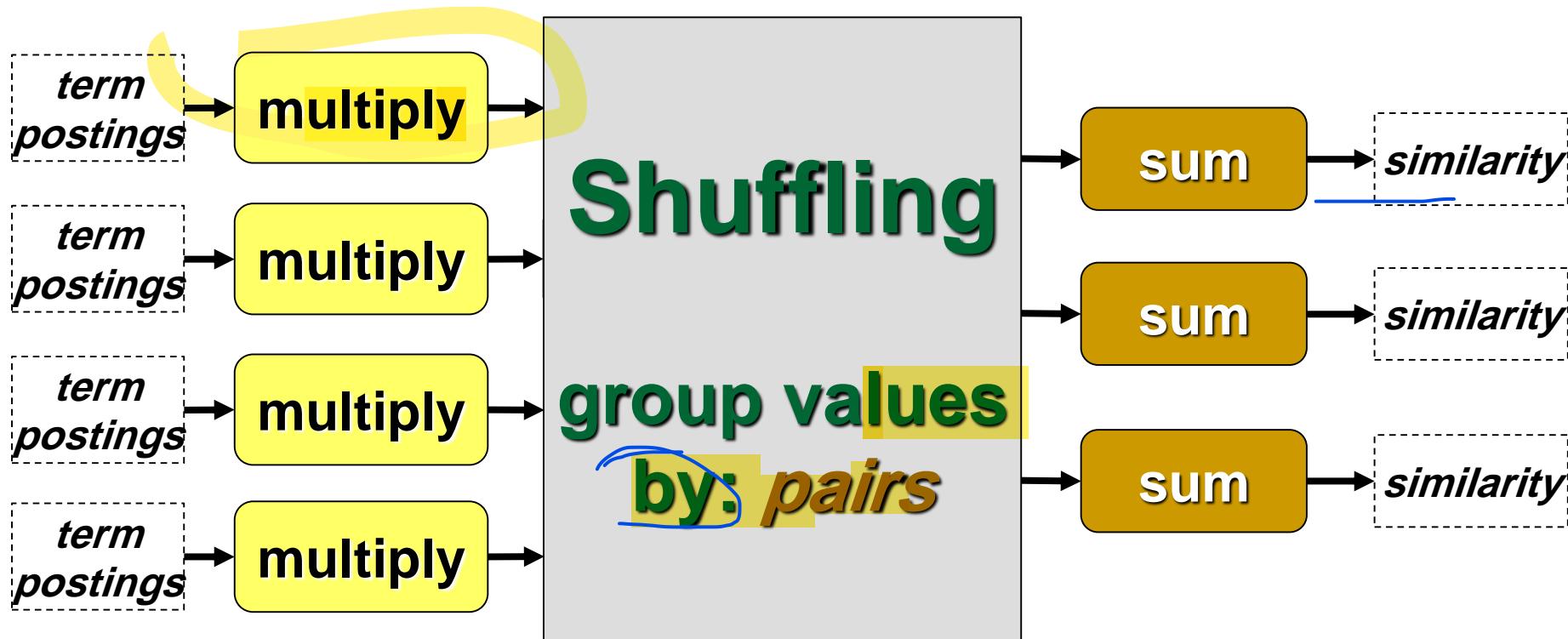


Step 2: Pairwise Similarity

(a) Generate pairs

(b) Group pairs

(c) Sum pairs

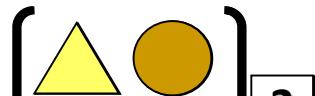
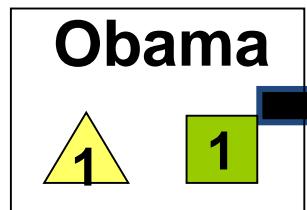
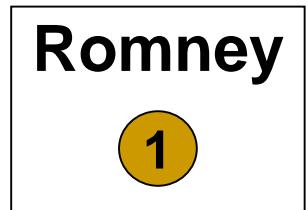
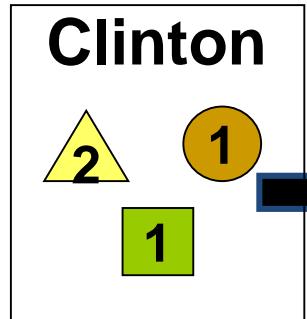


Step 2: Pairwise Similarity

(a) Generate pairs

(b) Group pairs

(c) Sum pairs



2



2



1



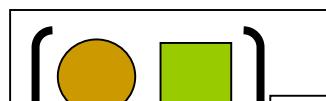
2



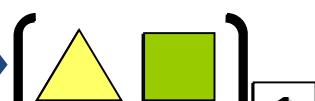
2



1



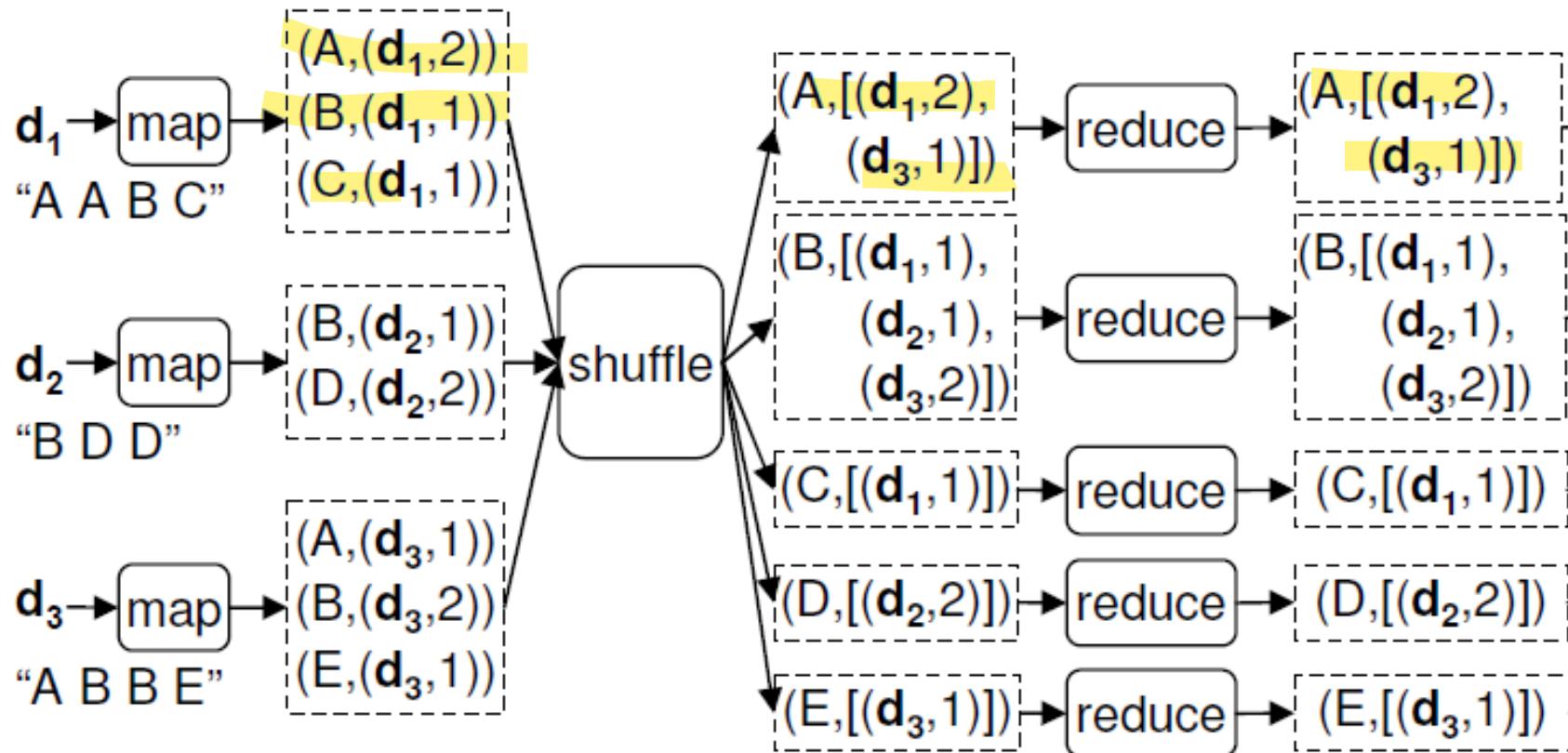
1



1

3

Another Ex.: Indexing Step



Another Ex: Pairwise Sim. Step

