

## *Inverted Indexing*



# Outline

- **Definition of an Inverted index**
- **Examples of Inverted Indices**
- **Representing an Inverted Index**
- **Processing a Query on a Linked Inverted Index**
- **Skip Pointers to Improve Merging**
- **Phrase Queries**
- **biwords**
- **Grammatical Tagging**
- **N-Grams**
- **Distributed Indexing**



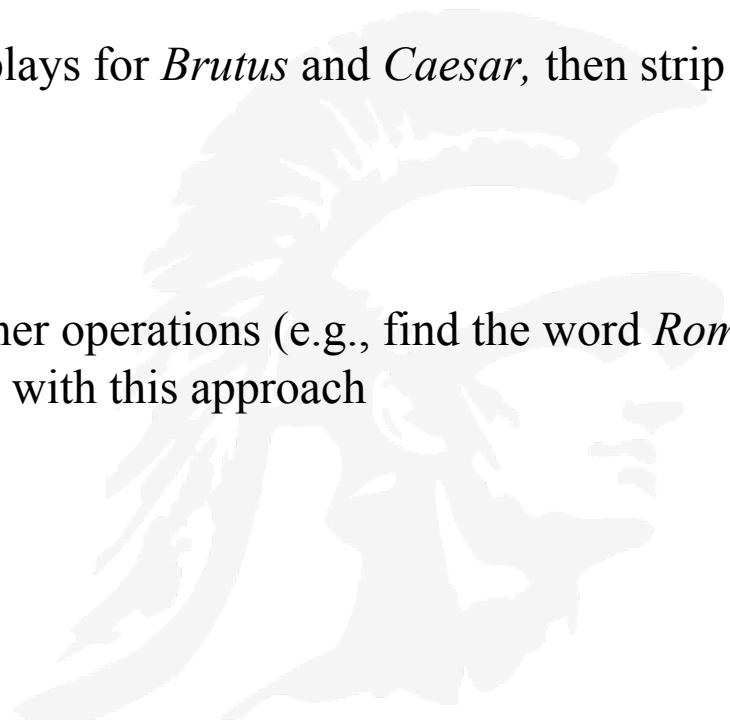
# Creating an Inverted Index

- An inverted index is typically composed of a vector containing all distinct words of the text collection in lexicographical order (which is called the **vocabulary**) and for each word in the vocabulary, a list of all documents (and text positions) in which that word occurs
  - This is nothing more than an index that one finds at the back of a book
- Terms in the inverted file index ***may be*** refined:
  - ***Case folding***: converting all uppercase letters to lower case
  - ***Stemming***: reducing words to their morphological roots
  - ***Stop words***: removing words that are so common they provide no information

# Processing a Query

## An Example

- The Query
  - Which plays of Shakespeare contain the words *Brutus* AND *Caesar* but NOT *Calpurnia*?
- One Possible Solution
  - One could grep all of Shakespeare's plays for *Brutus* and *Caesar*, then strip out lines containing *Calpurnia*?
    - Too Slow (for large corpora)
    - Requires lots of space
    - This method doesn't allow for other operations (e.g., find the word *Romans* near *countrymen*) are not feasible with this approach



# Term-Document Incidence Matrix

One way to think about an inverted index is to consider it as a sparse matrix where rows represent terms and columns represent documents

documents →	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	
terms ↗	Antony	1	1	0	0	0	1
terms ↗	Brutus	1	1	0	1	0	0
terms ↗	Caesar	1	1	0	1	1	1
terms ↗	Calpurnia	0	1	0	0	0	0
terms ↗	Cleopatra	1	0	0	0	0	0
terms ↗	mercy	1	0	1	1	1	1
terms ↗	worser	1	0	1	1	1	0

*The Query:*  
*Brutus AND Caesar but NOT Calpurnia*

1 if the play contains word, 0 otherwise

# Incidence Vectors

- So we have a 0/1 vector for each term.
- To answer the previous query: take the vectors for *Brutus*, *Caesar* and *Calpurnia* (complemented) and do a bitwise *AND*.
- $110100 \text{ AND } 110111 \text{ AND } 101111 = 100100$
- So the two plays matching the query are:  
“*Anthony and Cleopatra*”, “*Hamlet*”

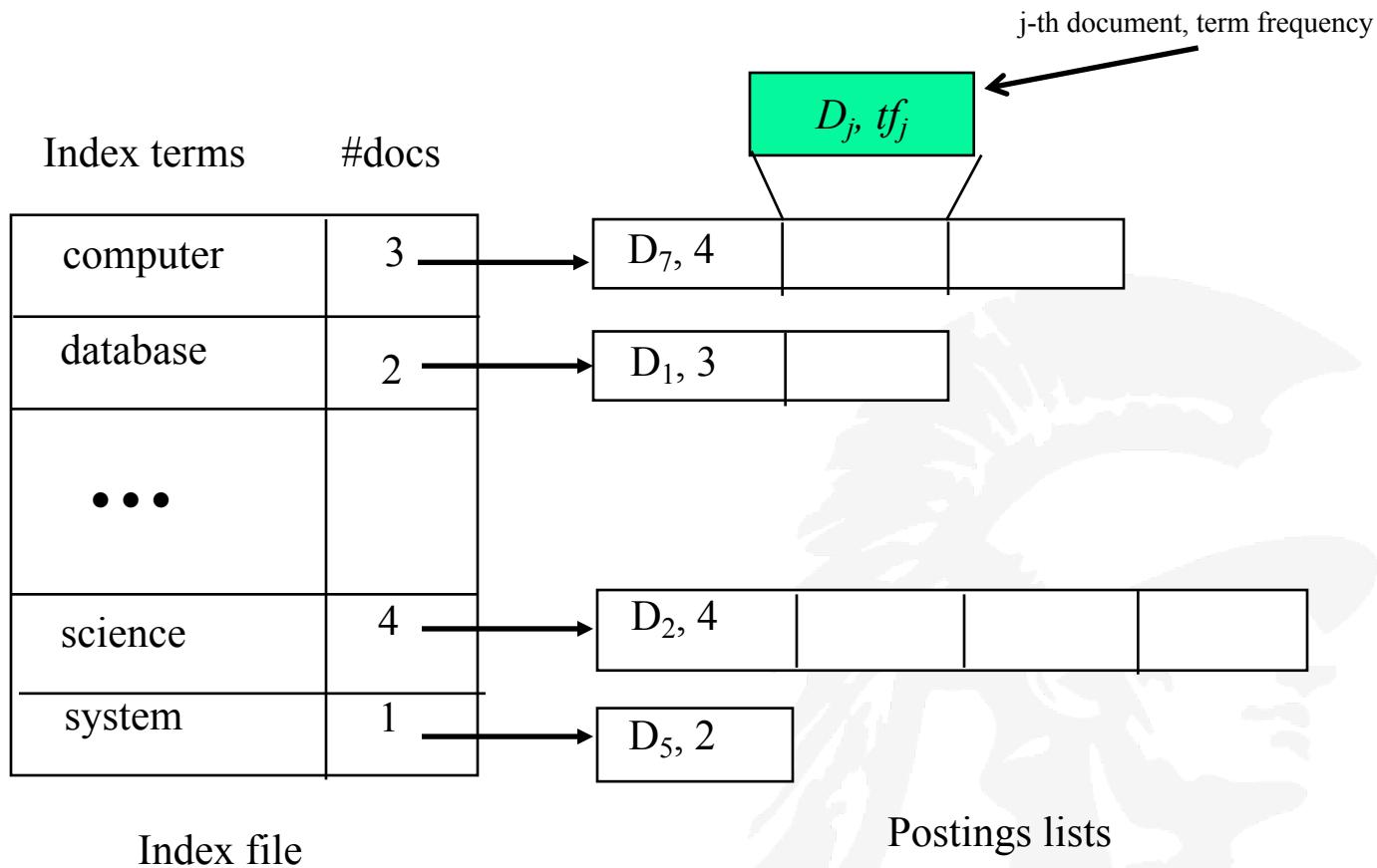
# Actual Answers to the Query

- **Antony and Cleopatra, Act III, Scene ii**
- ***Agrippa [Aside to DOMITIUS ENOBARBUS]:***
  - Why, Enobarbus,
  - When Antony found Julius Caesar dead,
  - He cried almost to roaring; and he wept
  - When at Philippi he found *Brutus* slain.
- **Hamlet, Act III, Scene ii**
- ***Lord Polonius:*** I did enact Julius Caesar I was killed i' the Capitol;  
*Brutus* killed me.

# Term-Document Incident Matrices are Naturally Sparse

- Given 1 million documents and 500,000 terms
- The (term x Document) matrix in this case will have size 500K x 1M or half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
  - So the matrix is extremely sparse, only 0.1% of the elements are 1
- So instead we use a data structure for an inverted index that exploits sparsity and then devise algorithms for query processing

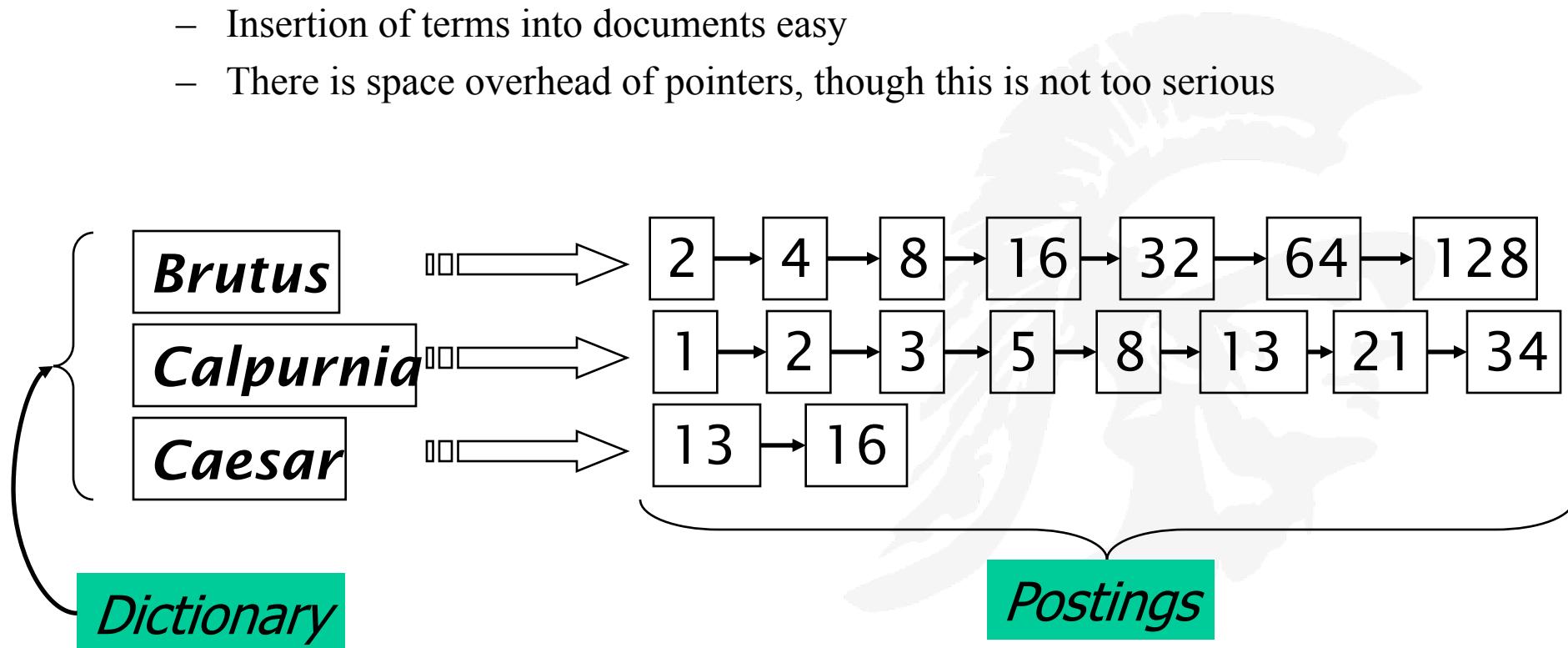
## Inverted Index Example



The two parts of an inverted index. The **dictionary** (Index file) is usually kept in memory, with pointers to each **postings list**, which is stored on disk. The dictionary has been sorted alphabetically and the postings list is sorted by document ID

# Inverted Index Stored In Two Parts

- For each term  $T$ , we must store a list of all documents that contain  $T$ .
- Linked lists are generally preferred to arrays, why . . .
  - Dynamic space allocation
  - Insertion of terms into documents easy
  - There is space overhead of pointers, though this is not too serious



# Parsing Documents To Create an Inverted Index

- Documents are parsed to extract words and these are saved with the document ID i.e a sequence of (possibly modified token, Document ID) pairs

Doc 1

I did enact Julius  
Caesar I was killed  
i' the Capitol;  
Brutus killed me.

Doc 2

So let it be with  
Caesar. The noble  
Brutus hath told you  
Caesar was ambitious



Term	Doc #
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

- If the corpus is known in advance, then after all documents have been parsed the inverted file is sorted by terms

## Initial capture of terms

- However, on the Web, documents are constantly being added and the terms are constantly increasing

Term	Doc #
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

## Refined list of terms

Term	Doc #
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

- Multiple term entries in a single document are merged.
  - Frequency information is added.

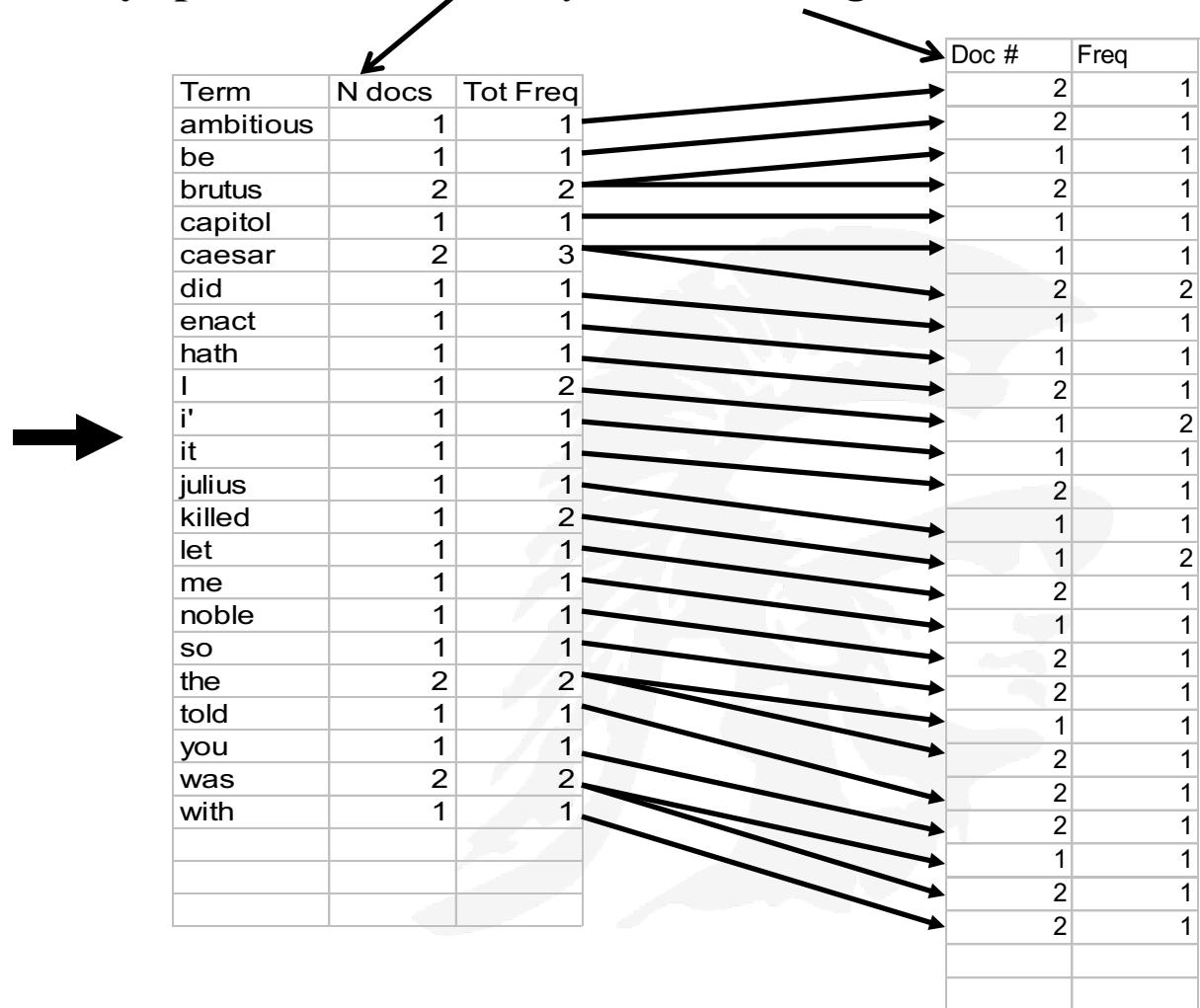
Term	Doc #
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2



Term	Doc #	Freq
ambitious	2	1
be	2	1
brutus	1	1
brutus	2	1
capitol	1	1
caesar	1	1
caesar	2	2
did	1	1
enact	1	1
hath	2	1
I	1	2
i'	1	1
it	2	1
julius	1	1
killed	1	2
let	2	1
me	1	1
noble	2	1
so	2	1
the	1	1
the	2	1
told	2	1
you	2	1
was	1	1
was	2	1
with	2	1

e.g. Caesar Occurs in Documents 1 and 2, With Total Frequency 3

- The file is commonly split into a *Dictionary* and a *Postings* file

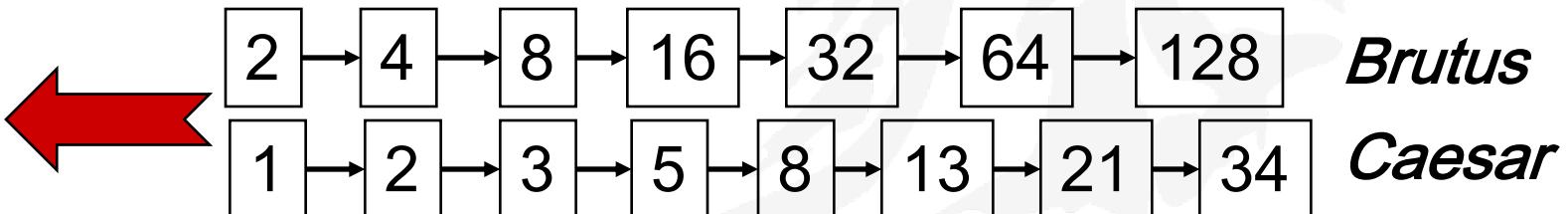


# Query Processing Across the Postings List

- Consider processing the query:

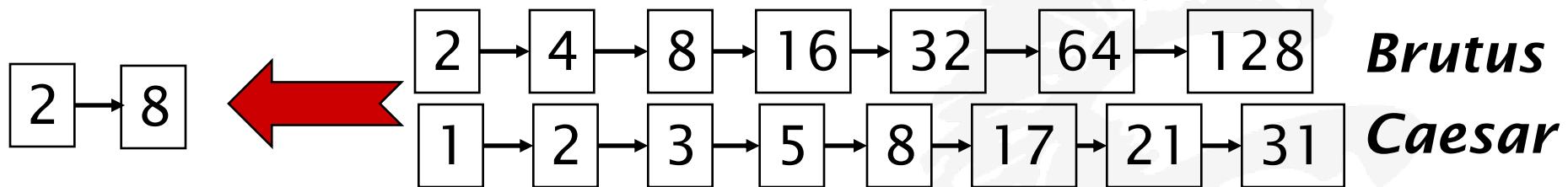
*Brutus AND Caesar*

- Locate *Brutus* in the Dictionary;
  - Retrieve its postings.
- Locate *Caesar* in the Dictionary;
  - Retrieve its postings.
- “Merge” the two postings and select the ones in common (postings are document ids):



# Basic Merge

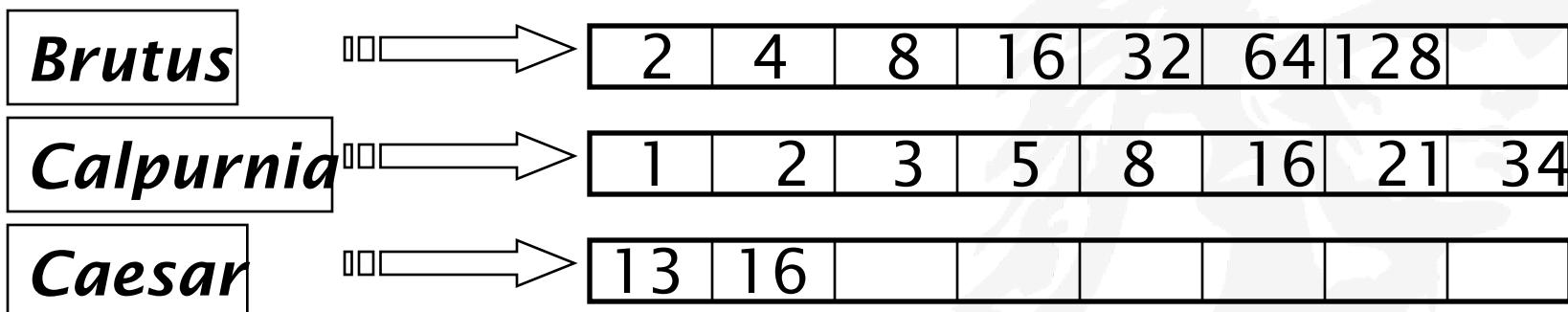
- Walk through the two postings simultaneously, in time linear in the total number of postings entries



If the list lengths are  $m$  and  $n$ , the merge takes  $O(m+n)$  operations.

# Query Optimization

- What is the best order for query processing?
- Consider a query that is an *AND* of  $t$  terms.
- For each of the  $t$  terms, get its postings, then *AND* together.

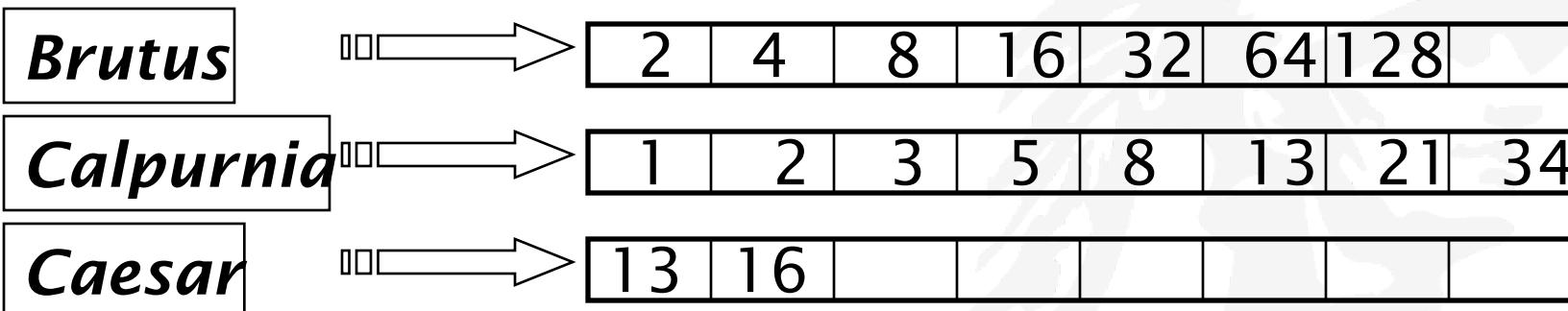


**Query: *Brutus AND Calpurnia AND Caesar***

# Query Optimization Example

- Process in order of increasing frequency of occurrence:
  - *start with smallest set, then keep cutting further.*

This is why we kept  
freq in dictionary

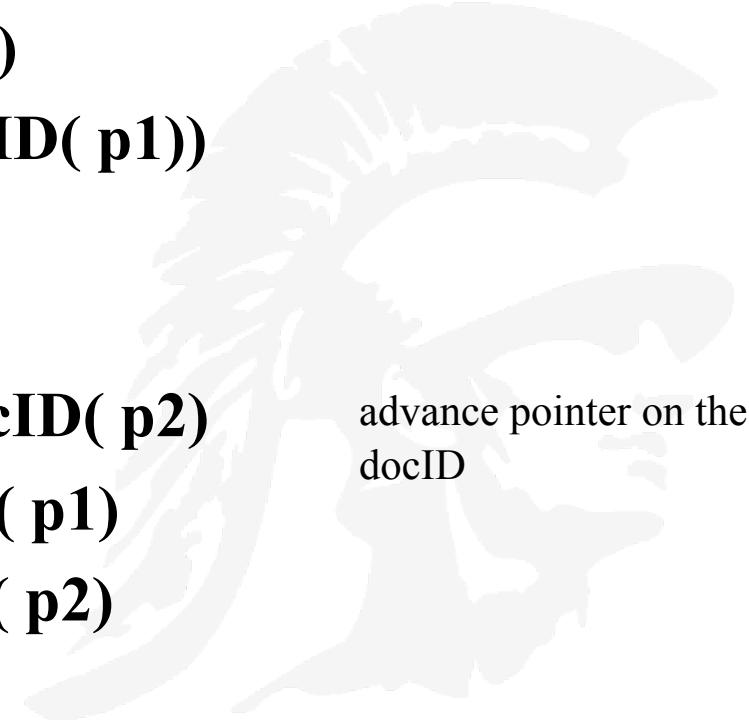


Execute the query as (*Caesar AND Brutus*) AND *Calpurnia*.

# Algorithm for the intersection of two postings lists p1 and p2

**INTERSECT( p1, p2)**

```
1 answer ← ()  
2 while p1 ≠ NIL and p2 ≠ NIL  
3 do if docID( p1) = docID( p2)  
4     then ADD(answer, docID( p1))  
5             p1 ← next( p1)  
6             p2 ← next( p2)  
7     else if docID( p1) < docID( p2)  
8         then p1 ← next( p1)  
9     else p2 ← next( p2)  
10    return answer
```



advance pointer on the smaller docID

## Algorithm for conjunctive queries that returns the set of documents containing each term in the input list of terms

**INTERSECT(< $t_1, \dots, t_n$ >)**

- 1 terms  $\leftarrow$  **SORTBYINCREASINGFREQUENCY(< $t_1, \dots, t_n$ >)**
- 2 result  $\leftarrow$  **postings(first(terms))**
- 3 terms  $\leftarrow$  **rest(terms)**
- 4 while terms  $\neq$  **NIL** and result  $\neq$  **NIL**
- 5 do result  $\leftarrow$  **INTERSECT(result, postings(first(terms)))**
- 6 terms  $\leftarrow$  **rest(terms)**
- 7 return result

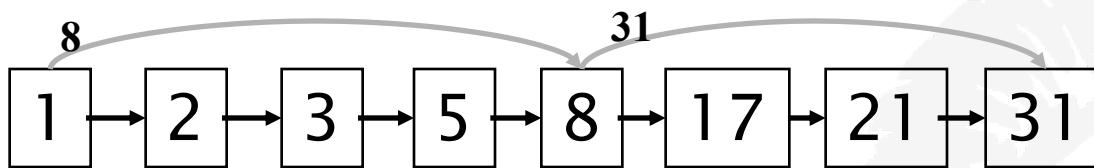
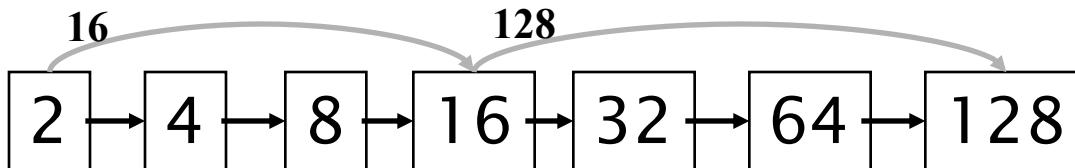
### Algorithm Strategy:

intersect each retrieved postings list with the current intermediate result in memory, where we initialize the intermediate result by loading the postings list of the least frequent term

To speed up the merging of postings we  
use the technique of *Skip Pointers*

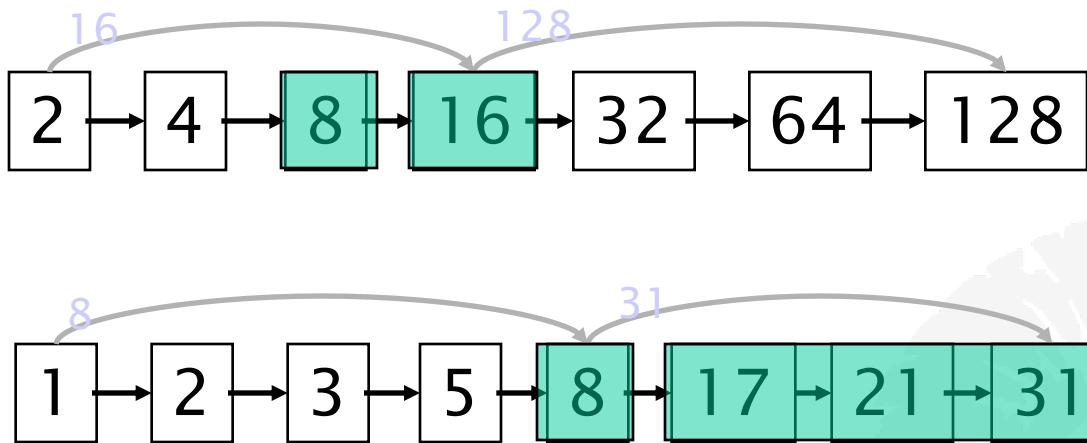
# The Technique of Skip Pointers

Augment postings with skip pointers (at indexing time)



- Why?
- To skip postings that will not figure in the search results.
- How?
- Where do we place skip pointers?

## Query Processing With Skip Pointers



Suppose we've stepped through the lists until we process 8 on each list.

When we get to 16 on the top list, we see that its successor is 32.

But the skip successor of 8 on the lower list is 31, so we can skip ahead past the intervening postings 17 and 21.

# Facts on Skip Pointers

- **Skip pointers are added at indexing time**; they are shortcuts, and they only help for AND queries and they are useful when the corpus is relatively static
- there are two questions that must be answered:
  - 1. where should they be placed?
  - 2. how do the algorithms change?
- **The Argument:** More skips means shorter skip spans, and that we are more likely to skip. But it also means lots of comparisons to skip pointers, and lots of space storing skip pointers. Fewer skips means few pointer comparisons, but then long skip spans which means that there will be fewer opportunities to skip.
- **The Solution:** A simple heuristic for placing skips, which has been found to work well in practice, is that *for a postings list of length  $P$ , use  $\sqrt{P}$  evenly-spaced skip pointers*. This heuristic possibly can be improved upon as it ignores any details of the distribution of query terms. **[Moffat and Zobel 1996]**
- See the YouTube video <http://www.youtube.com/watch?v=tPsCQOsA7j0> (15 min)

```

INTERSECTWITHSKIPS( $p_1, p_2$ )
1  answer  $\leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4    then ADD( $\text{answer}, \text{docID}(p_1)$ )
5       $p_1 \leftarrow \text{next}(p_1)$ 
6       $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8    then if  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2))$ 
9      then while  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2))$ 
10        do  $p_1 \leftarrow \text{skip}(p_1)$ 
11        else  $p_1 \leftarrow \text{next}(p_1)$ 
12    else if  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1))$ 
13      then while  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1))$ 
14        do  $p_2 \leftarrow \text{skip}(p_2)$ 
15        else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return  $\text{answer}$ 

```

# Faster Postings List Intersection via Skip Pointers

- Skip pointers will only be available for the original postings lists.
- For an intermediate result in a complex query, the call  $\text{hasSkip}(p)$  will always return false.
- Finally, note that the presence of skip pointers only helps for AND queries, not for OR queries.

# Phrase Queries

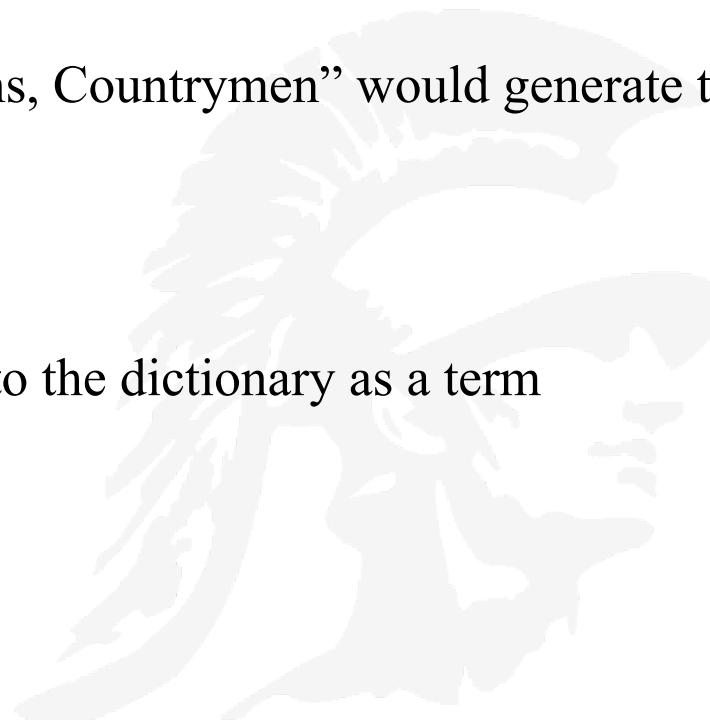


# Phrase queries

- We want to answer queries such as *stanford university* – as a phrase
  - Many search engines allow users to specify a phrase using double quotes, "Stanford University", which most people find easy to understand
  - As many as 10% of web search engine queries are phrase queries and many more are implied phrase queries, e.g. George Bush, Bernie Sanders
  - If we only store terms and documents in our inverted index,  $\langle term : docs \rangle$  then how we can answer a phrase query
- 
- There are two approaches to solving this problem
    1. Bi-word indexes (also called 2-grams or 2 shingles)
    2. Positional indexes

# Using Biword Indexes for Phrase Searching

- Definition: A *bi-word* (or a 2-gram) is a consecutive pair of terms in some text
- To improve phrase searching one approach is to index every bi-word in the text
- For example the text “Friends, Romans, Countrymen” would generate the bi-words
  - *friends romans*
  - *romans countrymen*
- Each of these bi-words is now added to the dictionary as a term



# Handling Longer Phrase Queries

- Consequences
  - Bi-words will cause an explosion in the vocabulary database
  - Queries longer than 2 words will have to be broken into bi-word segments
- Example: suppose the query is the 4 word phrase  
*stanford university palo alto*

The query can be broken into the Boolean query on bi-words:

*stanford university AND university palo AND palo alto*

- Matching the query to terms in the index will work, but may also produce false positives (i.e. occurrences of the bi-words, but not the full 4 word query)
- A Bi-word index that is extended to longer sequences, or even variable length sequences is called a *phrase index*

# Part-of-Speech Tagging

- **Many two word phrases have embedded stop words, e.g.**
  - the abolition of slavery
  - negotiation of the constitution
- **To salvage the bi-word indexing method on these examples one can use part-of-speech tagging**
  - Part-of-speech taggers classify words as nouns, verbs, etc. – or, in practice, often as finer grained classes like “plural proper noun”.
  - “Negotiation of the constitution” is transformed into “N X X N”
- **Many fairly accurate (c. 96% per-tag accuracy) part-of-speech taggers now exist, usually trained by machine learning methods on hand-tagged text**
  - See Manning and Schutze “*Foundations of Statistical Natural Language Processing*”
  - <https://nlp.stanford.edu/fsnlp/>
  - <https://www.issco.unige.ch/en/staff/robert/tatoo/tatoo.html>

## Alternate Solution - Using Positional Indexes

- Given the limitations of a bi-word index, (i.e. the enormous growth in the vocabulary) the alternate solution is most commonly used, called a **Positional Index**
- Store, for each *term* in the index, entries of the form:  
**<term, number of docs containing term;**  
***doc1:* position1, position2, position3 ... ;**  
***doc2:* position1, position2 ... ;**  
**etc.>**
- i.e. for each occurrence of the term in a document its position is stored

# Positional Index Example

for each term in the vocabulary, we store postings of the form

**docID: position1, position2, ...,**

where each position is a token index in the document.

Each posting will also usually record the term frequency

Adopting a positional index expands required postings storage significantly,  
even if we compress position values/offsets

<**be**: 993427;  
    ↑  
    Lots of documents  
**1**: 7, 18, 33, 72, 86, 231;  
    ↑  
    Lots of occurrences  
**2**: 3, 149;  
**4**: 17, 191, 291, 430, 434;  
**5**: 363, 367, ...>

- the term is "be", typically a stop word
- it occurs in 993,427 documents
- in document 1 "be" occurs at positions: 7, 18, 33, 72, 86 and 231

- this scheme expands postings storage *substantially (rather than the vocabulary)*

# Processing a Phrase Query

- Algorithm for matching a phrase query:
  1. Extract inverted index entries for each distinct term: e.g. *to*, *be*, *or*, *not*, *to*, *be*
  2. Merge their *doc:position* lists to enumerate all positions with “*to be or not to be*”.
  3. Match those documents that contain the terms in the adjacent positions
    - ***to:***
      - 2:1,17,74,222,551; 4:8,16,190,429,433; 7:13,23,191; ...
    - ***be:***
      - 1:17,19; 4:17,191,291,430,434; 5:14,19,101; ...
      - Same general method for proximity searches
    - In document 4 the word “*to*” appears in position 16 and the word “*be*” appears in position 17, so they are adjacent
    - **Note:** Figure 2.12 of our book has an algorithm for proximity searching

# Algorithm for Proximity Queries with $k$ words

## 2 The term vocabulary and postings lists

```

POSITIONALINTERSECT( $p_1, p_2, k$ )
1 answer  $\leftarrow \langle \rangle$ 
2 while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3 do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4   then  $l \leftarrow \langle \rangle$ 
5      $pp_1 \leftarrow \text{positions}(p_1)$ 
6      $pp_2 \leftarrow \text{positions}(p_2)$ 
7     while  $pp_1 \neq \text{NIL}$ 
8       do while  $pp_2 \neq \text{NIL}$ 
9         do if  $|\text{pos}(pp_1) - \text{pos}(pp_2)| \leq k$ 
10        then ADD( $l, \text{pos}(pp_2)$ )
11        else if  $\text{pos}(pp_2) > \text{pos}(pp_1)$ 
12          then break
13           $pp_2 \leftarrow \text{next}(pp_2)$ 
14        while  $l \neq \langle \rangle$  and  $|l[0] - \text{pos}(pp_1)| > k$ 
15          do DELETE( $l[0]$ )
16          for each  $ps \in l$ 
17            do ADD(answer,  $(\text{docID}(p_1), \text{pos}(pp_1), ps)$ )
18             $pp_1 \leftarrow \text{next}(pp_1)$ 
19             $p_1 \leftarrow \text{next}(p_1)$ 
20             $p_2 \leftarrow \text{next}(p_2)$ 
21        else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
22          then  $p_1 \leftarrow \text{next}(p_1)$ 
23        else  $p_2 \leftarrow \text{next}(p_2)$ 
24 return answer

```

The algorithm finds places where the two terms appear within  $k$  words of each other and returns a list of triples giving docID and the term position in  $p_1$  and  $p_2$ .

► **Figure 2.12** An algorithm for proximity intersection of postings lists  $p_1$  and  $p_2$ . The algorithm finds places where the two terms appear within  $k$  words of each other and returns a list of triples giving docID and the term position in  $p_1$  and  $p_2$ .

# Some High Frequency Noun Phrases from TREC and Patent DataSets

## *TREC*

<i>Frequency</i>	<i>Phrase</i>
<b>65824</b>	<b>united states</b>
<b>61327</b>	<b>article type</b>
<b>33864</b>	<b>Los Angeles</b>
<b>18062</b>	<b>Hong kong</b>
<b>17788</b>	<b>North Korea</b>
<b>17308</b>	<b>New York</b>
<b>15513</b>	<b>San Diego</b>
<b>15009</b>	<b>Orange county</b>

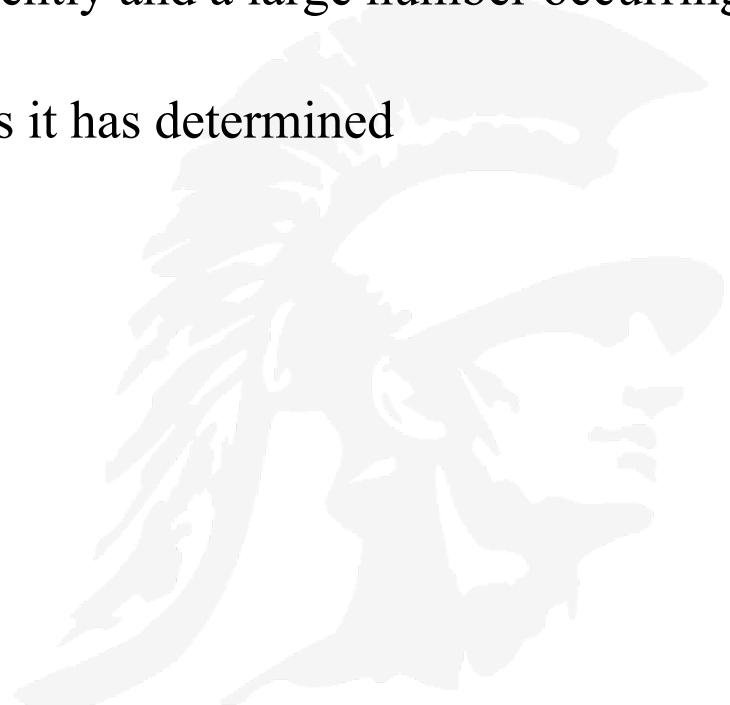
## *Patent*

<i>Frequency</i>	<i>Phrase</i>
<b>975362</b>	<b>present invention</b>
<b>191625</b>	<b>u.s. pat</b>
<b>147352</b>	<b>preferred embodiment</b>
<b>95097</b>	<b>carbon atoms</b>
<b>87903</b>	<b>group consisting</b>
<b>81809</b>	<b>room temperature</b>
<b>78458</b>	<b>seq id</b>
<b>75850</b>	<b>brief description</b>

The phrases above were identified by POS tagging; The data above shows that common phrases are used more frequently in patent data as patents have a very formal style; many of the TREC phrases are proper nouns, whereas patent phrases are those that occur in all patents

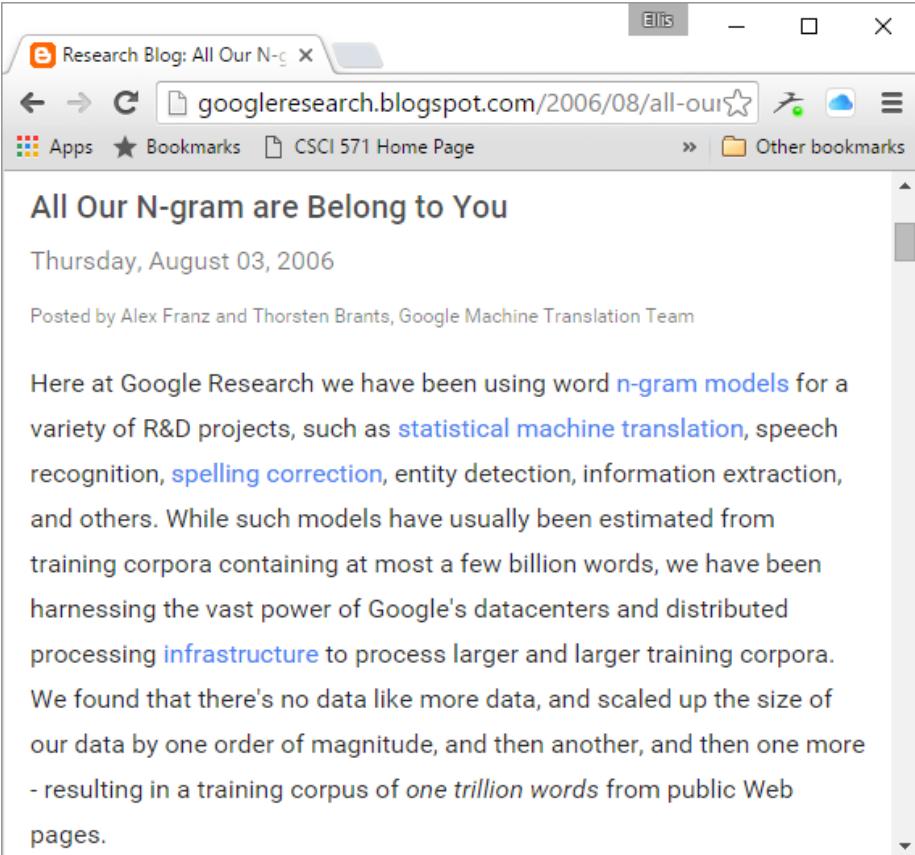
# Google Relies Upon *n*-gram Indexes

- Google has investigated the use of n-grams stored in its index,  $n \geq 2$ ;
- N-grams of all lengths form a **Zipf distribution (power law)** with a few common phrases occurring very frequently and a large number occurring with frequency 1
- Google has released the set of n-grams it has determined



# Google's N-Gram Database Facts

- Google made available a file of n-grams derived from the web pages it indexed
- <http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>
- Statistics for the Google n-gram sample
- Number of tokens 1,024,908,267,229  
(1 trillion, 24 billion, 908 million, . . . )
- Number of sentences 95,119,665,584
- Number of unigrams 13,588,391
- Number of bigrams 314,843,401
- Number of trigrams 977,069,902
- Number of four grams 1,313,818,354
- Number of five grams 1,176,470,663



The screenshot shows a web browser window with the title bar "Research Blog: All Our N-gram are Belong to You". The address bar contains the URL "googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html". The page content starts with the heading "All Our N-gram are Belong to You" and the date "Thursday, August 03, 2006". It is posted by "Alex Franz and Thorsten Brants, Google Machine Translation Team". The text discusses the use of word n-gram models for various R&D projects like statistical machine translation, speech recognition, spelling correction, entity detection, and information extraction. It mentions scaling up the data size by orders of magnitude and processing infrastructure to handle one trillion words from public Web pages.

# Examples of 3-Gram, 4-Gram Data

The following is an example of the 3-gram data contained this corpus:

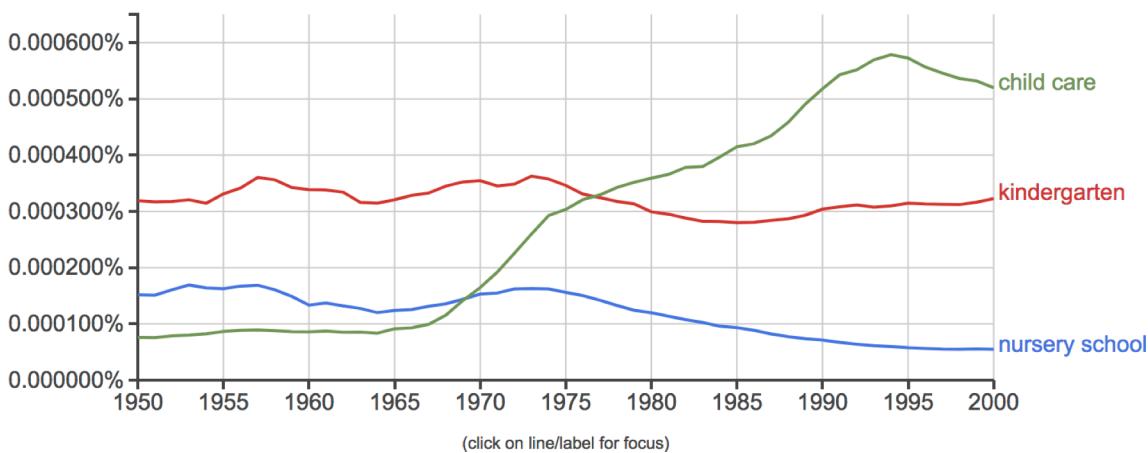
```
ceramics collectables collectibles 55
ceramics collectables fine 130
ceramics collected by 52
ceramics collectible pottery 50
ceramics collectibles cooking 45
ceramics collection , 144
ceramics collection . 247
ceramics collection </S> 120
ceramics collection and 43
ceramics collection at 52
ceramics collection is 68
ceramics collection of 76
ceramics collection | 59
ceramics collections , 66
ceramics collections . 60
ceramics combined with 46
ceramics come from 69
ceramics comes from 660
ceramics community , 109
ceramics community . 212
ceramics community for 61
```

The following is an example of the 4-gram data in this corpus:

```
serve as the incoming 92
serve as the incubator 99
serve as the independent 794
serve as the index 223
serve as the indication 72
serve as the indicator 120
serve as the indicators 45
serve as the indispensable 111
serve as the indispensable 40
serve as the individual 234
serve as the industrial 52
serve as the industry 607
serve as the info 42
serve as the informal 102
serve as the information 838
serve as the informational 41
serve as the infrastructure 500
serve as the initial 5331
serve as the initiating 125
serve as the initiation 63
serve as the initiator 81
----- -- --- ----- --
```

# Google's N-Gram Viewer

- The **Google N-Gram Viewer** or **Google Books N-Gram Viewer** is an online search engine that charts frequencies of any set of comma-delimited search strings using a yearly count of n -grams found in sources printed between 1500 and 2008 in **Google's** text corpora in English, Chinese (simplified), French, German, Hebrew, Italian, ...
  - <https://books.google.com/ngrams/>, or for examples see [books.google.com/ngrams/info](https://books.google.com/ngrams/info)
- The program can search for a single word or a phrase, including misspellings
- The n-grams are matched with the text within the selected corpus, and, if found in 40 or more books, are then plotted on a graph
- The data used for the search are composed of total \_ counts, 1-grams, 2-grams, 3-grams, 4-grams, and 5-grams files for each language



This shows trends in three n-grams from 1950 to 2000:  
 "nursery school" (a *2-gram* or *bigram*),  
 "kindergarten" (a *1-gram* or *unigram*),  
 and  
 "child care" (another bigram)

# Comparing N-Grams Across Languages

- S. Yang et al, N-gram statistics in English and Chinese: Similarities and differences, ICSC, 2007, Int'l Conf. on semantic computing, 454-460
- [http://static.googleusercontent.com/external\\_content/untrusted\\_dlcp/research.google.com/en/us/pubs/archive/33035.pdf](http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/pubs/archive/33035.pdf)
- The authors above analyzed 200 million randomly sampled English and Chinese Web pages and concluded:
  1. The distribution of the unique number of n-grams is similar between English and Chinese, though the Chinese distribution is shifted to larger N
  2. The distribution indicates that on average 1.5 Chinese characters correspond to 1 English word
  3. While frequency distributions of uni-grams and bi-grams are very different, the frequency distribution for 3-grams and 4-grams are strikingly similar

