

Introduction to Assignment Project Exam Help

Informa

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro
Lecture 1: Bool at

Introduction to Assignment Project Exam Help

Informa

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Term-document inc trices

Unstructured data in 1620

- Which plays of Shakespeare contain the words ***Brutus*** AND ***Caesar*** but NOT ***Calpurnia***?
- One could ~~Assignment Project Exam Help~~ **grep** all of Shakespeare's plays for ***Brutus*** and ***Caesar***, t ~~aining~~ <https://eduassistpro.github.io/> ~~aining~~ ***Calpurnia***?
- Why is that not the answer
 - Slow (for large corpora)
 - NOT ***Calpurnia*** is non-trivial
 - Other operations (e.g., find the word ***Romans*** near ***countrymen***) not feasible
 - Ranked retrieval (best documents to return)
 - Later lectures

Term-document incidence matrices

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

*Brutus AND Caesar BUT NOT
Calpurnia*

1 if play contains word, 0 otherwise

Incidence vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for *Brutus*, *Caesar* and *Calpurnia* (complemented) → bitwise AND.
- 110100 AND <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

	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

Answers to query

- Antony and Cleopatra, Act III, Scene ii

Agrippa [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,
When Antony found Julius **Caesar** dead,

He wept
Wh <https://eduassistpro.github.io/>
laid.

Add WeChat **edu_assist_pro**

- Hamlet, Act III, Sce

Lord Polonius: I did enact Julius **Caesar** I was killed i' the
Capitol; **Brutus** killed me.



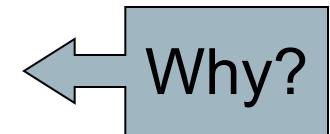
Bigger collections

- Consider $N = 1$ million documents, each with about 1000 words.
- Avg 6 bytes/word including spaces/punctuation
 - 6GB of data in <https://eduassistpro.github.io/>
- Say there are $M = 500K$ *dist* among these.
Add WeChat `edu_assist_pro`

Can't build the matrix

- 500K x 1M matrix has half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
 - matrix is extremely sparse.
- What's a bett <https://eduassistpro.github.io/>
 - We only record the 1 positions.

Assignment Project Exam Help



Add WeChat edu_assist_pro

Introduction to Assignment Project Exam Help

Informa

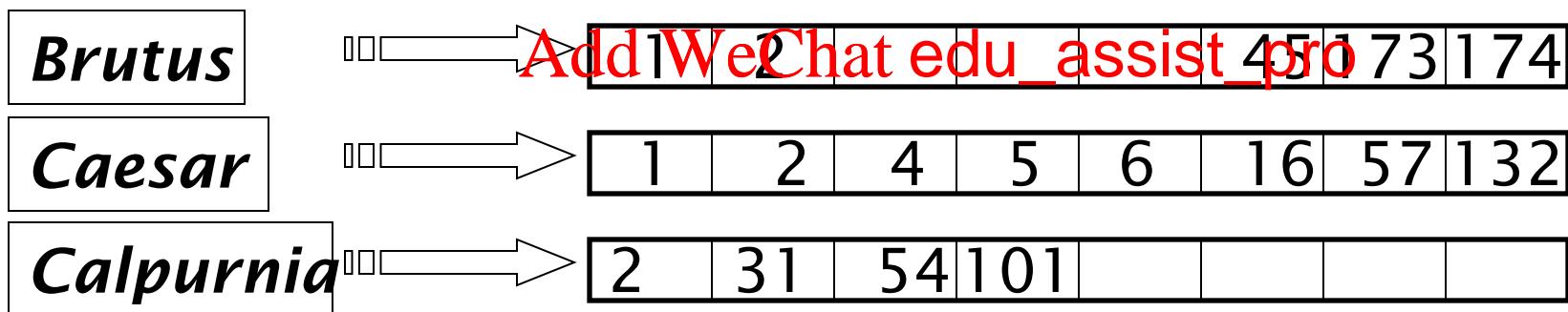
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
The Inverte

The key data structure underlying modern IR

Inverted index

- For each term t , we must store a list of all documents that contain t .
 - Identify each by a docID, a document serial number
- Can we use file <https://eduassistpro.github.io/>?



What happens if the word *Caesar* is added to document 14?

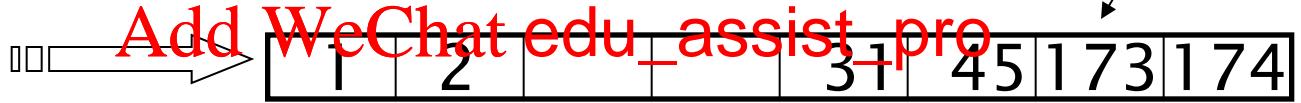
Inverted index

- We need variable-size postings lists
 - On disk, a continuous run of postings is normal and best
 - In memory, can use linked lists or variable length arrays
 - Some trade

<https://eduassistpro.github.io/>

Posting

Brutus



Caesar

Calpurnia

Dictionary

Postings

Sorted by docID (more later on why).

Inverted index construction

Documents to be indexed.



Friends, Romans, countrymen.

⋮

Assignment Project Exam Help

Token stream.

<https://eduassistpro.github.io/>

omans

Countrymen

More on these later.

Add WeChat edu_assist_pro
Linguistic modules

Modified tokens.

friend

roman

countryman

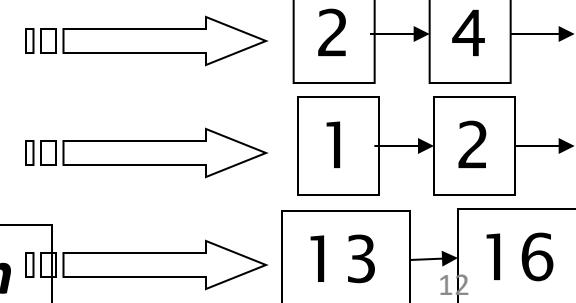
Inverted index.

Indexer

friend

roman

countryman



Indexer steps: Token sequence

- Sequence of (Modified token, Document ID) pairs.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Doc 1

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

Doc 2

Add WeChat edu_assist_pro →
So let it be with
Caesar. The noble
Brutus hath told you
Caesar was ambitious

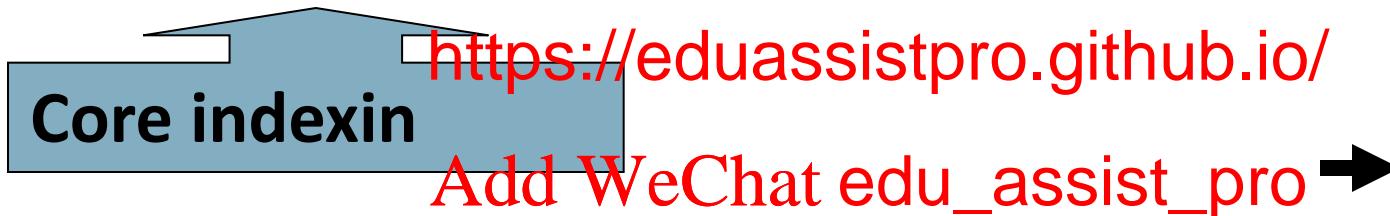
Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

Indexer steps: Sort

- ## ■ Sort by terms

- And then docID

Assignment Project Exam Help



Indexer steps: Dictionary & Postings

- Multiple term entries in a single document are merged.
- Split into Dictionary and Postings
- Doc. frequency information is added.

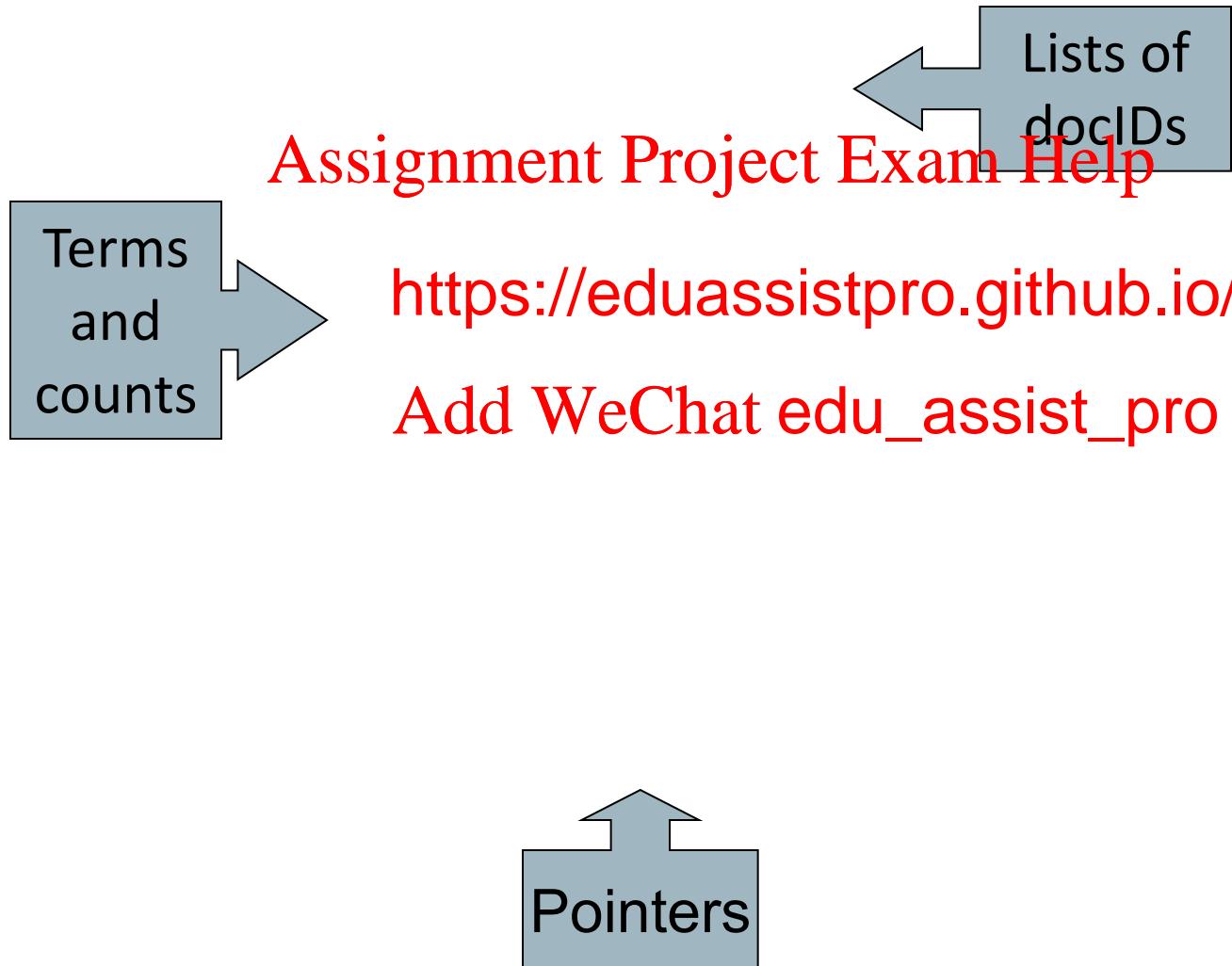
<https://eduassistpro.github.io/>



Add WeChat edu_assist_pro

Why frequency?
Will discuss later.

Where do we pay in storage?



Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Query processing with `ed_index`

The index we just built

- How do we process a query?
 - Later - what kinds of queries can we process?

Today's focus

Assignment Project Exam Help

<https://eduassistpro.github.io/>

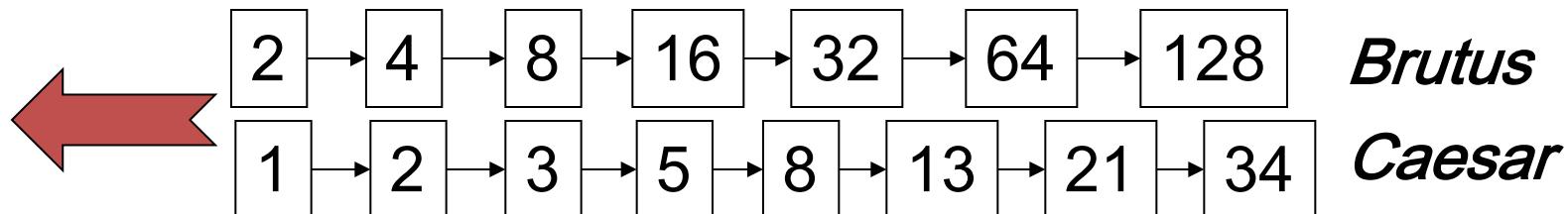
Add WeChat edu_assist_pro

Query processing: AND

- Consider processing the query:

Brutus AND Caesar

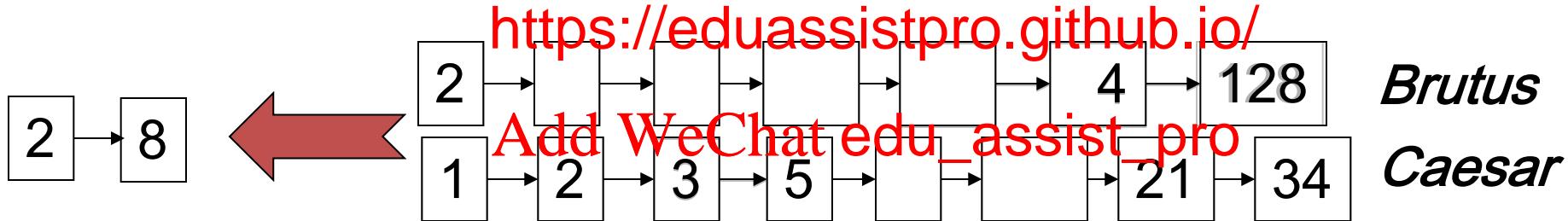
- Locate *Brutus* in the Dictionary;
▪ Retrieve its <https://eduassistpro.github.io/>
- Locate *Caesa*
▪ Retrieve its postings:
[Add WeChat edu_assist_pro](#)
- “Merge” the two postings:



The merge

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

Assignment Project Exam Help



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

Crucial: postings sorted by docID.

Intersecting two postings lists (a “merge” algorithm)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Boolean queries: Exact match

- The Boolean retrieval model is being able to ask a query that is a Boolean expression:
Assignment Project Exam Help
 - Boolean Queries are queries using *AND*, *OR* and *NOT* to join query <https://eduassistpro.github.io/>
 - Views each document as a s
 - Is precise: document match not:
Add WeChat edu_assist_pro
 - Perhaps the simplest model to build an IR system on
- Primary commercial retrieval tool for 3 decades.
- Many search systems you still use are Boolean:
 - Email, library catalog, Mac OS X Spotlight

Example: WestLaw

<http://www.westlaw.com/>

- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992; new federated search added 2010)
- Tens of ter <https://eduassistpro.github.io/> users
- Majority of users still us edu_assist_pro
- Example query:
 - What is the statute of limitations in cases involving the federal tort claims act?
 - LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM
 - foo! = foo*, /3 = within 3 words, /S = in same sentence

Example: WestLaw <http://www.westlaw.com/>

- Another example query:
 - Requirements for disabled people to be able to access a workplace
 - disabl! /p a <https://eduassistpro.github.io/> place (employment /3 place)
- Note that SPACE is disjunction!
- Long, precise queries; proximity operators; incrementally developed; not like web search
- Many professional searchers still like Boolean search
 - You know exactly what you are getting
- But that doesn't mean it actually works better...

Boolean queries:

More general merges

- Exercise: Adapt the merge for the queries:

Brutus AND NOT Caesar

Assignment Project Exam Help

Brutus OR N

<https://eduassistpro.github.io/>

Can we still run through the m Add WeChat edu_assist_pro $O(x+y)$?

What can we achieve?

Merging

What about an arbitrary Boolean formula?

(Brutus OR Caesar) AND NOT

(Antony OR Cleo)

- Can we always <https://eduassistpro.github.io/>?
 - Linear in what? [Add WeChat edu_assist_pro](#)
- Can we do better?

Query optimization

- What is the best order for query processing?
- Consider a query that is an AND of n terms.
- For each of the <https://eduassistpro.github.io/> terms, then
AND them to

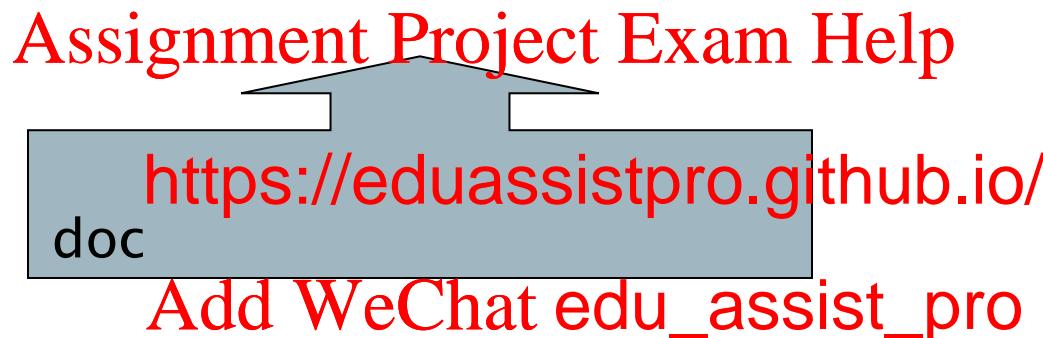
Add WeChat `edu_assist_pro`

Brutus	⇒	2	4	8	16	32	64	128
Caesar	⇒	1	2	3	5	8	16	21
Calpurnia	⇒	13	16					

Query: **Brutus AND Calpurnia AND Caesar**

Query optimization example

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



Brutus	⇒	2 4 8 16 32 64 128
Caesar	⇒	1 2 3 5 8 16 21 34
Calpurnia	⇒	13 16

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

More general optimization

- e.g., *(madding OR crowd) AND (ignoble OR strife) AND (light OR lord)*
- Get doc. freq
- Estimate the <https://eduassistpro.github.io/> sum of its doc. freq.'s (conservative)
- Process in increasing order of *OR* sizes.

Q: When will the conservative estimate fail?

Exercise

- Recommend a query processing order for

Assignment Project Exam Help

(tangerine OR tre
https://eduassistpro.github.io/
(marmalade OR skies) AND
(kaleidoscope OR Add WeChat edu_assist_pro

Q: Can we merge multiple lists (>2) simultaneously?

Introduction to Assignment Project Exam Help **Informa**

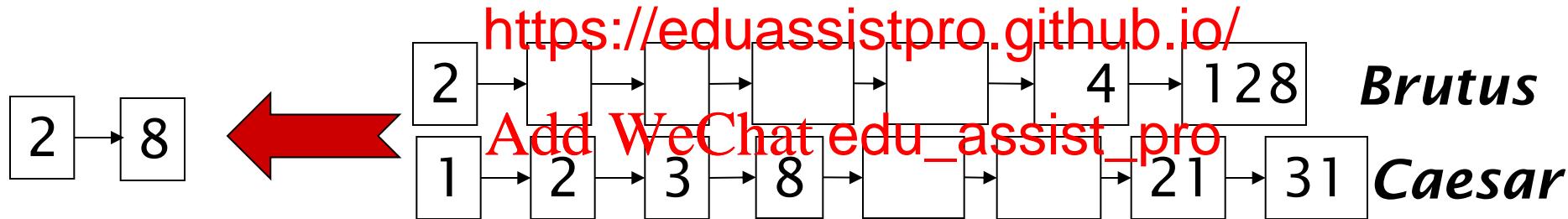
https://eduassistpro.github.io/

Add WeChat edu_assist_pro
Faster postings merges: ters and
Skip lists

Recall basic merge

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

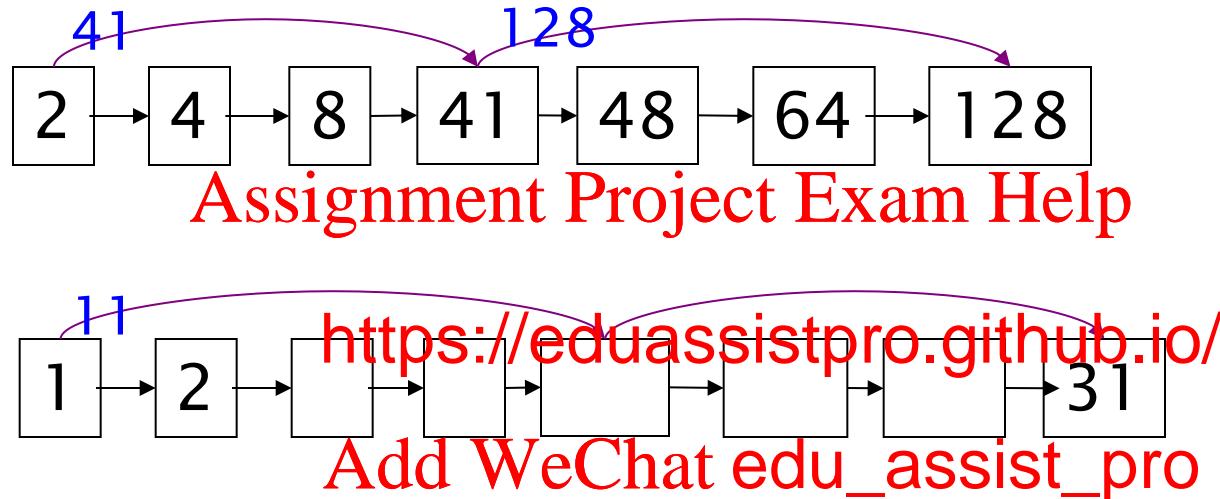
Assignment Project Exam Help



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

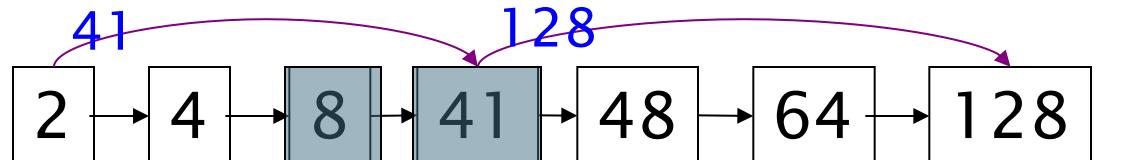
Can we do better?
Yes (if index isn't changing too fast).

Augment postings with skip pointers (at indexing time)

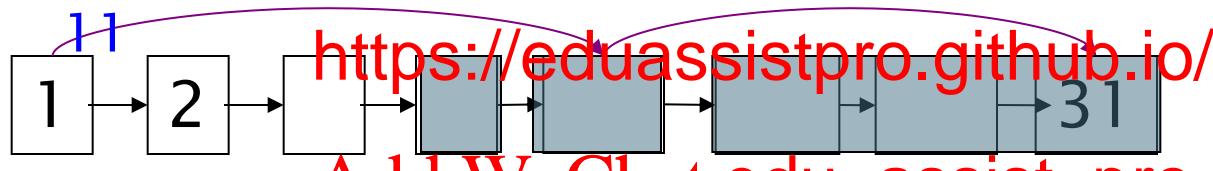


- Why?
- To skip postings that will not figure in the search results.
- How?
- Where do we place skip pointers?

Query processing with skip pointers



Assignment Project Exam Help



Add WeChat edu_assist_pro

Suppose we've stepped through until we process 8 on each list. We match it and advance.

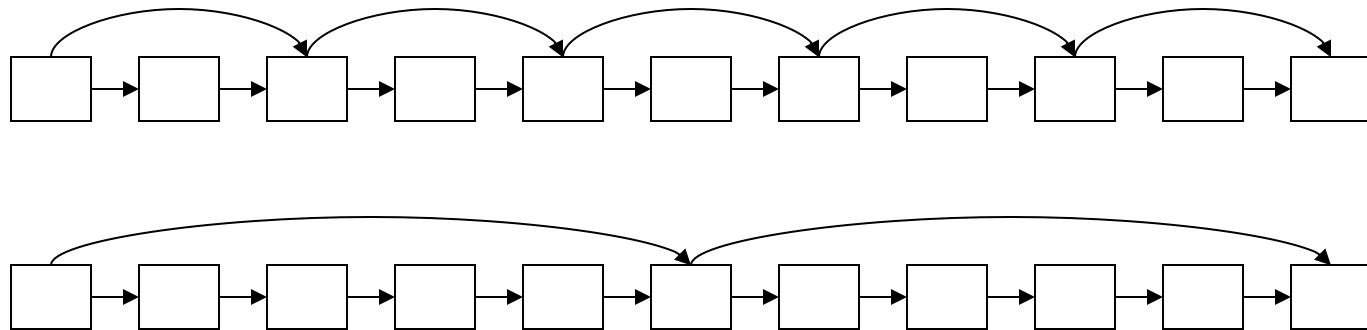
We then have 41 and 11 on the lower. 11 is smaller.

But the skip successor of 11 on the lower list is 31, so we can skip ahead past the intervening postings.

Where do we place skips?

- Tradeoff:
 - More skips → shorter skip spans → more likely to skip.
But lots of comparisons to skip pointers.
 - Fewer skips, but then long skip spans → few

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro



Placing skips

- Simple heuristic: for postings of length L , use $L^{1/2}$ evenly-spaced skip pointers.
- This ignores the distribution of query terms.
- Easy if the index is <https://eduassistpro.github.io/> harder if L keeps changing because of updates
[Add WeChat edu_assist_pro](#)
- This definitely used to help; with modern hardware it may not (Bahle et al. 2002) unless you're memory-based
 - The I/O cost of loading a bigger postings list can outweigh the gains from quicker in memory merging!

Introduction to Assignment Project Exam Help

Informa

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro
Phrase queries and indexes

Phrase queries

- Want to be able to answer queries such as “*stanford university*” – as a phrase
- Thus the sentence “*I went to university at Stanford*” is not a match <https://eduassistpro.github.io/>
 - The concept of phrase query is easily understood by users, one of the advanced search ideas that works
 - Many more queries are *implicit phrase queries*
- For this, it no longer suffices to store only $\langle \text{term} : \text{docs} \rangle$ entries

Solution 1: Biword indexes

- Index every consecutive pair of terms in the text as a phrase
- For example the text “Friends, Romans, Countrymen” <https://eduassistpro.github.io/iwords>
 - *friends romans*
 - *romans countrymen*
- Each of these biwords is now a dictionary term
- Two-word phrase query-processing is now immediate.

Longer phrase queries

- Longer phrases are processed as we did with wild-cards:
- *stanford university palo alto* can be broken into the Boolean query [https://eduassistpro.github.io/
stanford university AND univ AND palo alto
Add WeChat edu_assist_pro](https://eduassistpro.github.io/stanford_university_AND_univ_AND_palo_alto)

Without the docs, we cannot verify that the docs matching the above Boolean query do contain the phrase.

Can have false positives!

Issues for biword indexes

- False positives, as noted before
- Index blowup due to bigger dictionary
 - Infeasible for more than biwords, big even for them

<https://eduassistpro.github.io/>

- Biword indexes are not the solution (for all biwords) but can be part of a strategy

Add WeChat edu_assist_pro

Solution 2: Positional indexes

- In the postings, store, for each ***term*** the position(s) in which tokens of it appear:

Assignment Project Exam Help

<***term***, number<https://eduassistpro.github.io/>
doc1: position1, position2 ... ;
doc2: position1, position2 ... ;
etc.>

Add WeChat edu_assist_pro

Positional index example

<*be*: 993427;

1: 7, 18, 33, 72, 86, 231;

2: 3, 149;

4: 17, 191, 29 <https://eduassistpro.github.io/>

5: 363, 367, ...
Add WeChat edu_assist_pro

Which of docs 1,2,4,5
could contain “*to be*
or *not to be*”?

- For phrase queries, we use a merge algorithm recursively at the document level
- But we now need to deal with more than just equality

Processing a phrase query

- Extract inverted index entries for each distinct term:
to, be, or, not.
- Merge the ~~Assignment~~ ~~position~~ ~~Projects~~ ~~Exam~~ ~~Help~~ all positions with
<https://eduassistpro.github.io/>
 - *to:*
 - 2:1,17,74,222,551; 4:8,16,3,7:13,23,191; ...
 - *be:*
 - 1:17,19; 4:17,191,291,430,434; 5:14,19,101; ...
- Same general method for proximity searches

Proximity queries

- LIMIT! /3 STATUTE /3 FEDERAL /2 TORT
 - Again, here, $/k$ means “within k words of”.
- Clearly, positional indexes can be used for such queries; biwo <https://eduassistpro.github.io/>
- Exercise: Adapt the linear strings to make it work for any value of k ?
 - This is a little tricky to do correctly and efficiently
 - **See Figure 2.12 of IIR (Page 39)**
 - There's likely to be a problem on it!

Assignment Project Exam Help

<https://eduassistpro.github.io/>

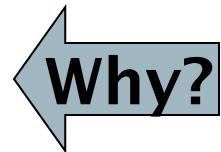
Add WeChat edu_assist_pro

Positional index size

- You can compress position values/offsets: we'll talk about that in lecture 5
- Nevertheless, a positional index expands postings storage *subst* <https://eduassistpro.github.io/>
- Nevertheless, a positional index is standardly used because of the power and proximity queries ... whether used explicitly or implicitly in a ranking retrieval system.

Positional index size

- Need an entry for **each occurrence**, not just once per document
- Index size depends on average document size
 - Average web <https://eduassistpro.github.io/>
 - SEC filings, books, even some ... easily 100,000 terms
- Consider a term with frequency 0.1%



Document size	Postings	Positional postings
1000	1	1
100,000	1	100

Rules of thumb

- A positional index is 2–4 as large as a non-positional index
- Positional index size 35–50% of volume of original text
<https://eduassistpro.github.io/>
- Caveat: all of this holds for “e” languages
Add WeChat edu_assist_pro

Combination schemes

- These two approaches can be profitably combined
 - For particular phrases (“*Michael Jackson*”, “*Britney Spears*”) it <https://eduassistpro.github.io/> postings list
 - Even more so for phrases like **Add WeChat edu_assist_pro**
- Williams et al. (2004) evaluate a more sophisticated mixed indexing scheme
 - A typical web query mixture was executed in $\frac{1}{4}$ of the time of using just a positional index
 - It required 26% more space than having a positional index alone

Resources for today's lecture

- *Introduction to Information Retrieval*, chapter 1
- Shakespeare:
 - [Assignment Project Exam Help](http://www.phyzone.com/shakespeare/)
 - Try the neat site <https://eduassistpro.github.io/>
- Add WeChat edu_assist_pro
- *Managing Gigabytes*, chapter 1
- *Modern Information Retrieval*, chapter 8.2

Resources for today's lecture

- Skip Lists theory: Pugh (1990)
 - Multilevel skip lists give same $O(\log n)$ efficiency as trees
- H.E. Williams, J. Zobel, and D. Bahle, 2004, “Fast Phrase Querying with Combined Indexes”, ACM Transactions on Information Systems, <https://eduassistpro.github.io/>
<http://www.seg.rmit.edu.au/~zobel/pubs/tis04.pdf>
- D. Bahle, H. Williams, and J. Zobel, Efficient querying with an auxiliary index. SIGIR 2002, pp. 215-2