

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`
Lecture 1: Bool at

Assignment Project Exam Help

Unstructured data

Add WeChat edu_assist_pro

- Which plays of Shakespeare contain the words ***Brutus*** AND ***Caesar*** but *NOT Calpurnia*?
- One could **grep** all of Shakespeare's plays for ***Brutus*** and ***Caesar***, t <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

Assignment Project Exam Help

Term-document in

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

Assignment Project Exam Help

Incidence vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for *Brutus*, *Caesar* and *Calpurnia* (complemented) \rightarrow bitwise AND.
- 110100 AND <https://eduassistpro.github.io/> 100100.

Add WeChat edu_assist_pro

Assignment Project Exam Help

Answers to query Add WeChat edu_assist_pro

- Antony and Cleopatra, Act III, Scene ii

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

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

Add WeChat edu_assist_pro

- Hamlet, Act III, Scene ii

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



Assignment Project Exam Help

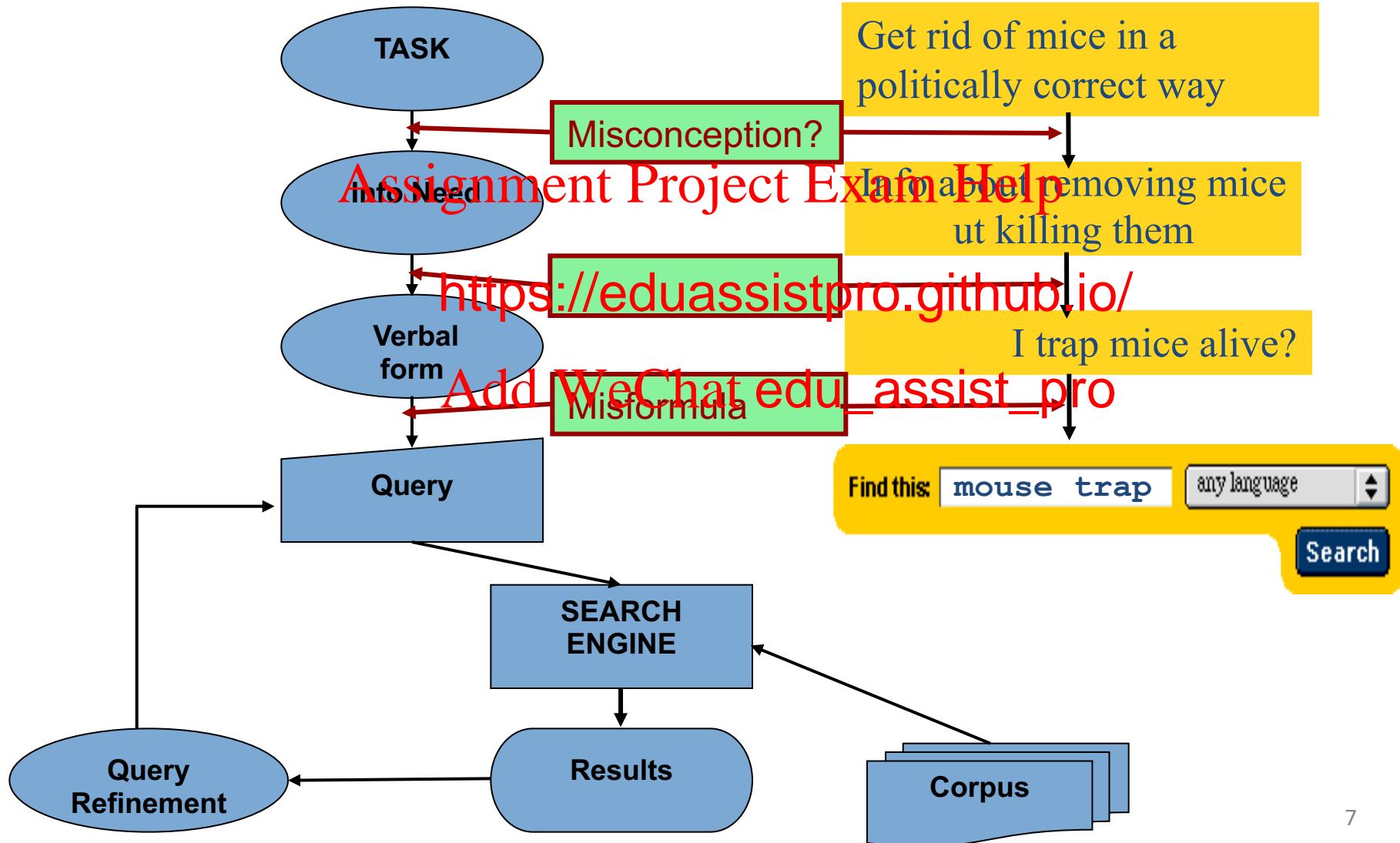
Basic assumptions of ~~Add WeChat edu_assist_pro~~ Retrieval

- **Collection:** Fixed set of documents
- **Goal:** Retrieve documents with information that is relevant to the user's information need and helps the user complete <https://eduassistpro.github.io/>

Add WeChat ~~edu_assist_pro~~

Assignment Project Exam Help

The classic search



Assignment Project Exam Help

How good are the docs?

- *Precision* : Fraction of retrieved docs that are relevant to user's information need
- *Recall* : Fraction of relevant docs in collection that are retrieved <https://eduassistpro.github.io/>
- More precise definitions a elements to follow in later lectures

Assignment Project Exam Help

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.

Assignment Project Exam Help

Can't build the ma

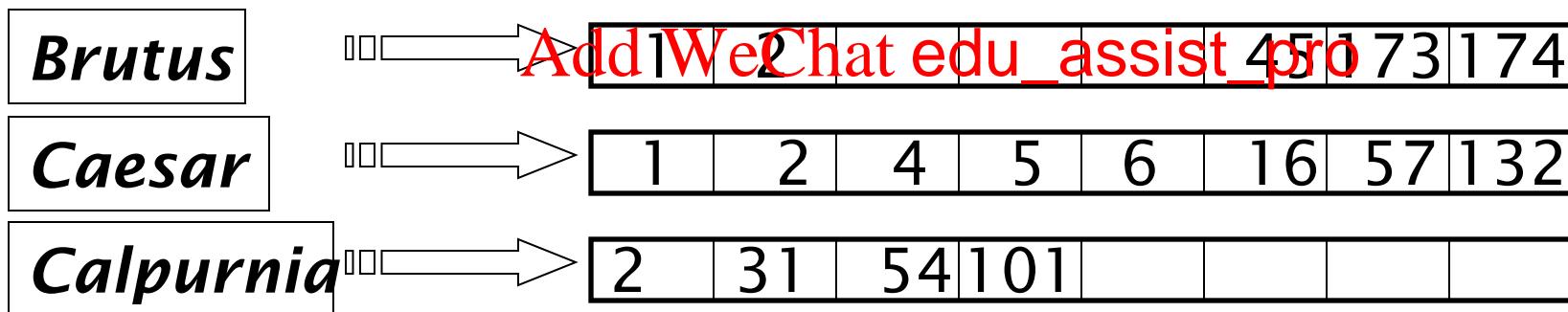
- 500K x 1M matrix has half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
 - matrix is extremely sparse.
- What's a bett <https://eduassistpro.github.io/>
 - We only record the 1 positions.



Assignment Project Exam Help

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



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

Assignment Project Exam Help

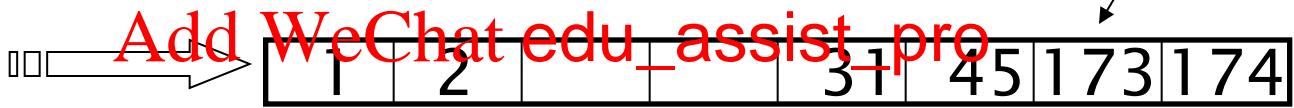
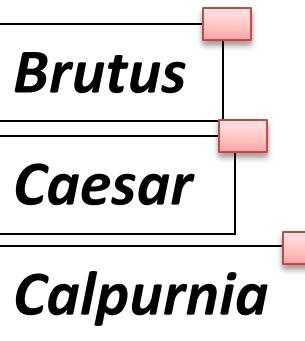
Inverted index

Add WeChat edu_assist_pro

- 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



Dictionary

Postings

Sorted by docID (more later on why).

Assignment Project Exam Help

Inverted index construction

Documents to be indexed.



Friends, Romans, countrymen.

Token stream.

Assignment Project Exam Help

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

omans

Countrymen

More on these later.

Add WeChat edu_assist_pro
Linguistic modules

Modified tokens.

friend

roman

countryman

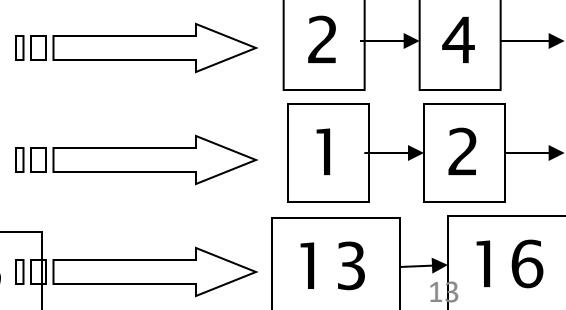
Inverted index.

Indexer

friend

roman

countryman



Assignment Project Exam Help

Indexer steps: Tokense

- 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

Assignment Project Exam Help

Indexer steps: Sort

- ## ■ Sort by terms

- And then docID

Assignment Project Exam Help

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

Core indexin

Add WeChat edu_assist_pro

Assignment Project Exam Help

Indexer steps: Dicti Add WeChat edu_assist_pro Postings

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

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

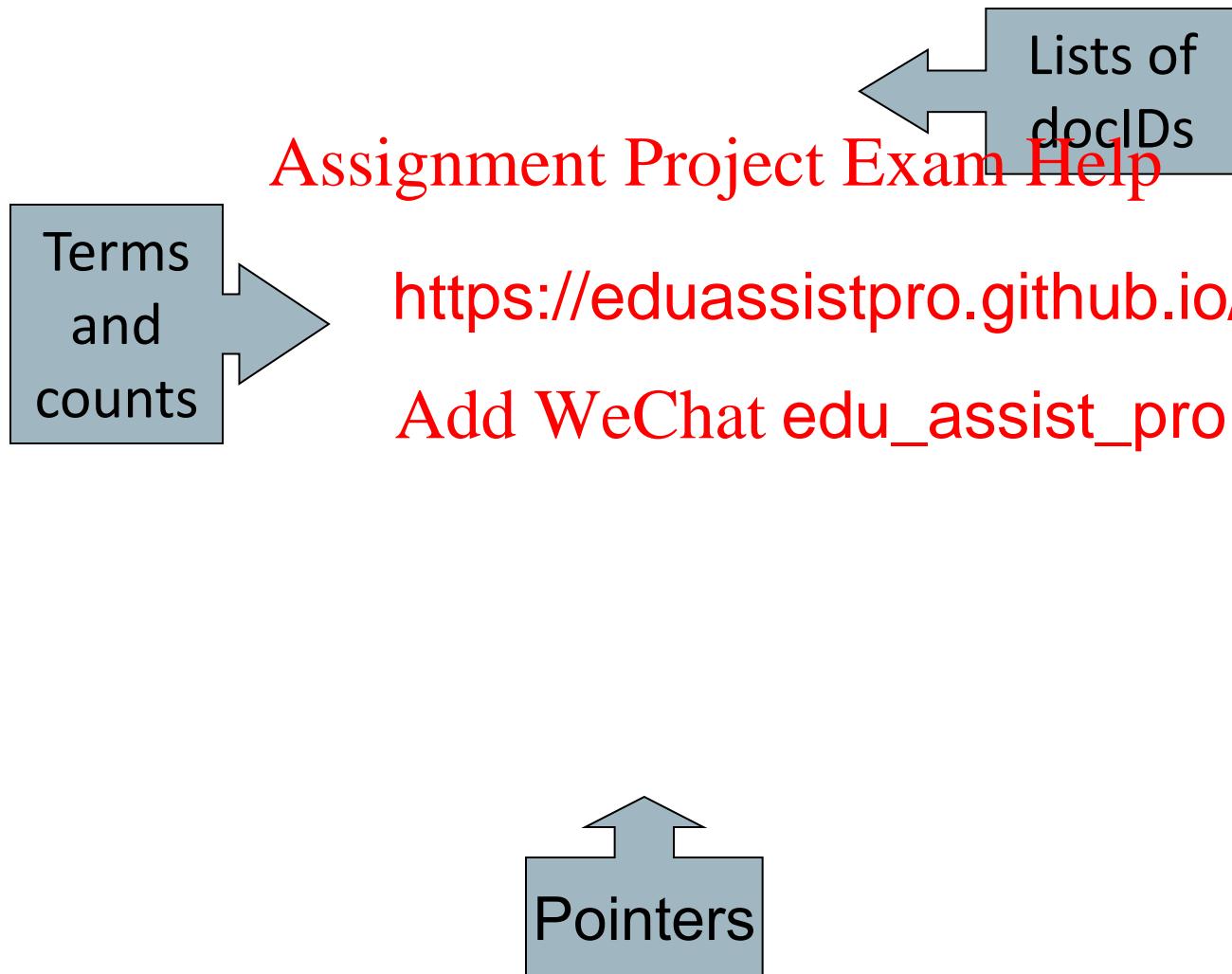


Add WeChat edu_assist_pro

Why frequency?
Will discuss later.

Assignment Project Exam Help

Where do we pay? Add WeChat | edu_assist_pro



Assignment Project Exam Help

More on Inverted Index | Add WeChat edu_assist_pro

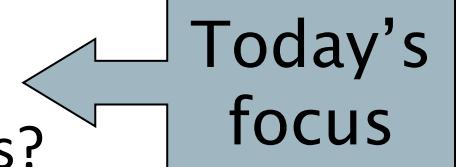
- Hashing-based construction methods are more efficient (though harder to implement)
- Inverted index can be compressed to
 - reduce index <https://eduassistpro.github.io/>
 - reduces transfer time between client and memory

Add WeChat edu_assist_pro

Assignment Project Exam Help

The index ~~Add WeChat~~ edu_assist_pro

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



Assignment Project Exam Help

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

Add WeChat edu_assist_pro

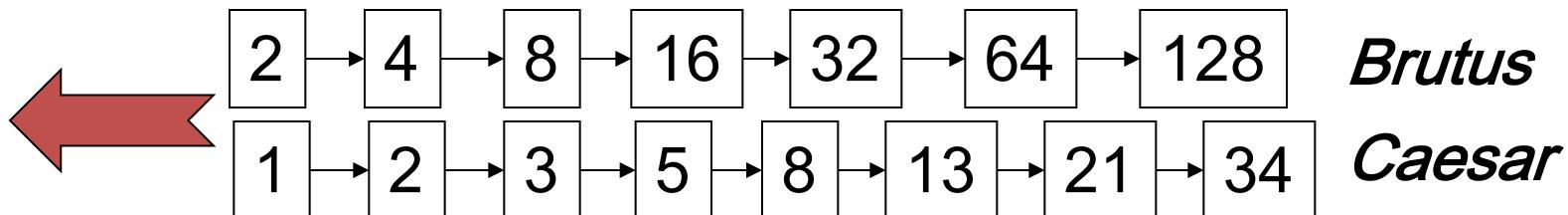
Assignment Project Exam Help

Query processing:

- Consider processing the query:

Brutus AND Caesar

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



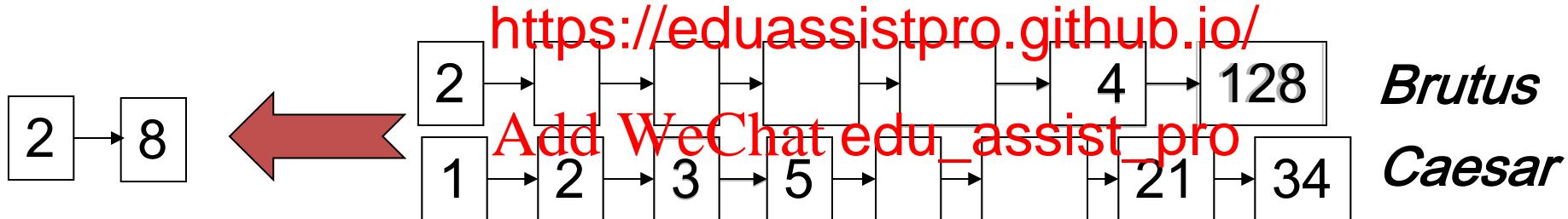
Assignment Project Exam Help

The merge

Add WeChat edu_assist_pro

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

Add WeChat edu_assist_pro



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

Assignment Project Exam Help

Boolean queries: E

- 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:
 - 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

Assignment Project Exam Help

Example: Add WeChat `edu_assist_pro` [.westlaw.com/](https://westlaw.com/)

- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992)
- Tens of ter <https://eduassistpro.github.io/> 00 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`

Assignment Project Exam Help

Example: Add WeChat `edu_assist_pro` [.westlaw.com/](https://westlaw.com/)

- Another example query:
 - Requirements for disabled people to be able to access a workplace
 - `disabl! /p a` <https://eduassistpro.github.io/> 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...

Assignment Project Exam Help
Boolean queries:

More general merge Add WeChat edu_assist_pro

- 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?

Assignment Project Exam Help

Merging Add WeChat edu_assist_pro

What about an arbitrary Boolean formula?

(Brutus OR Caesar) AND NOT

(Antony OR Cleo) Assignment Project Exam Help

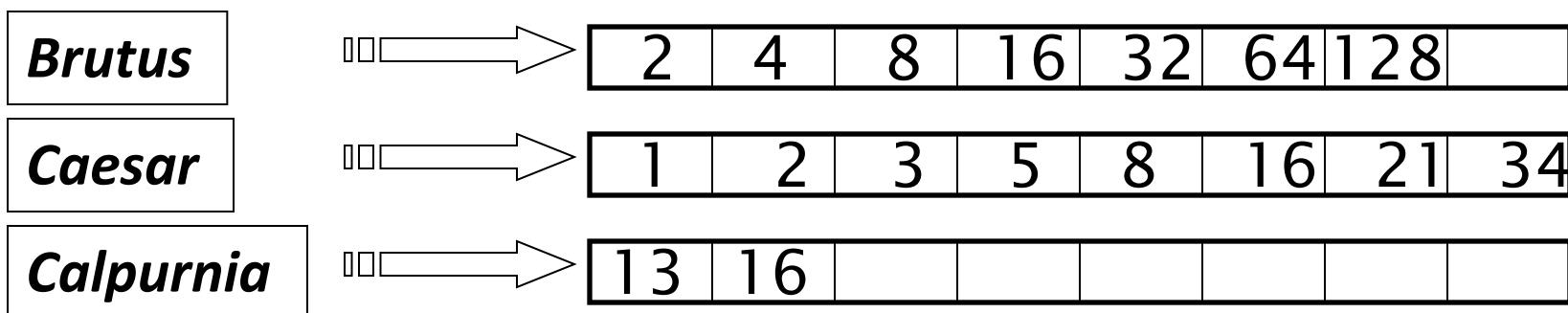
- Can we always <https://eduassistpro.github.io/>?
 - Linear in what? Add WeChat edu_assist_pro
- Can we do better?

Assignment Project Exam Help

Query optimization

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

Add WeChat edu_assist_pro



Query: **Brutus AND Calpurnia AND Caesar**

Assignment Project Exam Help

Query optimization

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

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

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

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

Assignment Project Exam Help

More general opti Add WeChat edu_assist_pro

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

Assignment Project Exam Help

Exercise Add WeChat edu_assist_pro

- Recommend a query processing order for

Assignment Project Exam Help

*(tangerine OR tree) AND https://eduassistpro.github.io/
(marmalade OR skies) AND
(kaleidoscope OR eyes)*

Q: Any more accurate way to estimate the cardinality of intermediate results?

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

Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

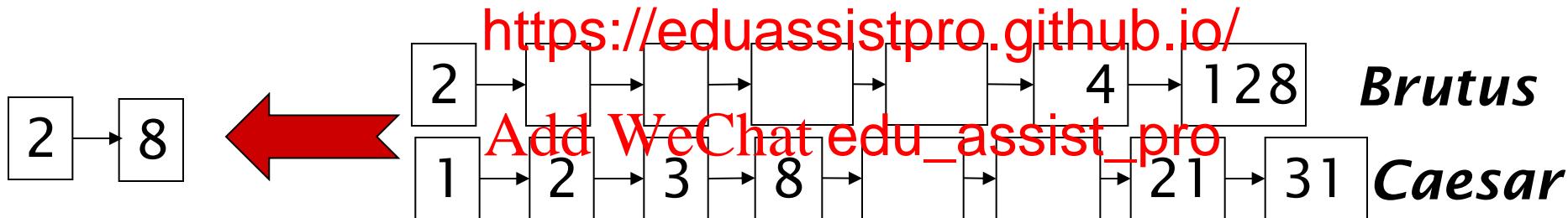
**FASTER POSTINGS MERGES:
SKIP POINTERS/SKIP LISTS**

Assignment Project Exam Help

Recall basic merge

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

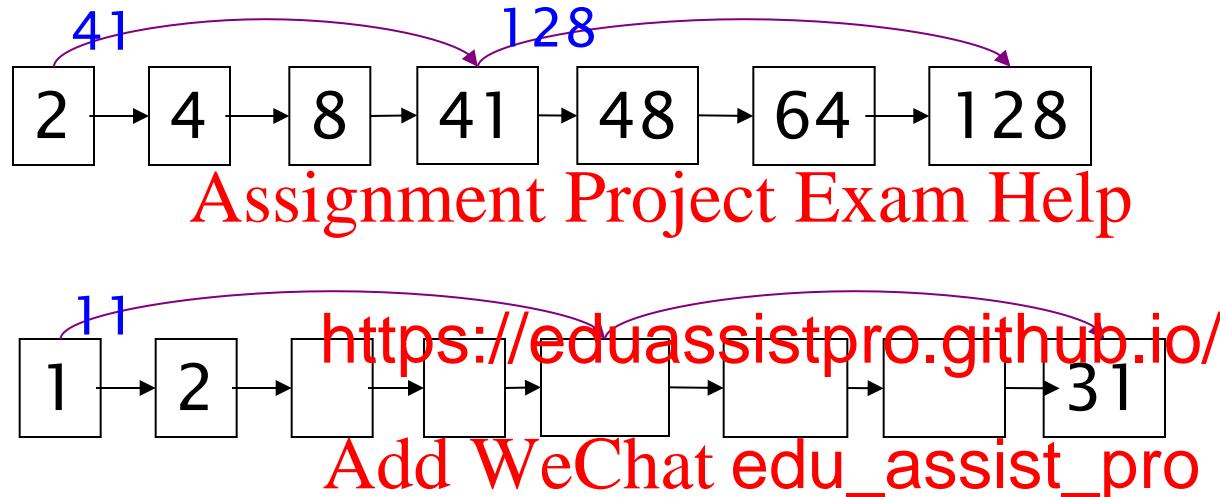
Assignment Project Exam Help



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)

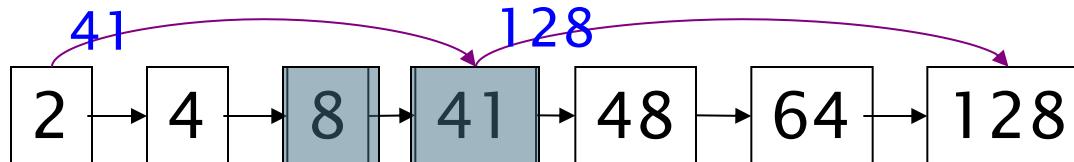


- Why?
- To skip postings that will not figure in the search results.
- How?
- Where do we place skip pointers?

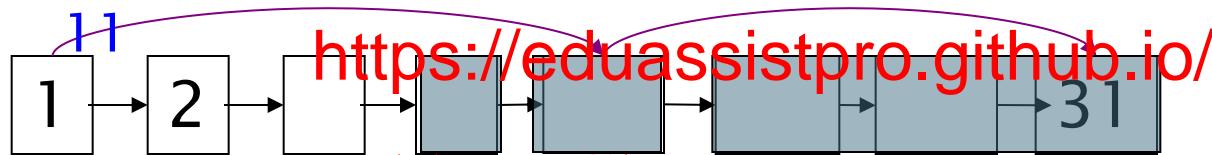
Assignment Project Exam Help

Query processing

Add WeChat edu_assist_pro



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.

Assignment Project Exam Help

Can we skip w/o skip pointers?

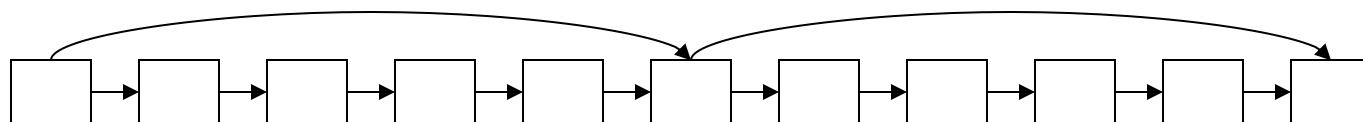
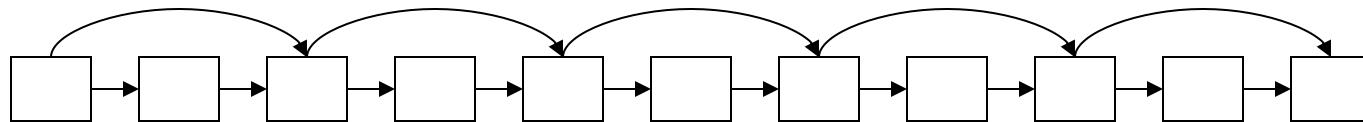
Where do we place skip pointers?

- Tradeoff:
 - More skips → shorter skip spans → more likely to skip.

But lots of comparisons to skip pointers.

- Fewer skips → few long skip spans → few comparisons to skip pointers.

Add WeChat edu_assist_pro



Assignment Project Exam Help

Placing skips

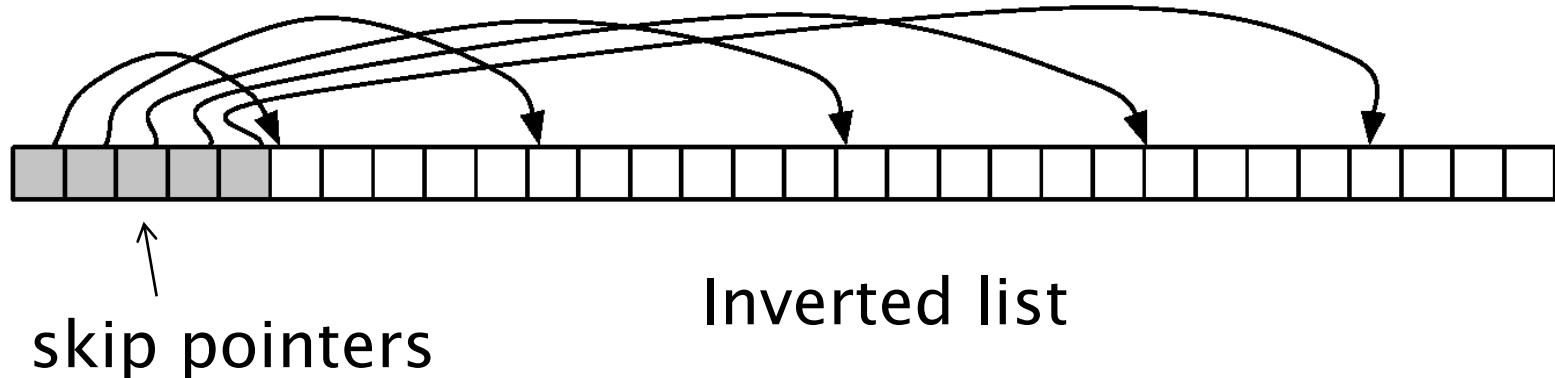
Add WeChat edu_assist_pro

- 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 static, 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!

Skip Pointers

- A skip pointer (d, p) contains a document number d and a byte (or bit) position p
 - Means there is an inverted list posting that starts at position p , a [as for document \$d\$](https://eduassistpro.github.io/)
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Assignment Project Exam Help

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

PHRASE QUERIES AND POSITIONAL INDEXES

Assignment Project Exam Help

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

Assignment Project Exam Help

Solution 1: Biword

- 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.

Assignment Project Exam Help

Longer phrase que

- Longer phrases are processed as we did with wild-cards:
- *stanford university palo alto* can be broken into the Boolean quer <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!

Assignment Project Exam Help

Extended biwords

- Parse the indexed text and perform part-of-speech-tagging (POST).
- Bucket the terms into (say) Nouns (N) and articles/prepositions (X).
- Call any string o <https://eduassistpro.github.io/> an extended biword.
 - Each such ext de a term i dictionary.
- Example: *catcher in the rye*

N	X	X	N
---	---	---	---
- Query processing: parse it into N's and X's
 - Segment query into enhanced biwords
 - Look up in index: *catcher rye*

Assignment Project Exam Help

Issues for biword indexing

- 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 hybrid strategy

Assignment Project Exam Help

Solution 2: Position

- 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.>

Assignment Project Exam Help

Positional index ex

<*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

Assignment Project Exam Help

Processing a phrase query *Add WeChat edu_assist_pro*

- Extract inverted index entries for each distinct term:
to, be, or, not.
- Merge the *Assignment Project Exam Help* terms into 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

Assignment Project Exam Help

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 handle proximity queries.
 - 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

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

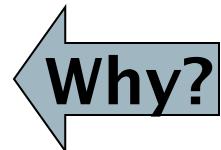
Positional index size

- 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.

Assignment Project Exam Help

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

Assignment Project Exam Help

Rules of thumb

- A positional index is 2–4 times 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

Assignment Project Exam Help

Combination scheme

- These two approaches can be profitably combined

Assignment Project Exam Help

- For particular phrases (“*Michael Jackson*”, “*Britney Spears*”) it <https://eduassistpro.github.io/> provides positional postings lists
 - Even more so for phrases like “*Michael Jackson*”
- 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

[Optional]

Assignment Project Exam Help

\$ < \text{any char}

Solution 3: Suffix T

- BANANA\$

- BANANA\$ pos:0
- ANANA\$ pos:1
- NANA\$ pos:2
- ANA\$ pos:3
- NA\$ pos:4 <https://eduassistpro.github.io/>
- A\$ pos:5



Sort on the strings

Add WeChat edu_assist_pro

- A\$ pos:5
- ANA\$ pos:3
- ANANA\$ pos:1
- BANANA\$ pos:0
- NA\$ pos:4
- NANA\$ pos:2

[Optional]

Assignment Project Exam Help

\$ < any char

Suffix Array

Add WeChat edu_assist_pro

- BANANA\$

- BANANA\$ pos:0
- ANANA\$ pos:1
- NANA\$ pos:2
- ANA\$ pos:3
- NA\$ pos:4
- A\$ pos:5
- \$ pos:6



- \$ pos:6
- A\$ pos:5
- ANA\$ pos:3
- ANANA\$ pos:1
- BANANA\$ pos:0
- NA\$ pos:4
- NANA\$ pos:2

- If the original string is available, each suffix can be completely specified by the first character

<https://eduassistpro.github.io/>: 5n

Add WeChat edu_assist_pro

Sort on the strings



B	A	N	A	N	A	\$
6	5	3	1	0	4	2

← Suffix array

Binary search (using offsets to fetch the 'key')

Assignment Project Exam Help

Resources for today

- *Introduction to Information Retrieval*, chapter 1
- Shakespeare:
 - <http://www.rhymezone.com/shakespeare/>
 - Try the neat site <https://eduassistpro.github.io/>
- *Managing Gigabytes*, chapter 1
- *Modern Information Retrieval*, chapter 8.2

Assignment Project Exam Help

Resources for today

- 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/phrse.pdf>
- D. Bahle, H. Williams and J. Zobel. Efficient querying with an auxiliary index. SIGIR 2002, pp. 215-2