

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 4: Index on

# Assignment Project Exam Help

## Plan

Add WeChat `edu_assist_pro`

---

- Last lecture:

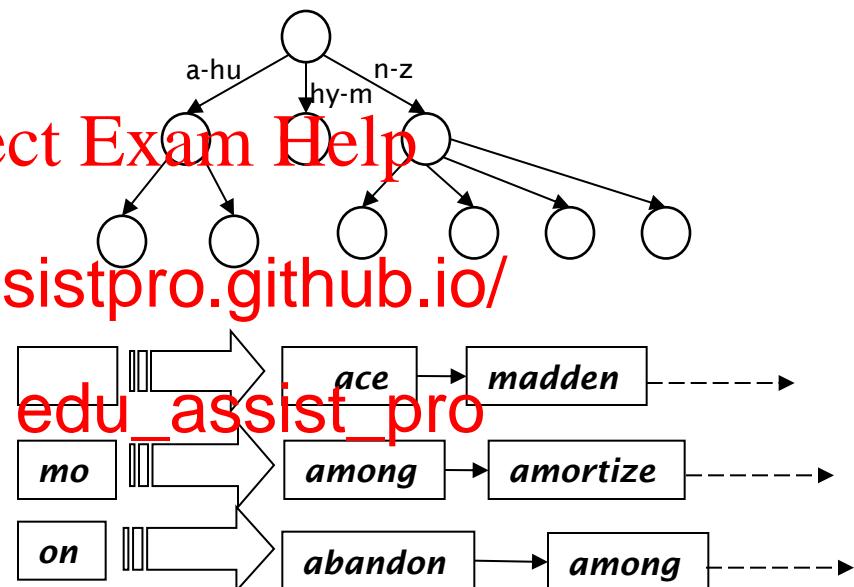
- Dictionary data structures
- Tolerant retrieval
  - Wildcards
  - Spell correction
  - Soundex

Assignment Project Exam Help

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

- This time:

- Index construction



# Assignment Project Exam Help

## Index construction

---

- How do we construct an index?
- What strategies can we use with limited main memory?  
Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Hardware basics

---

- Many design decisions in information retrieval are based on the characteristics of hardware
- We begin by reviewing hardware basics

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Hardware basics

- Access to data in memory is ***much*** faster than access to data on disk.
- Disk seeks: No data is transferred from disk while the disk head is b <https://eduassistpro.github.io/>
- Therefore: Transferring on disk to memory is faster th nk of data from disk to memory is faster th rring many small chunks.
- Disk I/O is block-based: Reading and writing of entire blocks (as opposed to smaller chunks).
- Block sizes: 8KB to 256 KB.

# Assignment Project Exam Help

## Hardware basics

---

- Servers used in IR systems now typically have several GB of main memory, sometimes tens of GB.
- Available disk space is several (2–3) orders of magnitude larger <https://eduassistpro.github.io/>
- Fault tolerance is very expensive to use many regular machines, much cheaper to use one fault tolerant machine.

# Assignment Project Exam Help

## Hardware Add WeChat edu\_assist\_pro

---

symbol	statistic	value
s	average seek time	$5 \text{ ms} = 5 \times 10^{-3} \text{ s}$
b	transf	$2 \mu\text{s} = 2 \times 10^{-8} \text{ s}$
	proces	$\text{https://eduassistpro.github.io/}$
p	low-level operation (e.g., compare & swap a word)	$\mu\text{s} = 10^{-8} \text{ s}$
	size of main memory	several GB
	size of disk space	1 TB or more

# Assignment Project Exam Help

## RCV1: Our collection is lecture

---

- Shakespeare's collected works definitely aren't large enough for demonstrating many of the points in this course. [Assignment Project Exam Help](#)
- The collection <https://eduassistpro.github.io/> either, but it's at least a more plausible example. [Add WeChat edu\\_assist\\_pro](#)
- As an example for applying scalable index construction algorithms, we will use the Reuters RCV1 collection.
- This is one year of Reuters newswire (part of 1995 and 1996)

# Assignment Project Exam Help

A Reuters RCV1 do Add WeChat edu\_assist\_pro

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Reuters RCV1 statistics (grounded)

symbol	statistic	value
N	documents	800,000
L	avg. # tokens per doc	200
M	avg. # bytes per token (incl. spaces/punct.)	4.5
	avg. # bytes per term	7.5
	non-positional postings	100,000,000

# Assignment Project Exam Help

## Recall IIR 1 index

Add WeChat edu\_assist\_pro

- Documents are parsed to extract words and these are saved with the Document ID.

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

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

Add WeChat edu\_assist\_pro

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

# Assignment Project Exam Help

## Key step Add WeChat edu\_assist\_pro

- After all documents have been parsed, the inverted file is sorted by terms.

Assignment Project Exam Help

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

We focus on this

We have 100M items to sort.

Add WeChat edu\_assist\_pro

Term	Doc #
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
it	1
killed	1
the	1
the	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
caesar	2
was	2
was	1
was	2
with	2

# Hash based in-memory index construction

- Another in-memory index construction method is to use hash-tables
  - Append  $(docid, pos)$  to the existing (partial) postings list of the token; crucially necessary
- Generally, fast
- Further optimizations
  - Dealing with collision: insert-at-back and move-to-front heuristics
  - Saving space and time: use ArrayList to implement postings lists

# Assignment Project Exam Help

## Scaling index const

- In-memory index construction does not scale.
- How can we construct an index for very large collections?  
Assignment Project Exam Help
- Taking into account <https://eduassistpro.github.io/> learned about  
Add WeChat edu\_assist\_pro
- Memory, disk, speed, etc.

# Assignment Project Exam Help

## Sort-based index

- As we build the index, we parse docs one at a time.
  - While building the index, we cannot easily exploit compression tricks (you can, but much more complex)
- The final posting lists are not yet complete until the end.
- At 12 bytes per posting list (term, doc, freq), demands a lot of space for large collections.
- $T = 100,000,000$  in the case of R
  - So ... we can do this in memory in 2009, but typical collections are much larger. E.g. the *New York Times* provides an index of >150 years of newswire
- Thus: We need to store intermediate results on disk.

# Assignment Project Exam Help

## Use the same algo Add WeChat edu\_assist\_pro on disk?

---

- Can we use the same index construction algorithm for larger collections, but by using disk instead of memory? Assignment Project Exam Help
- No: Sorting T <https://eduassistpro.github.io/> on disk is too slow – too ma
- We need an external sorting algo Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Bottleneck

Add WeChat edu\_assist\_pro

- Parse and build postings entries one doc at a time
- Now sort postings entries by term (then by doc within each term)  
Assignment Project Exam Help
- Doing this with <https://eduassistpro.github.io/> could be too slow
  - must sort  $T=100M$  record

Add WeChat edu\_assist\_pro



If every comparison took 2 disk seeks, and  $N$  items could be sorted with  $N \log_2 N$  comparisons, how long would this take?

# BSBI: Blocked sort-based Indexing

## (Sorting with few weeks)

- 12-byte (4+4+4) records (*term, doc, freq*).
- These are generated as we parse docs.
- Must now sort records by *term*.
- Define a Block <https://eduassistpro.github.io/>
  - Can easily fit a couple into memory
  - Will have 10 such blocks to sort
- Basic idea of algorithm:
  - Accumulate postings for each block, sort, write to disk.
  - Then merge the blocks into one long sorted order.

# 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

## Sorting 10 blocks of records

- First, read each block and sort within:
  - Quicksort takes  $2N \ln N$  expected steps
  - In our case  $2 \times (10M \ln 10M)$  steps
- *Exercise: estimate the time to sort each block from disk and add them together.*
- 10 times this estimate – giving 10 sorted runs of 10M records each.
- Done straightforwardly, need 2 copies of data on disk
  - But can optimize this

# 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

## Example Add WeChat edu\_assist\_pro

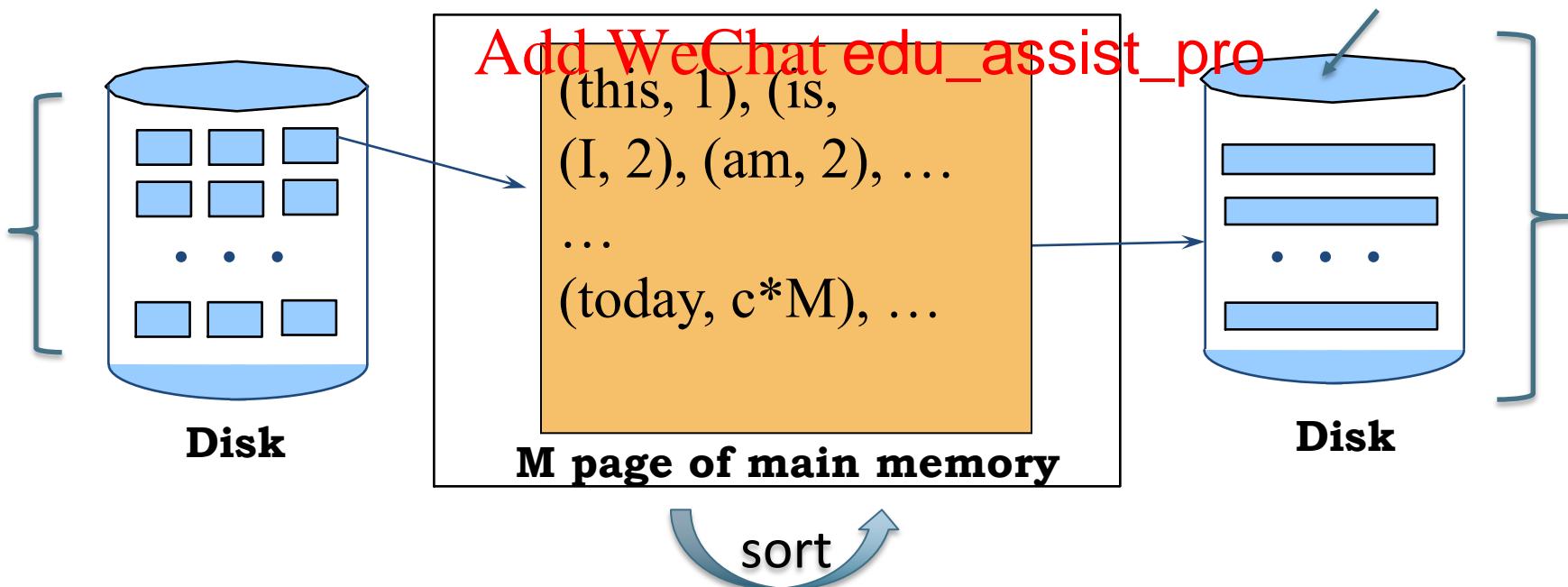
---

- Settings
  - **B**: Block/page size
  - **M**: Size of  $m$  (., = 10 blocks)
  - **N**: Number <https://eduassistpro.github.io/>
  - **R**: Size of the document emitted.
- Simplifying assumptions:
  - **R**: the same for all documents
  - **B** =  $c * R$ , for some integer  $c$  (e.g.,  $c = 5$ )
  - All I/Os have the same cost

# Assignment Project Exam Help

## External Merge-~~Sort~~ Add WeChat edu\_assist\_pro

- Phase I: load the (term, docID) pairs from  $(M*B)/R$  documents (*at a time*) into M buffer pages; sort
  - Result: ~~(Assignment Project Exam Help)~~ runs of length M pages
  - # of runs = 2 <https://eduassistpro.github.io/>

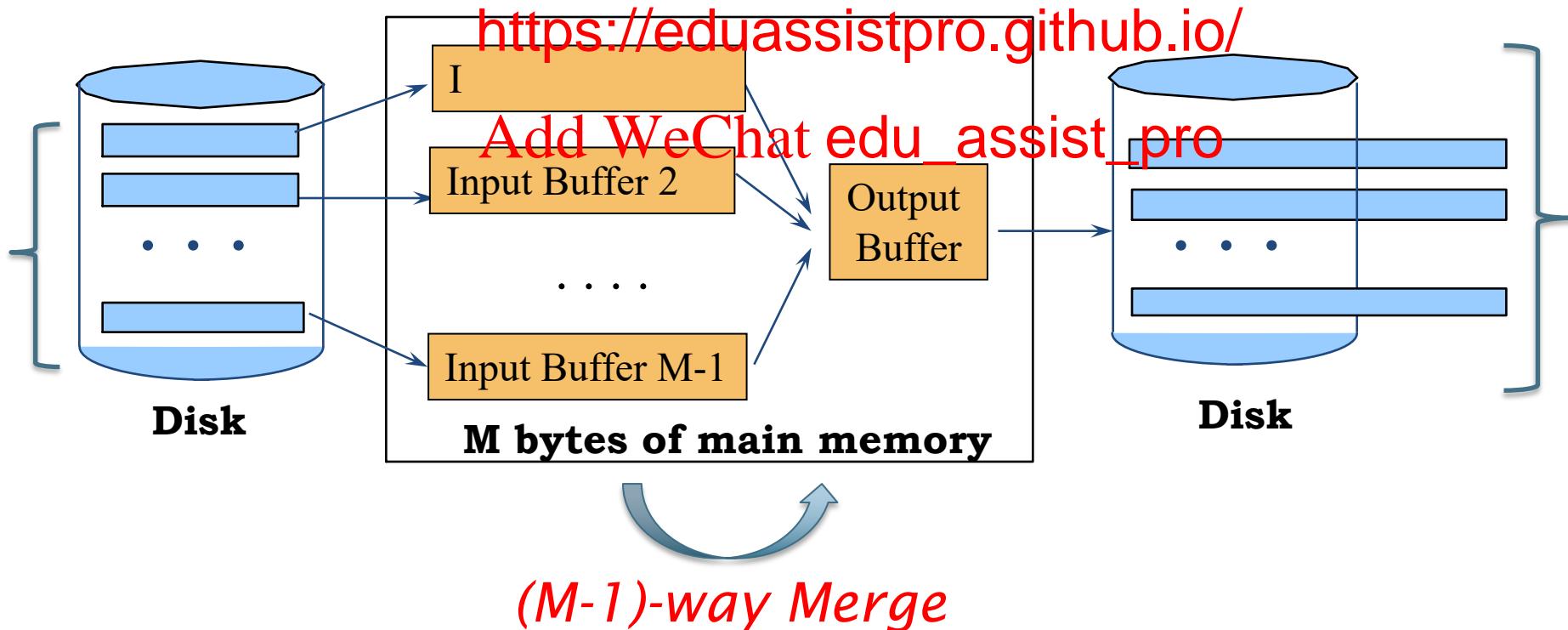


# Assignment Project Exam Help

## Phase II /1 Add WeChat edu\_assist\_pro

- Recursively merge (up to)  $M - 1$  runs into a new run
- Result: runs of length  $M$  ( $M - 1$ ) pages

Assignment Project Exam Help

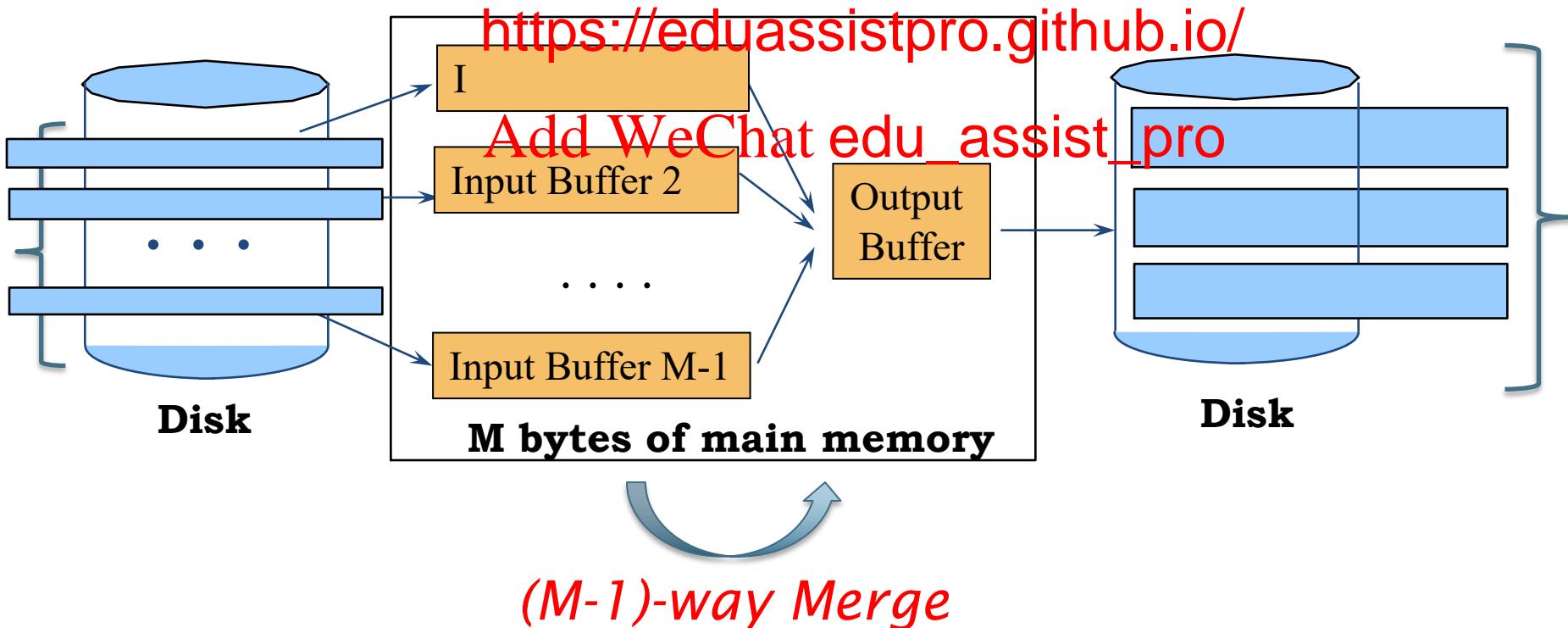


# Assignment Project Exam Help

## Phase II /2 Add WeChat edu\_assist\_pro

- Recursively merge (up to)  $M - 1$  runs into a new run
- Result: runs of length  $M$   $(M - 1)^2$  pages

Assignment Project Exam Help

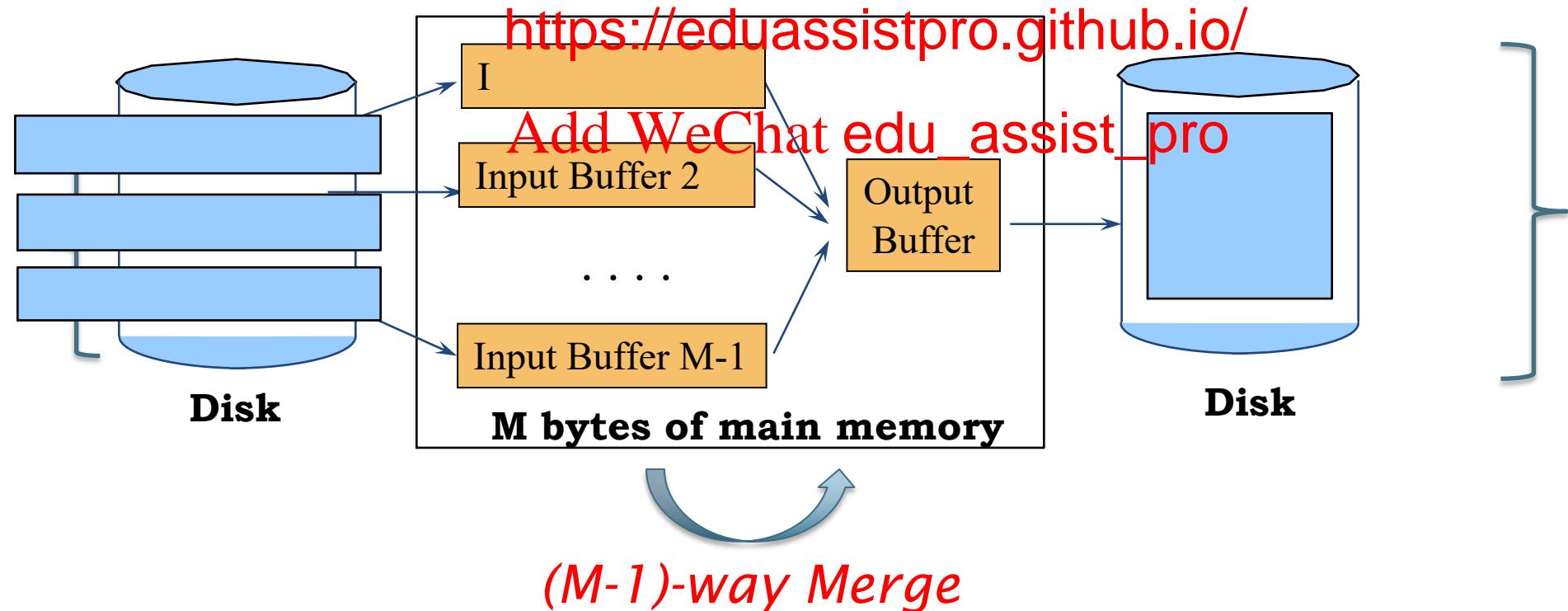


# Assignment Project Exam Help

## Phase II /3 Add WeChat edu\_assist\_pro

- Recursively merge (up to)  $M - 1$  runs into a new run
- Result: a single run

Assignment Project Exam Help



# Assignment Project Exam Help

## Cost of External M

- Number of passes:  $1 + \left\lceil \log_{M-1} \left[ \frac{NR}{MB} \right] \right\rceil$
- Total I/O cost:  $2 \cdot \binom{NR}{B} \cdot \left( 1 + \left\lceil \log_{M-1} \left[ \frac{NR}{MB} \right] \right\rceil \right)$  blocks/pages

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

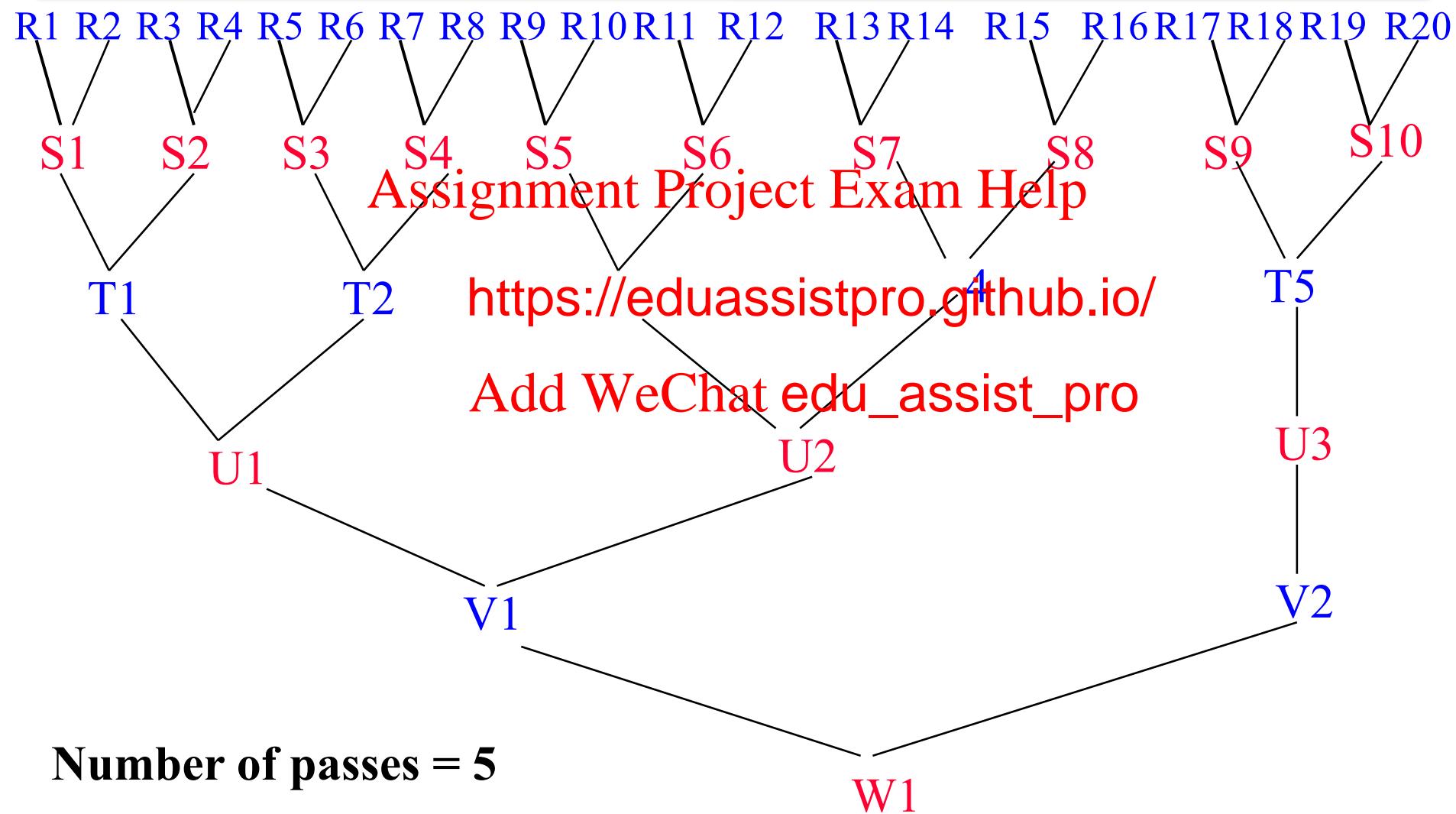
- How much data can we sort B RAM?
  - Assume B = 4KB
  - 1 pass → 10MB “data”
  - 2 passes →  $\approx 25GB$  “data” ( $M-1 = 2559$ )
  - 3 passes ?
- Can sort most reasonable inputs in 2 or 3 passes !

Another Example

# Assignment Project Exam Help

Example: 2-Way M r 20 Runs

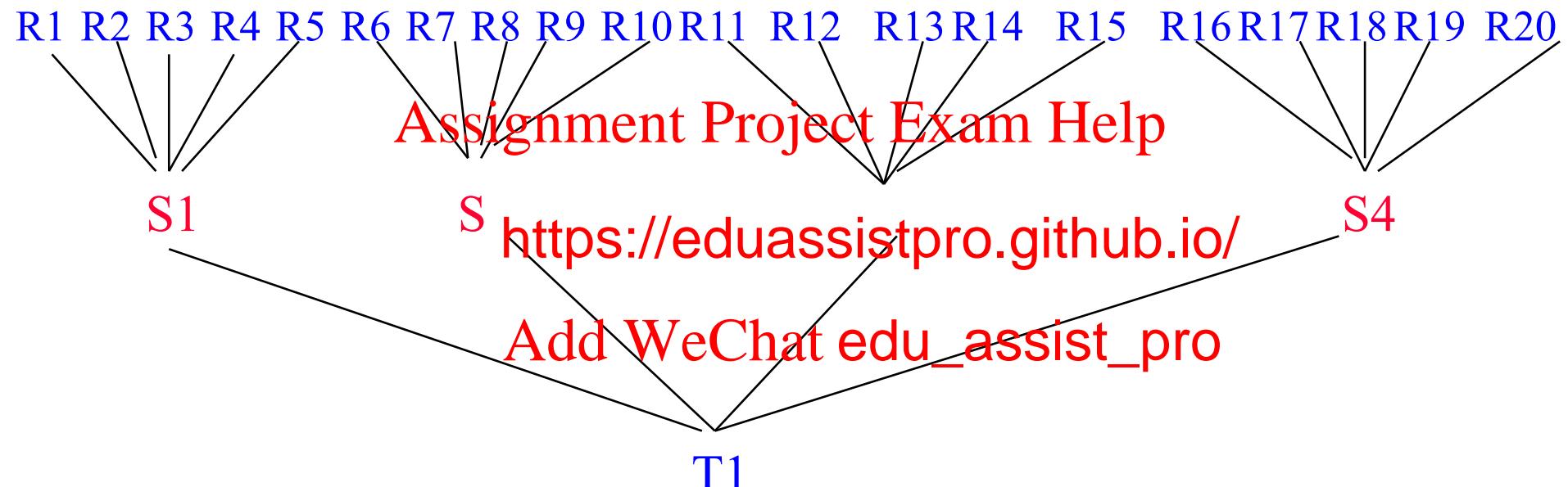
Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)



# Assignment Project Exam Help

Example: 5-Way Merge 20 Runs

---



**Number of passes = 2**

K-way Merge Add WeChat ~~Add WeChat~~ edu\_assist\_pro



Assignment Project Exam Help



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



Add WeChat ~~Add WeChat~~ edu\_assist\_pro

99



# Remaining problem with sort-based algorithm

[Assignment](#) [Project](#) [Exam](#) [Help](#)

Add WeChat edu\_assist\_pro

- Our assumption was: we can keep the dictionary in memory.
- We need the dictionary (which grows dynamically) in order to impl <https://eduassistpro.github.io/> D mapping.
- Actually, we could work wi cID postings instead of termID,docID po
- . . . but then intermediate files become very large. (We would end up with a scalable, but very slow index construction method.)

# SPIMI: Assignment Project Exam Help

## Single-pass ~~Add WeChat edu\_assist\_pro~~ in-memo

- Key idea 1: Generate separate dictionaries for each block – no need to maintain term-termID mapping across blocks.
- Key idea 2: D <https://eduassistpro.github.io/> postings lists
- With these two ideas we can generate a complete inverted index for each block.
- These separate indexes can then be merged into one big index.

# Assignment Project Exam Help

SPIMI-Invert Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

---

Assignment Project Exam Help

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

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

- Merging of blocks is analogous to BSBI.

# Assignment Project Exam Help

## SPIMI: Compression

- Compression makes SPIMI even more efficient.
  - Compression of terms
  - Compression of postings
- See next lecture <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Distributed indexing

- For web-scale indexing (don't try this at home!):  
must use a distributed computing cluster
- Individual machines are fault-prone
  - Can unpredictably fail
- How do we exploit such a problem?

For those interested in the topic, read the textbook for distributed indexing using the Map-Reduce paradigm.

Also check out:

[http://terrier.org/docs/v3.5/hadoop\\_indexing.html](http://terrier.org/docs/v3.5/hadoop_indexing.html)

# Assignment Project Exam Help

## Dynamic indexing

- Up to now, we have assumed that collections are static.
- They rarely are:
  - Documents change to be inserted.
  - Documents are deleted and removed from the collection.
- This means that the dictionary and postings lists have to be modified:
  - Postings updates for terms already in dictionary
  - New terms added to dictionary

# Assignment Project Exam Help

## Simplest approach

- Maintain “big” main index
- New docs go into “small” auxiliary index
  - Merge immediately with the big main index when memory is full  
<https://eduassistpro.github.io/>
- Search across
- Deletions
  - Invalidation bit-vector for deleted docs
  - Filter docs output on a search result by this invalidation bit-vector
- Periodically, re-index into one main index

# Assignment Project Exam Help

## Issues with main a Add WeChat edu\_assist\_pro lary indexes

---

- Problem of frequent merges – you touch stuff a lot
- Poor performance during merge
- Actually: **Assignment Project Exam Help**
  - Merging of the dex is efficient if we keep a separate <https://eduassistpro.github.io/>
  - Merge is the same as a simple app
  - But then we would need a lot of fil t for O/S.
- Assumption for the rest of the lecture: The index is one big file.
- In reality: Use a scheme somewhere in between (e.g., split very large postings lists, collect postings lists of length 1 in one file etc.)

# Assignment Project Exam Help

## Another Extreme

- Whenever memory is full, write the sub-index to the disk
  - Never merge sub-indexes
- Pros:
  - High indexing performance<https://eduassistpro.github.io/>
- Cons:
  - Slow query performance
    - Require  $\Omega(|C|/M)$  seeks to fetch the inverted list for a term

# Assignment Project Exam Help

## Compromise Log Add WeChat edu\_assist\_pro c-merge

---

- Comprise of the previous two extremes
- Generation of a sub-index
  - The one directly created from in-memory index has generation = <https://eduassistpro.github.io/>
  - Merge of multiple sub-indexes of same generation = g gives a new index with generation = g+1
- Invariant: no two sub-indexes can have the same generation
- When memory is full, create  $I_0$ 
  - If we have two sub-indexes of generation g, merge them to form a single sub-index of generation g+1

# Assignment Project Exam Help

Illustration Add WeChat edu\_assist\_pro

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Algorithm

Add WeChat edu\_assist\_pro

---

- Maintain a series of indexes, each twice as large as the previous one.
- Keep smallest ( $Z_0$ ) in memory
- Larger ones (I <https://eduassistpro.github.io/>)
- If  $Z_0$  gets too big ( $> n$ ), write  $I_0$
- or merge with  $I_0$  (if  $I_0$  already exists)  $\rightarrow Z_1$
- Either write merge  $Z_1$  to disk as  $I_1$  (if no  $I_1$ )
- or merge with  $I_1$  to form  $Z_2$
- etc.

# 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

## Logarithmic merge

- Auxiliary and main index: index construction time is  $O(C^2/M)$  as each posting is touched in each merge.
  - $C = \text{Collection size}$ , and  $M = \text{memory size}$
- Logarithmic <https://eduassistpro.github.io/> time  $O(\log C/M)$  (C log (C/M))
- So logarithmic merge is much more efficient for index construction
- But query processing now requires the merging of  $O(\log T)$  indexes
  - Whereas it is  $O(1)$  if you just have a main and auxiliary index

# Assignment Project Exam Help

Lucene      Add WeChat [edu\\_assist\\_pro](#)

---

- `Mergefactor` does not have to be 2
  - <http://lucene.sourceforge.net/talks/inktomi/>
- In animation ~~Assignment Project Exam Help~~
  - <http://blog.m-labs.ca/2012/02/visualizing-lucenes-seg.html>
  - <https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

# Assignment Project Exam Help

## Further issues with ~~Add WeChat edu\_assist\_pro~~ indexes

---

- Collection-wide statistics are hard to maintain
- E.g., when we spoke of spell-correction: which of several corrected alternatives do we present to the user?  
<https://eduassistpro.github.io/>
  - We said, pick the one with the most frequent spelling
- How do we maintain the total frequency across multiple indexes and invalidation bit vectors?
  - One possibility: ignore everything but the main index for such ordering
- Will see more such statistics used in results ranking

# Assignment Project Exam Help

## Dynamic indexing engines

- All the large search engines now do dynamic indexing
- Their indices have frequent incremental changes
  - News items, <https://eduassistpro.github.io/>
    - Sarah Palin, ...
- But (sometimes/typically) they periodically reconstruct the index from scratch
  - Query processing is then switched to the new index, and the old index is then deleted

# Assignment Project Exam Help

Add WeChat `edu_assist_pro`

---

Assignment Project Exam Help

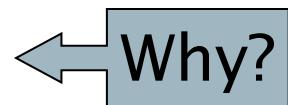
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

# Assignment Project Exam Help

## Other sorts of inde

- Positional indexes
  - Same sort of sorting problem ... just larger
- Building character n-gram indexes:
  - As text is par <https://eduassistpro.github.io/>
  - For each  $n$ -gram, need point tionary terms containing it – the “postings”
  - Note that the same “postings entry” will arise repeatedly in parsing the docs – need efficient hashing to keep track of this.
    - E.g., that the trigram you occurs in the term **deciduous** will be discovered on each text occurrence of **deciduous**
    - Only need to process each term once



# Assignment Project Exam Help

## Resources for today

- Chapter 4 of IIR
- MG Chapter 5
- Original publication by Dean and Ghemawat (2003) <https://eduassistpro.github.io/>
- Original publication on SPIE by Zobel (2003)