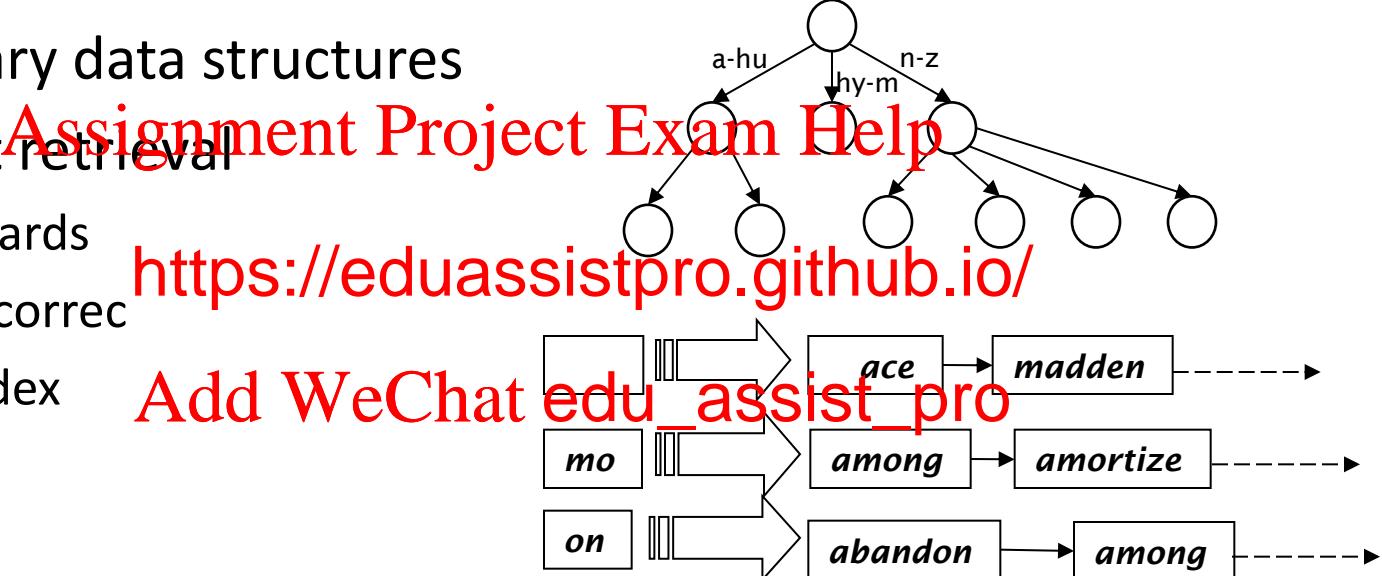


Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Lecture 4: Index on

Plan

- Last lecture:
 - Dictionary data structures
 - Tolerant retrieval
 - Wildcards
 - Spell correction
 - Soundex
- This time:
 - Index construction



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

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

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 Add WeChat edu_assist_pro}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.

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.

Hardware assumptions

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

RCV1: Our collection for this 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)

A Reuters RCV1 document

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Reuters RCV1 statistics (Rounded)

symbol	statistic	value
N	documents	800,000
L	avg. # tokens per doc	200
M	ter https://eduassistpro.github.io/	400,000
	avg. # bytes per token (incl. spaces/punct.)	Add WeChat edu_assist_pro
	avg. # bytes per token (without spaces/punct.)	4.5
	avg. # bytes per term	7.5
	non-positional postings	100,000,000

Recall IIR 1 index construction

- 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

Key step

- 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 #	Term	Doc #
I	1	ambitious	2
did	1	be	2
enact	1	brutus	1
julius	1	brutus	2
caesar	1	capitol	1
I	1	caesar	1
was	1	caesar	2
killed	1	caesar	2
it	1	did	1
the	1	enact	1
	1	hath	1
	1	I	1
	1	I	1
	1	i'	1
	2	it	2
be	2	julius	1
with	2	killed	1
caesar	2	killed	1
the	2	let	2
noble	2	me	1
brutus	2	noble	2
hath	2	so	2
told	2	the	1
you	2	the	2
caesar	2	told	2
was	2	you	2
ambitious	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; cr
<https://eduassistpro.github.io/> necessary me
- Generally, fas
- Further optimizations
 - Dealing with collision: ~~insert-at-back and move-to-front heuristics~~
 - Saving space and time: use ArrayList to implement postings lists

Scaling index construction

- 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/> constraints we just learned about
[Add WeChat](#) [edu_assist_pro](#)
- Memory, disk, speed, etc.

Sort-based index construction

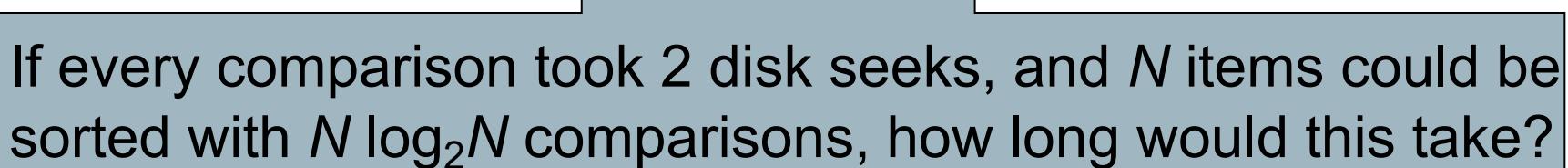
- As we build the index, we parse docs one at a time.
 - While building the index, we cannot easily exploit compression tricks (you can't much more complex)
Assignment Project Exam Help
- The final posting list is not complete until the end.
- At 12 bytes per posting (term, doc, freq), demands a lot of space for large collections.
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- $T = 100,000,000$ 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.

Use the same algorithm for 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 many steps
- We need an external sorting algo

Bottleneck

- 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 fewer disk seeks)

- 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

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

Add WeChat edu_assist_pro

Sorting 10 blocks of 10M 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 read each block from disk and sort it.* Add WeChat `edu_assist_pro`
- 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

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

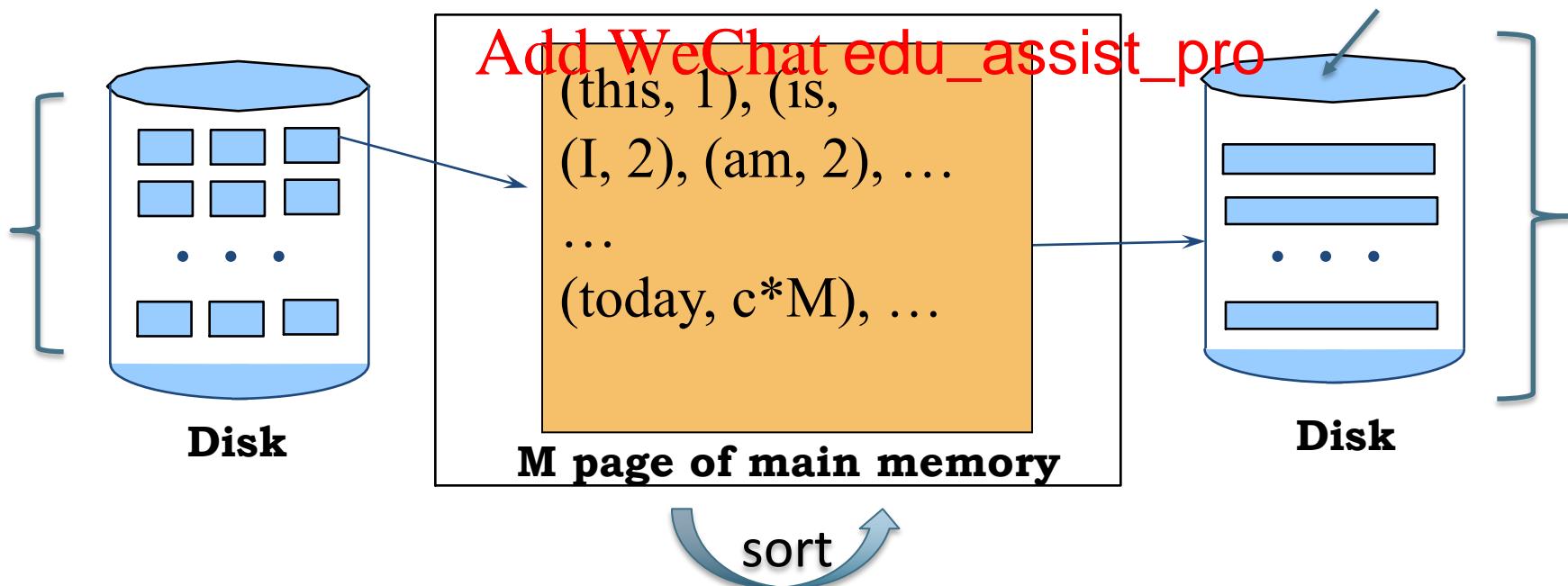
Add WeChat edu_assist_pro

Example

- Settings
 - **B:** Block/page size
 - **M:** Size of m
 - **N:** Number of documents
 - **R:** Size of the term-document vector emitted by the document.
- 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

External Merge-Sort: Phase I

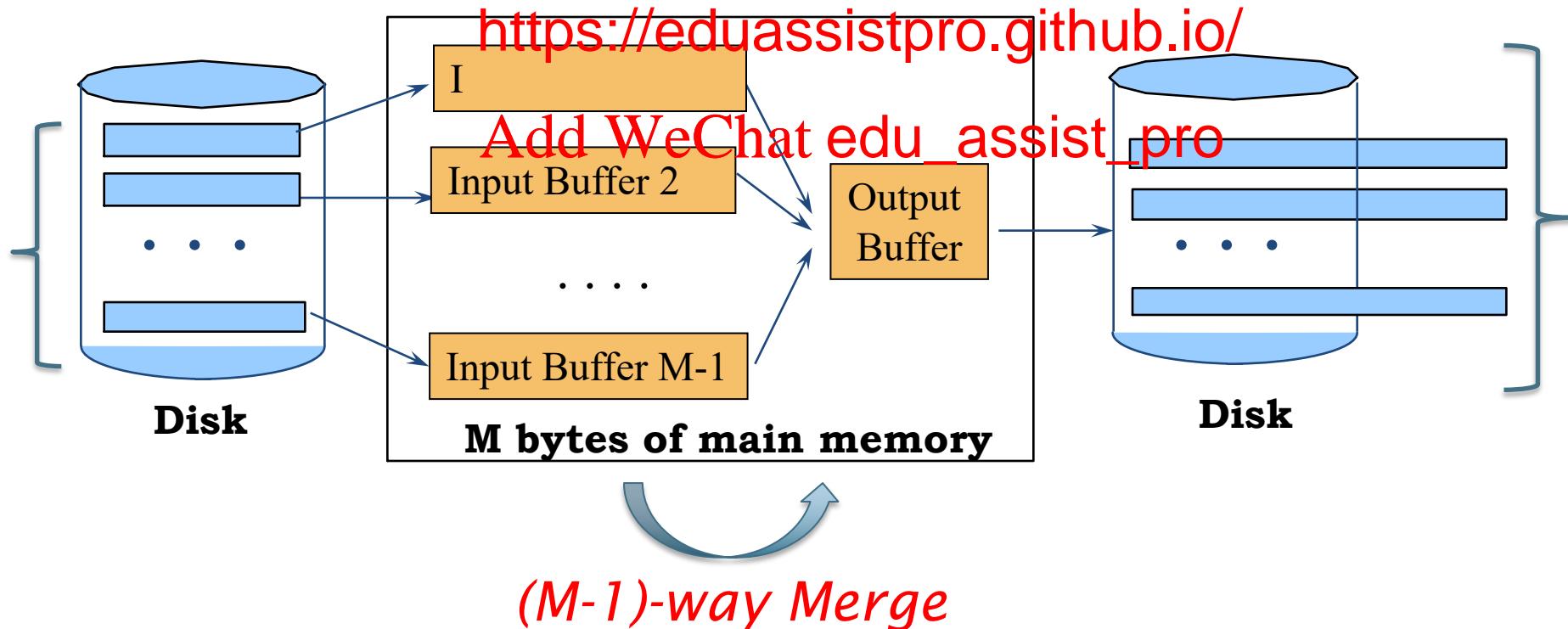
- Phase I: load the (term, docID) pairs from $(M*B)/R$ documents (*at a time*) into M buffer pages; sort
 - Result: ~~(Assignment Project Exam Help)~~ runs of length M pages
 - # of runs = 2 <https://eduassistpro.github.io/>



Phase II /1

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

Assignment Project Exam Help



K-way Merge: Using a **min-heap**



Assignment Project Exam Help



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



Add WeChat edu_assist_pro



K-way Merge: Using a **min-heap**

Heap.pop_min()



Assignment Project Exam Help



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



Add WeChat edu_assist_pro



K-way Merge: Using a **min-heap**



Assignment Project Exam Help



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



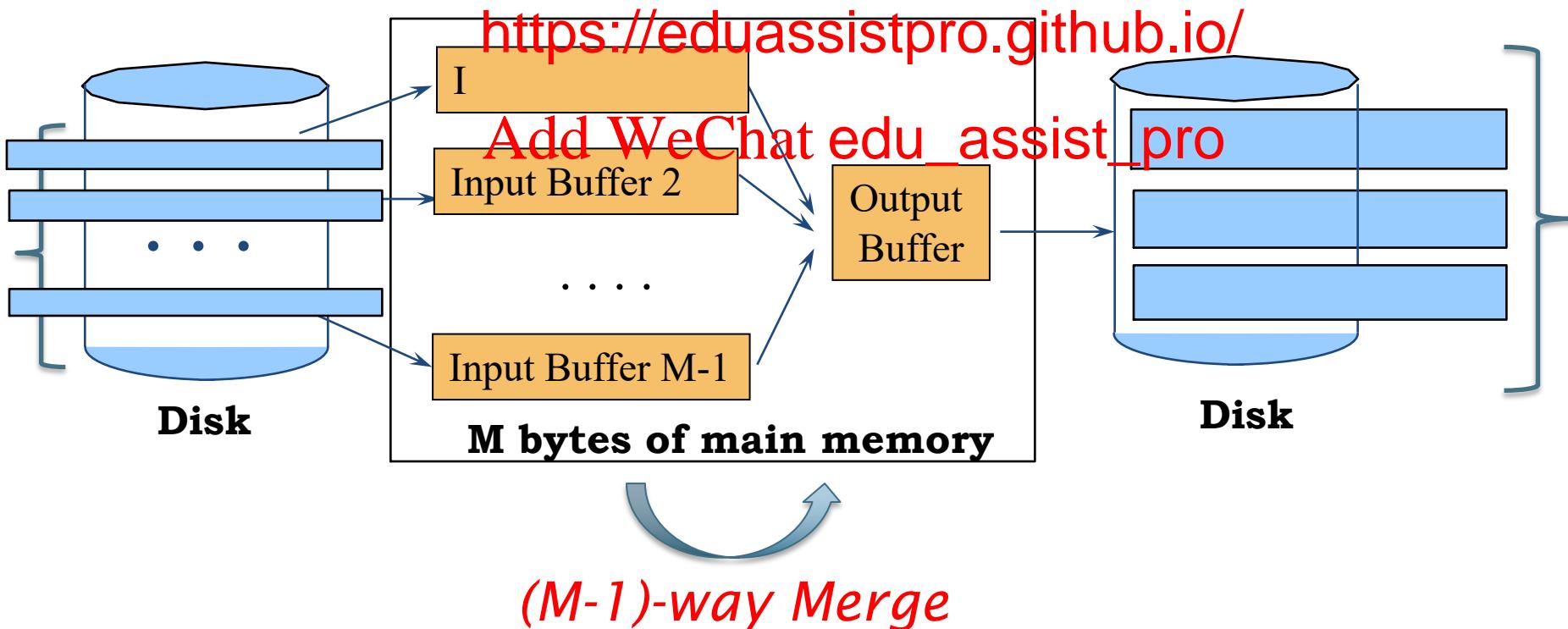
Add WeChat edu_assist_pro



99

Phase II /2

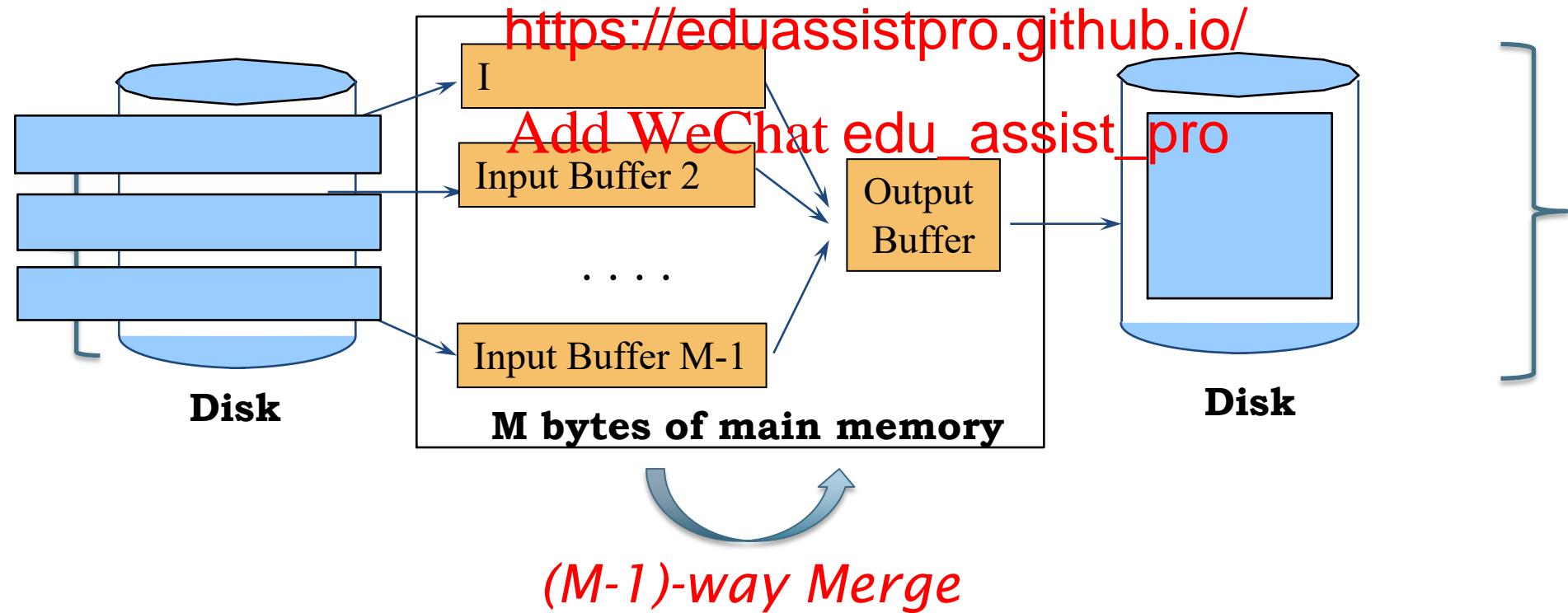
- **Recursively** merge (up to) $M - 1$ runs into a new run
- Result: runs of length M ($M - 1$)² pages
Assignment Project Exam Help



Phase II /3

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

Assignment Project Exam Help



Cost of External Merge Sort

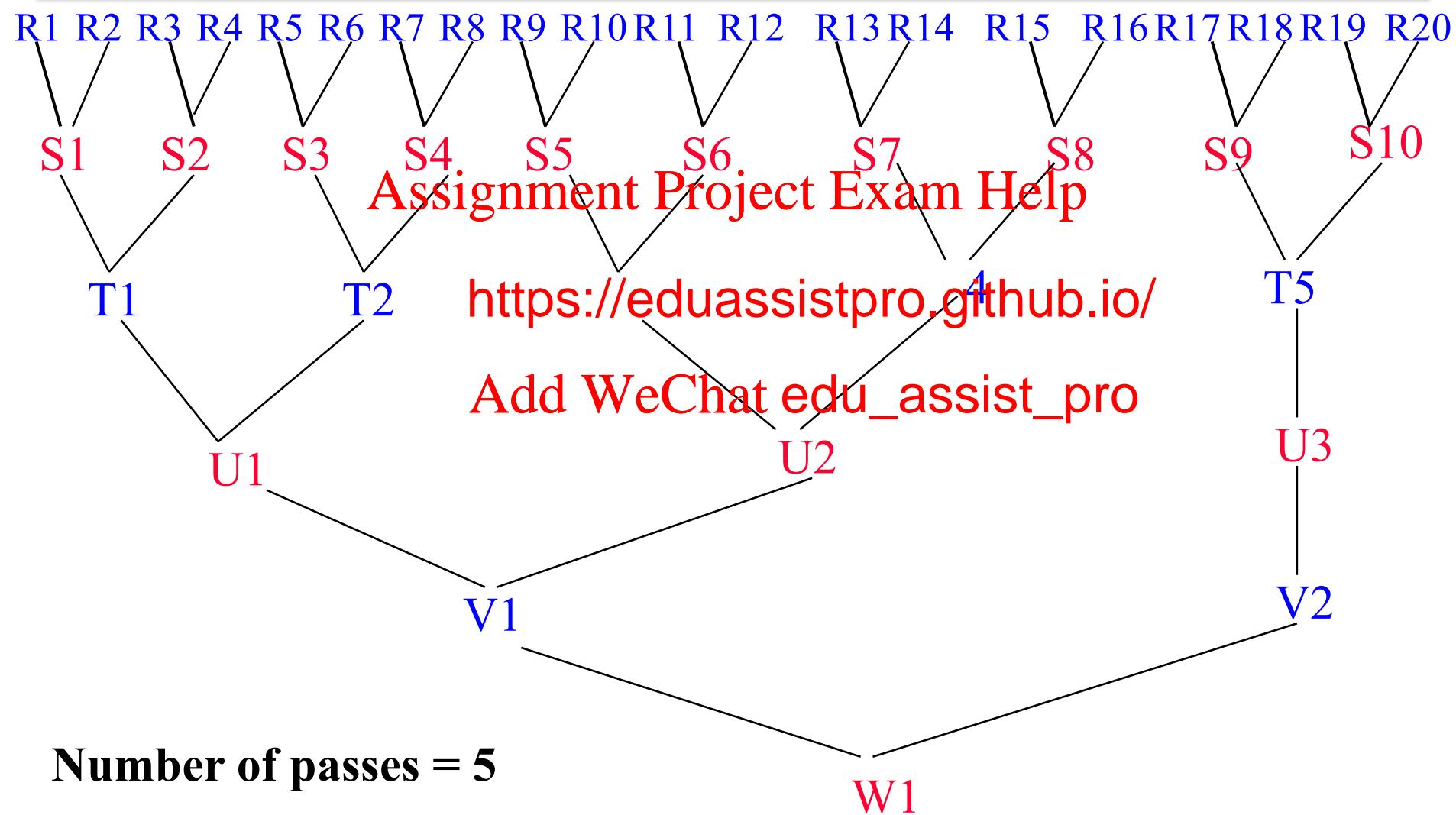
- Number of passes: $1 + \left\lceil \log_{M-1} \left\lceil \frac{NR}{MB} \right\rceil \right\rceil$
- Total I/O cost: $2 \cdot \binom{NR}{B} \cdot \left(1 + \left\lceil \log_{M-1} \left\lceil \frac{NR}{MB} \right\rceil \right\rceil \right)$ blocks/pages

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

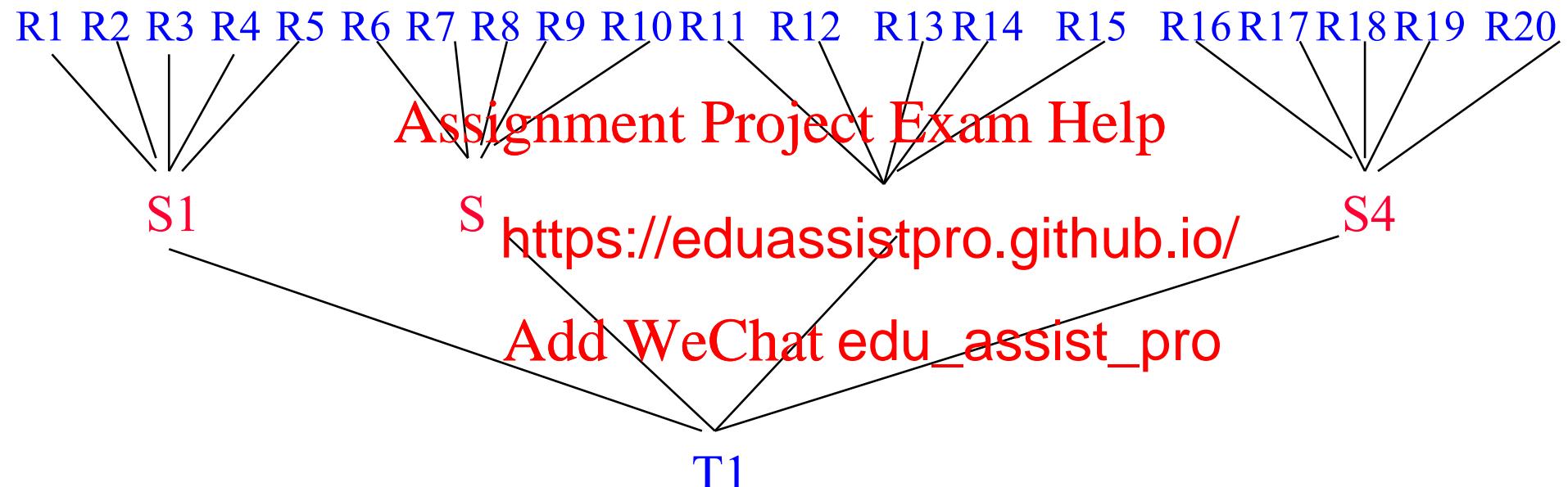
- How much data can we sort B RAM?
Add WeChat edu_assist_pro
 - Assume B = 4KB
 - 1 pass → 10MB “data”
 - 2 passes → $\approx 25GB$ “data” ($M-1 = 2559$)
 - 3 passes ?
- Can sort most reasonable inputs in 2 or 3 passes !

Another Example

Example: 2-Way Merge for 20 Runs



Example: 5-Way Merge for 20 Runs



Number of passes = 2

Remaining problem with sort-based algorithm

- 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 with [Add WeChat edu_assist_pro](#) 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:

Single-pass in-memory indexing

- 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 add WeChat `edu_assist_pro` a complete inverted index for each block.
- These separate indexes can then be merged into one big index.

SPIMI-Invert

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

- Merging of blocks is analogous to BSBI.

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

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?

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Dynamic indexing

- Up to now, we have assumed that collections are static.
- They rarely are:
 - Documents change <https://eduassistpro.github.io/> to be inserted.
 - Documents are deleted and [Add WeChat edu_assist_pro](#)
- This means that the dictionary strings lists have to be modified:
 - Postings updates for terms already in dictionary
 - New terms added to dictionary

Simplest approach – Immediate Merge

- 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
[Add WeChat edu_assist_pro](#)
- 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

Issues with main and auxiliary 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 **Add WeChat edu_assist_pro**
 - 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.)

Another Extreme – No Merge

- Whenever memory is full, write the sub-index to the disk
 - Never merge sub-indexes
- Pros:
 - High indexing performance

Assignment Project Exam Help
https://eduassistpro.github.io/
Add WeChat edu_assist_pro
- Cons:
 - Slow query performance
 - Require $\Omega(|C|/M)$ seeks to fetch the inverted list for a term

Compromise – Logarithmic 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 generation = g gives a new index with generation = [Add WeChat edu_assist_pro](#)
- 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

Illustration

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Algorithm

- 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 [Add WeChat edu_assist_pro](#)
- or merge with I_0 (if I_0 already $s Z_1$)
- 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

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

Add WeChat edu_assist_pro

Logarithmic merge

- Auxiliary and main index: index construction time is $O(C^2/M)$ as each posting is touched in each merge.
 - C = Collection size, and M = memory size
- Logarithmic <https://eduassistpro.github.io/> merged $O(\log C/M)$ time ($C \log (C/M)$)
 - Add WeChat edu_assist_pro
- So logarithmic merge is 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

Further issues with multiple indexes

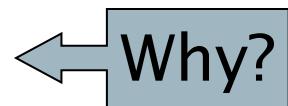
- Collection-wide statistics are hard to maintain
- E.g., when we spoke of spell-correction: which of several corrected alternatives do we present to the user?
<https://eduassistpro.github.io/>
 - We said, pick the one with the most frequent word
- How do we maintain the total statistics across multiple 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

Dynamic indexing at search 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

Other sorts of indexes

- 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 dictionary terms containing it – the “postings”
Add WeChat `edu_assist_pro`
 - 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



Resources for today's lecture

- Chapter 4 of IIR
- MG Chapter 5
- Original publication by Dean and Ghemawat (2003) <https://eduassistpro.github.io/>
- Original publication on SPIE [Add WeChat edu_assist_pro](#)