



Search

Mini Project

Building a Search Engine  
Phase - II

# Project Task

- Data: Wikipedia English Dump ~ 9.91 GB compressed
  - 🟡 Data link - <https://moodle.iiit.ac.in/mod/forum/discuss.php?d=22209>
- Search time < 1 sec (200-500ms) (C++/Java) ; ( < 5 sec; Python)
- Index size ~ Should be less than 1/4 of data size
- Support for field queries
- External tools and libraries like Lucene, WikiXMLj, elasticsearch, redis, etc not allowed.

# System capabilities

- Given a query / Field query output top 10 results(title of wiki document) sorted by relevance of document with respect to give query.
- Relevant results should come within expected time limit.
  - 🟡 < 1s for C++/Java
  - 🟡 < 5s for Python

# Phase II

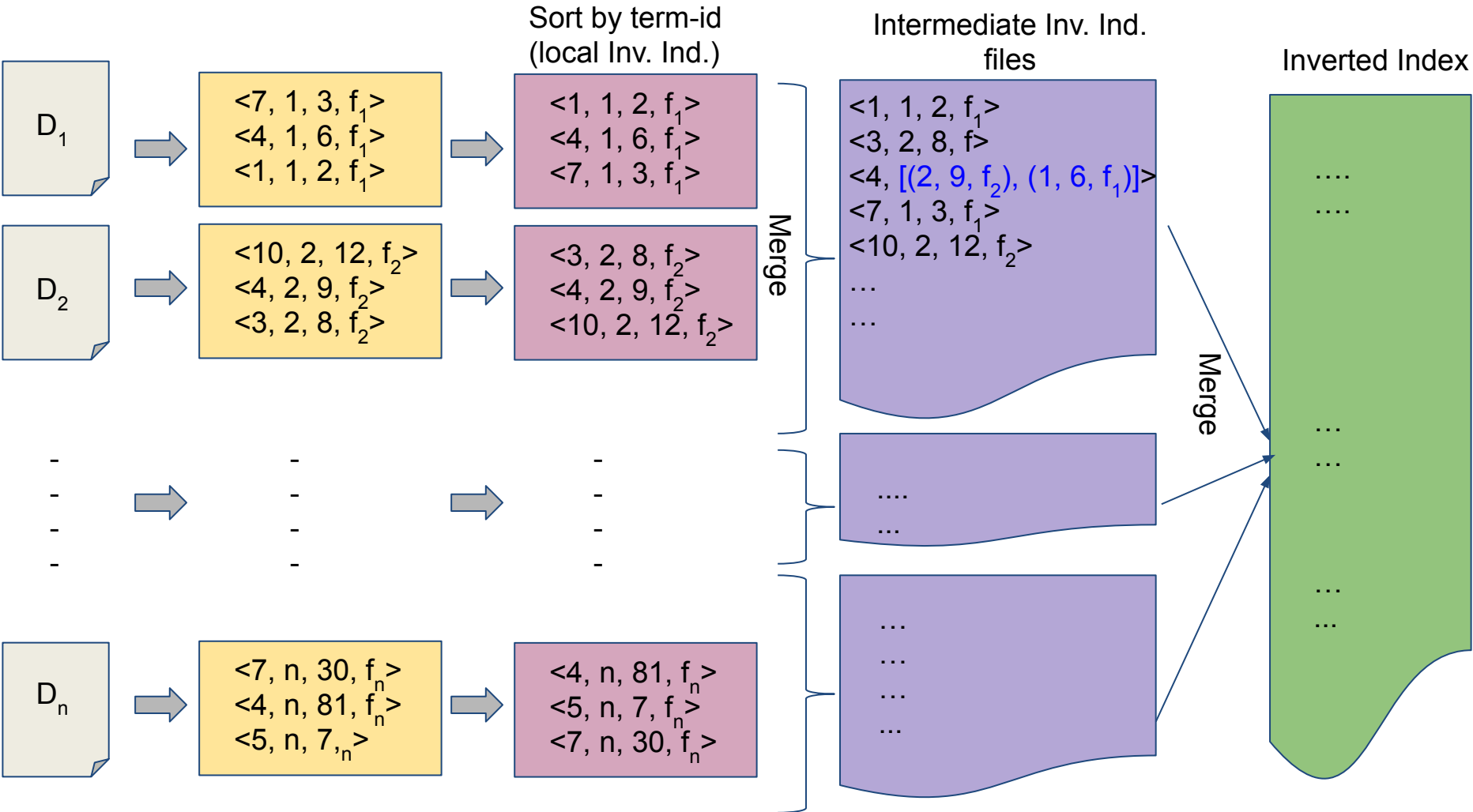
- Inverted index creation on whole wiki dump (~ 46 GB)
- Implement Ranking mechanism
- End to End search system

# Scalable Inverted Index Creation

- Main challenge is to build a huge index with limited memory.
- Sort-based Method
  - Build local inverted index
  - Merge local inverted index
  - Obtain large inverted index

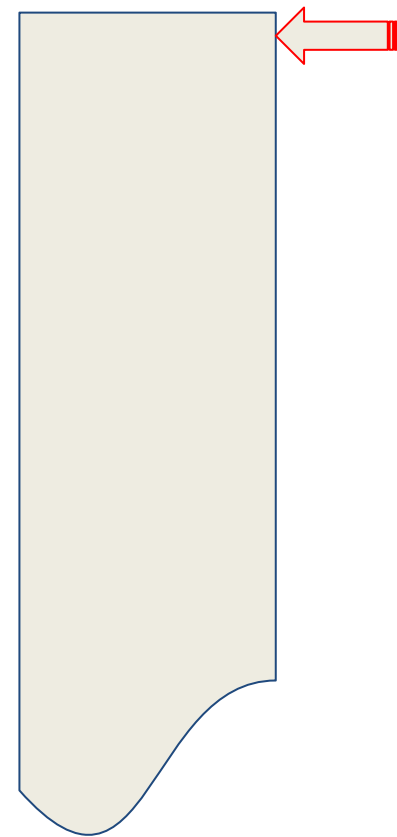
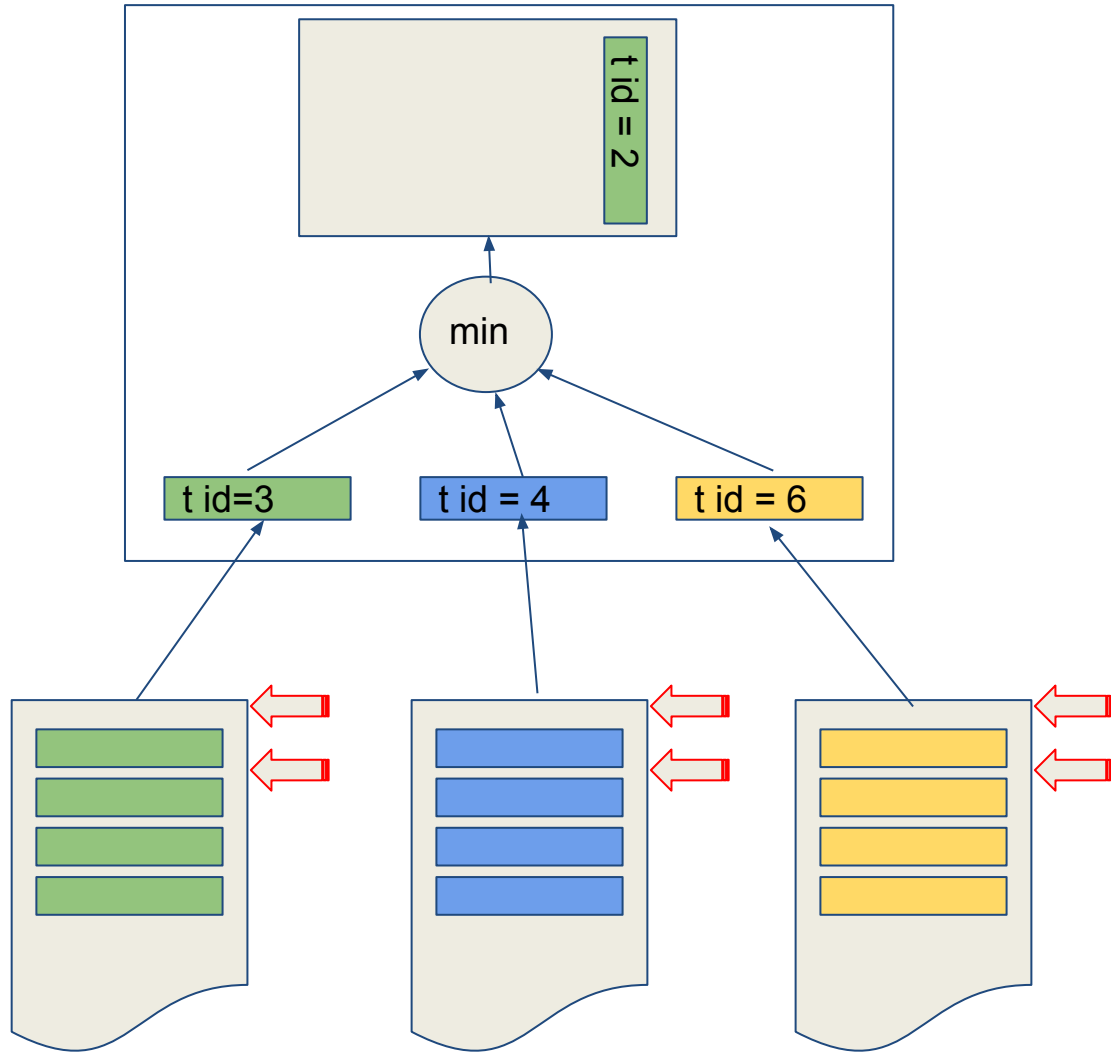
# What information is needed about a term?

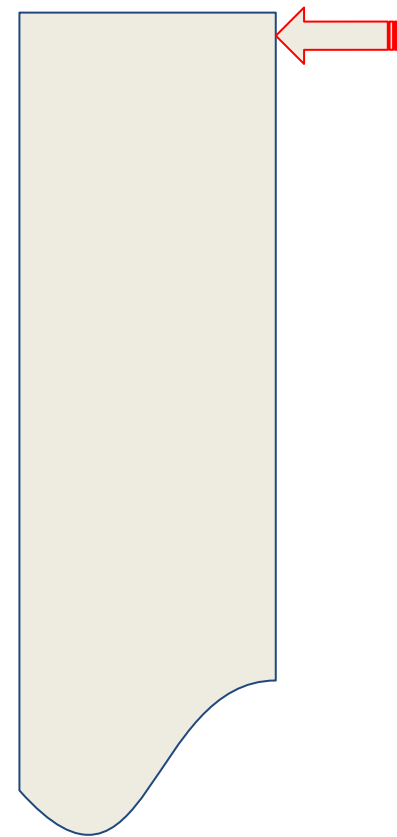
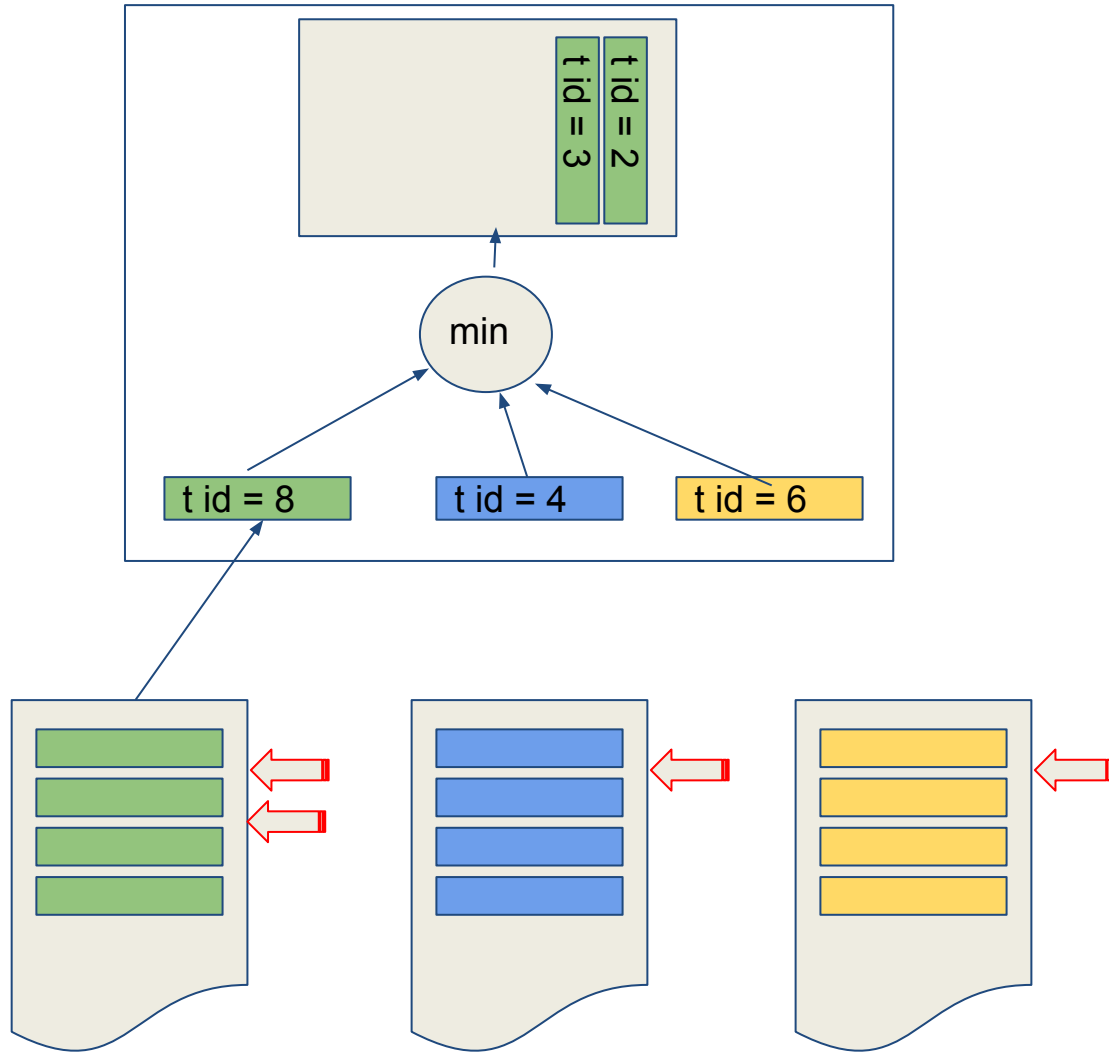
- <term id/word, doc id, term freq, field info>  
Eg: <2/ant, 12, 354, t:1b:4:c2>

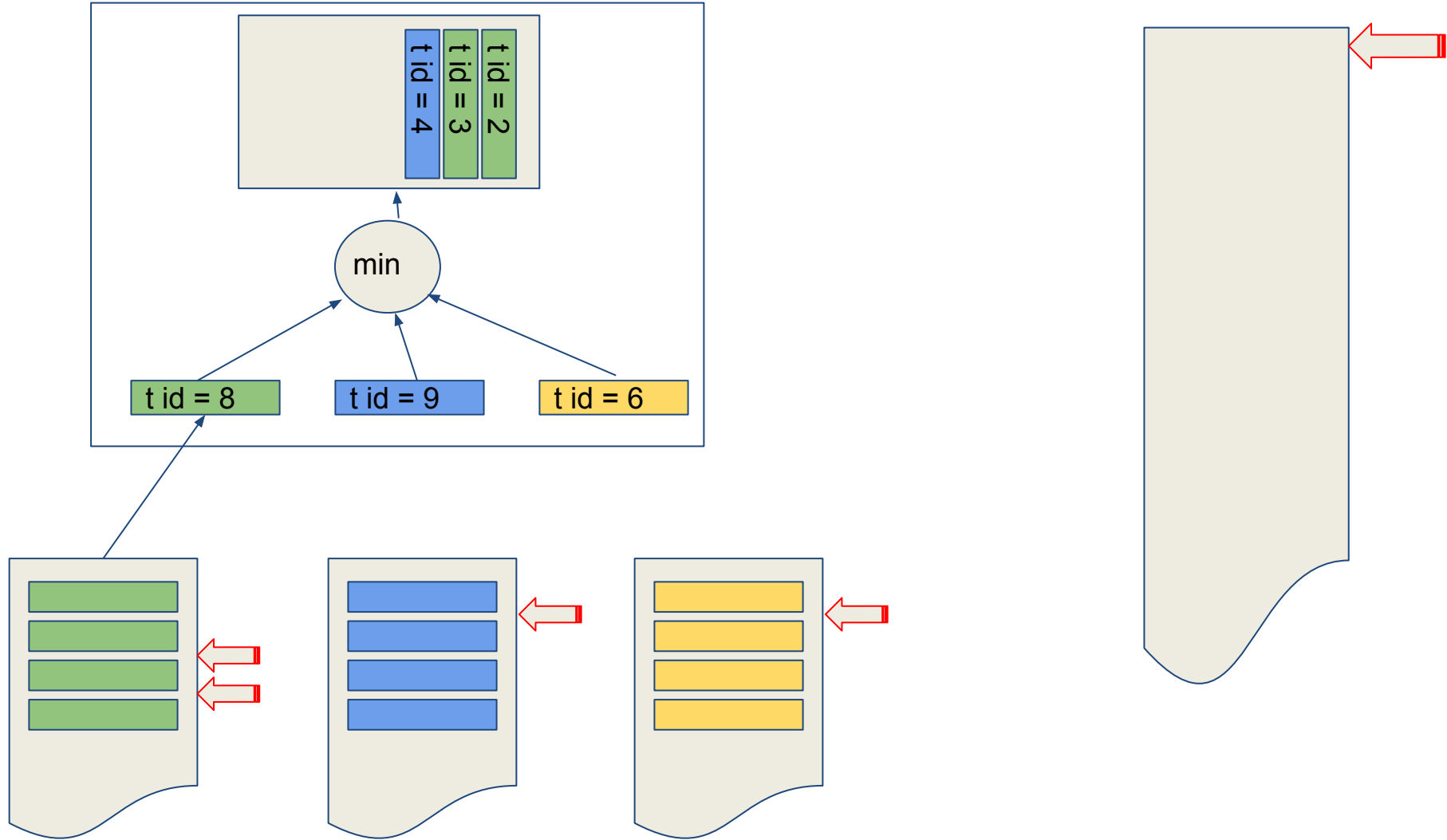


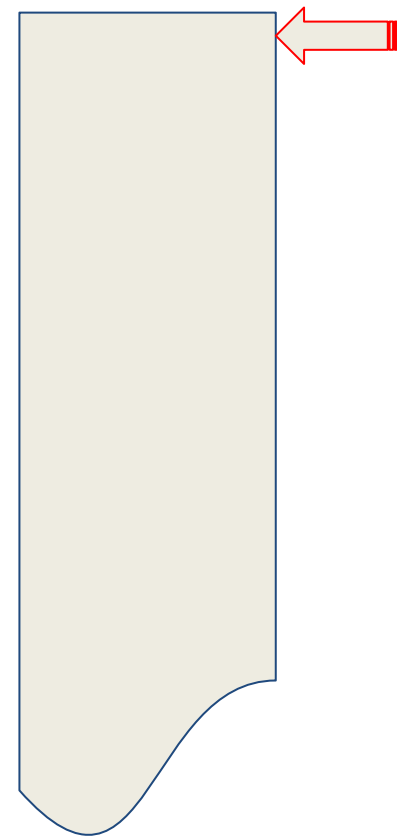
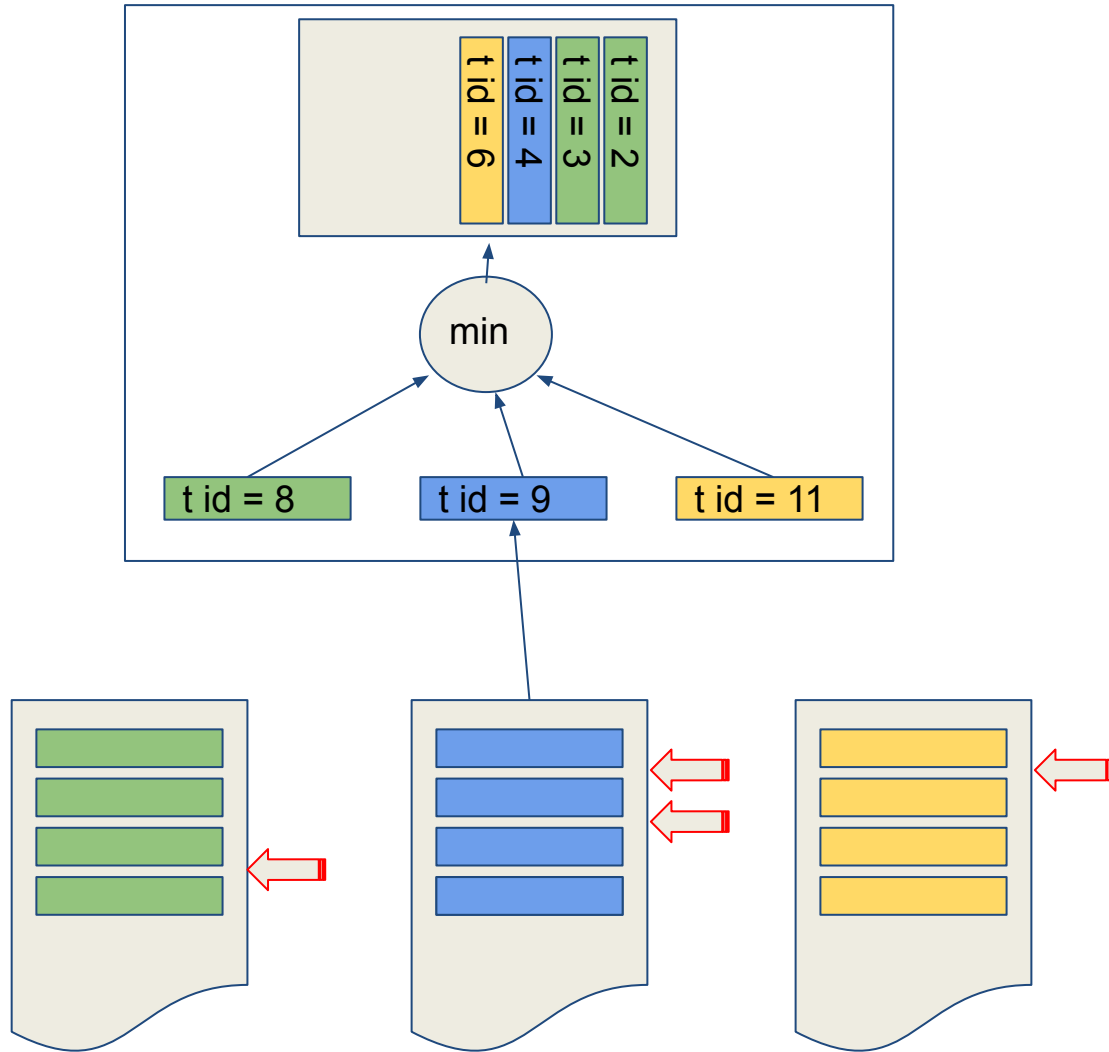
Merging sorted files

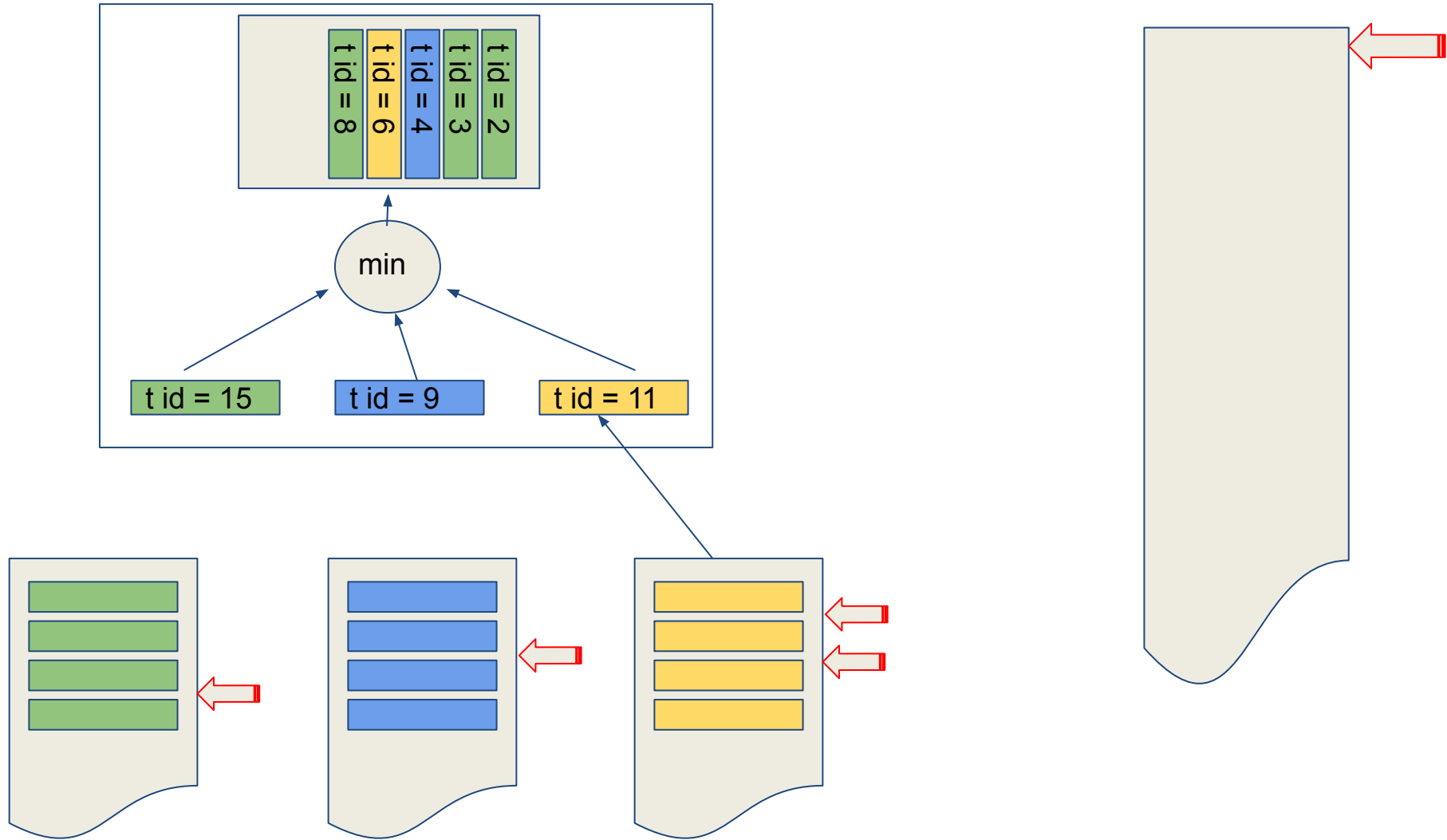


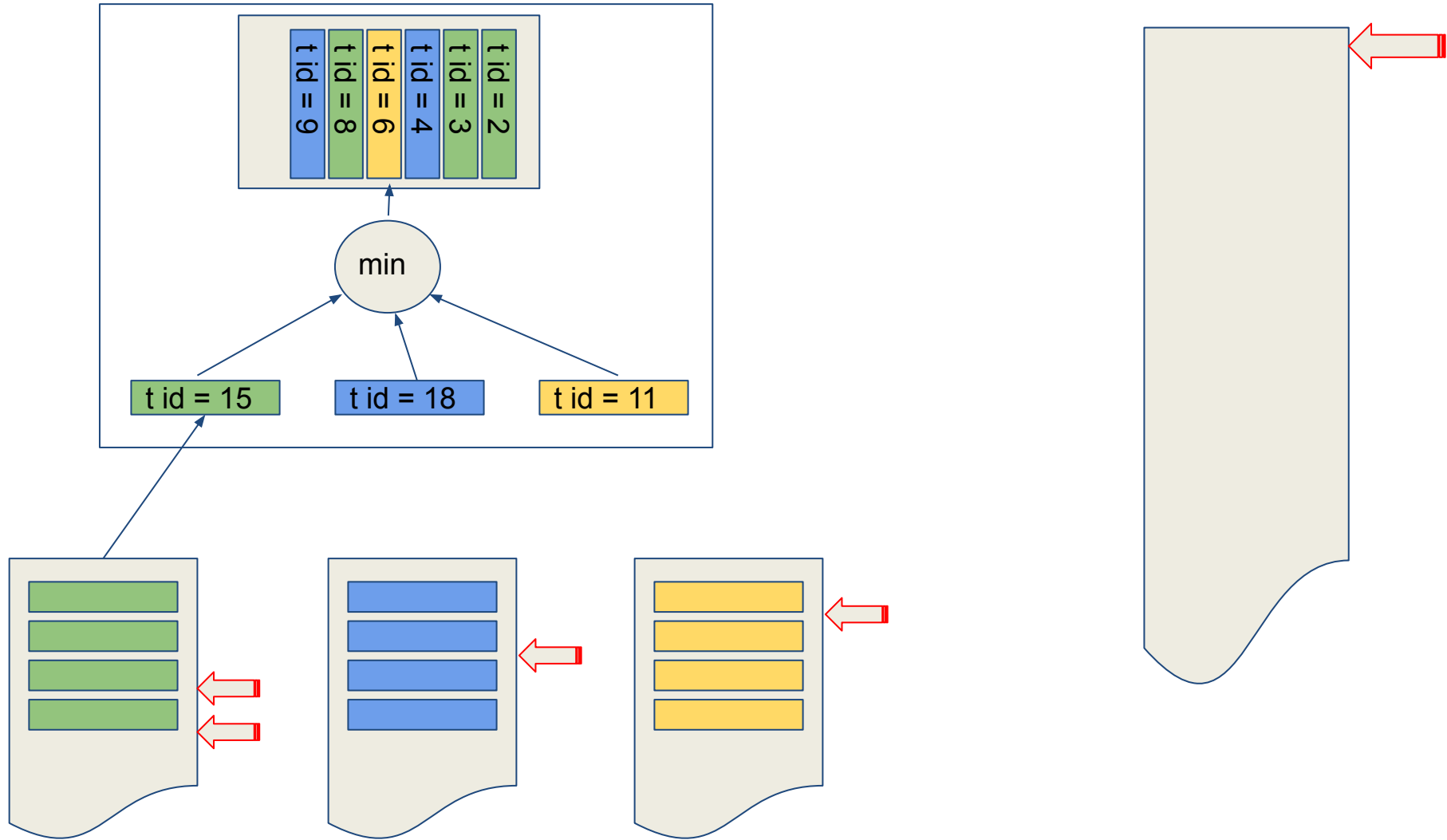


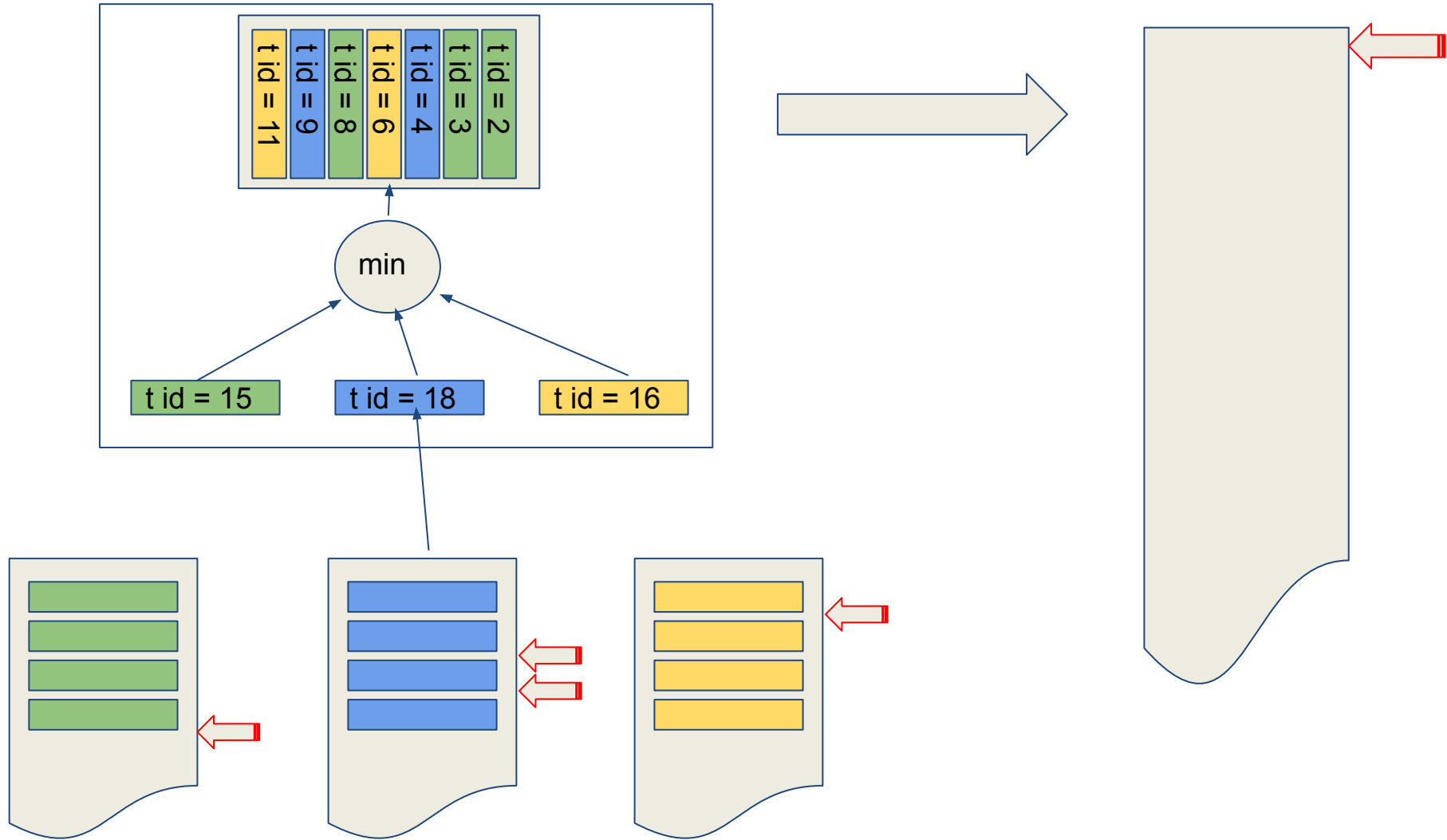


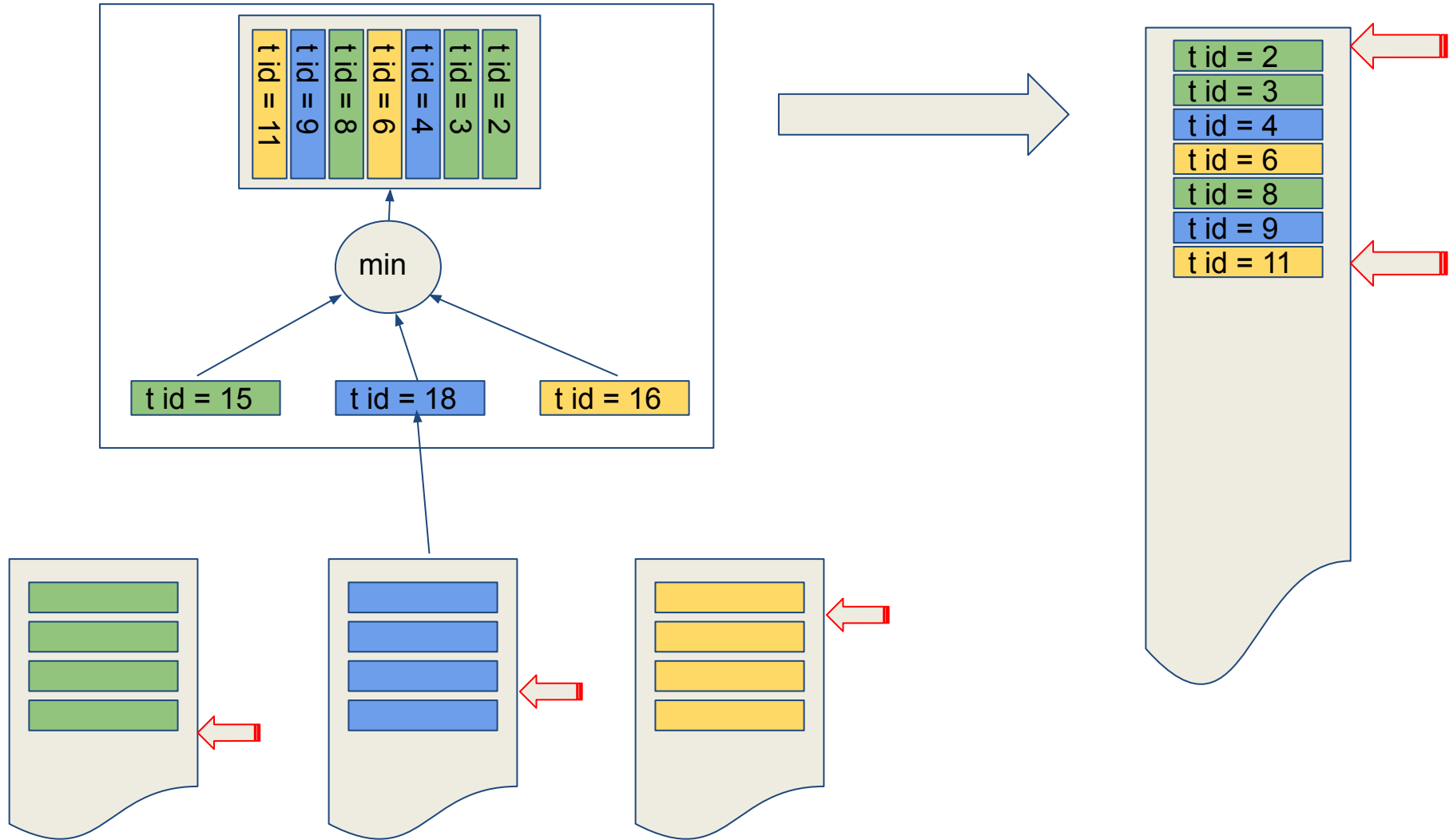








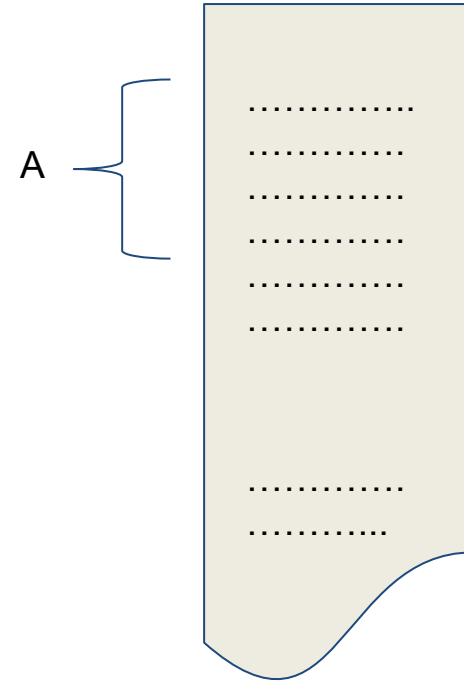






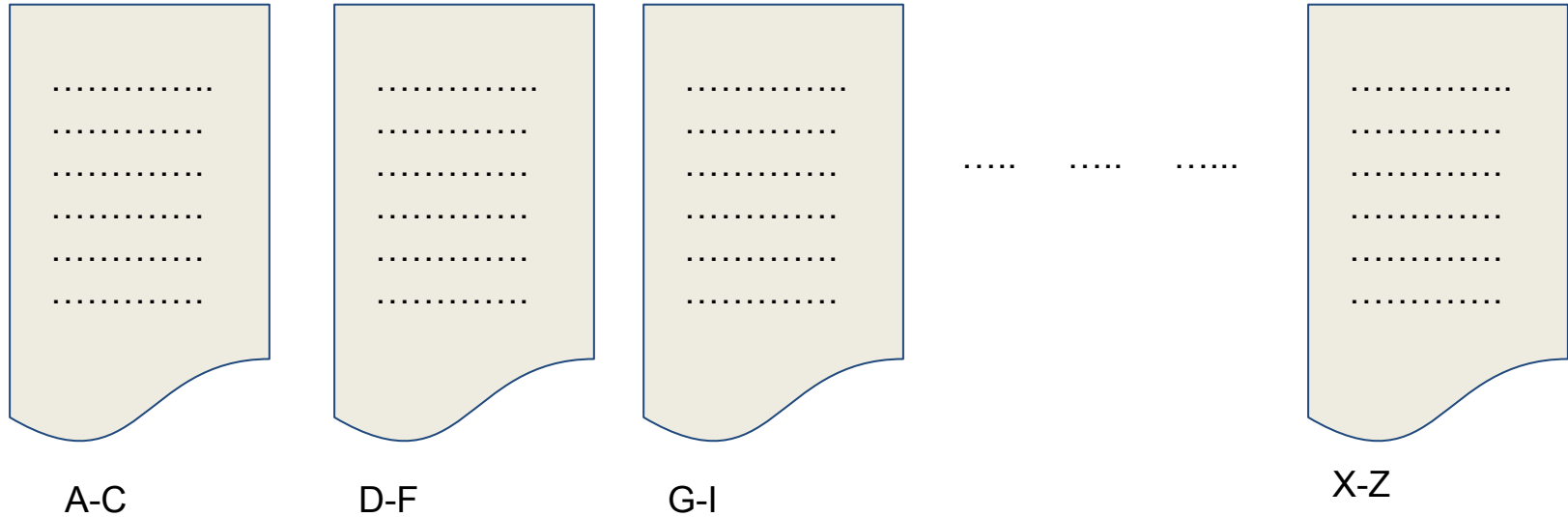
# Levels of Indexing...

First letter of Term	Range (line no.)
A	1-1000
B	1001-3050
C	3051 - 4800
....	...
Z	10000-10080



Inv. Index file

# Levels of Indexing..



Mapping

A-C	D-F	G-I	X-Z
0.txt	1.txt	2.txt	8.txt

Thank You!