

Constructing Inverted Index

- The main difficulty is to build a huge index with limited memory
- Memory-based methods: not usable for large collections
- Sort-based methods:
 - Step 1: Collect local (termID, docID, freq) tuples
 - Step 2: Sort local tuples (to make “runs”)
 - Step 3: Pair-wise merge runs
 - Step 4: Output inverted file