

Lecture 21

Multikey Files

Last Day: Query Processing with Inverted Lists

- Query Language
- Parse Trees
- Algorithms

Today: Threaded Files

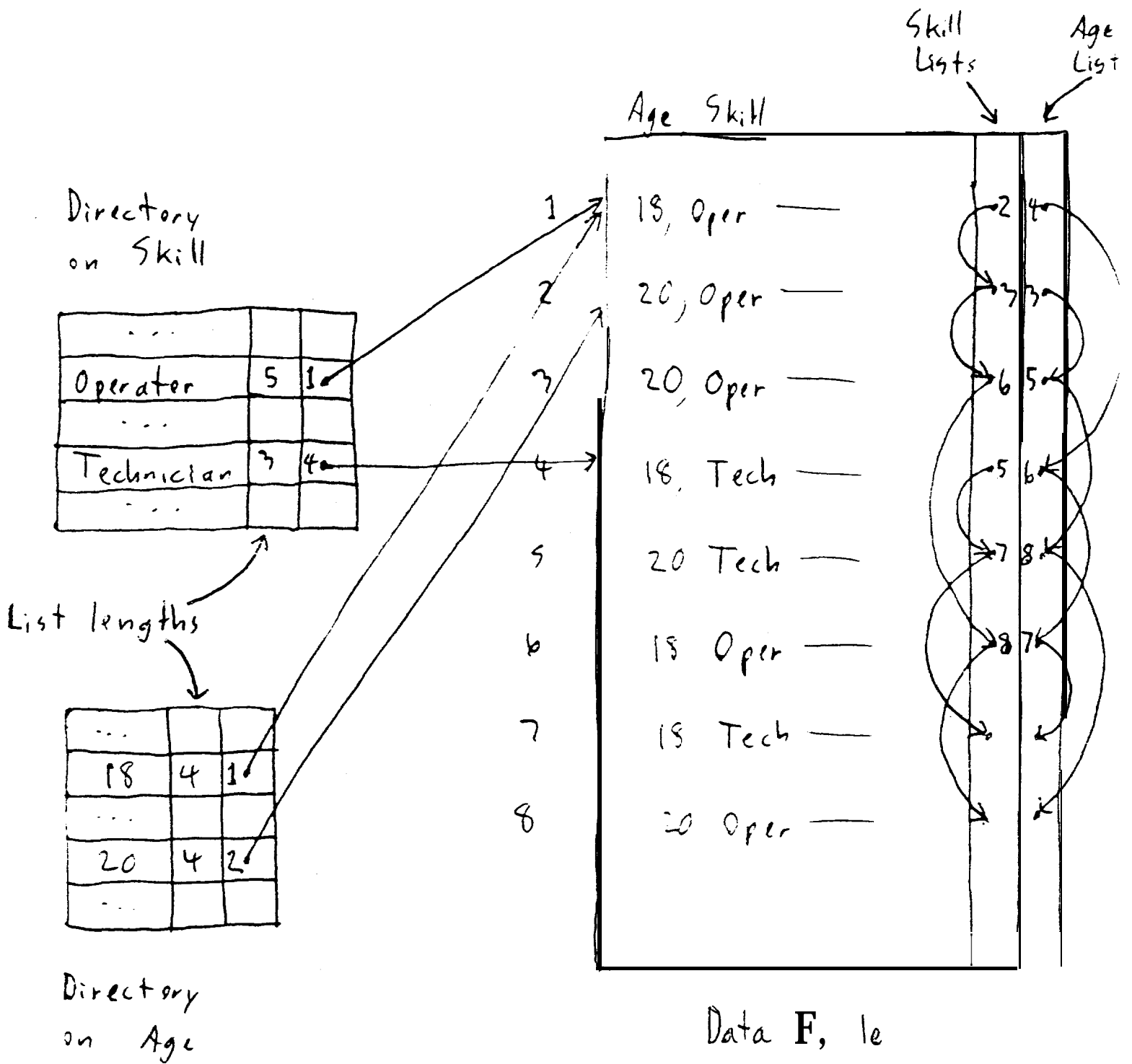
- Answering queries
 - Parse trees
 - Algorithms
-
-

(2) Threaded Files — which we shall
see next

Threaded Files

- Unlike inverted lists, the secondary indices for a threaded file do not have accession lists.
 - Instead, each data record has a number of extra pointer fields, one extra field for each secondary index.
 - Data records with the same value for a given secondary key are linked together into a list (a "thread") using one of the extra fields.
-
-

Example of a Threaded File



Data F, le

- This data record is the head of a list of a same data records with value v for field F

- The linked lists running through the data file are called "threads" (one thread for each secondary key field, F , and each value, v).

Query Processing with Threaded Files

Example 1: Retrieve all operators

- Find operator entry in skill directory.
- Follow the pointer to record 1 in the data file
- Print record 1, & follow its ^{skill} ptr to record 2.
- Print record 2, & follow its skill ptr to record 3.
- Print record 3, & follow its skill ptr to record 6.
- Print record 6, & follow its skill ptr to record 8.
- Print record 8, & terminate, since its skill ptr is empty.

Records 1, 2, 3, 6, 8 have been retrieved & printed.

Example 2

Retrieve all 18-year olds.

- Find the entry for 18 in the Age directory
 - Follow its pointer to record 1 in the data file.
 - Follow the Age pointers to retrieve records 1, 4, 6, 7 (the age 'thread' for 18-year olds)
 - Print each record as it is retrieved.
-
-
-

Example 3

Retrieve all 18-year old operators.

- Find the entry for 18 in the age directory.
 - Note that the length of its thread is 4.
 - Find the entry for operator in the skill director
 - Note that the length of its thread is 5
 - Follow the Age thread, since it is shorter
 - i.e., Retrieve records 1, 4, 6, 7.
 - Print those records with Skill = Operator.
 - i.e., Print records 1, 6.
-
-

Example 4

Retrieve all 18-year old technicians.

- Find the entry for 18 in the age direct
 - Note that the length of its thread is 4.
 - Find the entry for technician in Skill direct
 - Note that the length of its thread is 3
 - Follow the technician thread since it is shorter.
 - i.e., Retrieve records 4, 5, 7.
 - Print those records with Age = 18.
 - i.e., Print records 4, 7.
-
-
-
-

Example 5

Retrieve all records for 18-year olds
OR technicians

- Find the entry for 18 in the age dir
- Follow its thread, & retrieve each record
- Find the entry for technician in the skill dir
- Follow its thread, & retrieve each record
- Print each record in the union of the retrieved record sets.

- Note: Number of data-file accesses = $4 + 3$

ie, length of thread for 18 plus
 length of thread for technician

Query Processing For Threaded Files

In general, the approach is similar to that used for inverted lists.

- First, construct a parse tree for the query.
 - Assign to each node, X , of the parse tree a label, $SIZE(X)$.
 - $SIZE(X)$ represents the total length of all threads needed to answer the subquery at node X .
 - i.e., $SIZE(X)$ represents the number of data file accesses.
 - $SIZE(X) = \infty$ represents the entire file
-

Main Questions

Given a query and a File with indices on some fields

(1) do the indices help to answer the query (or do we have to scan the entire File)?

(2) If the indices help, how should we use them?

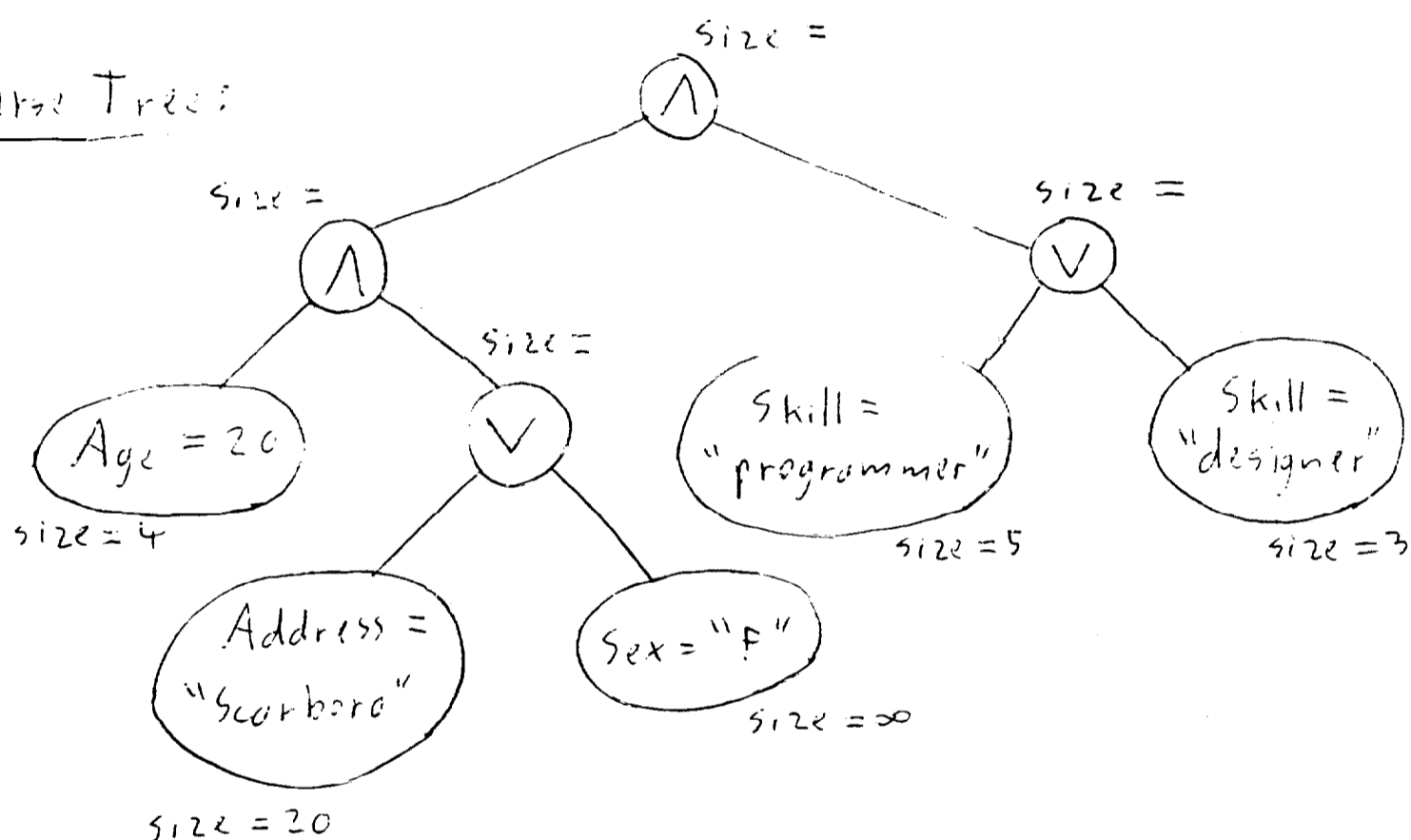
The indices help iff $SIZE(nof) \neq \infty$

Example: Computing SIZE(X)

Query:

(Age = 20 AND (Address = "Scarboro" OR sex = "F"
AND (Skill = "programmer" OR Skill = "designer"))

Parse Tree:



- Given indexes on Age, Address & Skill,
do they help? Yes, in this case.

To evaluate this query,

- Follow the thread for Age = 20.
- Keep those records with Address = "Scarboro"
or Sex = "F".
- Of these records, keep those with
Skill = "Programmer" or Skill = "designer".

Algorithm to Compute SIZE(X)

- (1) IF node X is a leaf of the form " $F=v$," then
- IF F is indexed, then $\text{SIZE}(X)$ is the length of the thread for value v + direct for F
 - IF F is not indexed, then $\text{SIZE}(X) = \infty$

(2) IF X has the form $(Y \text{ AND } Z)$, then
 $\text{SIZE}(X) = \min[\text{SIZE}(Y), \text{SIZE}(Z)]$

(3) IF X has the form $(Y \text{ OR } Z)$, then
 $\text{SIZE}(X) = \text{SIZE}(Y) + \text{SIZE}(Z)$

Here, $\infty + c = \infty$ and $\min(\infty, c) = c$.

Answering a Query

- Let X_R be the set of data records (not points) satisfying the subquery at node X .
(Eventually we want to compute X_R for $X = \text{root}$.)
 - If $\text{SIZE}(X) = \infty$, then the indices do not help and we must scan the entire file and check each record to see if it satisfies the query at X .
 - If $\text{SIZE}(X) \neq \infty$, then we can use the indices to answer the query more quickly.
-

Algorithm to compute XR (when $\text{SIZE}(X) \neq \infty$)

(1) IF X is a leaf, then XR is computed directly by following a thread from an index entry

(2) IF X has the form $(Y \text{ and } Z)$

and $\text{SIZE}(Y) \leq \text{SIZE}(Z)$

then compute YR

$XR := \emptyset$

For each record $r \in YR$,

if r satisfies the query at Z ,

then $XR := XR \cup \{r\}$

ie follow the shortest thread.

(3) IF X has the form $(Y \text{ or } Z)$

then compute YR, ZR

$$XR := YRUZR$$



R , Follow both threads