

Query Optimization

Query Tree

Consider the following schema and query:

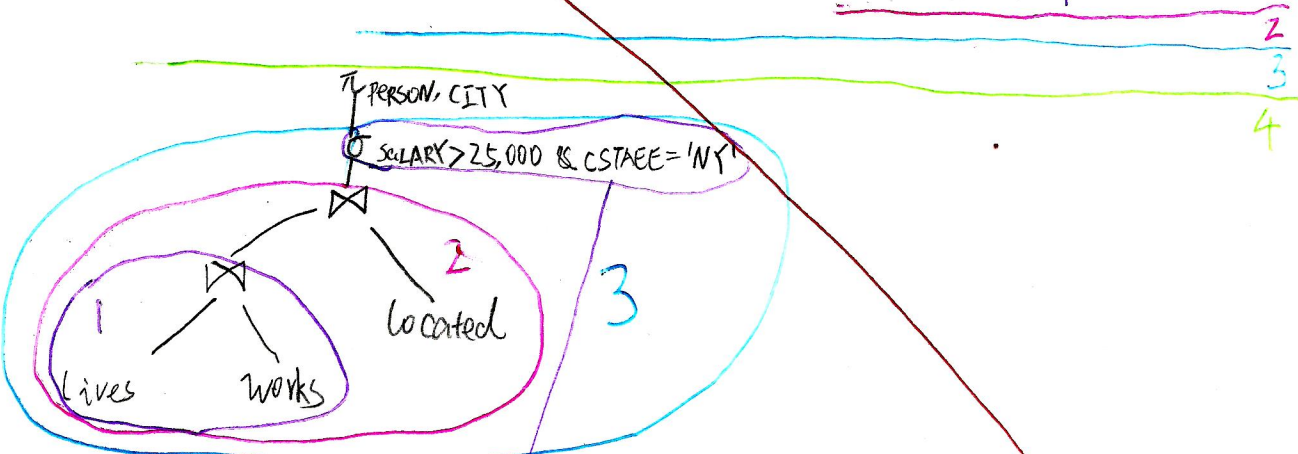
lives (PERSON, STREET, CITY, STATE, ZIP)

works (PERSON, COMPANY, SALARY, POSITION)

located (COMPANY, CCITY, CSTATE, CZIP)

Find the name and city of all people who earn more than 25,000 and work for a company located in the state of NY.

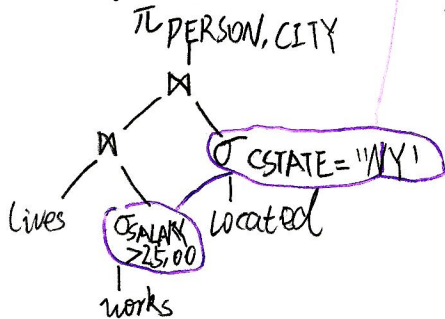
$\pi_{\text{PERSON, CITY}} \sigma_{\text{SALARY} > 25,000 \wedge \text{CSTATE} = \text{'NY'}} ((\text{lives} \bowtie \text{works}) \bowtie \text{located})$



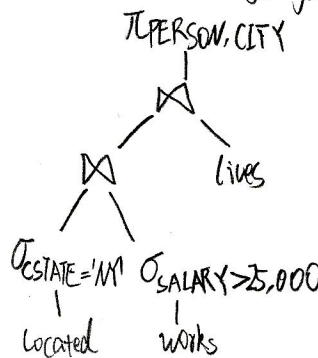
HEURISTIC ALGEBRAIC OPTIMIZATION ALGORITHM

1. Cascade selection (rule 1)
2. Move selection as far down the tree as possible (rules 2, 4, 6, 10)
3. Rearrange leaf node to get smaller intermediate relations (rule 9)
4. Combine a \bowtie with σ to yield \bowtie , if possible
5. Cascade projections and push down the tree (rules 3, 4, 7, 11)
6. Identify common subexpressions.

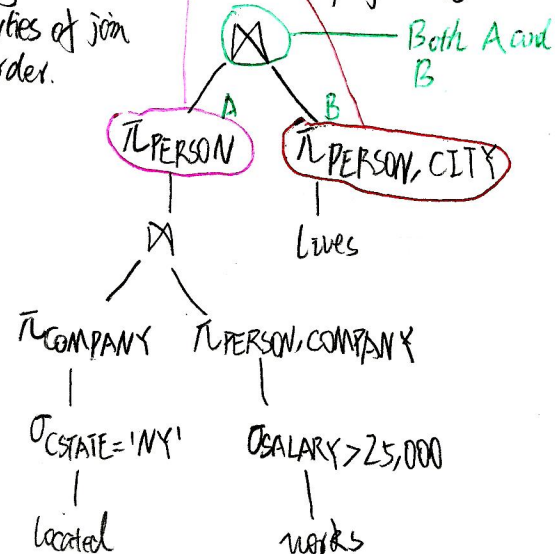
T₁: Cascade selections and push down tree



T₂: Using commutativity & associativity properties of join rearrange join order.



T₃: Introduce projections



Query Optimization

Consider the following relational schema and relational algebra query.

customer (cName, street, cCity)

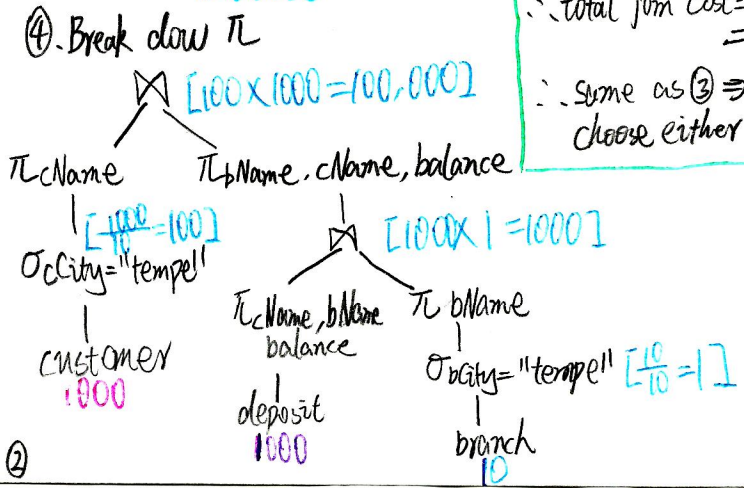
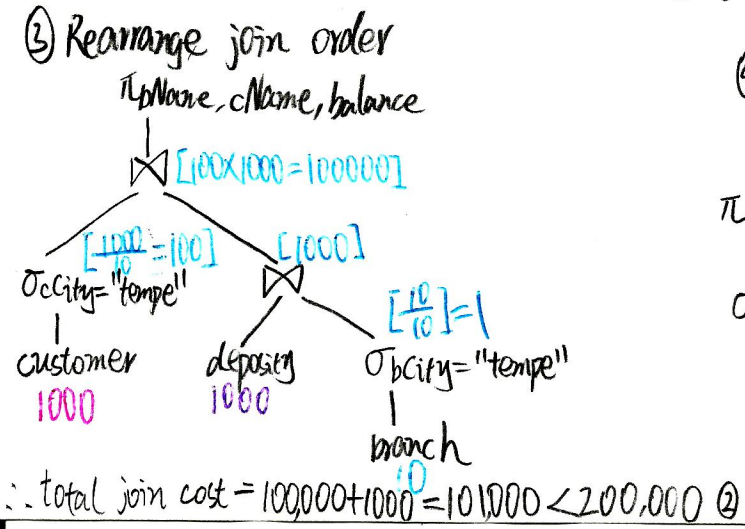
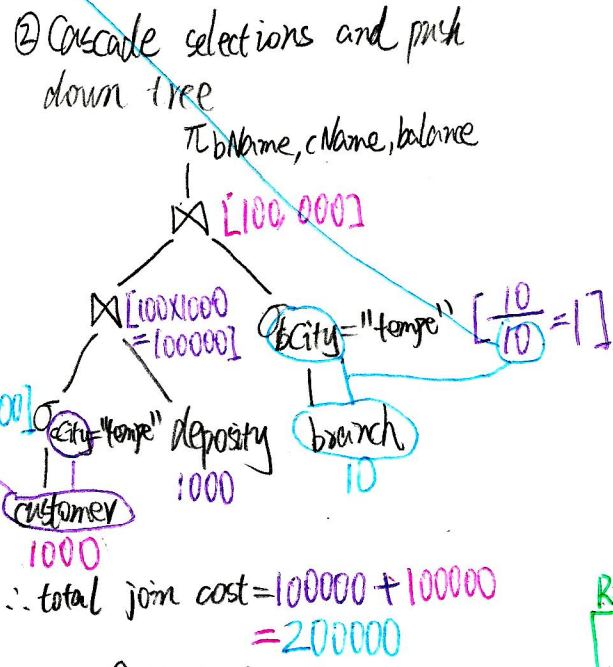
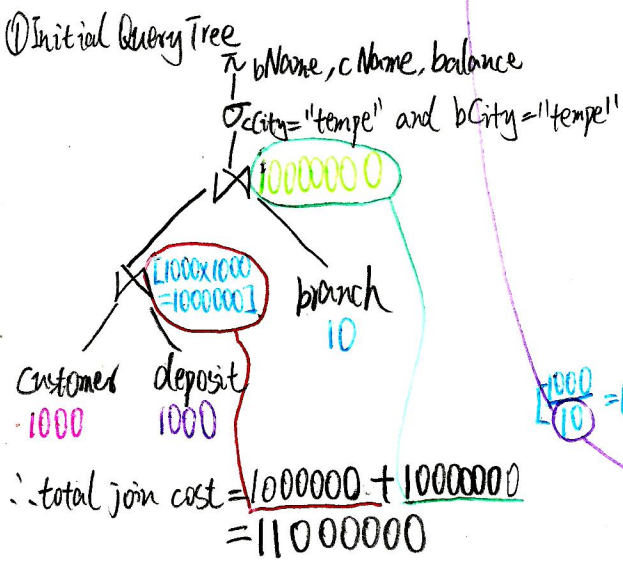
deposit (cName, bName, accountNumber, balance)

branch (bName, assets, bCity)

Identify a relational algebra tree that reflects the order of operations a decent query optimizer would choose. Use the following statistics to estimate the sizes of intermediate tables and total join costs.

- Assume, |customer| = 1000,
- |deposit| = 10000,
- |branch| = 10,
- $V(cCity, customer) = 10$,
- $V(bCity, branch) = 10$,
- Max (balance) = 100K

Query: $\pi_{bName, cName, balance} (\sigma_{cCity="tempe" \text{ and } bCity="tempe"} ((customer \bowtie deposit) \bowtie branch))$



Result

\therefore total join cost = $100000 + 1000 = 101000$

\therefore same as ③ \Rightarrow We could choose either ③ or ④.