

Relational algebra

1.

```
-- Yunyu Guo
-- Part 1: For a neighborhood entered by the user, show the popularity of tree species planted in a given year,
-- including the minimum and maximum tree species, associated planting zone factors,
-- filtering the neighborhoods with at least 2 species planted in the given year.
SET @neighborhoodName = ?;
SET @startYear = ?;
SET @endYear = ?;
SET @plantingZoneFactor = ?;

SELECT
  tr.neighborhood,
  COUNT(DISTINCT ts.treeID) AS species_count,
  MIN(ts.commonName) AS min_species,
  MAX(ts.commonName) AS max_species,
  MIN(tp.plantDate) AS firstPlantDate,
  MAX(tp.plantDate) AS lastPlantDate,
  GROUP_CONCAT(DISTINCT tpz.plantingZoneFactor) AS planting_zone_factors -- concatenate values from a
group into a single string, separated by commas
FROM treeRequests tr
JOIN treePlantings tp ON tr.requestID = tp.requestID
JOIN treeSpecies ts ON tp.treePlanted = ts.treeID
LEFT JOIN treeToPlantingZones tpz ON ts.treeID = tpz.treeID
-- WHERE tr.neighborhood = 'Bushrod'
WHERE tr.neighborhood = @neighborhoodName
  AND tr.neighborhood IN (
    SELECT DISTINCT tr2.neighborhood
    FROM treeRequests tr2
    JOIN treePlantings tp2 ON tr2.requestID = tp2.requestID
    -- WHERE YEAR(tp2.plantDate) BETWEEN 2023 AND 2025
    WHERE YEAR(tp2.plantDate) BETWEEN @startYear AND @endYear
    GROUP BY tr2.neighborhood
    HAVING COUNT(DISTINCT tp2.treePlanted) >= 2
  )
-- AND YEAR(tp.plantDate) BETWEEN 2023 AND 2025
  AND YEAR(tp.plantDate) BETWEEN @startYear AND @endYear
GROUP BY tr.neighborhood

UNION

-- Part 2: filter neighborhoods with a specific planting zone factor entered by the user
SELECT
  tr.neighborhood,
  COUNT(DISTINCT ts.treeID) AS species_count,
  MIN(ts.commonName) AS min_species,
  MAX(ts.commonName) AS max_species,
  MIN(tp.plantDate) AS firstPlantingDate,
```

```

MAX(tp.plantDate) AS lastPlantingDate,
GROUP_CONCAT(DISTINCT tpz.plantingZoneFactor) AS planting_zone_factors
FROM treeRequests tr
JOIN treePlantings tp ON tr.requestID = tp.requestID
JOIN treeSpecies ts ON tp.treePlanted = ts.treeID
LEFT JOIN treeToPlantingZones tpz ON ts.treeID = tpz.treeID
-- WHERE tr.neighborhood = 'Bushrod'
WHERE tr.neighborhood = @neighborhoodName
-- AND tpz.plantingZoneFactor = 'Highly urbanized zones'
AND tpz.plantingZoneFactor = @plantingZoneFactor
-- AND YEAR(tp.plantDate) BETWEEN 2023 AND 2025
AND YEAR(tp.plantDate) BETWEEN @startYear AND @endYear
GROUP BY tr.neighborhood;

```

Subquery: (// Subquery **validneighborhood** for IN clause

$$\prod \text{TR2.neighborhood} ($$

$$\sigma_{\text{TP2.treePlanted} \geq 2}$$

$$\gamma_{\text{TR2.neighborhood}, \text{COUNT}(\text{DISTINCT treePlanted})}$$

$$\sigma_{\text{YEAR}(\text{TP2.plantDate}) \geq @startYear \wedge \text{YEAR}(\text{TP2.plantDate}) \leq @endYear}$$

$$\text{TR2} \bowtie \text{TR2.requestID} = \text{TP2.requestID TP2})$$

Main query part 1:

$$\prod \text{TR.neighborhood}, \text{COUNT}(\text{DISTINCT TS.treeID}) \text{ AS species_count}, \text{MIN}(\text{TS.commonName}) \text{ AS}$$

$$\text{min_species}, \text{MAX}(\text{TS.commonName}) \text{ AS max_species}, \text{MIN}(\text{TP.plantDate}) \text{ AS firstPlantDate},$$

$$\text{MAX}(\text{TP.plantDate}) \text{ AS lastPlantDate}, \text{GROUP_CONCAT}(\text{DISTINCT TPZ.plantingZoneFactor}) \text{ AS}$$

$$\text{planting_zone_factors}$$

$$\gamma_{\text{COUNT}(\text{DISTINCT TS.treeID}), \text{MIN}(\text{TS.commonName}), \text{MAX}(\text{TS.commonName}),}$$

$$\text{MIN}(\text{TP.plantDate}), \text{MAX}(\text{TP.plantDate}), \text{GROUP_CONCAT}(\text{DISTINCT TPZ.plantingZoneFactor})}$$

$$\sigma_{\text{TR.neighborhood} = @neighborhoodName \wedge \text{YEAR}(\text{TP.plantDate}) \geq @startYear \text{ AND}$$

$$\text{YEAR}(\text{TP.plantDate}) \leq @endYear}$$

$$(\text{TR} \bowtie \text{TR.requestID} = \text{TP.requestID TP} \bowtie \text{TP.treePlanted} = \text{TS.treeID TS}) \bowtie \text{TS.treeID} = \text{TPZ.treeID TPZ})$$

$$\bowtie \text{neighborhood} = \text{validneighborhood}$$

U

Main query part 2:

$$\prod \text{TR.neighborhood}, \text{COUNT}(\text{DISTINCT TS.treeID}) \text{ AS species_count}, \text{MIN}(\text{TS.commonName}) \text{ AS}$$

$$\text{min_species}, \text{MAX}(\text{TS.commonName}) \text{ AS max_species}, \text{MIN}(\text{TP.plantDate}) \text{ AS firstPlantingDate},$$

MAX(TP.plantDate) AS lastPlantingDate, GROUP_CONCAT(DISTINCT TPZ.plantingZoneFactor) AS
planting_zone_factors

γTR.neighborhood

σTR.neighborhood = @neighborhoodName ∧ TPZ.plantingZoneFactor = @plantingZoneFactor

∧ YEAR(TP.plantDate) >= @startYear ∧ YEAR(TP.plantDate) <= @endYear(

(TR ⋈_{TR.requestID=TP.requestID} TP ⋈_{TP.treePlanted=TS.treeID} TS) ⋈_{TS.treeID=TPZ.treeID} TPZ)

2.

```
// get all assigned tree requests and plantings infor for the volunteer
// suppose the volunteer should see the address, phone number, plant date, and tree common name
public void seeTreeRequestAndTreePlantingAndTreeSpecies(int vid) {
```

```
    String sql = "SELECT tr.address, tr.phone, tp.plantDate, ts.commonName, vp.requestID " +
        "FROM volunteerPlants vp " +
        "JOIN treePlantings tp ON tp.plantID = vp.plantID " +
        "JOIN treeRequests tr ON tr.requestID = tp.requestID " +
        "JOIN siteVisits sv ON sv.requestNum = tr.requestID " +
        "JOIN recommendedTrees rt ON rt.visitID = sv.siteVisitID " +
        "JOIN treeSpecies ts ON ts.treeID = rt.treeID " +
        "WHERE vp.vid = ?";
```

π tr.address, tr.phone, tp.plantDate, ts.commonName, vp.requestID (

σ vp.vid = @volunteerID (

((((vp \bowtie vp.plantID = tp.plantID treePlantings tp)

\bowtie tp.requestID = tp.requestID treeRequests tr)

\bowtie tr.requestID = sv.requestNum siteVisits sv)

\bowtie sv.siteVisitID = rt.visitID recommendedTrees rt)

\bowtie rt.treeID = ts.treeID treeSpecies ts

)

)

3.1

```
SELECT tr.neighborhood,
       ROUND(AVG(vp.workHour), 2) AS AverageWorkingHour,
       COUNT(DISTINCT tp.requestID) AS treePlantedCount,
       SUM(vp.workloadFeedback = 'overload') AS OverloadCount,
       ROUND(SUM(vp.workloadFeedback = 'overload')/COUNT(*), 2) AS OverloadRate
FROM volunteerPlants vp
  INNER JOIN treePlantings tp ON vp.requestID = tp.requestID
  INNER JOIN treeRequests tr ON tr.requestID = tp.requestID
  INNER JOIN neighborhoods n ON tr.neighborhood = n.name
WHERE n.district = (
  SELECT n2.district
    FROM neighborhoods n2
   INNER JOIN treeRequests tr2 ON tr2.neighborhood = n2.name
   INNER JOIN treePlantings tp2 ON tr2.requestID = tp2.requestID
   INNER JOIN volunteerPlants vp2 ON vp2.requestID = tp2.requestID
  GROUP BY n2.district
  HAVING SUM(vp2.workHour) >= ALL(
    SELECT SUM(vp3.workHour)
      FROM neighborhoods n3
     INNER JOIN treeRequests tr3 ON tr3.neighborhood = n3.name
     INNER JOIN treePlantings tp3 ON tr3.requestID = tp3.requestID
     INNER JOIN volunteerPlants vp3 ON vp3.requestID = tp3.requestID
    GROUP BY n3.district
  )
)
GROUP BY tr.neighborhood
HAVING COUNT(DISTINCT tp.requestID) > 1;
```

Score of the query: 6 points complex

- Tables joined (1–2:0 points, ≥3:1 point) 1 point
- Non-inner/natural join? (no:0 points, yes:1 point)
- # of subqueries (0:0 points, 1:1 point, >2:2 points) 2 points
- # queries comprising result via union/intersect (0:0 points, ≥1:1 point)
- Aggregate function(s) and grouping rows? (no:0 points, yes:1 point) 1 point
- # WHERE/HAVING conditions not for joins (≤1:0 points, >1:1 point) 1 point
- Non-aggregation functions or expressions in SELECT/WHERE? (no:0 points, yes:1 point)
- Strong motivation/justification for the query in the domain? (no:0 points, yes:1 point) 1 point

subquery_sum_of_workhour_each_district:

$\Pi_{\text{SUM}(vp3.\text{workHour}) \ \forall \ n3.\text{district}, \text{SUM}(vp3.\text{workHour})}$
(
 neighborhoods n3
 $\bowtie_{tr3.\text{neighborhood}=n3.\text{name}}$ treeRequests tr3
 $\bowtie_{tr3.\text{requestID}=tp3.\text{requestID}}$ treePlantings tp3

```

    ⋈ vp3.requestID=tp3.requestID volunteerPlants vp3
)

```

subquery_find_max_workhour_district:

```

Πn2.district
(σSUM(vp2.workHour) ≥ ALL(subquery_sum_of_workhour_each_district)
  ⋈n2.district,SUM(vp2.workHour)
  (neighborhoods n2
    ⋈n2.name=tr2.neighborhood treeRequests tr2
    ⋈tr2.requestID=tp2.requestID treePlantings tp2
    ⋈vp2.requestID=tp2.requestID volunteerPlants vp2
  )
)

```

main_query:

```

Πtr.neighborhood,ROUND(AVG(vp.workHour),2)-> AverageWorkingHour,
COUNT(DISTINCT tp.requestID)-> treePlantedCount,
SUM(vp.workloadFeedback='overload') -> OverloadCount,
ROUND(SUM(vp.workloadFeedback='overload')/COUNT(*),2)-> OverloadRate
σtreePlantedCount>1
⋈tr.neighborhood,ROUND(AVG(vp.workHour),2),COUNT(DISTINCT tp.requestID),
SUM(vp.workloadFeedback='overload'),
ROUND(SUM(vp.workloadFeedback='overload')/COUNT(*),2
  σn.district=(subquery_find_max_workhour_district)
  (volunteerPlants vp
    ⋈vp.requestID=tp.requestID treePlantings tp
    ⋈tr.requestID=tp.requestID treeRequests tr
    ⋈tr.neighborhood=n.name neighborhoods n
  )
)

```

3.2 equivalent RA & SQL query

```

SELECT tr.neighborhood,
       ROUND(AVG(vp.workHour), 2) AS AverageWorkingHour,
       COUNT(DISTINCT tp.requestID) AS treePlantedCount,
       SUM(vp.workloadFeedback = 'overload') AS OverloadCount,
       ROUND(SUM(vp.workloadFeedback = 'overload') * 1.0 / COUNT(*), 2) AS OverloadRate
FROM volunteerPlants vp
  INNER JOIN treePlantings tp ON vp.requestID = tp.requestID
  INNER JOIN treeRequests tr ON tr.requestID = tp.requestID
  INNER JOIN neighborhoods n ON tr.neighborhood = n.name
WHERE n.district IN (
  SELECT district
  FROM (
    SELECT district, SUM(workHour) AS totalWork
    FROM volunteerPlants vp2
      INNER JOIN treePlantings tp2 ON vp2.requestID = tp2.requestID
      INNER JOIN treeRequests tr2 ON tr2.requestID = tp2.requestID
      INNER JOIN neighborhoods n2 ON tr2.neighborhood = n2.name
    GROUP BY district
  ) AS subquery_totalwork_per_district
 WHERE totalWork = (
   SELECT MAX(totalWork)
   FROM (
    SELECT district, SUM(workHour) AS totalWork
    FROM volunteerPlants vp3
      INNER JOIN treePlantings tp3 ON vp3.requestID = tp3.requestID
      INNER JOIN treeRequests tr3 ON tr3.requestID = tp3.requestID
      INNER JOIN neighborhoods n3 ON tr3.neighborhood = n3.name
    GROUP BY district
   ) AS subquery_totalwork_per_district
  )
)
GROUP BY tr.neighborhood
HAVING COUNT(DISTINCT tp.requestID) > 1;

```

subquery_totalwork_per_district:

$\Pi_{\text{district}, \text{SUM}(\text{workHour}) \rightarrow \text{totalWork}}$

$\gamma_{\text{district}, \text{SUM}(\text{workHour})}$

(volunteerPlants vp3

⋈_{vp3.requestID=tp3.requestID} treePlantings tp3

⋈_{tr3.requestID=tp3.requestID} treeRequests tr3

⋈_{tr3.neighborhood=n3.name} neighborhoods n3

)

)

subquery_max_totalwork_hour:

$\Pi_{\text{MAX}(\text{totalWork})}$ subquery_totalwork_per_district

subquery_max_totalwork_district:

Π_{district} subquery_totalwork_per_district

main_query:

$\Pi_{\text{tr.neighborhood}, \text{ROUND}(\text{AVG}(\text{vp.workHour}), 2) \rightarrow \text{AverageWorkingHour},$
 $\text{COUNT}(\text{DISTINCT tp.requestID}) \rightarrow \text{treePlantedCount},$
 $\text{SUM}(\text{vp.workloadFeedback} = \text{'overload'}) \rightarrow \text{OverloadCount},$
 $\text{ROUND}(\text{SUM}(\text{vp.workloadFeedback} = \text{'overload'}) / \text{COUNT}(*), 2) \rightarrow \text{OverloadRate}$
 $\sigma_{\text{treePlantedCount} > 1}$
 $\gamma_{\text{tr.neighborhood}, \text{ROUND}(\text{AVG}(\text{vp.workHour}), 2), \text{COUNT}(\text{DISTINCT tp.requestID}),$
 $\text{SUM}(\text{vp.workloadFeedback} = \text{'overload'}),$
 $\text{ROUND}(\text{SUM}(\text{vp.workloadFeedback} = \text{'overload'}) / \text{COUNT}(*), 2)$
 $\sigma_{\text{n.district IN (subquery_max_totalwork_district)}}$
 $(\text{volunteerPlants vp}$
 $\quad \bowtie_{\text{vp.requestID} = \text{tp.requestID}} \text{treePlantings tp}$
 $\quad \bowtie_{\text{tr.requestID} = \text{tp.requestID}} \text{treeRequests tr}$
 $\quad \bowtie_{\text{tr.neighborhood} = \text{n.name}} \text{neighborhoods n}$
 $)$

Justification:

Version 2 is typically more efficient than Version 1 because:

- It avoids using the ALL clause (which is hard to optimize)
- Uses a MAX() aggregate (which is more efficient)
- Enables better subquery reuse (maybe using a WITH will improve this significantly)

4.1

```

SELECT n.name AS neighborhood, t2.commonName AS mostRecommendedTree
FROM recommendedTrees rt2
INNER JOIN treeSpecies t2 ON rt2.treeID = t2.treeID
INNER JOIN siteVisits sv2 ON rt2.requestID = sv2.requestID
INNER JOIN treeRequests tr2 ON tr2.requestID = sv2.requestID
INNER JOIN neighborhoods n ON n.name = tr2.neighborhood
GROUP BY n.name, t2.commonName
HAVING COUNT(*) = (
    SELECT MAX(tree_count)
    FROM (
        SELECT n2.name AS neighborhood, rt3.treeID, COUNT(*) AS tree_count
        FROM recommendedTrees rt3
        INNER JOIN siteVisits sv3 ON rt3.requestID = sv3.requestID
        INNER JOIN treeRequests tr3 ON tr3.requestID = sv3.requestID
        INNER JOIN neighborhoods n2 ON n2.name = tr3.neighborhood
        GROUP BY n2.name, rt3.treeID
    ) AS subquery
    WHERE subquery.neighborhood = n.name
)
ORDER BY n.name;

```

Score of the query: 6 points, complex

- Tables joined (1–2:0 points, ≥3:1 point) 1 point
- Non-inner/natural join? (no:0 points, yes:1 point)
- # of subqueries (0:0 points, 1:1 point, >2:2 points) 2 points
- # queries comprising result via union/intersect (0:0 points, ≥1:1 point)
- Aggregate function(s) and grouping rows? (no:0 points, yes:1 point) 1 point
- # WHERE/HAVING conditions not for joins (≤1:0 points, >1:1 point) 1 point
- Non-aggregation functions or expressions in SELECT/WHERE? (no:0 points, yes:1 point)
- Strong motivation/justification for the query in the domain? (no:0 points, yes:1 point) 1 point

$\Pi_{n.name, t2.commonName}$

$\sigma_{COUNT(*)=}$

$\Pi_{MAX(tree_count)}$

$\sigma_{subquery.neighborhood=n.name}$

(
 $\Pi_{n2.name, rt3.treeID, COUNT(*) \rightarrow tree_count}$
 $\gamma_{n2.name, rt3.treeID, COUNT(*)}$
 (recommendedTrees rt3
 $\bowtie_{sv3.requestID=rt3.requestID}$ siteVisits sv3
 $\bowtie_{tr3.requestID=sv3.requestID}$ treeRequests tr3
 $\bowtie_{tr3.neighborhood=n2.name}$ neighborhoods n2
)

)-> subquery

)

Yn.name,t2.commonName,COUNT(*)

(recommendedTrees rt3

⌘ sv3.requestID=rt3.requestID siteVisits sv3

⌘ tr3.requestID=sv3.requestID treeRequests tr3

⌘ tr3.neighborhood=n2.name neighborhoods n2

)

4.2 equivalent RA & query

```

SELECT n.name AS neighborhood, t.commonName AS mostRecommendedTree
FROM recommendedTrees rt
INNER JOIN treeSpecies t ON rt.treeID = t.treeID
INNER JOIN siteVisits sv ON rt.requestID = sv.requestID
INNER JOIN treeRequests tr ON tr.requestID = sv.requestID
INNER JOIN neighborhoods n ON n.name = tr.neighborhood
INNER JOIN (
    SELECT n2.name AS neighborhood, rt2.treeID, COUNT(*) AS tree_count
    FROM recommendedTrees rt2
    INNER JOIN siteVisits sv2 ON rt2.requestID = sv2.requestID
    INNER JOIN treeRequests tr2 ON tr2.requestID = sv2.requestID
    INNER JOIN neighborhoods n2 ON n2.name = tr2.neighborhood
    GROUP BY n2.name, rt2.treeID
) AS TreeCounts ON TreeCounts.neighborhood = n.name AND TreeCounts.treeID = rt.treeID
INNER JOIN (
    SELECT SubQuery.neighborhood AS neighborhood, MAX(tree_count) AS max_tree_count
    FROM (
        SELECT n3.name AS neighborhood, rt3.treeID, COUNT(*) AS tree_count
        FROM recommendedTrees rt3
        INNER JOIN siteVisits sv3 ON rt3.requestID = sv3.requestID
        INNER JOIN treeRequests tr3 ON tr3.requestID = sv3.requestID
        INNER JOIN neighborhoods n3 ON n3.name = tr3.neighborhood
        GROUP BY n3.name, rt3.treeID
    ) AS SubQuery
    GROUP BY SubQuery.neighborhood
) AS MaxTreeCounts ON MaxTreeCounts.neighborhood = n.name
WHERE TreeCounts.tree_count = MaxTreeCounts.max_tree_count
GROUP BY n.name, t.commonName
ORDER BY n.name;

```

$\Pi_{n.name, t.commonName}$

$\sigma_{TreeCounts.tree_count = MaxTreeCounts.max_tree_count}$

$\gamma_{n.name, t.commonName}$

(recommendedTrees rt

$\bowtie_{rt.treeID = t.treeID}$ treeSpecies t

$\bowtie_{rt.requestID = sv.requestID}$ siteVisits sv

$\bowtie_{tr.requestID = sv.requestID}$ treeRequests tr

$\bowtie_{tr.neighborhood = n.name}$ neighborhoods n

$\bowtie_{TreeCounts.neighborhood = n.name \wedge TreeCounts.treeID = rt.treeID}$

(

$\Pi_{n.name \rightarrow neighborhood, rt2.treeID, COUNT(*) \rightarrow tree_count}$

$\gamma_{2.name, rt2.treeID, COUNT(*)}$

(recommendedTrees rt2

```

        ⋈r2.treeID=sv2.treeID siteVisits sv2
        ⋈tr2.treeID=sv2.treeID treeRequests tr2
        ⋈n2.name=tr2.neighborhood neighborhoods n2
    )
)AS TreeCounts
⋈MaxTreeCounts.neighborhood=n.name
(
    ⋈SubQuery.neighborhood->neighborhood,MAX(tree_count)->max_tree_count
    ⋈SubQuery.neighborhood,MAX(tree_count)
    (
        ⋈n3.name->neighborhood,rt3.treeID,COUNT(*)->tree_count
        (
            recommendedTrees rt3
            ⋈rt3.treeID=sv3.treeID siteVisits sv3
            ⋈tr3.treeID=sv3.treeID treeRequests tr3
            ⋈n3.name=tr3.neighborhood neighborhoods n3
        )
    )AS SubQuery
)AS MaxTreeCounts
)

```

Justification:

Version 2 is typically more efficient than Version 1 because:

- It avoids correlated subqueries (which are costly and hard to optimize)
- Uses joins with pre-aggregated data (enabling better performance and reuse)
- Allows the optimizer to apply indexes and efficient join strategies
- Performs a fixed number of aggregations, regardless of group count
- Scales better with large datasets due to reduced redundant computation

5.

```
-- justification/motivation: we want to track tree species and volunteer number impact in actively planting neighborhood.
-- PS: the current result is empty table as no valid row yet; we temporarily set given year is 2023.
-- Show the number of distinct tree species planted by neighborhood in a given year,
-- only for neighborhoods where total planting events that year exceeded 2,
-- including the number of volunteers involved in each neighborhood's plantings.
-- time complexity:

SELECT n.name AS neighborhood, COUNT(DISTINCT t.commonName) AS species_count, COUNT(DISTINCT
vp.vid) AS total_volunteers
FROM treePlantings tp
  JOIN treeRequests tr ON tp.requestRefNum = tr.referenceNum
  JOIN neighborhoods n ON tr.neighborhood = n.name
  JOIN trees t ON TRIM(LEADING ' ' FROM tp.treePlanted) = t.commonName
  LEFT JOIN volunteerPlants vp ON vp.plantID = tp.plantID -- some tree planting might have no volunteer
WHERE YEAR(tp.plantDate) = 2023 AND tr.neighborhood IN (
  SELECT tr2.neighborhood
  FROM treePlantings tp2
    JOIN treeRequests tr2 ON tp2.requestRefNum = tr2.referenceNum
  WHERE YEAR(tp2.plantDate) = 2023
  GROUP BY tr2.neighborhood
  HAVING COUNT(tp2.plantID) >= 2
)
GROUP BY n.name;
```

1. **Tables joined (1–2: 0 points, ≥3: 1 point):**
 - Main query joins 5 tables. Score: +1 point
2. **Non-inner/natural join? (no: 0 points, yes: 1 point):**
 - there is a LEFT JOIN volunteerPlants vp. Score: +1 point
3. **# of subqueries (0: 0 points, 1: 1 point, >1: 2 points):**
 - There is one subquery (used in the IN clause). Score: +1 point
4. **Aggregate function(s) and grouping rows? (no: 0 points, yes: 1 point):**
 - the main query uses COUNT(DISTINCT ...) aggregate functions and has a GROUP BY n.name. Score: +1 point
5. **# WHERE/HAVING conditions not for joins (≤1: 0 points, >1: 1 point):**
 - Main query WHERE: YEAR(tp.plantDate) = 2023 (1), tr.neighborhood IN (...) (2). Score: +1 point
6. **Non-aggregation functions or expressions in SELECT/WHERE? (no: 0 points, yes: 1 point):**
 - TRIM(LEADING ' ' FROM tp.treePlanted) is used in a JOIN condition. Score: +1 point
7. **Strong motivation/justification for the query in the domain? (no: 0 points, yes: 1 point):**

- The query identifies neighborhoods with a minimum level of planting activity (≥ 2 plantings) in a specific year (2023) and then calculates the species diversity and volunteer involvement for those same active neighborhoods within that year. This is a plausible analytical question for urban forestry or community programs. Score: +1 point

Total Score **7 points**

RA expression:

Subquery IN clause: As **QualifyingNeighborhoods**

QualifyingNeighborhoods = Π TR.neighborhood σ planting_count ≥ 2

γ TR.neighborhood, COUNT(TP2.plantID) (σ YEAR(TP.plantDate) = 2023 (TP \bowtie TP.requestRefNum = TR.referenceNum TR))

Main query:

Π N.name AS neighborhood, species_count, total_volunteers

γ N.name, COUNT(DISTINCT T.commonName) AS species_count, COUNT(DISTINCT VP.vid) AS total_volunteers ((σ YEAR(TP.plantDate) = 2023 ((((TP \bowtie TP.requestRefNum = TR.referenceNum TR) \bowtie TR.neighborhood = N.name N) \bowtie TRIM(LEADING ' ' FROM TP.treePlanted) = T.commonName T) \bowtie TP.plantID = VP.plantID VP)) \bowtie TR.neighborhood = **QualifyingNeighborhoods**))

Equivalent RA expressions:

1. Apply filters as early as possible to reduce the size of intermediate relations before joins: Applying the $\text{YEAR}(\text{tp.plantDate}) = 2023$ filter (σ) to the TP relation before joining it would reduce the number of rows processed in subsequent joins. Applying the Semijoin (\bowtie) filter earlier also reduces the data flowing into the later joins.

Π N.name AS neighborhood, species_count, total_volunteers (

γ N.name, COUNT(DISTINCT T.commonName) AS species_count, COUNT(DISTINCT VP.vid) AS
total_volunteers (

(

(

(

(

(σ YEAR(TP.plantDate) = 2023(TP)) //Filter TP early for main query

\bowtie TP.requestRefNum = TR.referenceNum TR // Join pre-filtered TP with TR

)

\bowtie TR.neighborhood = QualifyingNeighborhoods // Apply IN clause filter Semijoin

early

(

// Start definition of QualifyingNeighborhoods

Π TR.neighborhood (// Project result of subquery

σ planting_count >= 2 (// Apply Having

```

    γTR.neighborhood, COUNT(TP.plantID) AS planting_count ( // Group/Count for
subquery

    ( σYEAR(TP.plantDate) = 2023(TP) )      // Filter TP early for subquery

    ⋈ TP.requestRefNum = TR.referenceNum TR    // Compute subquery
intermediate

    )

    )

    )

    // end definition of QualifyingNeighborhoods

    )

    )

    ⋈ TR.neighborhood = N.name N                // Join with Neighborhoods

    )

    ⋈ TRIM(LEADING ' ' FROM TP.treePlanted) = T.commonName T    // Join with Trees

    )

    ⋈ TP.plantID = VP.plantID VP                // Left Join with Volunteers

    )

) // Grouping and Aggregation applied before final projection

```

2. Replacing IN with JOIN

Computes the qualifying neighborhoods first and then uses a standard JOIN instead of a semijoin (\bowtie) to filter the main query path, because database optimizers often have more sophisticated strategies available for optimizing standard joins compared to semijoins or IN subqueries.


```

ΠN.name AS neighborhood, species_count, total_volunteers (

    γN.name, COUNT(DISTINCT T.commonName) AS species_count, COUNT(DISTINCT
VP.vid) AS total_volunteers ( // Grouping/Aggregation

    (

        (

            (

                (

                    (

                        ( σYEAR(TP.plantDate) = 2023(TP) )           // Filter TP early for main query

                        ⋈ TP.requestRefNum = TR.referenceNum TR    // Join pre-filtered TP with TR

                    )

                        ⋈ TR.neighborhood = QualifyingNeighborhoods // JOIN with
QualifyingNeighborhoods instead of IN

                )

                    // Start Definition of QualifyingNeighborhoods

                    ΠTR.neighborhood (

                        σ planting_count >= 2 (                      // Apply Having

                            γTR.neighborhood, COUNT(TP.plantID) AS planting_count ( //
Group/Count for subquery

                                ( σYEAR(TP.plantDate) = 2023(TP) )    // Filter TP early for subquery

```

```

        ⋈ TP.requestRefNum = TR.referenceNum TR // Compute subquery
intermediate
    )
    )
    )
    // End Definition of QualifyingNeighborhoods
    )
    )
    ⋈ TR.neighborhood = N.name N // Join with Neighborhoods
    )
    ⋈ TRIM(LEADING ' ' FROM TP.treePlanted) = T.commonName T // Join with Trees
    )
    ⋈ TP.plantID = VP.plantID VP //Left Join with Volunteers
    )
    )
    )

```

Execution plan and visualization

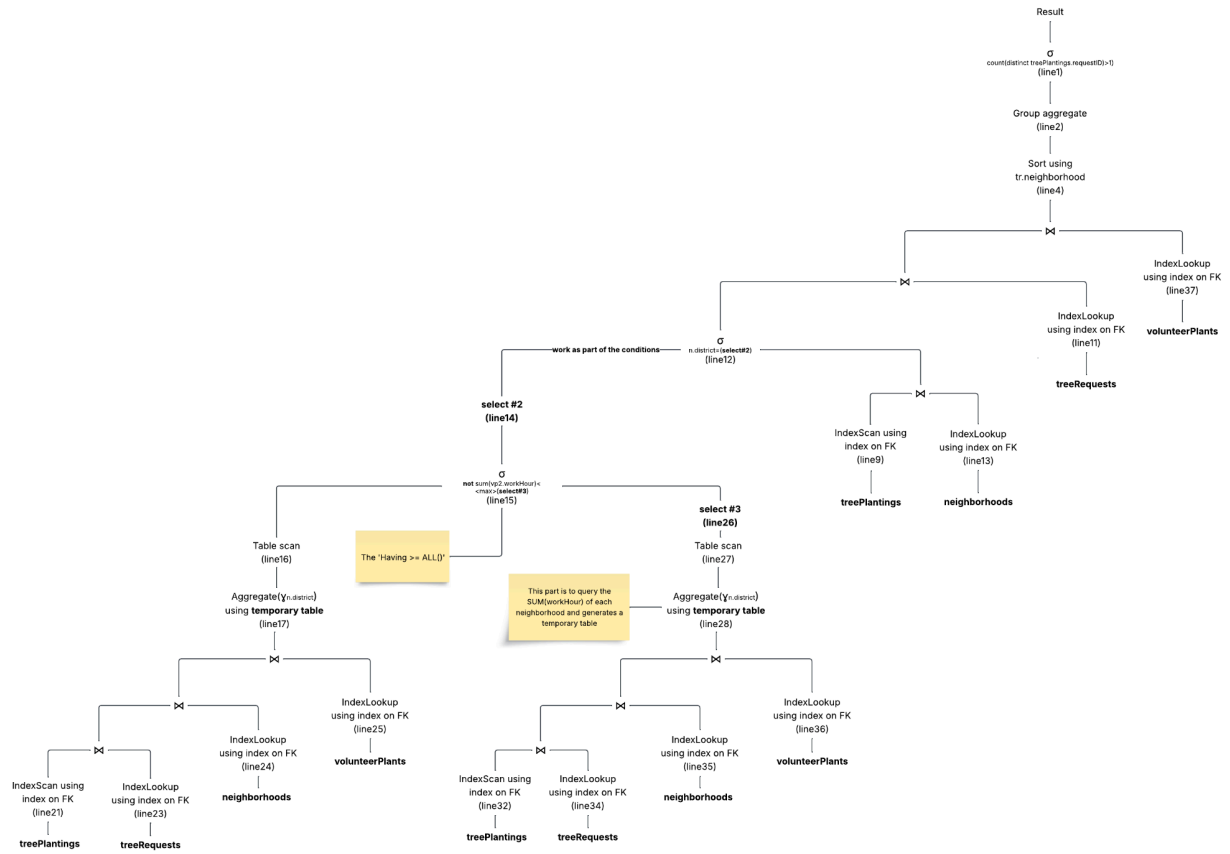
This is the execution plan that is generated by running `EXPLAIN ANALYZE` on the report query about the workload of volunteers in the district that has the maximum total work hours. For the query please see 3.1 or the first query in reports.sql.

```
1  -> Filter: (count(distinct treePlantings.requestID) > 1)
2  -> Group aggregate: count(distinct treePlantings.requestID), count(0), sum(tmp_field),
3    avg(volunteerPlants.workHour), count(distinct treePlantings.requestID), sum(tmp_field)
4  -> Sort: tr.neighborhood
5  -> Stream results
6    -> Nested loop inner join
7      -> Nested loop inner join
8        -> Nested loop inner join
9          -> Covering index scan on tp using FK_planting_aid
10         -> Filter: (tr.neighborhood is not null)
11         -> Single-row index lookup on tr using PRIMARY (requestID = tp.requestID)
12         -> Filter: (n.district = (select #2))
13         -> Single-row index lookup on n using PRIMARY (name = tr.neighborhood)
14         -> Select #2 (subquery in condition; run only once)
15           -> Filter: <not>((`sum(vp2.workHour)` < <max>(select #3)))
16           -> Table scan on <temporary>
17             -> Aggregate using temporary table
18               -> Nested loop inner join
19                 -> Nested loop inner join
20                   -> Nested loop inner join
21                     -> Covering index scan on tp2 using FK_planting_aid
22                     -> Filter: (tr2.neighborhood is not null)
23                     -> Single-row index lookup on tr2 using PRIMARY (requestID = tp2.requestID)
24                     -> Single-row index lookup on n2 using PRIMARY (name = tr2.neighborhood)
25                     -> Index lookup on vp2 using PRIMARY (requestID = tp2.requestID)
26           -> Select #3 (subquery in condition; run only once)
27             -> Table scan on <temporary>
28               -> Aggregate using temporary table
29                 -> Nested loop inner join
30                   -> Nested loop inner join
31                     -> Nested loop inner join
32                       -> Covering index scan on tp3 using FK_planting_aid
33                       -> Filter: (tr3.neighborhood is not null)
34                       -> Single-row index lookup on tr3 using PRIMARY (requestID = tp3.requestID)
35                       -> Single-row index lookup on n3 using PRIMARY (name = tr3.neighborhood)
36                     -> Index lookup on vp3 using PRIMARY (requestID = tp3.requestID)
37           -> Index lookup on vp using PRIMARY (requestID = tp.requestID)
```

And below is the visualization of the execution plan. To review the diagram easily, please use this link to Lucid chart:

https://lucid.app/lucidchart/1cdf3ac8-bced-4d66-9733-e5d723fca36c/edit?page=0_0&invitationId=inv_e595e45b-df1f-4694-8fff-eaf8c16fa2be#

We have already sent the invitation to you using the email a.monge@northeastern.edu, if you have not received the invitation, please contact us.



Stored procedure

This stored procedure aims to delete one treeRequest given the requestID.

Analysed steps for the procedure:

1. Verify if this treeRequest exists
2. Find the related siteVisit using the requestID
 - a. Find the related recommendedTree data using the requestID
 - b. Delete all the related rows in table recommendedTrees
3. Delete the related row in table siteVisits(if exists)
4. Find the related treePlanting using the requestID
 - a. Find the related volunteerPlant data using the requestID
 - b. Delete all the related rows in table volunteerPlants
5. Delete the related row in table treePlantings(if exists)
6. Delete the treeRequest data in table treeRequests

Since we used ON DELETE CASCADE in the FK constraint of table recommendedTrees and table volunteerPlants, we actually only need to delete the related row in siteVisits and treePlantings.(In table siteVisits and treePlantings the rule is ON DELETE NO ACTION, so we need to delete the data manually)

The code: (you can find this piece of code in the end of file ddl.sql)

```
DROP PROCEDURE IF EXISTS delete_a_treeRequest;
DELIMITER //

CREATE PROCEDURE delete_a_treeRequest(
    IN treeRequestID INT, -- id of the tree request
    OUT status VARCHAR(500) -- message providing status of request
)
BEGIN
    DECLARE row_count INT DEFAULT 0;
    DECLARE siteVisits_deleted_count INT DEFAULT 0;
    DECLARE treePlantings_deleted_count INT DEFAULT 0;
    DECLARE recommendedTrees_deleted_count INT DEFAULT 0;
    DECLARE volunteerPlants_deleted_count INT DEFAULT 0;

    -- verify if the tree request exists
    SELECT COUNT(requestID) INTO row_count
        FROM treeRequests WHERE requestID = treeRequestID;

    SET status = '';

    IF row_count = 1 THEN -- this tree request exists
        SELECT COUNT(requestID) INTO row_count
            FROM siteVisits WHERE requestID = treeRequestID;

        IF row_count = 1 THEN -- the related site visit exists, delete it
            SET siteVisits_deleted_count = 1;
            -- count the deleted rows in the recommendedTrees
            SELECT COUNT(*) INTO recommendedTrees_deleted_count
```

```

        FROM recommendedTrees WHERE requestID = treeRequestID;
        DELETE FROM siteVisits WHERE requestID = treeRequestID;
    end if;

    SELECT COUNT(requestID) INTO row_count
        FROM treePlantings WHERE requestID = treeRequestID;

    IF row_count = 1 THEN -- the related tree planting event exists, delete
it
        SET treePlantings_deleted_count = 1;
        -- count the deleted rows in the volunteerPlants
        SELECT COUNT(*) INTO volunteerPlants_deleted_count
            FROM volunteerPlants WHERE requestID = treeRequestID;
        DELETE FROM treePlantings WHERE requestID = treeRequestID;
    end if;

    DELETE FROM treeRequests WHERE requestID = treeRequestID;

    SET status = CONCAT(status, 'Deleted the tree request(id): ',
treeRequestID,
                        ', deleted count of rows in siteVisits: ',
siteVisits_deleted_count,
                        ', deleted count of rows in recommendedTrees: ',
recommendedTrees_deleted_count,
                        ', deleted count of rows in treePlantings: ',
treePlantings_deleted_count,
                        ', deleted count of rows in volunteerPlants: ',
volunteerPlants_deleted_count
                        );

    ELSE
        SET status = 'This tree request does not exist!';
    end if;
end //

DELIMITER ;

```

For the Java code that called this stored procedure, please see in the presentation video or task running screenshots in the final report.