# Topic 1: Data models and query languages
# Unit 3: Relational Algebra (RA)
# Lecture 8

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp23)

https://northeastern-datalab.github.io/cs7240/sp23/

2/3/2023

# Pre-class conversations

- Online today: please use your cameras!

- Last class summary

- Please do look through the slides between lectures

- It is time to start to hand in your first scribes (some ideas today)

- Project discussions (in 2 weeks: Fri 2/17: project ideas)

- today:
  - End of Algebra / Codd's theorem
  - Datalog (recursion), and a new tool: Souffle

# Algebra and the connection to logic and queries

- Algebra

- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations

- Equivalence RA and safe RC (Codd's theorem)
  - RA → RC
  - RC → RA

# DRC → RA: Intuition

Proof (Sketch):

- Show first that for every relational database schema **S**, there is a relational algebra expression **E** such that for every database instance **D**, we have that ADom(**D**) = **E**(**D**).

- Use the above fact and induction on the construction of RC formulas to obtain a translation of RC under the active domain interpretation to RA.

# DRC → RA: Intuition          $E(A,B)$

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra

  uses the logical equivalence:  $\forall y . \phi \equiv$          **?**

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra

  uses the logical equivalence:   $\forall y.\, \phi \equiv \quad \neg \exists y.\, \neg \phi$

- As an illustration, consider:    $\forall y.\, E(x, y) \equiv$    **?**

# DRC → RA: Intuition $\qquad$ E(A,B)

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra

  uses the logical equivalence: $\forall y. \phi \equiv \quad \neg \exists y. \neg \phi$

- As an illustration, consider: $\forall y. E(x, y) \equiv \neg \exists y. \neg E(x, y)$

  and recall: $ADom(D) = \quad$ **?**

# DRC → RA: Intuition                                    $E(A,B)$

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra

  uses the logical equivalence:  $\forall y. \phi \equiv \quad \neg \exists y. \neg \phi$

- As an illustration, consider:  $\forall y. E(x,y) \equiv \neg \exists y. \neg E(x,y)$

  and recall:  $\boxed{ADom(D) = \pi_A(E) \cup \pi_B(E)}$

| DRC formula $\phi$ | RA expression for $\phi^{adom}$ |
|---|---|
| $\neg E(x,y)$ | |
| $\exists y. \neg E(x,y)$ | ? |
| $\neg \exists y. \neg E(x,y)$ | |

# DRC → RA: Intuition                                    $E(A,B)$

- In this translation, the most interesting part is the simulation of the universal quantifier ∀ in relational algebra

  uses the logical equivalence:   $\forall y.\, \phi \equiv \quad \neg \exists y.\, \neg \phi$

- As an illustration, consider:   $\forall y.\, E(x,y) \equiv \neg \exists y.\, \neg E(x,y)$

  and recall:   $\boxed{\text{ADom(D)} = \pi_A(E) \cup \pi_B(E)}$

| DRC formula $\phi$ | RA expression for $\phi^{\text{adom}}$ |
|:---:|:---:|
| $\neg E(x,y)$ | $(\, \text{ADom(D)} \times \text{ADom(D)} \,) - E$ |
| $\exists y.\, \neg E(x,y)$ | $\pi_{-B}\;\; \pi_A[(\, \text{ADom(D)} \times \text{ADom(D)} \,) - E]$ |
| $\neg \exists y.\, \neg E(x,y)$ | $\text{ADom(D)} - \pi_A[(\, \text{ADom(D)} \times \text{ADom(D)} \,) - E]$ |

# Entire Story in One Slide (repeated slide)

1. RC = FOL over DB

2. RC can express "bad queries" that depend not only on the DB, but also on the domain from which values are taken  (domain dependence)

3. We cannot test whether an RC query is "good," but we can use a "good" subset of RC that captures all "good" queries  (safety)

4. "Good" RC and RA can express the same queries! (equivalence = Codd's theorem)

# Discussion

- What is the monotone fragment of RA ?

    ?

- What are the safe queries in RA ?

    ?

- Where do we use RA (applications) ?

    ?

# Discussion

- What is the monotone fragment of RA ?
  - Basic except difference (–): ∪, $\sigma$, $\pi$, ⋈

- What are the safe queries in RA ?

  **?**

- Where do we use RA (applications) ?

  **?**

# Discussion

- What is the monotone fragment of RA ?
  - Basic except difference (–): $\cup$, $\sigma$, $\pi$, $\bowtie$

$A(\vec{x}) - B(\vec{x})$

$\neg B(\vec{x})$

- What are the safe queries in RA ?
  - All RA queries are safe

- Where do we use RA (applications) ?

?

# Discussion

- ## What is the monotone fragment of RA ?
  - Basic except difference (–): ∪, $\sigma$, $\pi$, ⋈

- ## What are the safe queries in RA ?
  - All RA queries are safe

- ## Where do we use RA (applications) ?
  - Translating SQL (from WHAT to HOW)
  - Directly as query languages (e.g. Pig-Latin)

  *See next pages*

EXAMPLE 1. *Suppose we have a table* urls: (url, category, pagerank). *The following is a simple SQL query that finds, for each sufficiently large category, the average pagerank of high-pagerank urls in that category.*

```
SELECT category, AVG(pagerank)
FROM urls WHERE pagerank > 0.2
GROUP BY category HAVING COUNT(*) > 10^6
```

*An equivalent Pig Latin program is the following. (Pig Latin is described in detail in Section 3; a detailed understanding of the language is not required to follow this example.)*

```
good_urls = FILTER urls BY pagerank > 0.2;
groups = GROUP good_urls BY category;
big_groups = FILTER groups BY COUNT(good_urls)>10^6;
output = FOREACH big_groups GENERATE
            category, AVG(good_urls.pagerank);
```

EXAMPLE 1. *Suppose we have a table* urls: (url, category, pagerank). *The following is a simple SQL query that finds, for each sufficiently large category, the average pagerank of high-pagerank urls in that category.*

```
SELECT category, AVG(pagerank)
FROM urls WHERE pagerank > 0.2
GROUP BY category HAVING COUNT(*) > 10^6
```

*An equivalent Pig Latin program is the following. (Pig Latin is described in detail in Section 3; a detailed understanding of the language is not required to follow this example.)*

```
good_urls = FILTER urls BY pagerank > 0.2;
groups = GROUP good_urls BY category;
big_groups = FILTER groups BY COUNT(good_urls)>10^6;
output = FOREACH big_groups GENERATE
              category, AVG(good_urls.pagerank);
```
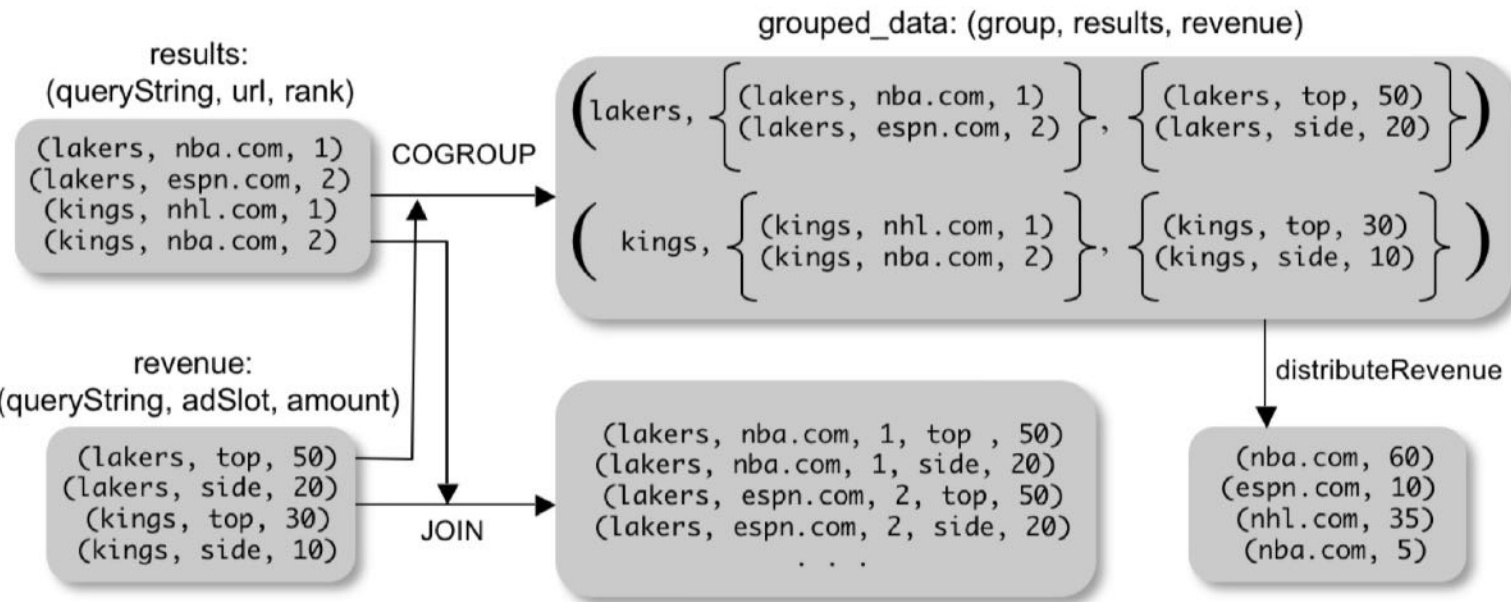
results:
(queryString, url, rank)

(lakers, nba.com, 1)
(lakers, espn.com, 2)
(kings, nhl.com, 1)
(kings, nba.com, 2)

grouped_data: (group, results, revenue)

(lakers, { (lakers, nba.com, 1) (lakers, espn.com, 2) }, { (lakers, top, 50) (lakers, side, 20) } )

( kings, { (kings, nhl.com, 1) (kings, nba.com, 2) }, { (kings, top, 30) (kings, side, 10) } )

COGROUP

revenue:
(queryString, adSlot, amount)

(lakers, top, 50)
(lakers, side, 20)
(kings, top, 30)
(kings, side, 10)

JOIN

(lakers, nba.com, 1, top , 50)
(lakers, nba.com, 1, side, 20)
(lakers, espn.com, 2, top, 50)
(lakers, espn.com, 2, side, 20)
. . .

distributeRevenue

(nba.com, 60)
(espn.com, 10)
(nhl.com, 35)
(nba.com, 5)

**Figure 2: COGROUP versus JOIN.**

grouped_data  =  COGROUP results BY queryString,
                          revenue BY queryString;

### 3.5.2  JOIN in Pig Latin

Not all users need the flexibility offered by COGROUP. In many cases, all that is required is a regular equi-join. Thus, Pig Latin provides a JOIN keyword for equi-joins. For example,

join_result  =  JOIN results BY queryString,
                        revenue BY queryString;

It is easy to verify that JOIN is only a syntactic shortcut for COGROUP followed by flattening. The above join command is equivalent to:

temp_var  =  COGROUP results BY queryString,
                     revenue BY queryString;
join_result  =  FOREACH temp_var GENERATE
                        FLATTEN(results), FLATTEN(revenue);