

Introduction to Database Systems

CSE 414

Lecture 15: SQL++ Wrapup

world

```
 {{ {"mondial":  
     {"country":  
         [{"-car_code": "AL", ...}  
          {"name": "Albania"}, ...  
        ], ...  
      }, ...  
    }}}
```

country

```
 {{ { "-car_code": "AL",  
      "gdp_total": 4100,  
      ...  
    }, ...  
  }}}
```

Find each country's GDP

```
SELECT x.mondial.country.name, c.gdp_total  
FROM world AS x, country AS c  
WHERE x.mondial.country.`-car_code` = c.`-car_code`;
```

x.mondial.country is an array of objects. No field as -car_code!

Error: Type mismatch!

Need to “unnest” the array

In General

Needs to be an array
or dataset
(i.e., iterable)

Object to be
iterated on

```
SELECT ...  
FROM R AS x, S AS y  
WHERE x.f1 = y.f2;
```

These cannot evaluate to an array or dataset!

Need to
“unnest”
the array

Unnesting collections

mydata

```
{"A": "a1", "B": [{"C": "c1", "D": "d1"}, {"C": "c2", "D": "d2"}]}  
{"A": "a2", "B": [{"C": "c3", "D": "d3"}]}  
{"A": "a3", "B": [{"C": "c4", "D": "d4"}, {"C": "c5", "D": "d5"}]}
```

```
SELECT x.A, y.C, y.D  
FROM mydata AS x, x.B AS y;
```

Form cross product between
each x and its x.B

Answer

```
{"A": "a1", "C": "c1", "D": "d1"}  
{"A": "a1", "C": "c2", "D": "d2"}  
{"A": "a2", "C": "c3", "D": "d3"}  
{"A": "a3", "C": "c4", "D": "d4"}  
{"A": "a3", "C": "c5", "D": "d5"}
```

Unnesting collections

mydata

```
{"A": "a1", "B": [{"C": "c1", "D": "d1"}, {"C": "c2", "D": "d2"}]}  
{"A": "a2", "B": [{"C": "c3", "D": "d3"}]}  
{"A": "a3", "B": [{"C": "c4", "D": "d4"}, {"C": "c5", "D": "d5"}]}
```

```
SELECT x.A, y.C, y.D  
FROM mydata AS x UNNEST x.B AS y;
```

Answer

Same as before

```
{"A": "a1", "C": "c1", "D": "d1"}  
{"A": "a1", "C": "c2", "D": "d2"}  
{"A": "a2", "C": "c3", "D": "d3"}  
{"A": "a3", "C": "c4", "D": "d4"}  
{"A": "a3", "C": "c5", "D": "d5"}
```

world

```
 {{ {"mondial":  
      {"country":  
        [{"-car_code": "AL", ...}  
         {"name": "Albania"}, ...  
        ], ...  
       }, ...  
     }}}
```

country

```
 {{ { "-car_code": "AL",  
      "gdp_total": 4100,  
      ...  
    }, ...  
  }}}
```

Find each country's GDP

```
SELECT y.name, c.gdp_total  
FROM world AS x, x.mondial.country AS y, country AS c  
WHERE y.-car_code = c.-car_code;
```

Answer

```
{ "name": "Albania", "gdp_total": "4100" }  
{ "name": "Greece", "gdp_total": "101700" }  
...
```

world

```
{{ {"mondial":  
    {"country": [{Albania}, {Greece}, ...],  
     "continent": [...],  
     "organization": [...],  
     ...  
     ...  
   }  
 }  
}}
```

Return province and city names

```
SELECT z.name AS province_name, u.name AS city_name  
FROM world x, x.mondial.country y, y.province z, z.city u  
WHERE y.name = "Greece";
```

The problem:

Error: Type mismatch!

```
"name": "Greece",  
"province": [ ...  
  {"name": "Attiki",  
   "city": [ {"name": "Athens"}, {"name": "Pireus"}, ...]  
   ...},  
  {"name": "Ipiros",  
   "city": {"name": "Ioannia"} }  
  ...], ...]
```

city is an array

city is an object

world

```
{{ {"mondial":  
    {"country": [{Albania}, {Greece}, ...],  
     "continent": [...],  
     "organization": [...],  
     ...  
     ...  
   }  
 }  
}}
```

Return province
and city names

```
SELECT z.name AS province_name, u.name AS city_name  
FROM world x, x.mondial.country y, y.province z,  
  
(CASE WHEN z.city IS missing THEN []  
      WHEN IS_ARRAY(z.city) THEN z.city  
      ELSE [z.city] END) AS u  
  
WHERE y.name="Greece";
```

Even better

Useful Functions

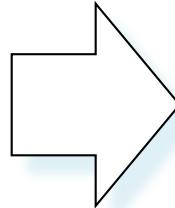
- `is_array`
- `is_boolean`
- `is_number`
- `is_object`
- `is_string`
- `is_null`
- `is_missing`
- `is_unknown = is_null or is_missing`

Other Features

- Unnesting
- Nesting
- Grouping and aggregate
- Joins
- Multi-value join

Nesting

C
[{A:a1, B:b1},
 {A:a1, B:b2},
 {A:a2, B:b1}]



We want:

[{A:a1, Grp:[{b1, b2}]},
 {A:a2, Grp:[{b1}]}}]

```
SELECT DISTINCT x.A,  
  (SELECT y.B FROM C AS y WHERE x.A = y.A) AS Grp  
FROM C AS x
```

Using LET syntax:

```
SELECT DISTINCT x.A, g AS Grp  
FROM C AS x  
LET g = (SELECT y.B FROM C AS y WHERE x.A = y.A)
```

Grouping and Aggregates

C

```
[{A:a1, F:[{B:b1}, {B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3}, {B:b4}, {B:null}], G:[ ]},  
 {A:a3, F:[{B:b6}]} , G:[{C:c2},{C:c3}]]]
```

Count the number of elements in the F array for each A

```
SELECT x.A, COLL_COUNT(x.F) AS cnt  
FROM C AS x
```

```
SELECT x.A, COUNT(*) AS cnt  
FROM C AS x, x.F AS y  
GROUP BY x.A
```

These are
NOT
equivalent!

Grouping and Aggregates

Function	NULL	MISSING	Empty Collection
COLL_COUNT	counted	counted	0
COLL_SUM	returns NULL	returns NULL	returns NULL
COLL_MAX	returns NULL	returns NULL	returns NULL
COLL_MIN	returns NULL	returns NULL	returns NULL
COLL_AVG	returns NULL	returns NULL	returns NULL
ARRAY_COUNT	not counted	not counted	0
ARRAY_SUM	ignores NULL	ignores NULL	returns NULL
ARRAY_MAX	ignores NULL	ignores NULL	returns NULL
ARRAY_MIN	ignores NULL	ignores NULL	returns NULL
ARRAY_AVG	ignores NULL	ignores NULL	returns NULL

Grouping and Aggregates

C

```
[{A:a1,  F:[{B:b1}, {B:b2}],          G:[{C:c1}]},  
 {A:a2,  F:[{B:b3}, {B:b4}, {B:null}], G:[ ]},  
 {A:a3,  F:[{B:b6}]],                  G:[{C:c2},{C:c3}]}]
```

Lesson:

Read the *\$@# manual!!

```
SELECT x.A, COUNT(*) AS cnt  
FROM C AS x, x.F AS y  
GROUP BY x.A
```

These are

NOT
equivalent!

Joins

Two flat collection

```
coll1 = [{A:a1, B:b1}, {A:a1, B:b2}, {A:a2, B:b1}]\ncoll2 = [{B:b1, C:c1}, {B:b1, C:c2}, {B:b3, C:c3}]
```

Answer

```
SELECT x.A, x.B, y.C\nFROM coll1 AS x, coll2 AS y\nWHERE x.B = y.B
```

```
[{A:a1, B:b1, C:c1},\n {A:a1, B:b1, C:c2},\n {A:a2, B:b1, C:c1},\n {A:a2, B:b1, C:c2}]
```

```
SELECT x.A, x.B, y.C\nFROM coll1 AS x JOIN coll2 AS y ON x.B = y.B
```

Outer Joins

Two flat collection

coll1 [{A:a1, B:b1}, {A:a1, B:b2}, {A:a2, B:b1}]

coll2 [{B:b1, C:c1}, {B:b1, C:c2}, {B:b3, C:c3}]

```
SELECT x.A, x.B, y.C  
FROM coll1 AS x RIGHT OUTER JOIN coll2 AS y  
    ON x.B = y.B
```

Answer

[{A:a1, B:b1, C:c1},
 {A:a1, B:b1, C:c2},
 {A:a2, B:b1, C:c1},
 {A:a2, B:b1, C:c2},
 {B:b3, C:c3}]

Ordering

coll1

```
[{A:a1, B:b1}, {A:a1, B:b2}, {A:a2, B:b1}]
```

```
SELECT x.A, x.B  
FROM coll AS x  
ORDER BY x.A
```

Data type matters!

"90" > "8000" but
90 < 8000 !

Multi-Value Join

river

```
[{"name": "Donau",    "-country": "SRB A D H HR SK BG RO MD UA"},  
 {"name": "Colorado", "-country": "MEX USA"},  
 ... ]
```

```
SELECT ...  
FROM country AS x, river AS y,  
     split(y.`-country`, " ") AS z  
WHERE x.`-car_code` = z
```

String

Separator

A collection

```
split("MEX USA", " ") = ["MEX", "USA"]
```

Behind the Scenes

Query Processing on NFNF data:

- Option 1: give up on query plans, use standard java/python-like execution
- Option 2: represent the data as a collection of flat tables, convert SQL++ to a standard relational query plan

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

Relational representation

coll:

id	A
1	a1
2	a2
3	a1

F

parent	B
1	b1
1	b2
2	b3
2	b4
2	b5
3	b6

G

parent	C
1	c1
3	c2
3	c3

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

SQL++

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y  
WHERE x.A = "a1"
```

Relational representation

coll:

id	A
1	a1
2	a2
3	a1

F

parent	B
1	b1
1	b2
2	b3
2	b4
2	b5
3	b6

G

parent	C
1	c1
3	c2
3	c3

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

SQL++

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y  
WHERE x.A = "a1"
```

Relational representation

coll:

id	A
1	a1
2	a2
3	a1

F

parent	B
1	b1
1	b2
2	b3
2	b4
2	b5
3	b6

G

parent	C
1	c1
3	c2
3	c3

SQL

```
SELECT x.A, y.B  
FROM coll AS x, F AS y  
WHERE x.id = y.parent AND x.A = "a1"
```

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

SQL++

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y  
WHERE x.A = "a1"
```

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y, x.G AS z  
WHERE y.B = z.C
```

Relational representation

coll:

id	A
1	a1
2	a2
3	a1

F

parent	B
1	b1
1	b2
2	b3
2	b4
2	b5
3	b6

G

parent	C
1	c1
3	c2
3	c3

SQL

```
SELECT x.A, y.B  
FROM coll AS x, F AS y  
WHERE x.id = y.parent AND x.A = "a1"
```

Flattening SQL++ Queries

A nested collection

```
coll =  
[ {A:a1, F:[{B:b1},{B:b2}], G:[{C:c1}]},  
 {A:a2, F:[{B:b3},{B:b4},{B:b5}], G:[ ]},  
 {A:a1, F:[{B:b6}], G:[{C:c2},{C:c3}]}]
```

SQL++

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y  
WHERE x.A = "a1"
```

```
SELECT x.A, y.B  
FROM coll AS x, x.F AS y, x.G AS z  
WHERE y.B = z.C
```

Relational representation

coll:

id	A
1	a1
2	a2
3	a1

F

parent	B
1	b1
1	b2
2	b3
2	b4
2	b5
3	b6

G

parent	C
1	c1
3	c2
3	c3

SQL

```
SELECT x.A, y.B  
FROM coll AS x, F AS y  
WHERE x.id = y.parent AND x.A = 'a1'
```

```
SELECT x.A, y.B  
FROM coll AS x, F AS y, G AS z  
WHERE x.id = y.parent AND x.id = z.parent  
AND y.B = z.C
```

Semistructured Data Model

- Several file formats: Json, protobuf, XML
- The data model is a tree
- They differ in how they handle structure:
 - Open or closed
 - Ordered or unordered
- Query language needs to take NFNF into account
 - Various “extra” constructs introduced as a result

Conclusion

- Semi-structured data best suited for *data exchange*
- “General” guidelines:
 - For quick, ad-hoc data analysis, use a “native” query language: SQL++, or AQL, or Xquery
 - Where “native” = how data is stored
 - Modern, advanced query processors like AsterixDB / SQL++ can process semi-structured data as efficiently as RDBMS
 - For long term data analysis: spend the time and effort to normalize it, then store in a RDBMS