# Exercise Session 11: FLWOR Expressions on Spark.
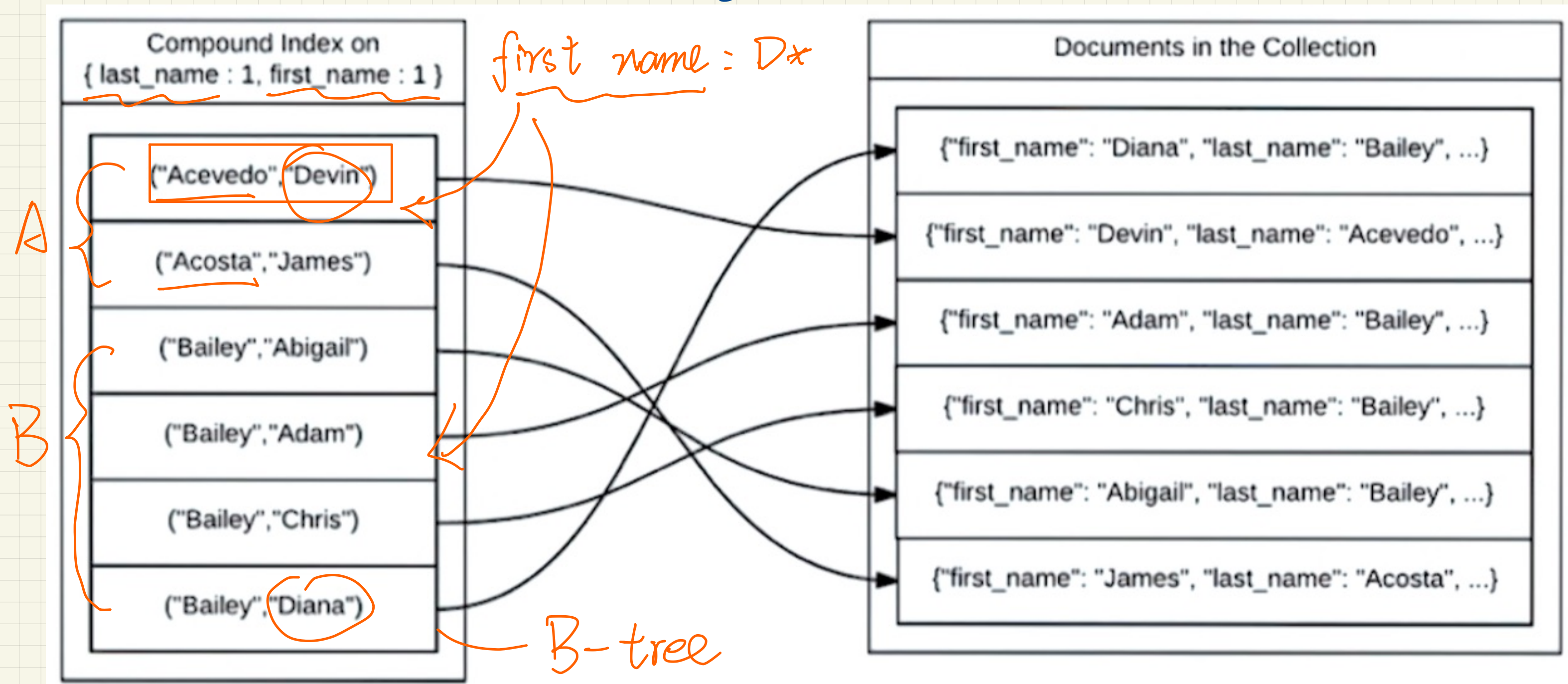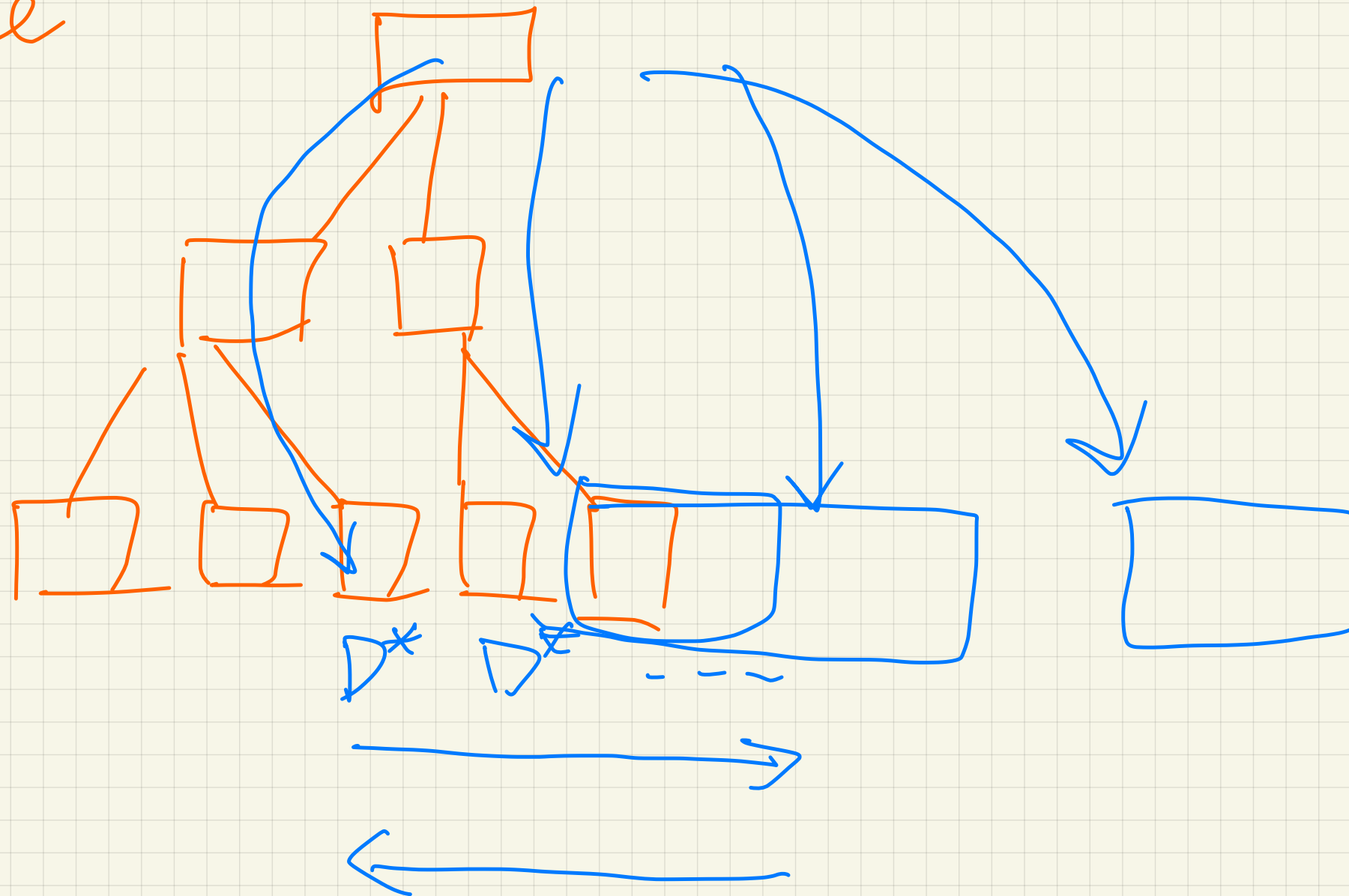
## Agenda:

① Compound index in MongoDB (Session 10)

② Inner workings of the FLOWR expressions

  ⎰ Tuple streams

  ⎱ Execution and implementation on Spark.

③ Tips on JSONiq.

# MongoDB Compound Indices ( session 10 )

- Compound keys are "flat" => 1D regardless of #fields
- Non-prefix queries can only use the index as a "filter"

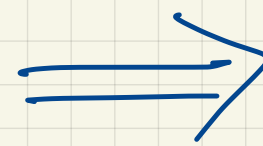first name : D*

**Compound Index on**
**{ last_name : 1, first_name : 1 }**

A {
("Acevedo","Devin")
("Acosta","James")
}

B {
("Bailey","Abigail")
("Bailey","Adam")
("Bailey","Chris")
("Bailey","Diana")
}

B-tree

**Documents in the Collection**

{"first_name": "Diana", "last_name": "Bailey", ...}

{"first_name": "Devin", "last_name": "Acevedo", ...}

{"first_name": "Adam", "last_name": "Bailey", ...}

{"first_name": "Chris", "last_name": "Bailey", ...}

{"first_name": "Abigail", "last_name": "Bailey", ...}

{"first_name": "James", "last_name": "Acosta", ...}

# B-tree

# Motivating Example of JSONiq

```
2 let $arr := ["a", "b", "b", "c", "c", "c"]
```

$\Longrightarrow$

```
{"c": 3}
{"a": 1}
{"b": 2}
```

**A:**

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in $arr[]
4 group by $x
5 return {
6     $x : count($arr)
7 }
```

**B:**

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : count($arr[][$$ = $x])
7 }
```

**C:**

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : count($x)
7 }
```

**D:**

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : 1
7 }
```

## Outputs:

**A** ✓

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in $arr[]
4 group by $x
5 return {
6     $x : count($arr)
7 }
```

```
Took: 0.033636609313964844 ms
{"c": 3}
{"a": 1}
{"b": 2}
```

**B** ✓

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : count($arr[][$$ = $x])
7 }
```

```
Took: 0.044761896133422851 ms
{"c": 3}
{"a": 1}
{"b": 2}
```

**C** ✗

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c", "c"]
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : count($x)
7 }
```

```
Took: 0.03741312026977539 ms
{"c": 1}
{"a": 1}
{"b": 1}
```

**D** ✗

```
1 %%jsoniq
2 let $arr := ["a", "b", "b", "c", "c",
3 for $x in distinct-values($arr[])
4 group by $x
5 return {
6     $x : 1
7 }
```

```
Took: 0.03769087791442871 ms
{"c": 1}
{"a": 1}
{"b": 1}
```

# FLOWR ~~Semantics~~ Effects

- `let` — Creates a <u>new column</u>

- `for` — <u>Reshapes</u> a list of values <u>into a column</u>

- `where` — Fliters out tuples given a criterion

- `group by` — Collects tuples of <u>the same key</u> to <u>a sequence</u>

- `return` — Emits streams of tuples <u>one by one</u>

    SELECT

# Example Breakdown

A single object

```
let $arr := ["a", "b", "b", "c", "c", "c"]
```

**①**

$arr

["a", "b", "b", "c", "c", "c"]

**②**

```
for $x in $arr[]
```

unpacking

seq. ⇒

```
["a", "b", "b", "c", "c", "c"]
```

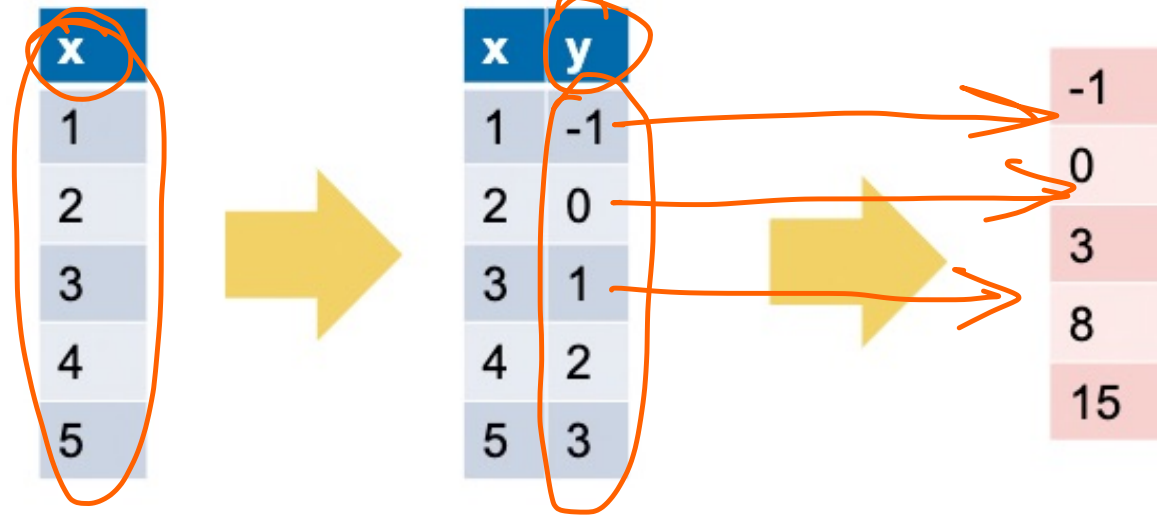| $x | $ arr |
|----|-------|
| a | ["a", "b", "b", "c", "c", "c"] |
| b | ["a", "b", "b", "c", "c", "c"] |
| b | ["a", "b", "b", "c", "c", "c"] |
| c | ["a", "b", "b", "c", "c", "c"] |
| c | ["a", "b", "b", "c", "c", "c"] |
| c | ["a", "b", "b", "c", "c", "c"] |

**③** group by $x

| $x | $ arr' |
|----|--------|
| a | [ ]  seq. |
| b | [ ], [ ], |
| c | [ ], [ ], [ ] |

tuple

aggregation on each tuple

**④**

```
return {
    $x : count($arr)
}
```

for $x in 1 to 5
let $y := $x - 2
return $x * $y

| x |
|---|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |

| x | y |
|---|---|
| 1 | -1 |
| 2 | 0 |
| 3 | 1 |
| 4 | 2 |
| 5 | 3 |

| |
|---|
| -1 |
| 0 |
| 3 |
| 8 |
| 15 |

for $x in 1 to 5
group by $y := $x mod 2
return {
  "foo" : $y,
  "bar" : count($x)
}

| x |
|---|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |

| x | y |
|---|---|
| 1 | 1 |
| 2 | 0 |
| 3 | 1 |
| 4 | 0 |
| 5 | 1 |

*a sequence* / *foo*

| x | y |
|---|---|
| (1, 3, 5) | 1 |
| (2, 4) | 0 |

| |
|---|
| {"foo":1,"bar":3} |
| {"foo":0,"bar":2} |

for $x in 1 to 10
where $x - 2 gt 5
return $x * $x

*output stream*

| x | (where) |
|---|---------|
| 1 | false |
| 2 | false |
| 3 | false |
| 4 | false |
| 5 | false |
| 6 | false |
| 7 | false |
| 8 | true |
| 9 | true |
| 10 | true |

| x |
|---|
| 8 |
| 9 |
| 10 |

| |
|---|
| 64 |
| 81 |
| 100 |

# Execution

## Streaming

1
6
4
3
2

**ReturnIterator**

↑

Stream of tuples

**ForIterator**

$i: { "field" : 1 }
$i: { "field" : 6 }
$i: { "field" : 4 }
$i: { "field" : 3 }
$i: { "field" : 2 }

## Parallel

**ReturnIterator**

Stream of tuples

**ForIterator**

$i: { "field" : 1 }
$i: { "field" : 6 }
$i: { "field" : 4 }
$i: { "field" : 3 }
$i: { "field" : 2 }

1
6

$i: { "field" : 1 }
$i: { "field" : 6 }

4
3
2

$i: { "field" : 4 }
$i: { "field" : 3 }
$i: { "field" : 2 }

RDD

# Implimentation atop Spark



**AST**

**Expression Tree**

**Iterator Tree**

## AST

- count()
  - FLWOR
    - for
      - $i
      - json-lines()
        - "file.json"
    - return
      - .field
        - $i

## Expression Tree

- CountFunctionCall
  - FLWORExpression
    - ForClause
      - VariableExpression
      - JSONLinesFunctionCall
        - StringLiteralExpression
    - ReturnClause
      - ObjectLookupExpression
        - StringLiteralExpression
        - VariableExpression

## Iterator Tree

Sequence of items

- CountIterator
  - ReturnIterator
    - ObjectLookupExpression
      - VariableExpression
      - StringLiteralExpression
    - ForIterator
      - VariableName
      - JSONFileIterator
        - StringIterator

Stream of tuples

# Spark Implementation

_new col_

for $person in _json-file_("persons")

...

⬇

Implementation:
mapping of the RDD to a single-column DataFrame

_serialized_
_binary format_

| person |
| --- |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } |
| { "Name" : { "First" : "Srinivasa", "Last" : "Ramanujan" }, "Countries" : [ "IN", "UK" ] } |
| { "Name" : { "First" : "Kurt", "Last" : "Gödel" }, "Countries" : [ "CZ", "A", "US" ] } |
| { "Name" : { "First" : "John", "Last" : "Nash" }, "Countries" : "US" } |
| { "Name" : { "First" : "Alan", "Last" : "Turing" }, "Countries" : "UK" } |
| { "Name" : { "First" : "Maryam", "Last" : "Mirzakhani" }, "Countries" : [ "IR", "US" ] } |
| { "Name" : "Pythagoras", "Countries" : [ "GR" ] } |
| { "Name" : { "First" : "Nicolas", "Last" : "Bourbaki" }, "Number" : 9, "Countries" : null } |

**for** $person **in** *json-file*("persons")
**for** $country **in** *flatten*($person.Countries[])[$$ ne null]
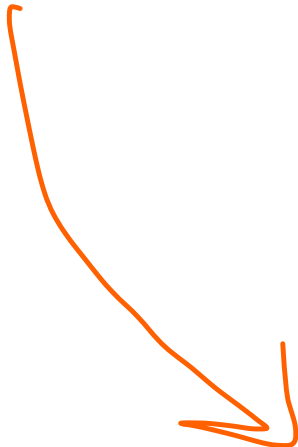...

*let*

*Cartesian Product*

Implementation:
Spark SQL + UDF

| person | country |
|---|---|
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | D |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | I |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | CH |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | A |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | BE |
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | US |
| { "Name" : { "First" : "Srinivasa", "Last" : "Ramanujan" }, "Countries" : [ "IN", "UK" ] } | IN |
| { "Name" : { "First" : "Srinivasa", "Last" : "Ramanujan" }, "Countries" : [ "IN", "UK" ] } | UK |
| { "Name" : { "First" : "Kurt", "Last" : "Gödel" }, "Countries" : [ "CZ", "A", "US" ] } | CZ |
| { "Name" : { "First" : "Kurt", "Last" : "Gödel" }, "Countries" : [ "CZ", "A", "US" ] } | A |

```
for $person in json-file("persons")
for $country in flatten($person.Countries[])[$$ ne null]
where $country eq "US"
...
```

Implementation:
Spark SQL + UDF

| person | country |
|---|---|
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | US |
| { "Name" : { "First" : "Kurt", "Last" : "Gödel" }, "Countries" : [ "CZ", "A", "US" ] } | US |
| { "Name" : { "First" : "John", "Last" : "Nash" }, "Countries" : "US" } | US |
| { "Name" : { "First" : "Maryam", "Last" : "Mirzakhani" }, "Countries" : [ "IR", "US" ] } | US |

```
for $person in json-file("persons")
for $country in flatten($person.Countries[])[$$ ne null]
where $country eq "US"
let $first-name eq $person.Name.First
...
```

Implementation:
Spark SQL + UDF

| person | country | first-name |
|---|---|---|
| { "Name" : { "First" : "Albert", "Last" : "Einstein" }, "Countries" : [ "D", "I", "CH", "A", "BE", "US" ] } | US | Albert |
| { "Name" : { "First" : "Kurt", "Last" : "Gödel" }, "Countries" : [ "CZ", "A", "US" ] } | US | Kurt |
| { "Name" : { "First" : "John", "Last" : "Nash" }, "Countries" : "US" } | US | John |
| { "Name" : { "First" : "Maryam", "Last" : "Mirzakhani" }, "Countries" : [ "IR", "US" ] } | US | Maryam |

```
for $person in json-file("persons")
for $country in flatten($person.Countries[])[$$ ne null]
where $country eq "US"
let $first-name eq $person.Name.First
return $first-name
```

Implementation:
Mapping back to an RDD

emit
one by one

"Albert"
"Kurt"
"John"
"Maryam"

# Final Tips

## ① Iterate over arrays:

```
1  %%jsoniq
2  for $i in [1, 2, 3]
3  return $i
```

Took: 0.03023815155029297 ms
[1, 2, 3]

→ a single item

```
1  %%jsoniq
2  for $i in [1, 2, 3][]
3  return $i
```

Took: 0.03351902961730957 ms
1
2
3

3 items (ints)

## ② Implicit existential Quantification:

```
1  %%jsoniq
2  10 eq 5
```

Took: 0.02959108352611328 ms
false

```
1  %%jsoniq
2  1 to 10 = 5
```

Took: 0.04393911361694336 ms
true

[1, 10]        [9, 20]

```
1  %%jsoniq
2  1 to 10 > 9 to 20
```

Took: 0.03199887275695801 ms
true

```
1  %%jsoniq
2  1 to 10 > 11 to 20
```

Took: 0.04635882377624512 ms
false