

# Incremental Data Transformations on Wide-Column Stores with NôtaQL

2014-06-29

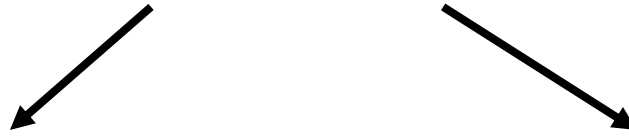
**M. Sc. Johannes Schildgen**

schildgen@cs.uni-kl.de



"A DBA walks into a NoSQL bar, but turns and leaves because he couldn't find a table"

# Column Families

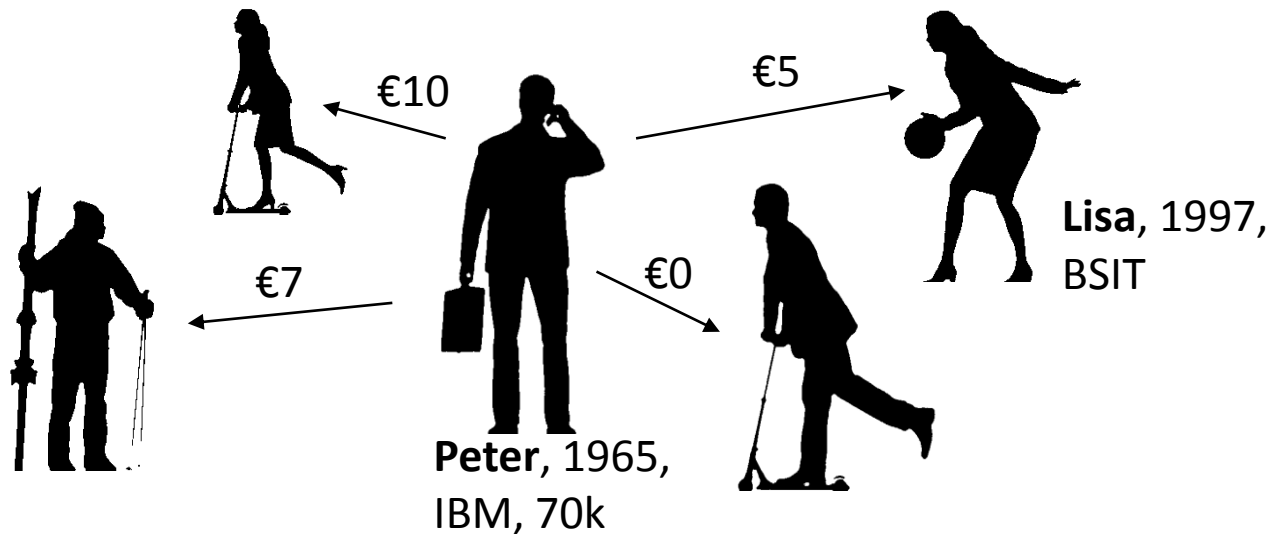


<u>RowId</u>	info	children

# Column Families



<u>RowId</u>	info			children			
Peter	<b>born</b>	<b>cmpny</b>	<b>salary</b>	<b>Lisa</b>	<b>Carl</b>	<b>Susi</b>	<b>Toni</b>
	1965	IBM	70k	€5	€0	€10	€7
Lisa	<b>born</b>	<b>school</b>					
	1997	BSIT					



# HBase API

```
put `pers`, `Carl`,  
  `info:born`, `1982`
```

```
put `pers`, `Carl`,  
  `info:school`, `BSIT`
```

```
put `pers`, `Carl`,  
  `info:school`, `BUIIT`
```

```
get `pers`, `Carl`
```

# Jaspersoft HBase QL



$\sigma_{\text{school}='BSIT'}\text{pers}$

```
{
  "tableName": "pers",
  "deserializerClass": "com.jaspersoft...DefaultDeserializer",
  "filter": {
    "SingleColumnValueFilter": {
      "family": „info“,
      "qualifier": „school“,
      "compareOp": "EQUAL",
      "comparator": { "SubstringComparator": { "substr":
        „BSIT" } }
    }
  }
}
```

# Phoenix

$\sigma_{\text{school}='BSIT'}$  pers

```
SELECT * FROM pers WHERE school = 'BSIT'
```

„Parent of  
each person?“

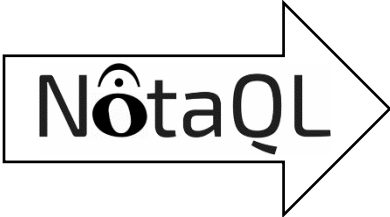
Nôta QL





Input Cell

	Column
RowID	Value

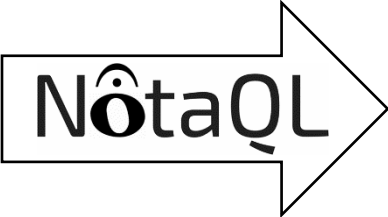


Output Cell

	Column
RowID	Value

Input Cell

	<b>_c</b>
<b>_r</b>	<b>_v</b>



Output Cell

	<b>_c</b>
<b>_r</b>	<b>_v</b>

$\pi_{born}$  pers

Input Cell

	<b>_c</b> <i>born</i>
<b>_r</b> <i>Lisa</i>	<b>_v</b> <i>1997</i>

Output Cell

	<b>_c</b> <i>born</i>
<b>_r</b> <i>Lisa</i>	<b>_v</b> <i>1997</i>

**NotaQL**

```
OUT._r <- IN._r,  
OUT.born <- IN.born;
```

$\pi_{born, school}$  pers

Input Cell

	<b>_c</b>
	<b><i>born</i></b>
<b>_r</b>	<b>_v</b>
<i>Lisa</i>	<i>1997</i>

Output Cell

	<b>_c</b>
	<b><i>born</i></b>
<b>_r</b>	<b>_v</b>
<i>Lisa</i>	<i>1997</i>

**NotaQL**

```
OUT._r <- IN._r,  
OUT.born <- IN.born,  
OUT.school <- IN.school;
```

$\sigma_{\text{school}='BSIT'}$  pers

Input Cell

	<b>_c</b>
	<i>born</i>
<b>_r</b>	<b>_v</b>
<i>Lisa</i>	<i>1997</i>

Output Cell

	<b>_c</b>
	<i>born</i>
<b>_r</b>	<b>_v</b>
<i>Lisa</i>	<i>1997</i>

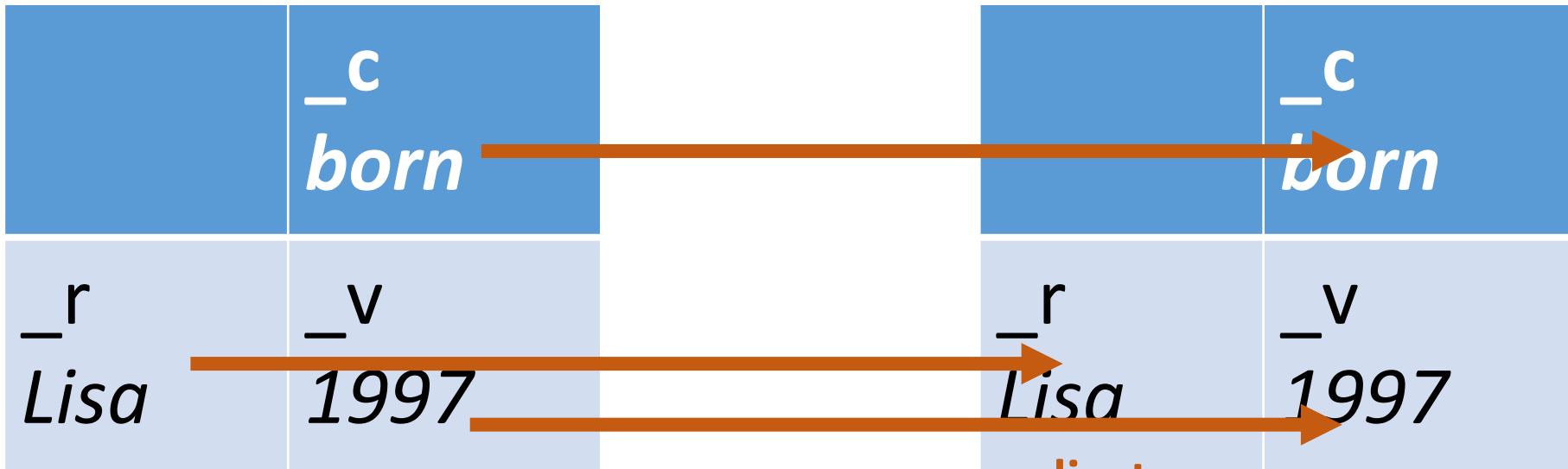
**NotaQL**

```
OUT._r <- IN._r,  
OUT.$(IN._c) <- IN._v;
```

$\sigma_{\text{school}='BSIT'}$  pers

Input Cell

Output Cell



row predicate

NotaQL

```
IN-FILTER: school='BSIT',  
OUT._r <- IN._r,  
OUT.$(IN._c) <- IN._v;
```

That was:

**Selection and Projection**



Now:

**Grouping**

Salary sum of  
each company.

Input Cell

	<b>_c</b> <i>cmpny</i>
<b>_r</b> <i>Peter</i>	<b>_v</b> <i>IBM</i>

Output Cell

	<b>_c</b> <i>salsum</i>
<b>_r</b> <i>IBM</i>	<b>_v</b> <i>645k</i>

NotaQL

```
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```

RowId	info		
Eve	born	cmpny	salary
	1965	IBM	70k
Carl	born	cmpny	job
	1966	IBM	intern
Julia	born	cmpny	salary
	1967	IBM	80k
Lisa	born	school	salary
	1997	BSIT	1k



	salsum
IBM	70k



	salsum
IBM	150k



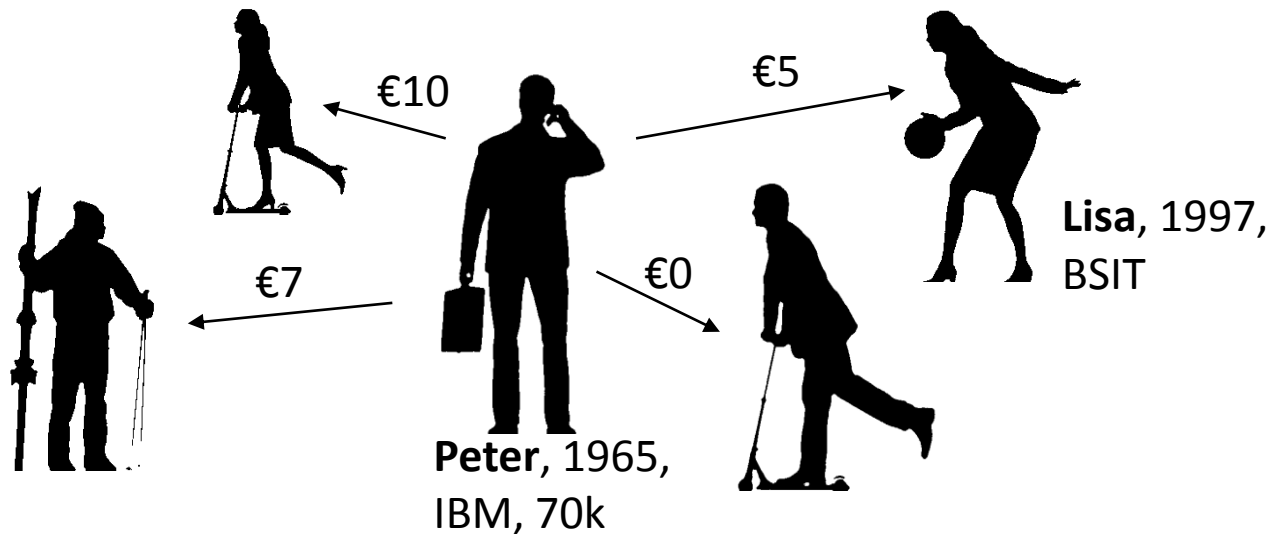
	salsum
IBM	80k

NotaQL

```
OUT._r <- IN.cmpny,
OUT.salsum <- SUM(IN.salary):
```

# Advanced Transformations: **More Filters**

RowId	info			children			
Peter	born	cmpny	salary	Lisa	Carl	Susi	Toni
	1965	IBM	70k	€5	€0	€10	€7
Lisa	born	school					
	1997	BSIT					



<u>RowId</u>	info			children			
Peter	born	cmpny	salary	Lisa	Carl	Susi	Toni
	1965	IBM	70k	€5	€0	€10	€7
Lisa	born	school					
	1997	BSIT					

**NotaQL** `OUT._r <- IN._r,`  
`OUT.$ (IN._c) <- IN._v;`

<u>RowId</u>	info			children			
Peter	<b>born</b>	<b>cmpny</b>	<b>salary</b>	<b>Lisa</b>	<b>Carl</b>	<b>Susi</b>	<b>Toni</b>
	1965	IBM	70k	€5	€0	€10	€7
Lisa	<b>born</b>	<b>school</b>					
	1997	BSIT					

**NotaQL**

```

IN-FILTER: COL_COUNT(children) > 0
OUT._r <- IN._r,
OUT.$(IN._c) <- IN._v;

```

<u>RowId</u>	info			children			
Peter	born	cmpry	salary	Lisa	Carl	Susi	Toni
	1965	IBM	70k	€5	€0	€10	€7
Lisa	born	school					
	1997	BSIT					

**NotaQL**

```

IN-FILTER: COL_COUNT(children) > 0
OUT._r <- IN._r,
OUT.$(IN.children._c) <- IN._v;

```



<u>RowId</u>	info			children			
Peter	born	cmpry	salary	Lisa	Carl	Susi	Toni
	1965	IBM	70k	€5	€0	€10	€7
Lisa	born	school					
	1997	BSIT					

NotaQL

```

IN-FILTER: COL_COUNT(children) > 0
OUT._r <- IN._r,
OUT.$(IN.children._c? (@>5))
                                     cell predicate
                                     <- IN._v;

```

<u>RowId</u>	info			children			
Peter	born	cmpry	salary	Lisa	Carl	Susi	Toni
	1965	IBM	70k	€5	€0	€10	€7
Lisa	born	school					
	1997	BSIT					

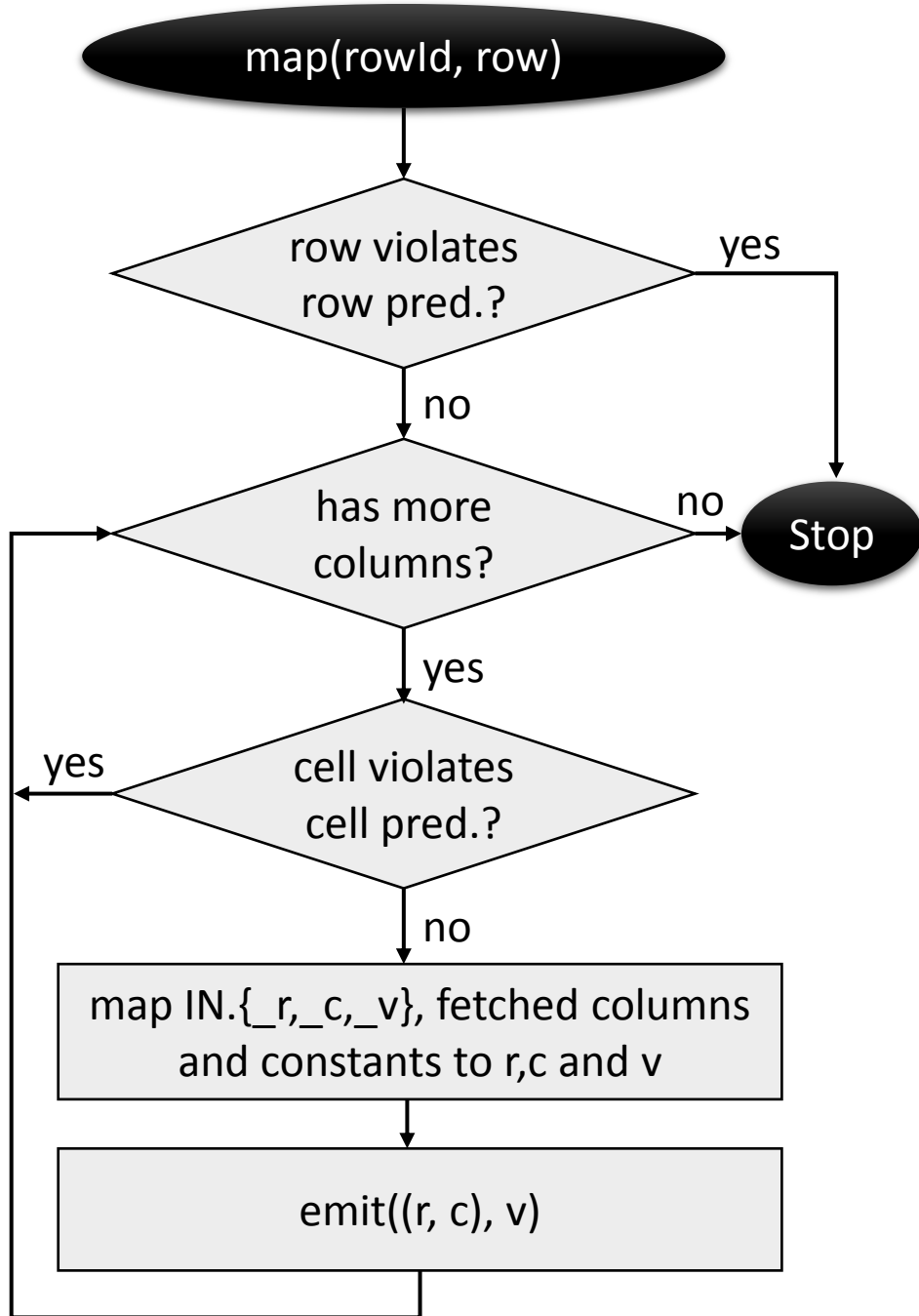
NôtaQL

```

IN-FILTER: COL_COUNT(children) > 0
OUT._r <- IN._r,
OUT.$(IN.children._c?(!Carl))
                                     cell predicate
                                     <- IN._v;

```

# NotaQL Transformation Platform: **MapReduce**

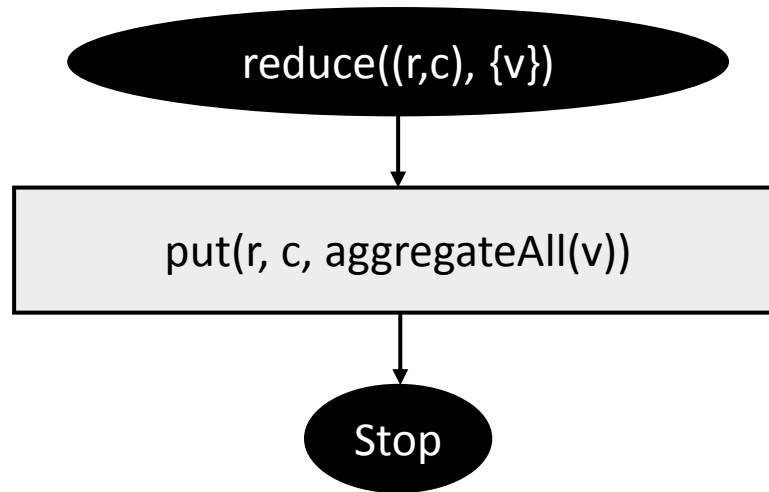


RowId	info		
Peter	born	cmpny	salary
	1965	IBM	70k

*Note: The entire table above is crossed out with a red X.*

	salsum
IBM	70k

`((IBM, salsum), 70k)`



`((IBM, salsum), {70k, 80k, 10k})`

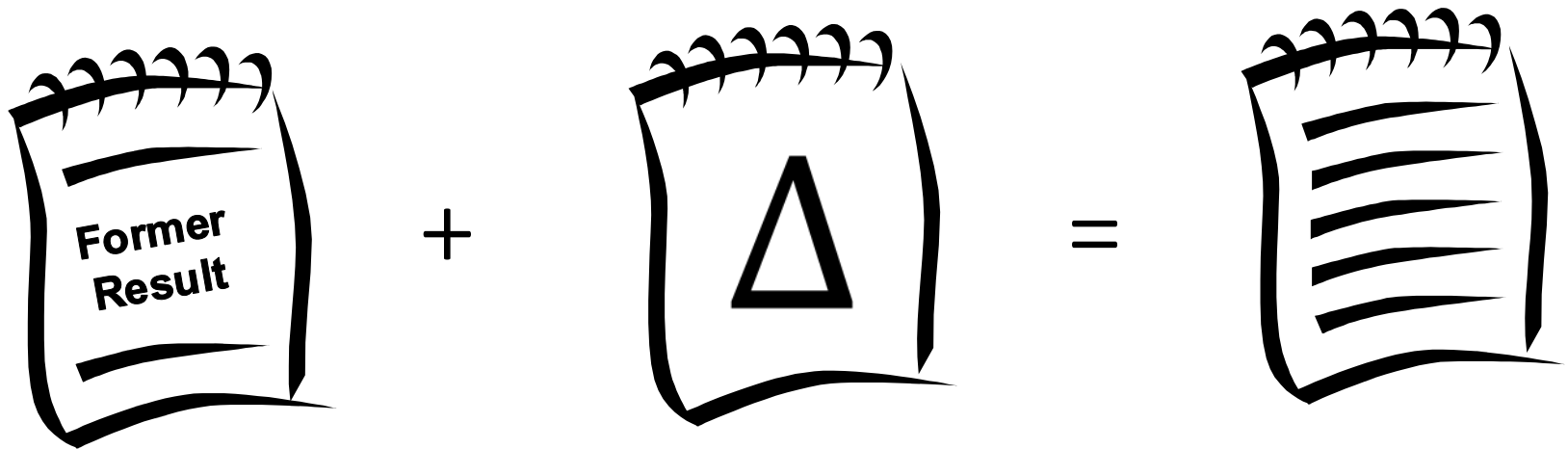
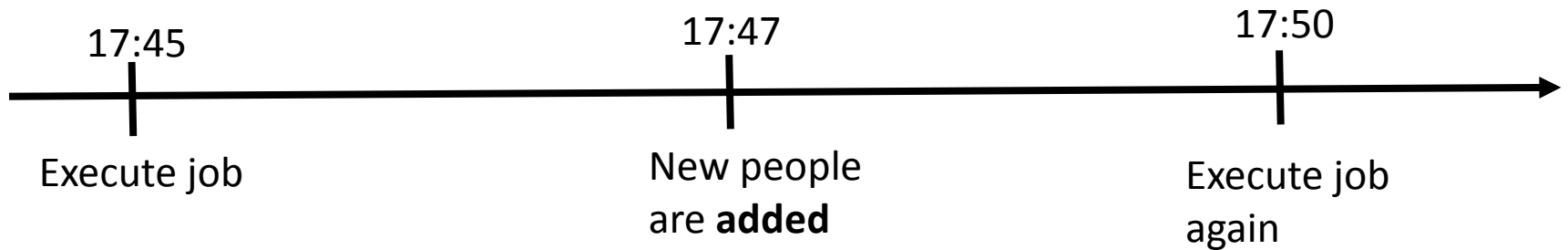
`((IBM, salsum), 160k)`

# Incremental Transformations: **Self-Maintainability**

Salary sum of all people born before 1980 per company.

NotaQL

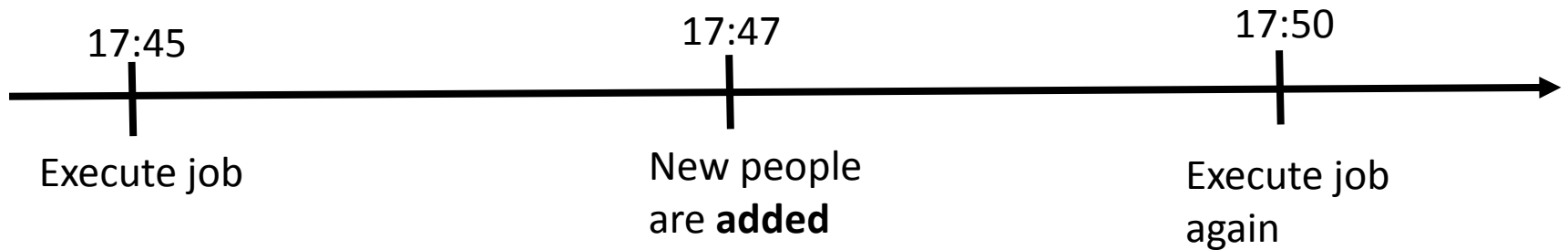
```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



Salary sum of all people born before 1980 per company.

NotaQL

```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



	salsum		born	cmpny	salary		salsum	
IBM	150k	Melissa	1989	IBM	50k	+	IBM	230k
		Nora	1977	IBM	80k			

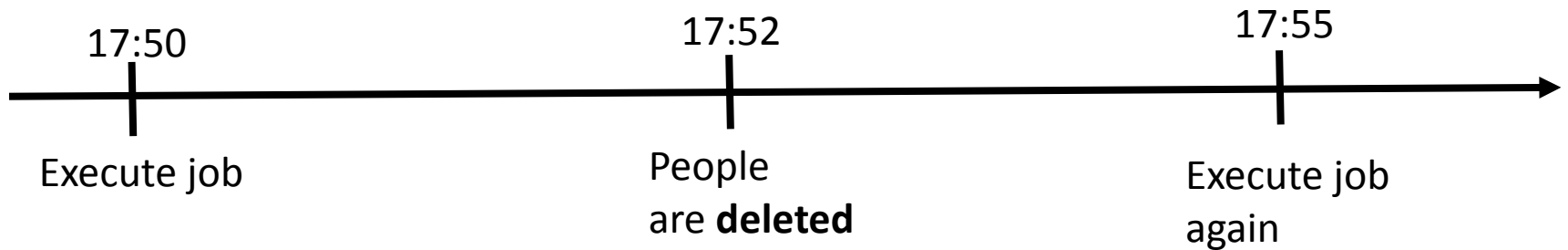
Diagram illustrating the state of data before and after a query execution. The initial state shows a table with one row: IBM, 150k. The second state shows a table with two rows: Melissa (born 1989, IBM, 50k) and Nora (born 1977, IBM, 80k). The final state shows the result of the query: IBM, 230k. Arrows indicate the flow of data from the initial state to the second state, and from the second state to the final state. The final state is the sum of the initial state and the new data added.



Salary sum of all people born before 1980 per company.

NotaQL

```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



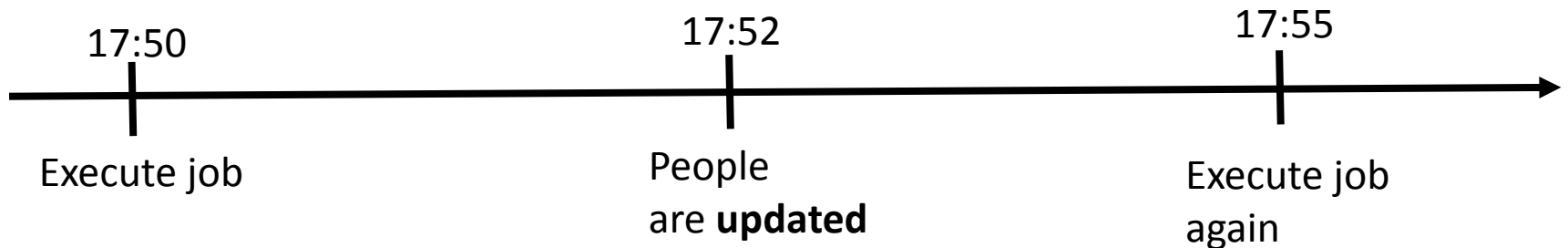
	salsum		born	cmpny	salary		salsum
IBM	230k	Melissa	1989	IBM	50k	- IBM	150k
		Nora	1977	IBM	80k		



Salary sum of all people born before 1980 per company.

NotaQL

```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



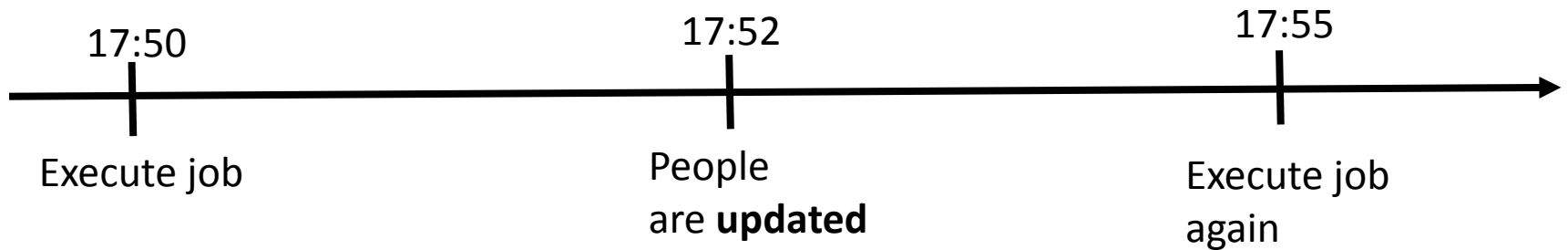
	born	cmpny	salary
Peter	<del>1965</del> 1990	IBM	70k

→ == 70k

Salary sum of all people born before 1980 per company.

NotaQL

```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



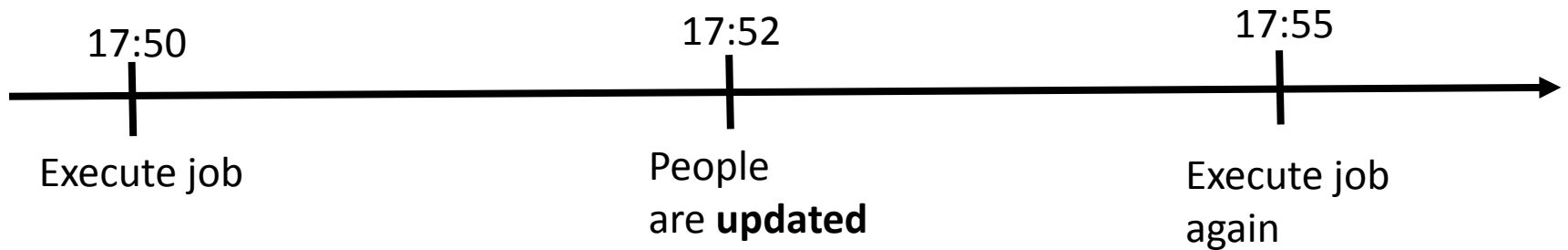
	born	cmpny	salary
Peter	<del>1965</del> <del>1990</del> 1965	IBM	70k

➔ += 70k

Salary sum of all people born before 1980 per company.

NotaQL

```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```



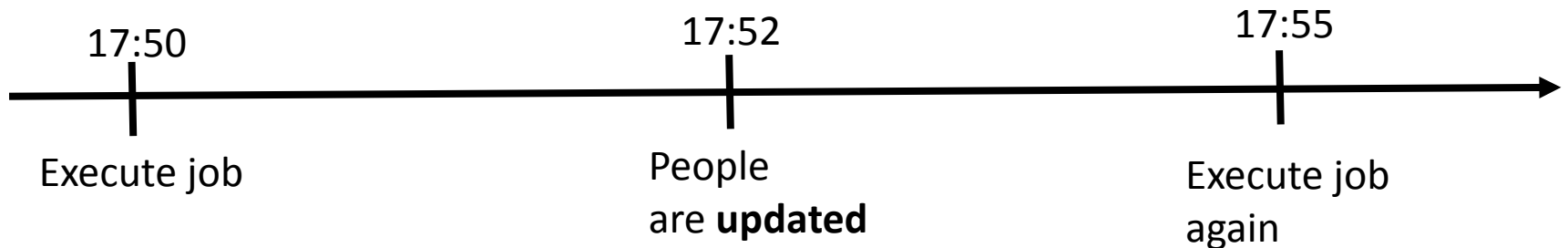
	born	cmpny	salary
Peter	1965	IBM	<del>70k</del> 75k

➔ += 75k-70k

Salary sum of all people born before 1980 per company.

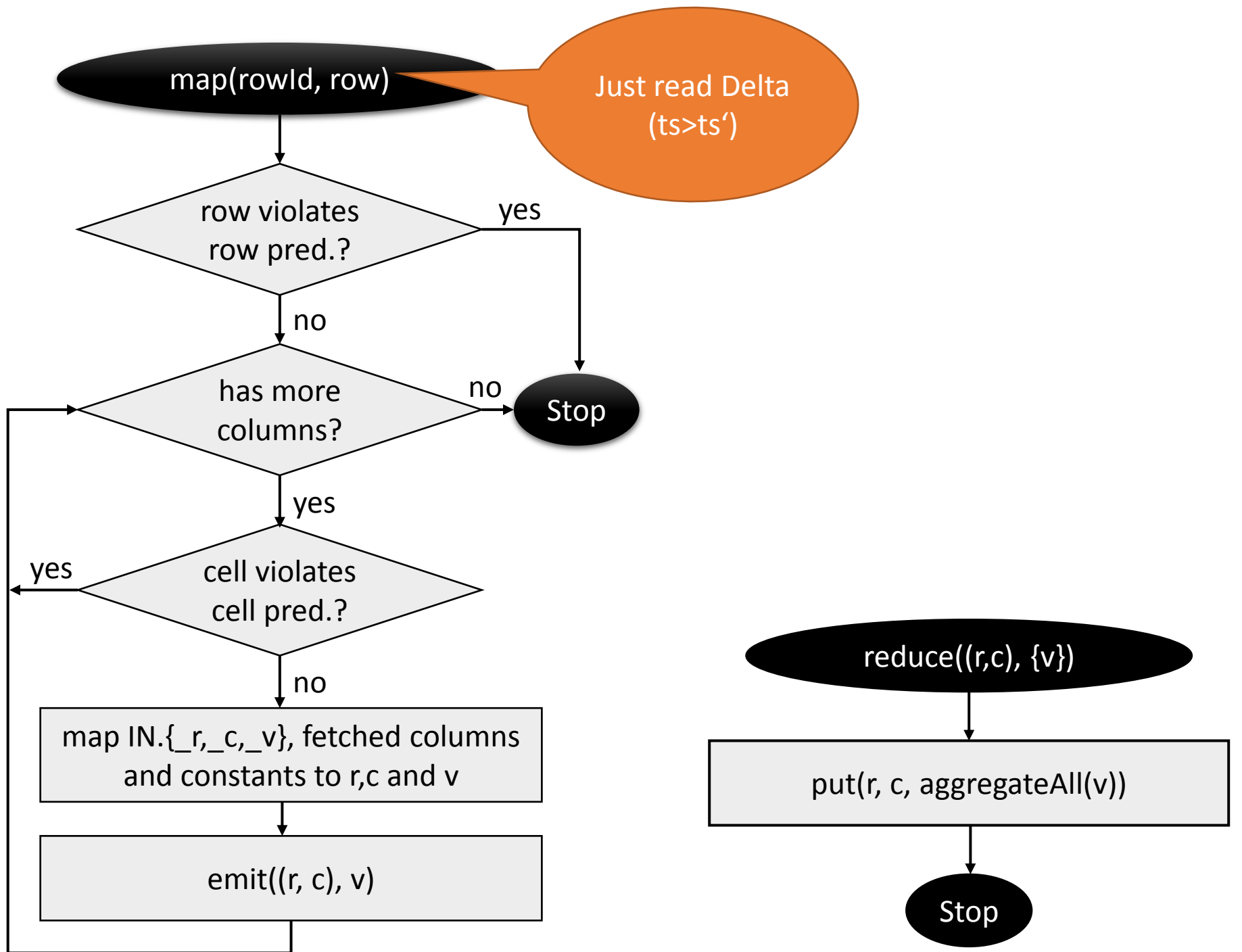
NotaQL

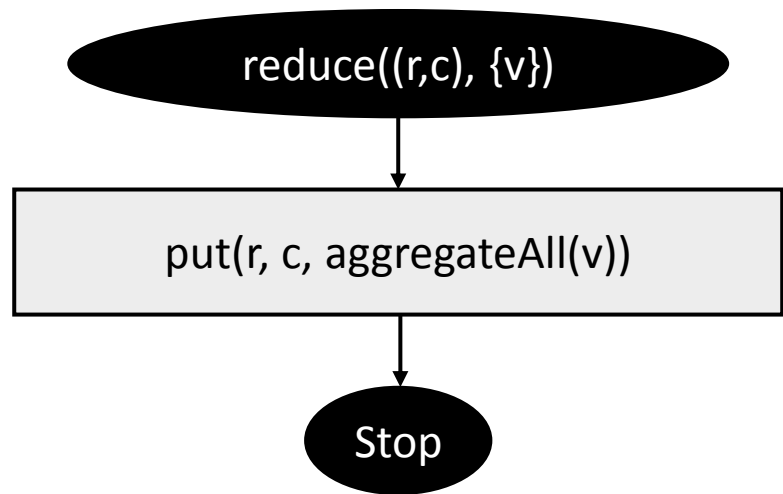
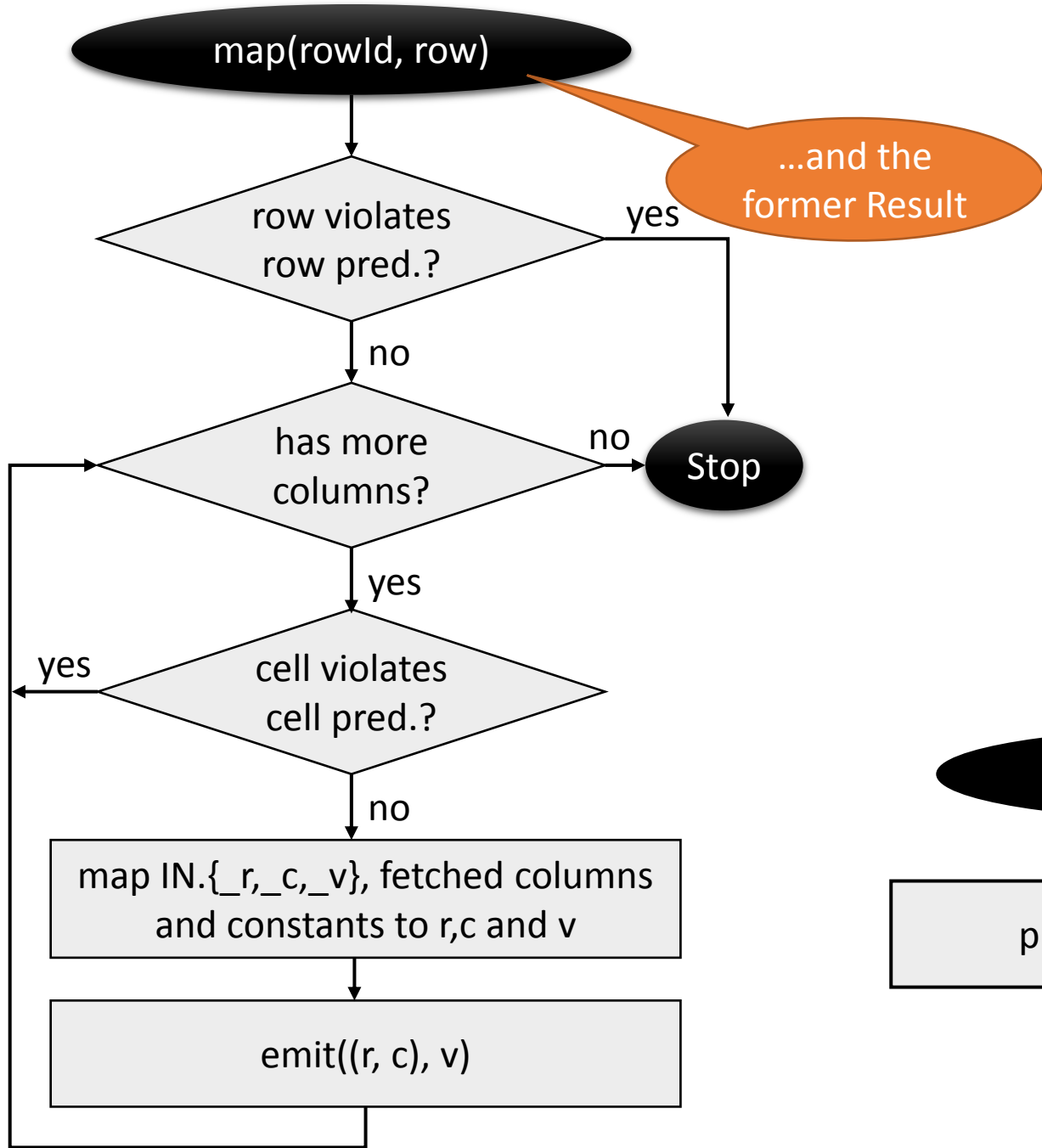
```
IN-FILTER: born<1980,  
OUT._r <- IN.cmpny,  
OUT.salsum <- SUM(IN.salary);
```

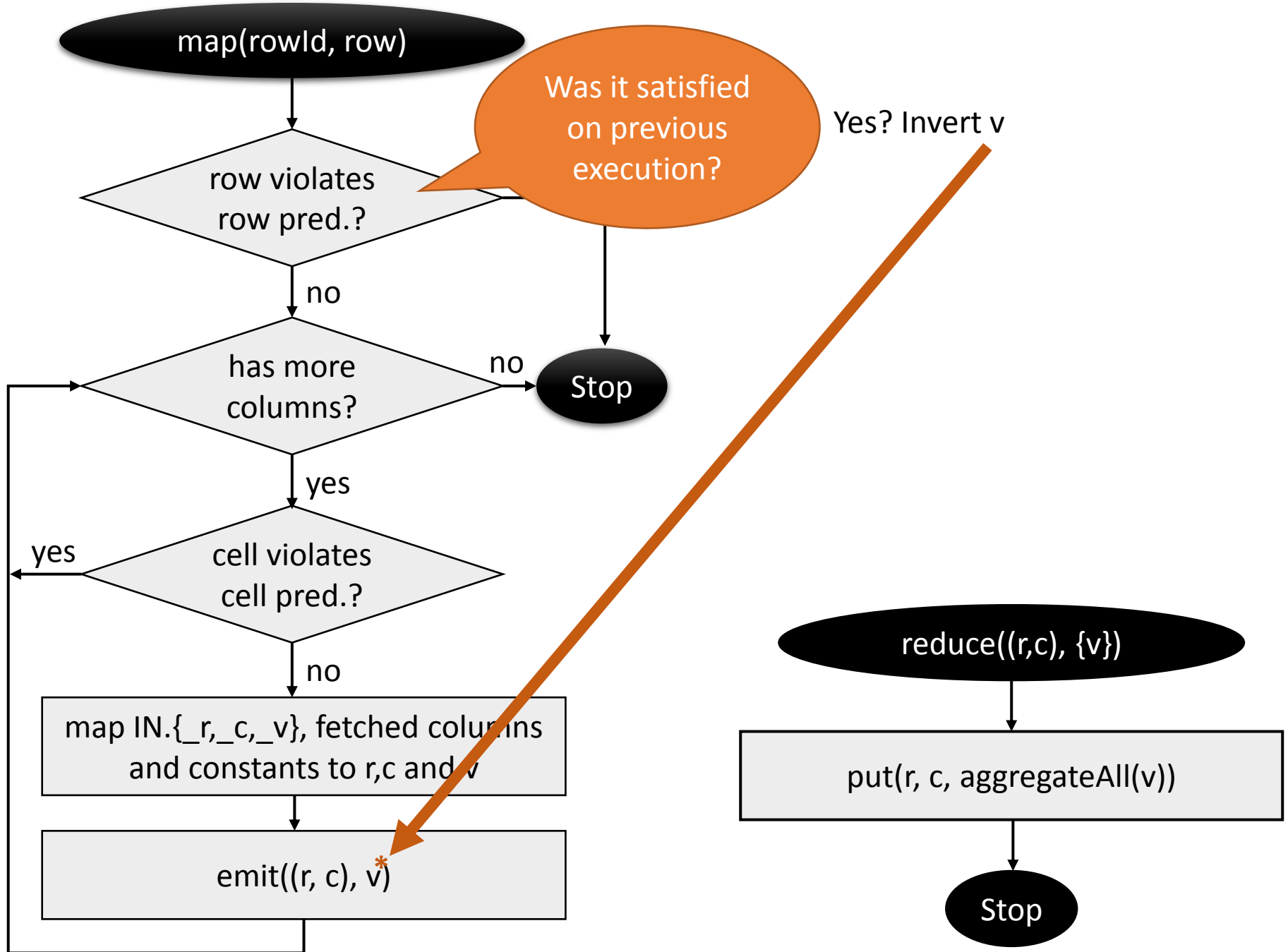


	born	cmpny	salary
Peter	1965	<del>IBM</del> SAP	70k

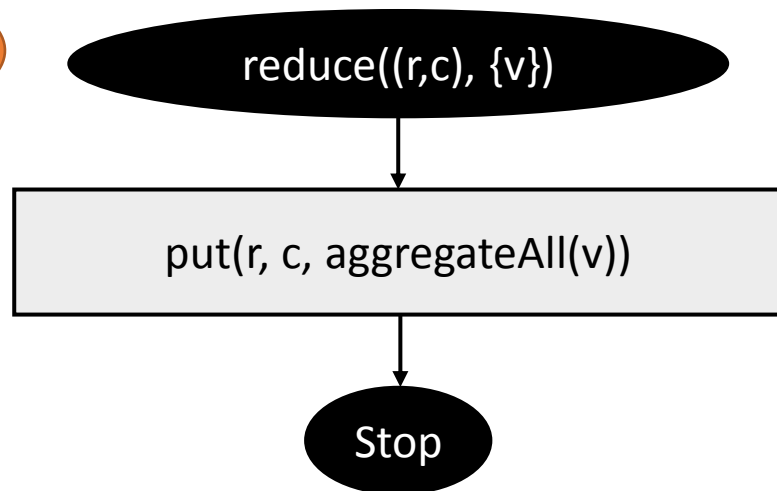
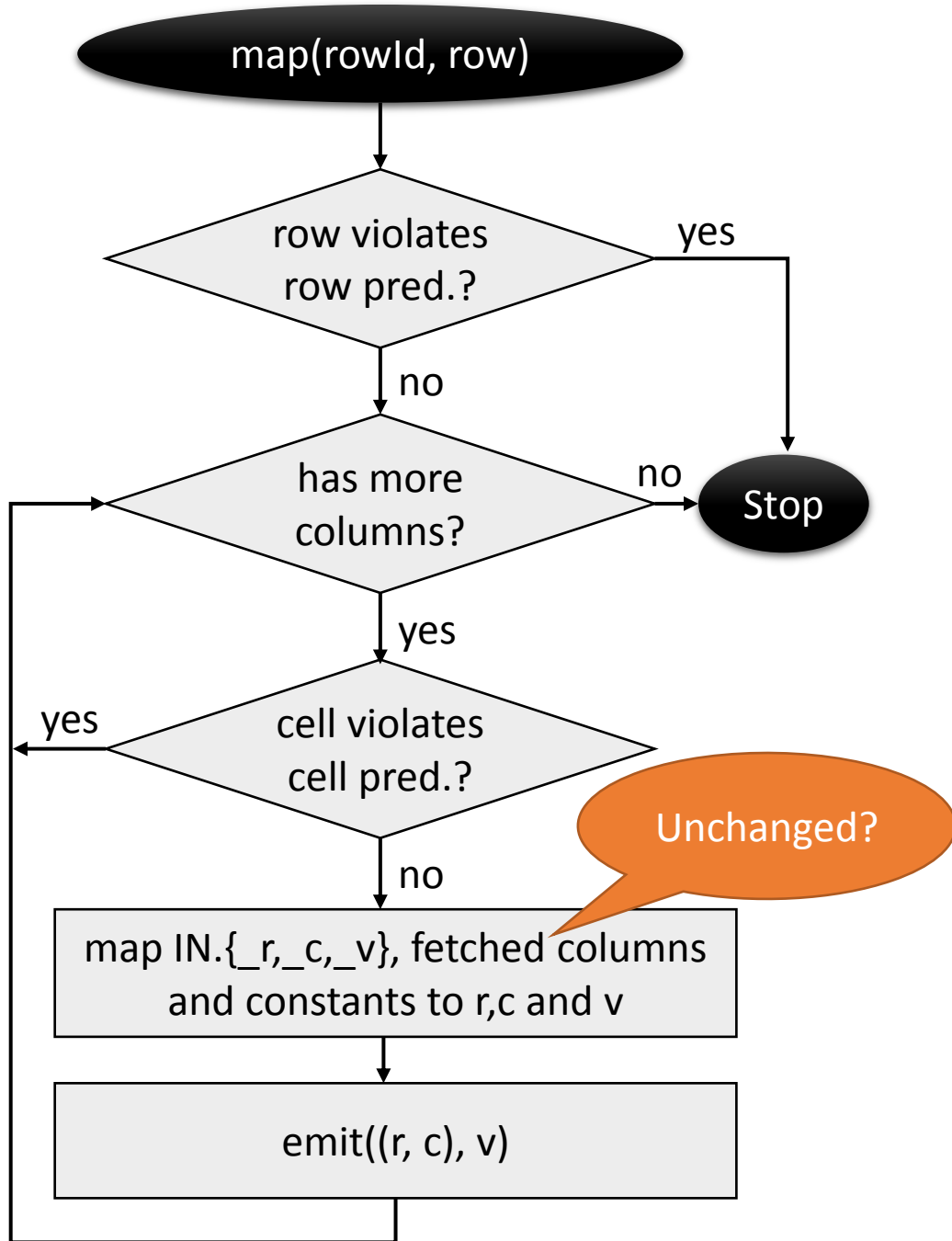
→ == 70k





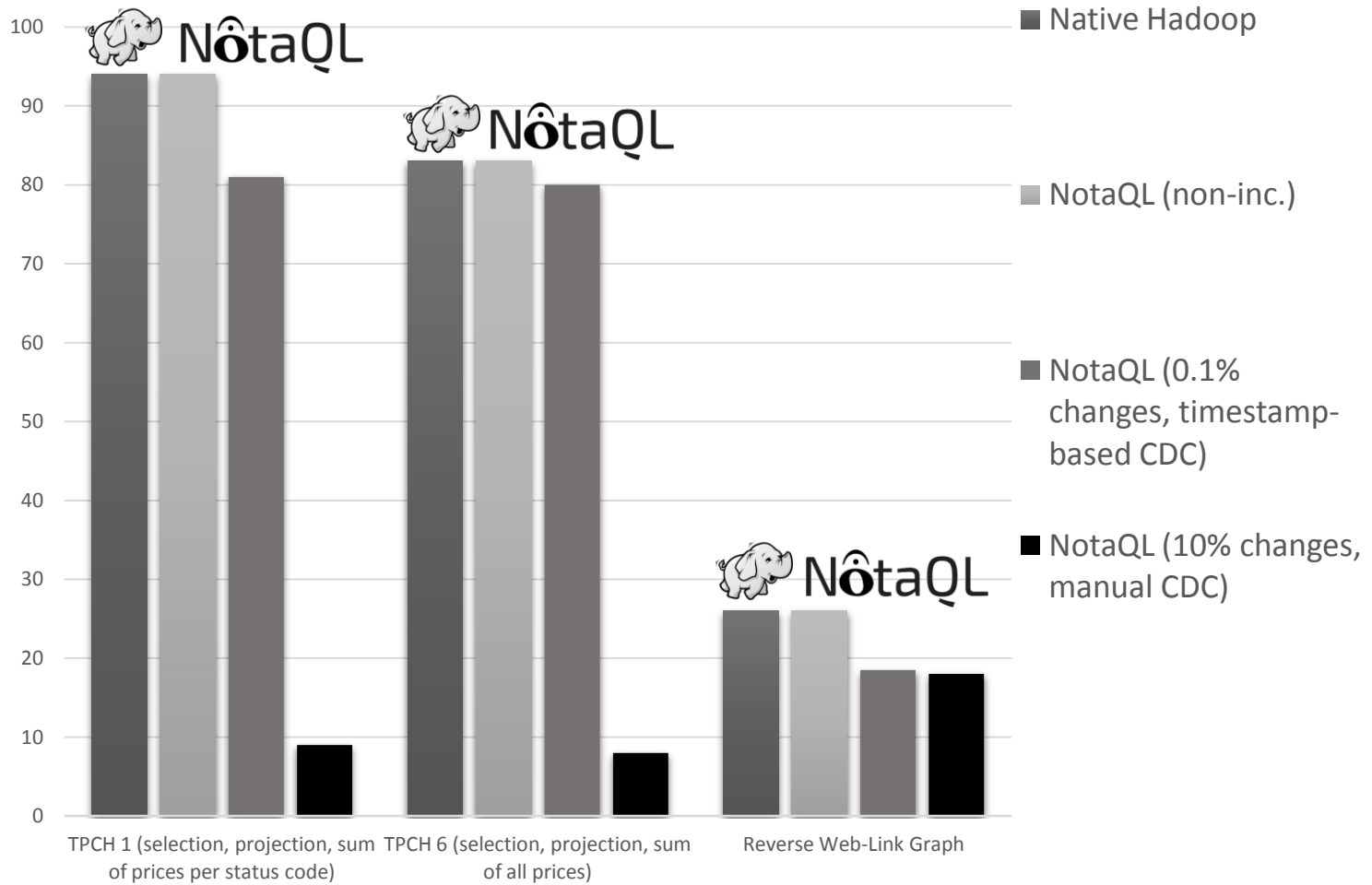






Evaluation:

**Performance**



# Conclusion

# NôtaQL

Incremental!

# ≠SQL

- Selection, Projection
- Grouping, Aggregation
- Schema-Flexible
- Horizontal Aggregation
- Metadata↔Data
- Graph Processing
- Text Processing



# Thank you!

$\langle \text{NotaQL} \rangle \models [\text{IN-FILTER: } \langle \text{predicate} \rangle, ] \langle \text{rowspec} \rangle, \langle \text{cellspec} \rangle (, \langle \text{cellspec} \rangle) * [; ]$

$\langle \text{rowspec} \rangle \models \text{OUT.}_{-r} <- \langle \text{vdata} \rangle$

$\langle \text{cellspec} \rangle \models \text{OUT.}(\langle \text{colname} \rangle \mid \$(\langle \text{input} \rangle)) <- (\langle \text{vdata} \rangle \mid \langle \text{aggfun} \rangle(\langle \text{vdata} \rangle))$

$\langle \text{input} \rangle \models (\text{IN.}_{-r} \mid \text{IN.}[\langle \text{colfamily} \rangle :](_{-c} \mid _{v}) \mid \text{IN.}\langle \text{colname} \rangle)[?(\langle \text{predicate} \rangle)]$

$\langle \text{vdata} \rangle \models \langle \text{input} \rangle \mid \langle \text{const} \rangle \mid \langle \text{vdata} \rangle (+ \mid - \mid * \mid /) \langle \text{vdata} \rangle$

$\langle \text{aggfun} \rangle \models \text{COUNT} \mid \text{SUM} \mid \text{MIN} \mid \text{MAX} \mid \text{AVG}$

$\langle \text{const} \rangle \models '(A \dots Z \mid a \dots z \mid 0 \dots 9) + ' \mid (0 \dots 9) +$

$\langle \text{colname} \rangle \models [\langle \text{colfamily} \rangle :](A \dots Z \mid a \dots z \mid 0 \dots 9) +$

$\langle \text{colfamily} \rangle \models (A \dots Z \mid a \dots z \mid 0 \dots 9) +$

$\langle \text{predicate} \rangle \models (\langle \text{colname} \rangle \mid @ \mid \text{col\_count}([\langle \text{colfamily} \rangle]))[\langle \text{op} \rangle(\langle \text{colname} \rangle \mid \langle \text{const} \rangle)]$   
 $\mid (\text{NOT} \mid !)\langle \text{predicate} \rangle \mid \langle \text{predicate} \rangle(\text{AND} \mid \text{OR})\langle \text{predicate} \rangle$

$\langle \text{op} \rangle \models = \mid ! = \mid < \mid < = \mid > \mid > =$