

Scala **Integrated** Query

Christopher Vogt



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

RWTHAACHEN
UNIVERSITY

Scala Integrated Query

short SIQ

- a type safe database query library
like Scala Query *by Stefan Zeiger* and SQueryl
- compiles a subset of Scala to SQL

Feature Comparison

	SQL	Scala Query	SQueryl	Microsoft's LINQ-to-SQL on .NET	Scala Integrated Query
Queries	✓	✓	✓	✓	✓
Type Safety		✓	✓	✓	✓
List Order					✓
Avalanche Safe Nesting					✓
...

AVALANCHE SAFE NESTING

Example (in pseudo code)

```
for nation <- nations  
  ( nation, customers of this nation )
```

```
Result: List[ ( Nation, List[Customer] ) ]  
List(  
  ("GERMANY", List("Martin","Chris" ) ),  
  ("ARGENTINA", List("Miguel",...) )  
)
```

How many SQL Queries?

```
for nation <- nations  
  ( nation, customers of this nation )
```

Depends on your code:

	SQL	Scala Query	SQueryl	Microsoft's LINQ-to-SQL on .NET	Scala Integrated Query
straight- forward	$n + 1$	$n + 1$	$n + 1$	$n + 1$	2
manually tweaked	2	2	2	2	

$n = \text{nations.length}$

Avalanche
Safe

Example: Scala Query

```
val nations = ( for( n <- nation ) yield n ).list

for( n <- nations ) yield {
  val cs =
    ( for( c <- customer;
      if c.nationkey === n.nationkey
    ) yield c
  ).list
  (n, cs)
} // : List[ ( Nation, List[Customer] ) ]
```

SQL Mapping in Scala Query

JVM

1 SQL Query

Example: 200 results

DBMS

```
val nations = ( for( n <- nation ) yield n ).list
```

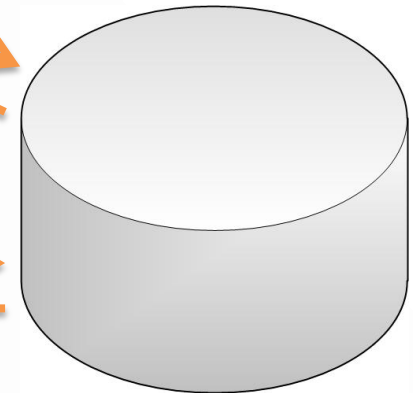
```
for( n <- nations ) yield {
```

```
  val cs =
```

```
    ( for( c <- customer;  
          if c.nationkey === n.nationkey  
        ) yield c
```

```
    (n, cs)
```

```
}
```



200 SQL queries
An „avalanche“

SQL does this better

Scala Integrated Query: No Avalanche

JVM

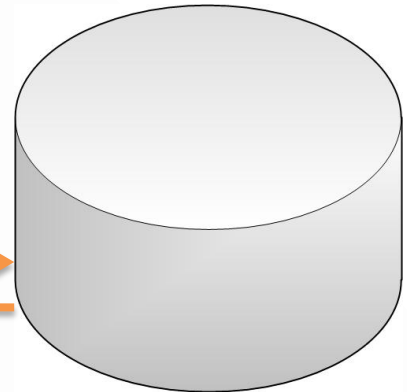
```
val nations = for( n <- nation ) yield n

(for( n <- nations ) yield {
  val cs =
    for( c <- customer;
      if c.nationkey == n.nationkey
    ) yield c
  (n, cs)
})
```

.fromdb

Mapped to SQL as a whole

DBMS



2 SQL queries
No „avalanche“

SQL Mapping in Scala Query

JVM

```
val nations = ( for( n <- nation ) yield n ).list
```

```
for( n <- nations ) yield {
```

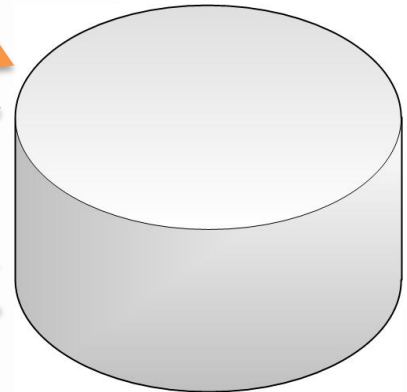
```
  val cs =
```

```
    ( for( c <- customer;  
          if c.nationkey === n.nationkey  
        ) yield c  
    ).list
```

```
  (n, cs)
```

```
}
```

DBMS



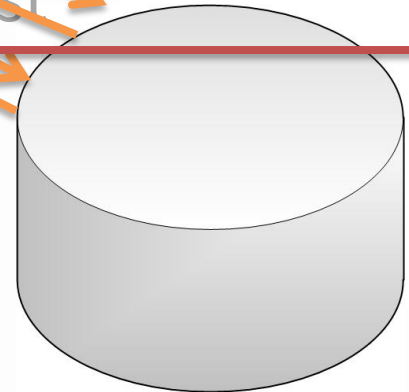
Avalanche avoidance in Scala Query

JVM

```
val nations = ( for( n <- nation ) yield n ).list
val customers = ( for( c <- customer ) yield c ).list
for( n <- nations ) yield {
  val cs =
    for( c <- customers;
      if c.nationkey == n.nationkey
    ) yield c
  (n, cs)
}
```

2 SQL queries
No „avalanche“

DBMS



In SQL you get this for free!

Scala Integrated Query: No Avalanche

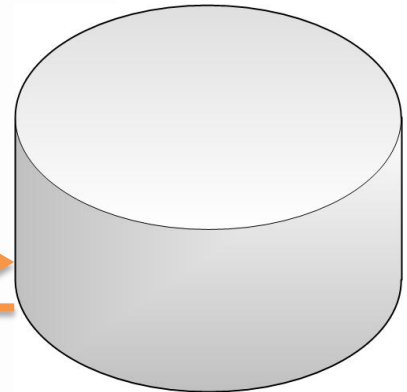
JVM

```
val nations = for( n <- nation ) yield n

(for( n <- nations ) yield {
  val cs =
    for( c <- customer;
      if c.nationkey == n.nationkey
    ) yield c
  (n, cs)
})
```

.fromdb

DBMS



Avalanche Safety
for free

SQ: Let's simplify

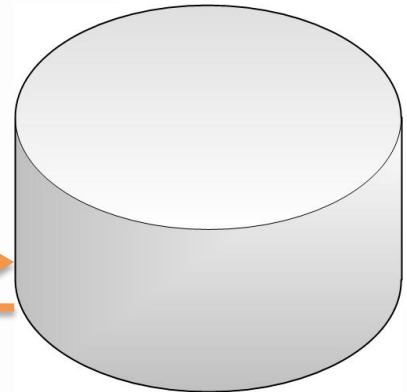
JVM

```
val nations = for( n <- nation ) yield n
```

```
(for( n <- nations ) yield {  
  val cs =  
    for( c <- customer;  
      if c.nationkey == n.nationkey  
    ) yield c  
  (n, cs)  
})
```

.fromdb

DBMS



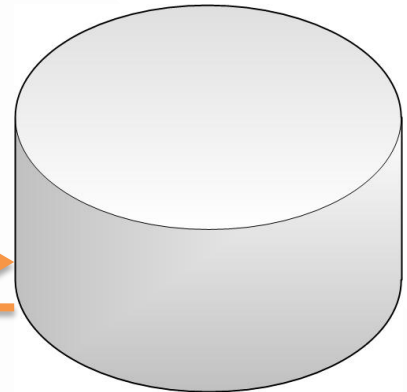
SQ: Let's simplify

JVM

```
(for( n <- nation ) yield {  
  val cs =  
    for( c <- customer;  
      if c.nationkey == n.nationkey  
    ) yield c  
  (n, cs)  
})
```

.fromdb

DBMS



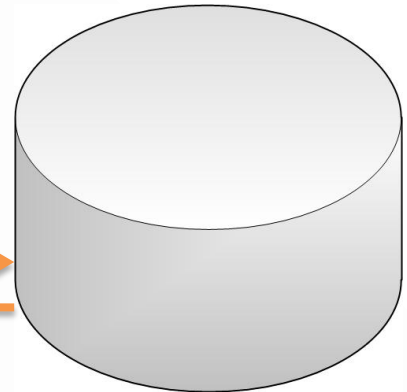
SQ: Let's simplify

JVM

```
nation.map( n =>
  ( n, customer.withFilter(
    _.nationkey == n.nationkey
  ))
)
```

.fromdb

DBMS



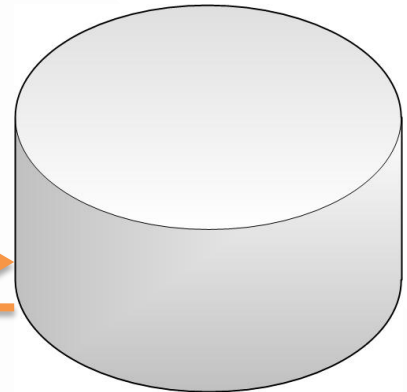
SQ: Query Decomposition

JVM

```
def customers( nationkey : Rep[Int] ) =  
  customer.withFilter(  
    _.nationkey == nationkey  
  )  
  
nation.map( n =>  
  (n, customers(n.nationkey)) )  
)
```

.fromdb

DBMS



still only
2 SQL queries

SQ: Data Model

- Seen so far
 - List[(Nation, List[Customer])]
- Supported: **Arbitrarily nested Lists and Tuples**
- you can mix database tables and Iterables

How many SQL queries exactly?

- Depends only on **result type**
- **Lower Bound**
Number of list constructors
- **Upper Bound**
Number of all type constructors excluding tuples

- Example

List[(A, List[B])]

count
word „List“

2 ≤ number of SQL queries ≤ 4

count
all words

HOW DOES AVALANCHE SAFETY WORK?

SIQ implements Ferry



- nested comprehensions \rightarrow relational algebra
- theory developed by Tom Schreiber
University of Tübingen
- (Great stuff, Tom 😊!)



Ferry SQL Mapping

for nation <- **nations**

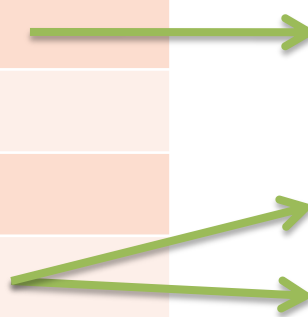
(nation, **customers** of this nation)

Query 1

name	inner
ARGENTINA	1
BRAZIL	2
...	...
GERMANY	7
...	...

Query 2

inner	name
1	Miguel
...	...
7	Martin
7	Christopher
...	...

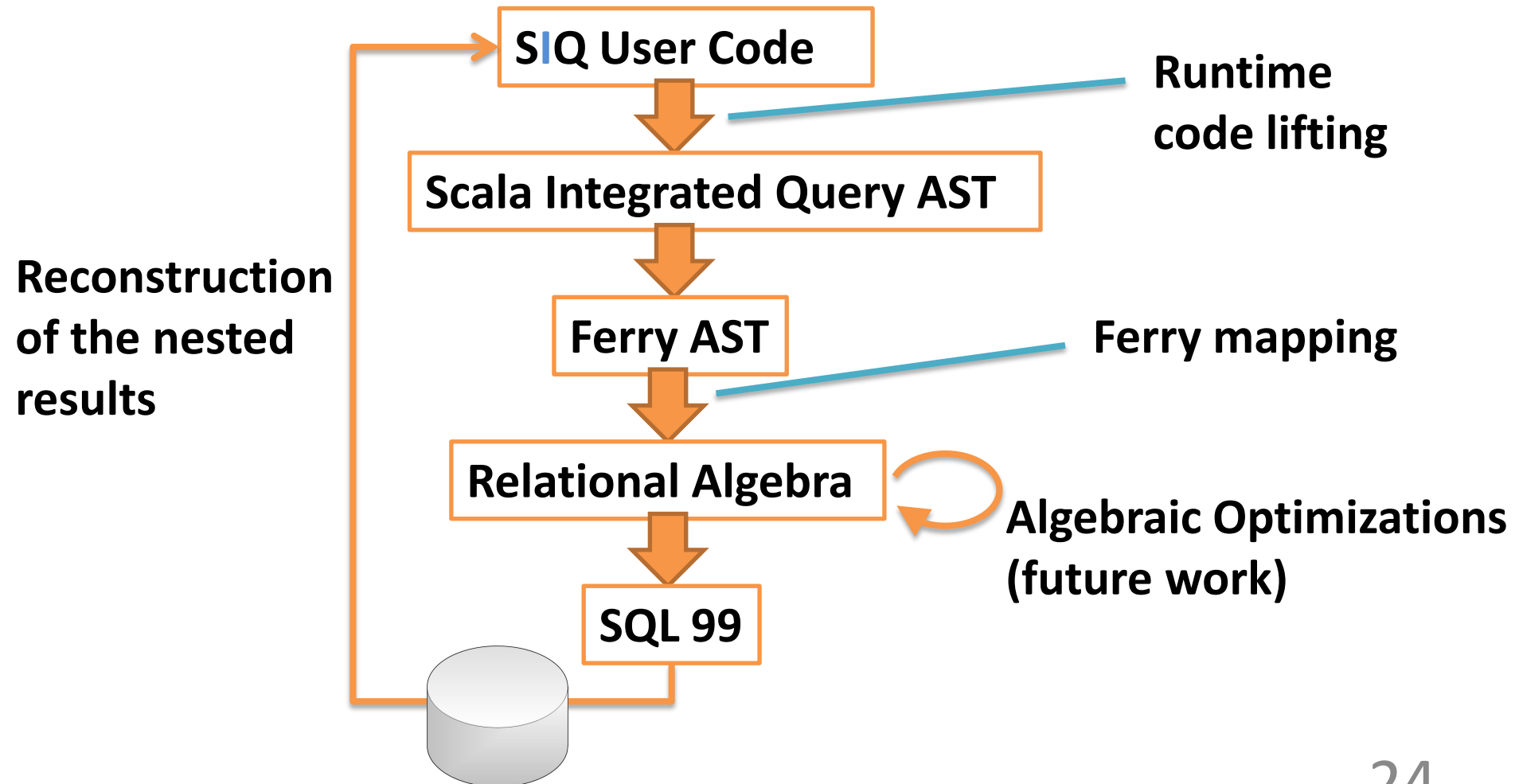


STAGING

SIQ Implementation

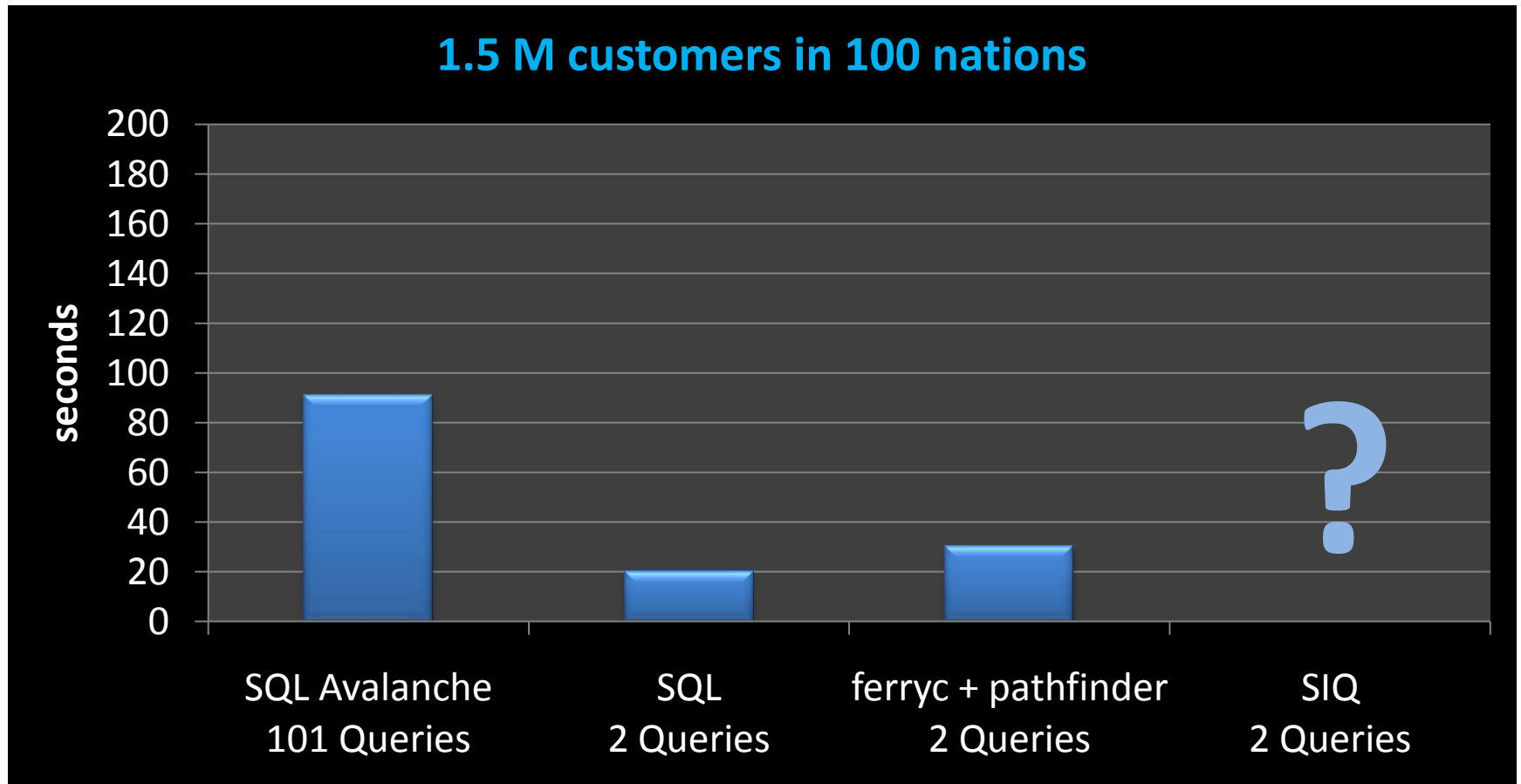
- **Staging** (EPFL / Stanford)
a technique used in
 - OptiML
 - Delite
 - Scala **Integrated** Query
 - ...
- light-weight modular staging framework
- **Scala Virtualized**

Staging in Scala Integrated Query



BENCHMARK

SQL execution time



10 GB TCP-H, Average of 10 runs, Postgres, Windows 7 32bit, Intel Core i5 M 580, 4GB RAM
Ferry list order requires SQL ROW_NUMBER, so for comparability used in all queries here

Supported Scala Subset

currently	coming up
<div>map</div> <div>flatMap</div> <div>withFilter</div> <div>length</div> <div>+ - * / %</div> <div>distinct</div> <div>flatten</div> <div>++</div> <div>== >= <= > <</div> <div>+</div> <div>&& </div> <div>unzip</div> <div>orderBy</div> <div>Tuple1-23</div> <div>zip</div> <div>groupBy</div>	<div>head</div> <div>db functions</div> <div>max</div> <div>min</div> <div>take</div> <div>apply</div> <div>mkString</div> <div>drop</div> <div>zipWithIndex</div> <div>sum</div> <div>avg</div> <div>more ...</div>

Future Work

- Make it production ready
 - I will be joining Typesafe in Autumn
- Algebraic optimizations
 - collaboration with University of Tübingen
 - taken from their „Pathfinder“ project
 - with help of Tom Schreiber, the inventor of Ferry
- Collaboration with Stefan Zeiger / Scala Query

Acknowledgements

- **Prof. Bernhard Rumpe**
- **Prof. Martin Odersky**
- **Miguel Garcia**
- **Tiark Rumpf**
- **Tom Schreiber**
- **Prof. Torsten Grust**

check out the prototype at

<http://code.google.com/p/scala-integrated-query/>

Other Ferry implementations exist for
.NET/LINQ, Haskell, Ruby, Links

Thank you!

<http://code.google.com/p/scala-integrated-query/>

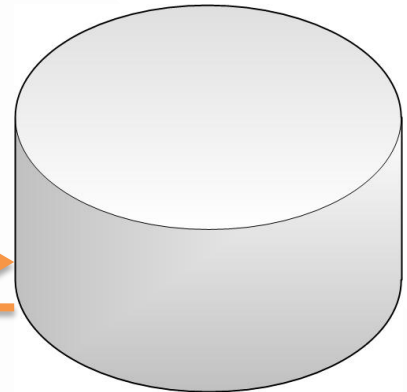
SQ: Element access before execution

JVM

```
def customers( nationkey : Rep[Int] ) =  
  customer.withFilter(  
    _.nationkey == nationkey  
  )  
  
nation.map( n =>  
  (n, customers(n.nationkey) )  
).withFilter( _._1.name == „EGYPT“ )  
  .map( _._2)
```

.fromdb

DBMS



only 1 SQL query
in total

BEHIND THE SCENES

- CODE LIFTING

Code lifting - internal view (simplified)

- code-lifting, query returns AST
- `for(c <- customer; if c.custkey == 2) yield c`

```
Comprehension(  
  Table("customer"),  
  Equals( Column("custkey"), Const( 2 ) ),  
  ( Column("custkey"), Column("name"), ... )  
)
```

- `SELECT * FROM customer WHERE custkey = 2`

Code lifting – How lifting works

- implicit conversions
- overloading
 - methods: flatten, head, +, -, ...
 - in Scala Virtualized also: ==, if-then-else, ...
 - for-comprehensions: map, flatMap, withFilter

```
for( c <- customer; if c.custkey == 2 ) yield c
```



```
customer.withFilter( c => c.custkey == 2 )
```

Code lifting - user view

- types wrapped by `Rep[...]`

type-driven lifting
methods redefined for
`Rep[...]`
return again `Rep[...]`

- `customer.withFilter(_.custkey == 2).fromdb`

`Rep[Seq[Employee]]`

`Rep[Employee]`

`Rep[Int]`

`Int`

`Rep[Boolean]`

`Rep[Seq[Employee]]`

`Seq[Employee]`

BEHIND THE SCENES - FERRY

Ferry Mapping

for **nation** <- nations

(**nation**, **customers** of this nation)

- done in 2 steps (all in the database)
 1. **loop lifting**: for each nation ALL customers
 2. **filter all at once**: keep only relevant customers

Ferry Mapping

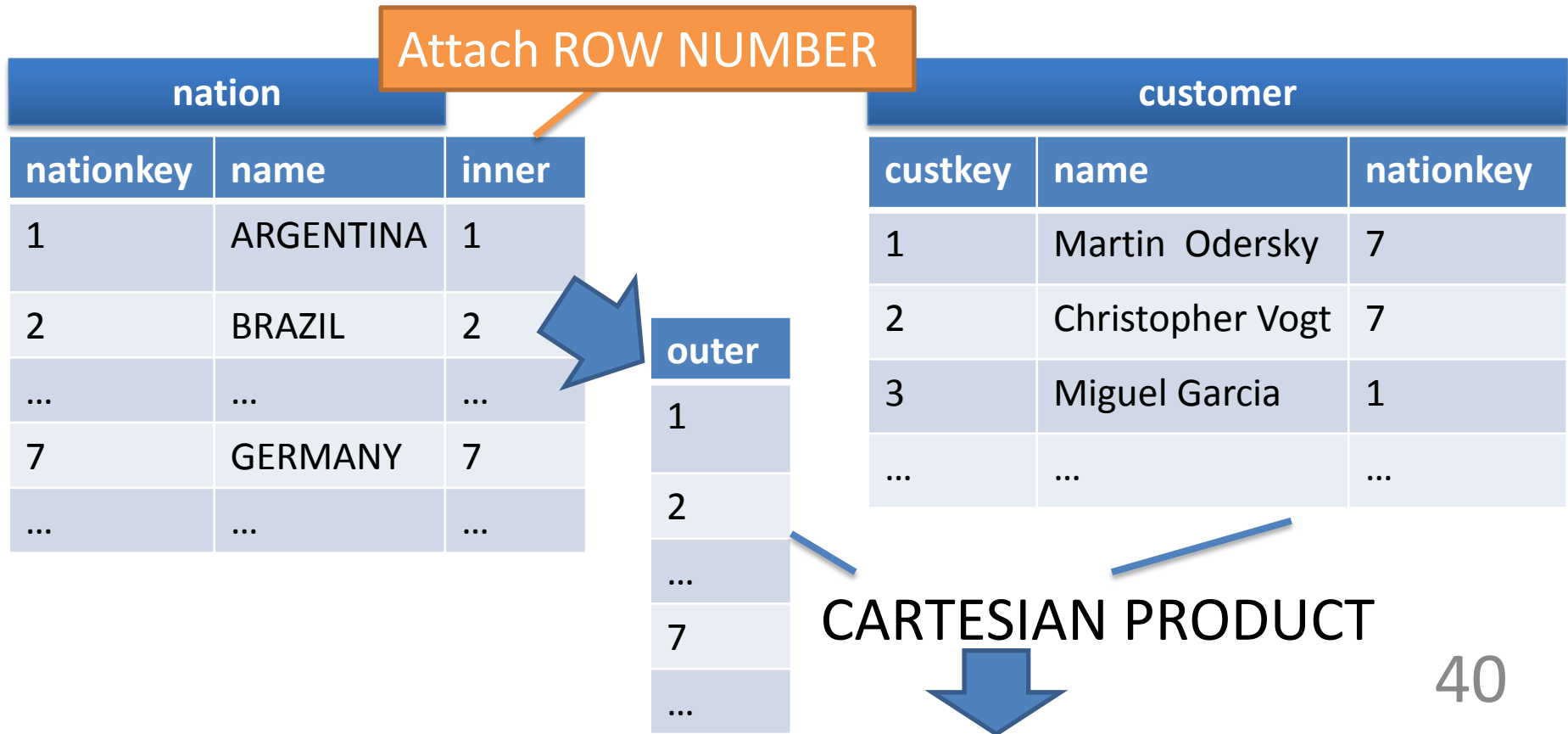
for **nation** <- nations
(**nation**, **customers** of this nation)

nation	
nationkey	name
1	ARGENTINA
2	BRAZIL
...	...
7	GERMANY
...	...

customer		
custkey	name	nationkey
1	Martin Odersky	7
2	Christopher Vogt	7
3	Miguel Garcia	1
...

Step 1: Loop Lifting

for each nation ALL customers
for(n <- nation) yield customer



After Loop Lifting: Nesting

for each nation ALL customers
for(n <- nation) yield customer



can be used to
construct nested results

Query 1		
nationkey	name	inner
1	ARGENTINA	1
2	BRAZIL	2
...
7	GERMANY	7
...

Query 2			
outer	custkey	name	nationkey
1	1	Martin Odersky	7
1	2	Christopher Vogt	7
1	3	Miguel Garcia	1
...
2	1	Martin Odersky	7
2	2	Christopher Vogt	7
2	3	Miguel Garcia	1

Step 2: Filter all at once

customers grouped by nation

for(n <- nation) yield

customer.withFilter(_.nationkey == n.nationkey)

Query 1: q1

nationkey	name	inner
1	ARGENTINA	1
2		
...		
7		
...

Query 2: q2

outer	custkey	name	nationkey
1	1	Martin Odersky	7
...			
2		Martin Odersky	7
2	2	Christopher Vogt	7
2	3	Miguel Garcia	1
...			

JOIN ON

inner = outer AND q1.nationkey = q2.nationkey

Step 2: Filter all at once

customers grouped by nation

for(n <- nation) yield

customer.withFilter(_.nationkey == n.nationkey)

Query 1: q1

nationkey	name	inner
1	ARGENTINA	1
2	BRAZIL	2
...
7	GERMANY	7
...

Query 2: q2

outer	custkey	name	nationkey
1	1	Martin Odersky	7
1	2	Christopher Vogt	7
1	3	Miguel Garcia	1
...
2	1	Martin Odersky	7
2	2	Christopher Vogt	7
2	3	Miguel Garcia	1

Nested Results

customers grouped by nation 

```
for( n <- nation ) yield
```

```
  customer.withFilter( _.nationkey == n.nationkey )
```

Query 1		
nationkey	name	inner
1	ARGENTINA	1
2	BRAZIL	2
...
7	GERMANY	7
...

Query 2'			
outer	custkey	name	nationkey
1	3	Miguel Garcia	1
...
7	1	Martin Odersky	7
7	2	Christopher Vogt	7
...

Scala Iterables

Correlate **main memory data** and **database tables**

```
for( i <- List(17,3,4).toddb; n <- nation;  
      if i == n.nationkey ) yield n
```

One query, no avalanche:

```
WITH t( i ) AS (VALUES (17),(3),(4))  
SELECT nation.* FROM t, ...
```

works for any iterable

List Order

Return items in order

```
for( i <- List(17,3,4).todb; n <- nation;  
      if i == n.nationkey ) yield n
```

Results are in order of the in-memory list:

```
List(Nation(17,PERU), Nation(3,CANADA), ... )
```

Mapping of List Order

- Get nations in order given ids
- `for(i <- List(17,3,4).toddb; n <- nation;
if i == n.nationkey) yield n`

Attach ROW NUMBER

VALUES (17), (3), (4)	
pos	i
1	17
2	3
3	4

JOIN ON
`i == nationkey`

ORDER BY pos

nation	
nationkey	name
...	...
3	CANADA
4	EGYPT
...	...
17	PERU
...	...