

Cassandra - CQL

Big Data Systems

Dr. Rubi Boim

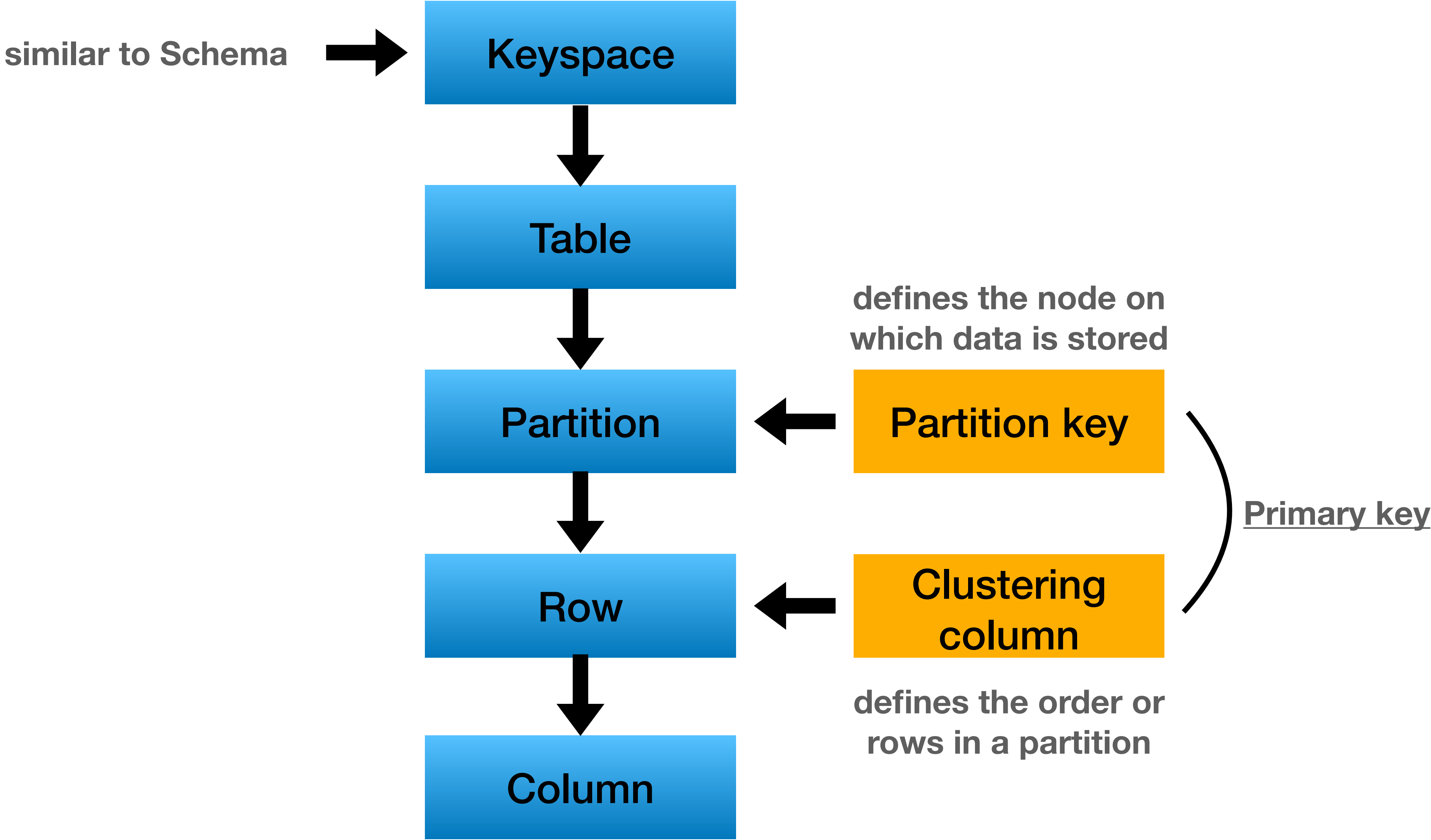
Cassandra CQL

- Terminology
- Keyspaces
- Tables
- Data types
- DDL / DML



Spoiler - most slides will be on SELECT

Terminology (Cassandra)



Keyspace

- High level container - AKA “schemas” from rDB
- **replication factor strategy**
 - “SimpleStrategy”: entire cluster
 - “NetworkTopologyStrategy”: different settings for each DS

Keyspace

```
CREATE KEYSPACE BigDataCourse WITH REPLICATION = {  
  'class'          : 'SimpleStrategy',  
  'replication_factor': 1  
};
```

```
CREATE KEYSPACE BigDataCourse WITH REPLICATION = {  
  'class'          : 'NetworkTopologyStrategy',  
  'israel'         : 3 , // Datacenter 1  
  'us'             : 2   // Datacenter 2  
};
```

Use & Describe

- **USE:** switch between key spaces in CQL

USE bigdatacourse

JAVA:

```
CassandraConnectionPool connectionPool.setKeyspace("bigdatacourse")
```

- **DESCRIBE:** display detailed information in CQL
(see manual for more options)

DESCRIBE KEYSACES/KEYSPACE/TABLES/TABLE/...

CREATE TABLE

```
CREATE TABLE students (  
  column1    TEXT,  
  column2    INT,  
  column3    UUID,  
  PRIMARY KEY (column1)  
);
```

```
CREATE TABLE [IF NOT EXISTS] [keyspace_name.] table_name (  
  column_definition [, ...]  
  PRIMARY KEY (column_name [, column_name ...])  
[WITH table_options  
  | CLUSTERING ORDER BY (clustering_column_name order)]  
  | ID = 'table_hash_tag'  
  | COMPACT STORAGE]
```

Data types (basic)

- TEXT utf8
- INT signed 32bits
- BIGINT signed 64bits
- TIMESTAMP 64bits
- FLOAT 32bits floating point
- DOUBLE 64bits floating point
- DECIMAL variable-precision decimal
- UUID universally unique identifier, 128bits
- TIMEUUID sortable UUID, embedded timestamp
- BLOB arbitrary bytes

Data types (basic)

- TEXT utf8
- INT signed 32bits
- BIGINT signed 64bits
- TIMESTAMP 64bits
- FLOAT 32bits floating point
- DOUBLE 64bits floating point
- DECIMAL variable-precision decimal
- **UUID** universally unique identifier, 128bits
- **TIMEUUID** **sortable UUID, embedded timestamp**
- BLOB arbitrary bytes

Unique across all nodes,
regardless of the number of nodes

Note on generating unique IDs

- Not trivial for distributed systems
- UUID / TIMEUUID are great
 - Downside - requires 128bit
what's the problem with java primitives?

Note on generating unique IDs

- Not trivial for distributed systems
- UUID / TIMEUUID are great
- Downside - requires 128bit
what's the problem with java primitives?



Max primitive is 64bit (long)

More data types

- COUNTER
- LIST
- SET
- MAP
- More on these later...

SELECT

```
SELECT * FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse  
WHERE column1 = "1234" LIMIT 100
```

```
SELECT count(*) FROM BigDataCourse
```

- “Limited” compared to RDBMS
sum / avg / min / max or only supported on new versions
no joins / having / union...

SELECT

```
SELECT * FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse  
WHERE column1 = "1234" LIMIT 100
```

```
SELECT count(*) FROM BigDataCourse
```

ANTI PATTERN

Can be very slow and expensive - when?

- “Limited” compared to RDBMS

sum / avg / min / max or only supported on new versions

no joins / having / union...

SELECT

```
SELECT * FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse
```

```
SELECT column1, column2 FROM BigDataCourse  
WHERE column1 = "1234" LIMIT 1
```

```
SELECT count(*) FROM BigDataCourse
```

Even if counting a single row, it can be expensive (on a really big wide row)

ANTI PATTERN

Can be very slow and expensive - when?

- “Limited” compared to RDBMS

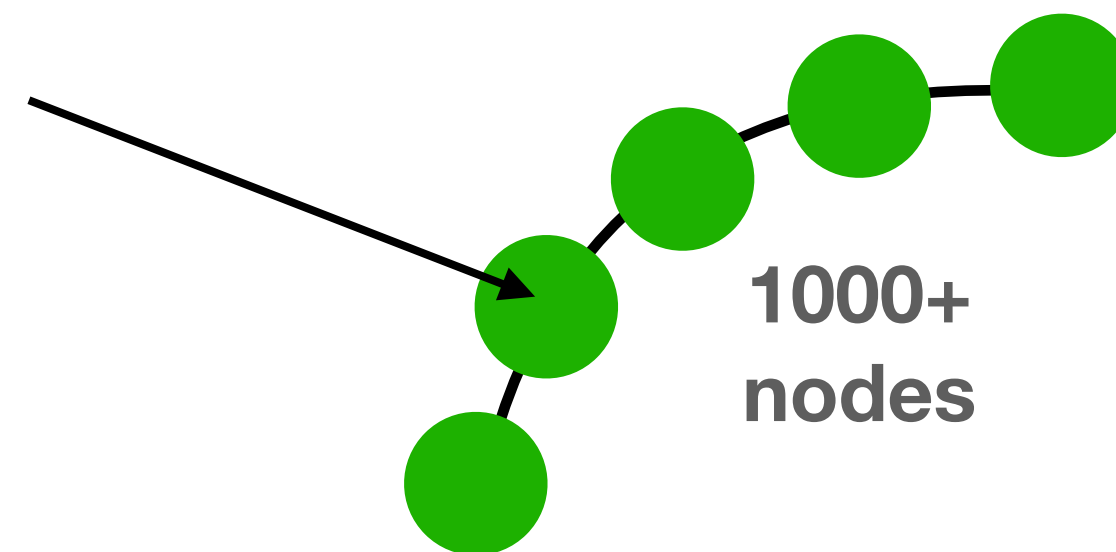
sum / avg / min / max or only supported on new versions

no joins / having / union...

SELECT - partitions and keys

- TLDR; **provide the partition key to the query**

```
SELECT * FROM users  
WHERE user_id = "1234"
```

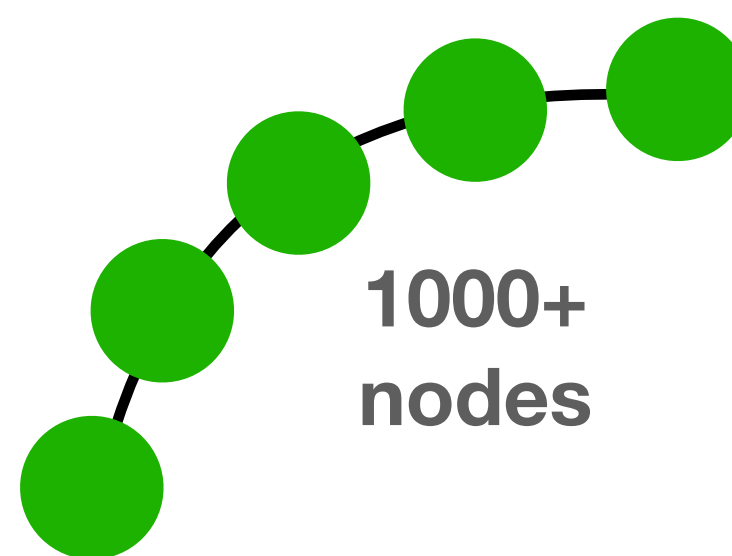


users	
user_id	K
name	
birth_year	
...	

SELECT - partitions and keys

- What happens if no partition is given?

```
SELECT * FROM users
```



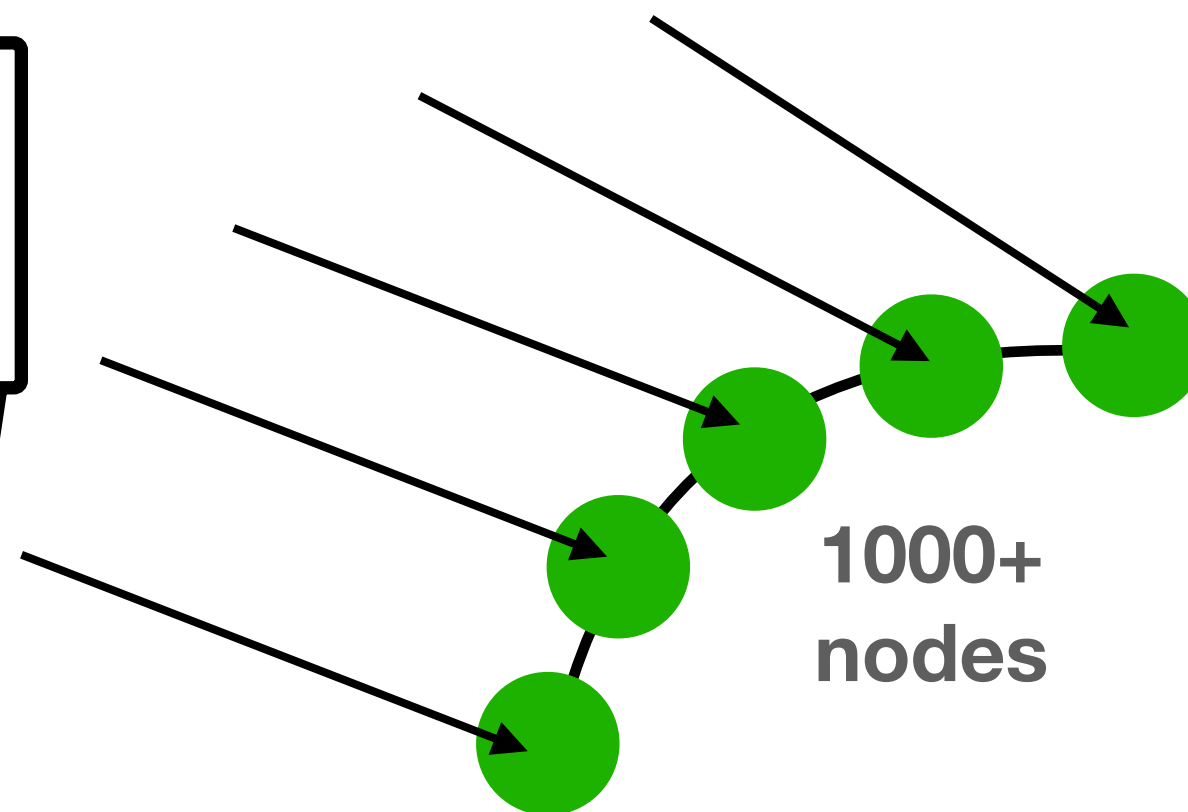
users	
user_id	K
name	
birth_year	
...	

SELECT - partitions and keys

- What happens if no partition is given?

```
SELECT * FROM users
```

We need to contact all servers
(as all partitions are valid)



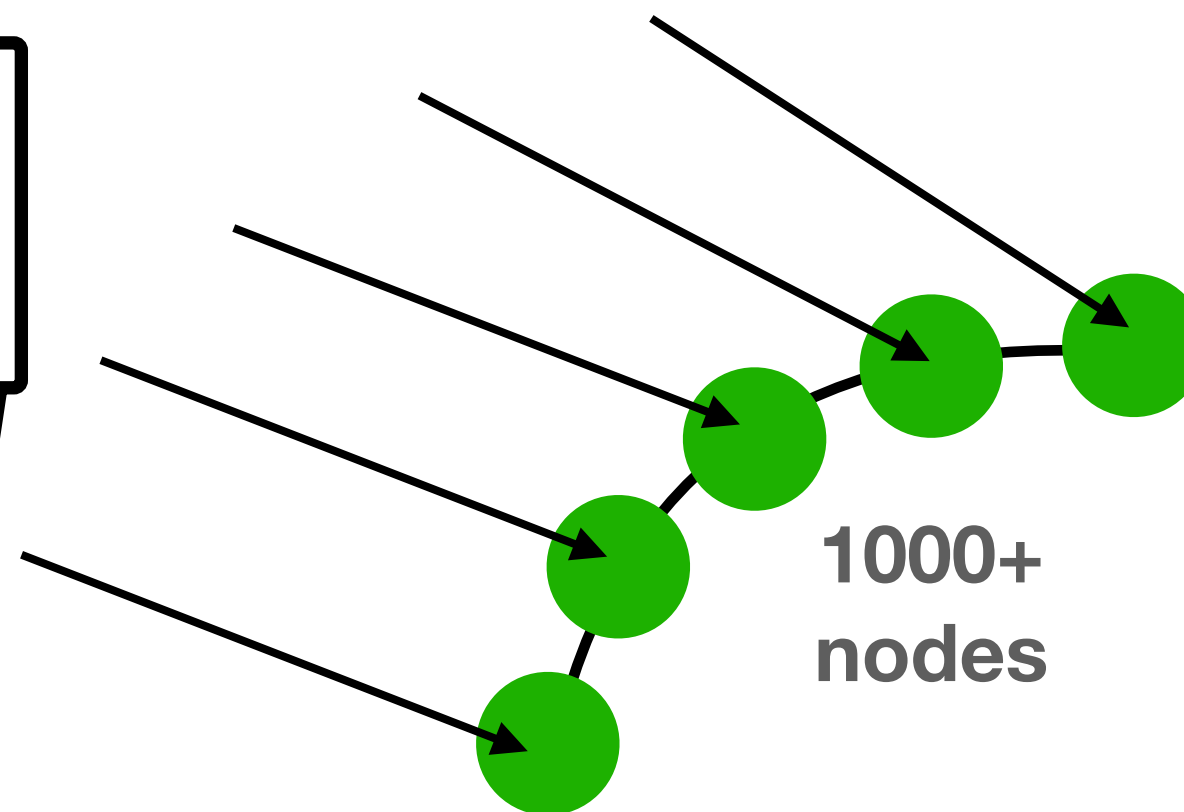
users	
user_id	K
name	
birth_year	
...	

SELECT - partitions and keys

- What happens if no partition is given?

```
SELECT * FROM users
```

We need to contact all servers
(as all partitions are valid)



This is valid!
Lets see some examples

users	
user_id	K
name	
birth_year	
...	

SELECT - partitions and keys

Each user “creates” a partition (user_id is partition_key)

Assume there are **10k nodes** in the cluster and **no replication**
- If there are **100k users**, would the query be optimal?
(that is, we would not check unnecessary nodes/partitions)

?

users
K
...

1000
nodes

SELECT - partitions and keys

Each user “creates” a partition (user_id is partition_key)

Assume there are **10k nodes** in the cluster and **no replication**
- If there are **100k users**, would the query be optimal?
(that is, we would not check unnecessary nodes/partitions)

YES - why?

?

users

K

ar

...

1000
nodes

SELECT - partitions and keys

Each user “creates” a partition (user_id is partition_key)

Assume there are **10k nodes** in the cluster and **no replication**

- If there are **100k users**, would the query be optimal?

(that is, we would not check unnecessary nodes/partitions)

YES - why?

There are 100k partitions which are distributed on 10k nodes

?

users

K

ar

...

1000
nodes

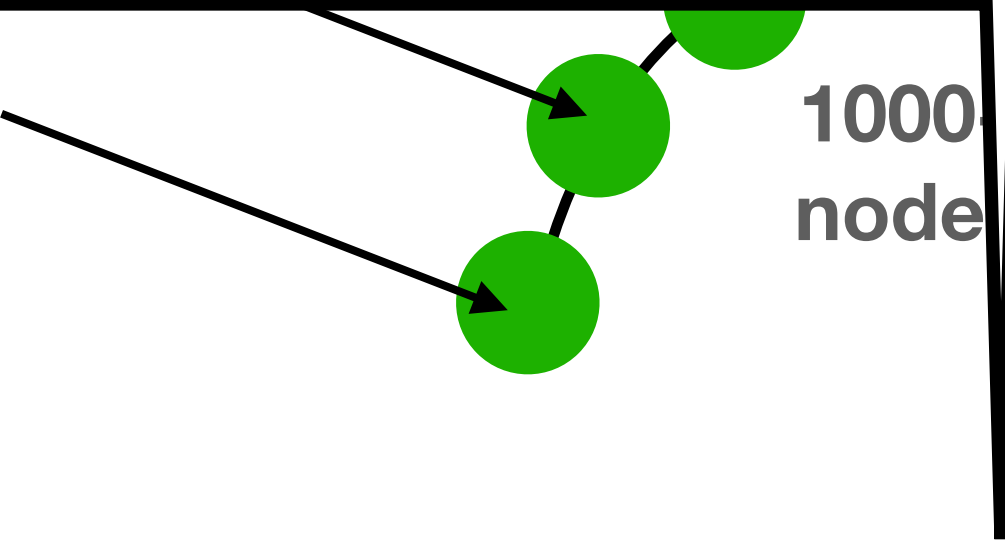
SELECT - partitions and keys

Each user “creates” a partition (user_id is partition_key)

Assume there are **10k nodes** in the cluster and **no replication**
 - If there are **10 users**, would the query be optimal?
 (that is, we would not check unnecessary nodes/partitions)

?

users	
	K



...

SELECT - partitions and keys

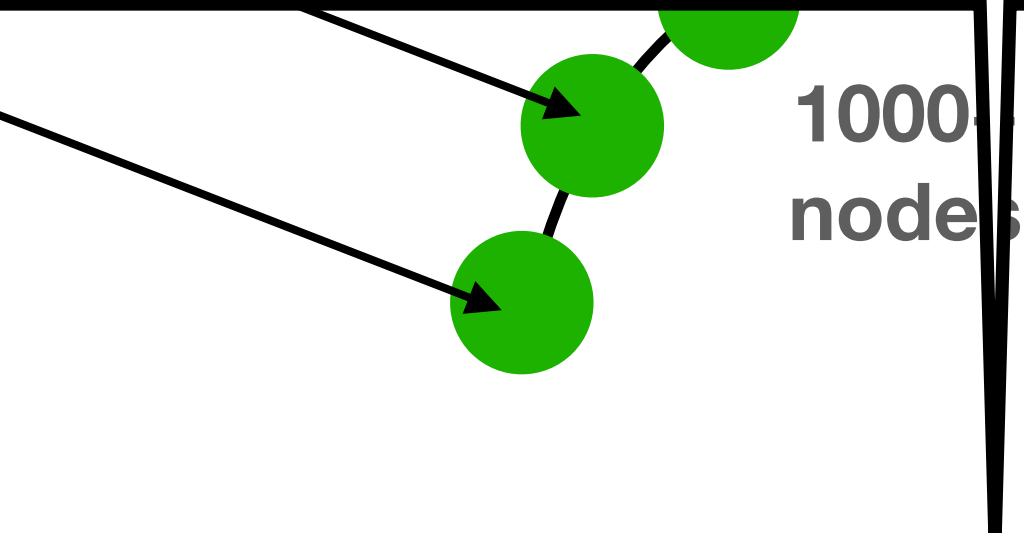
Each user “creates” a partition (user_id is partition_key)

Assume there are **10k nodes** in the cluster and **no replication**
- If there are **10 users**, would the query be optimal?
(that is, we would not check unnecessary nodes/partitions)

NO - why?

?

users	
	K



...

SELECT - partitions and keys

Each user “creates” a partition (user_id is partition_key)

?

Assume there are **10k nodes** in the cluster and **no replication**

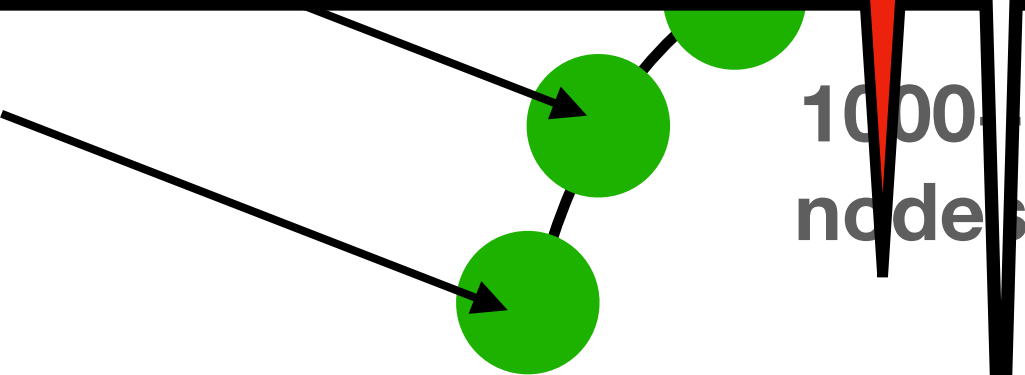
- If there are 10 users
(that is, we would have 10 partitions)

```
SELECT * from <TABLE> - Summary
```

NO - why?

The there are 10k nodes. We will have 10k partitions

Although this is allowed - this is in general anti pattern
Use with caution



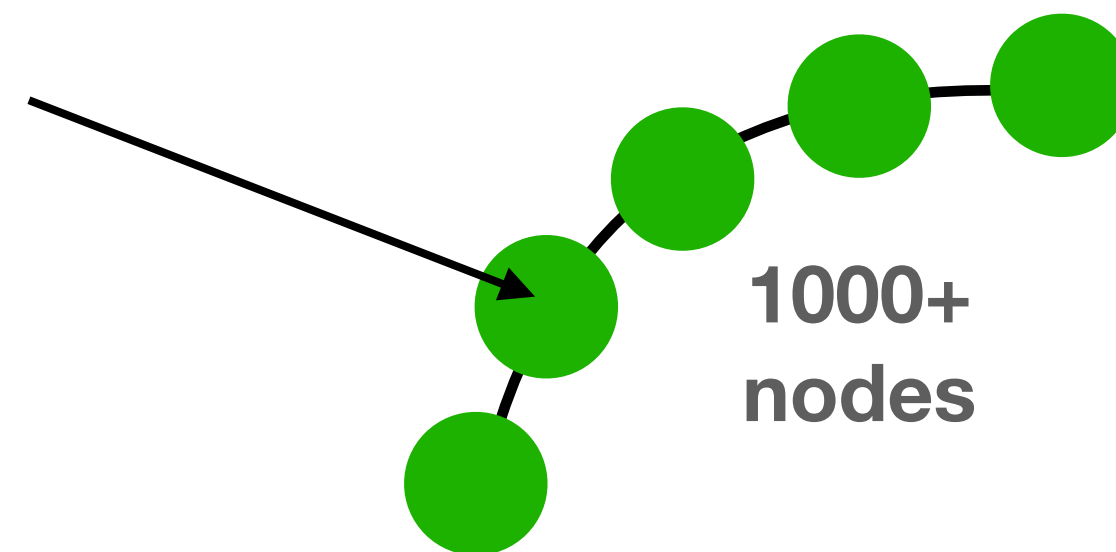
...

The right way for this scenario is to create a single partition for these 10 users, then read 1 partition

SELECT - partitions and keys

- Try a different model

```
SELECT * FROM users  
WHERE country = "israel"
```



Note
K is the partition key (**NOT the key**)
▼ C is the clustering column,
Together both are the key

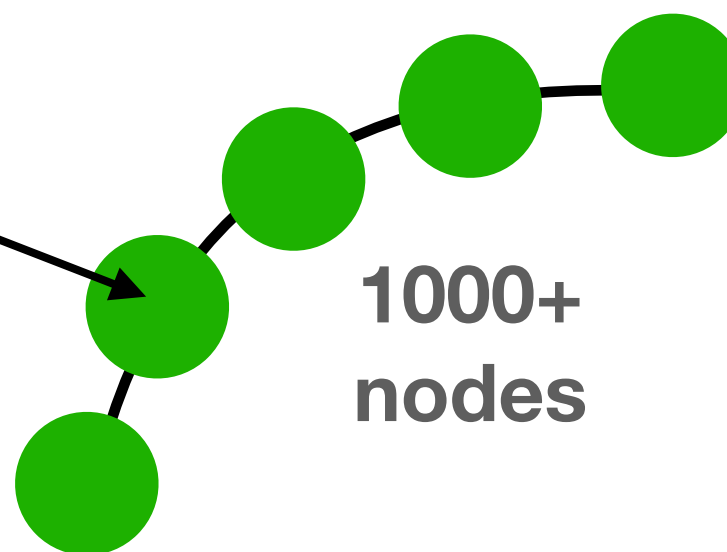
users	
country	K
user_id	▼ C
name	
birth_year	
...	

SELECT - partitions and keys

- Try a different model

```
SELECT * FROM users  
WHERE country = "israel"
```

Reading the users from Israel is fast



users	
country	K
user_id	▼C
name	
birth_year	
...	

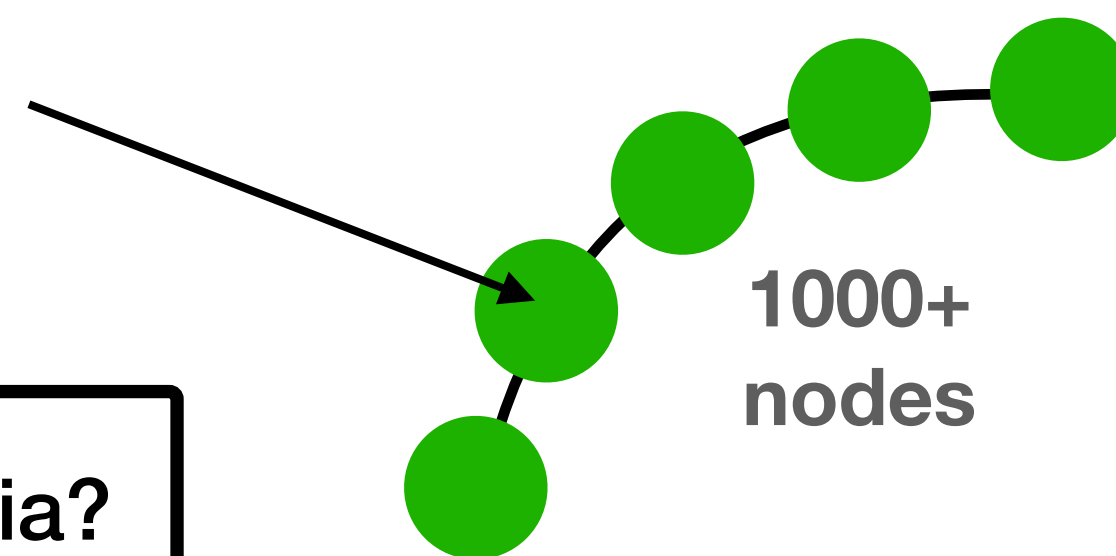
SELECT - partitions and keys

- Try a different model

```
SELECT * FROM users  
WHERE country = "israel"
```

users	
country	K
user_id	▼C
name	
birth_year	
...	

What happen if the country is India?



SELECT - partitions and keys

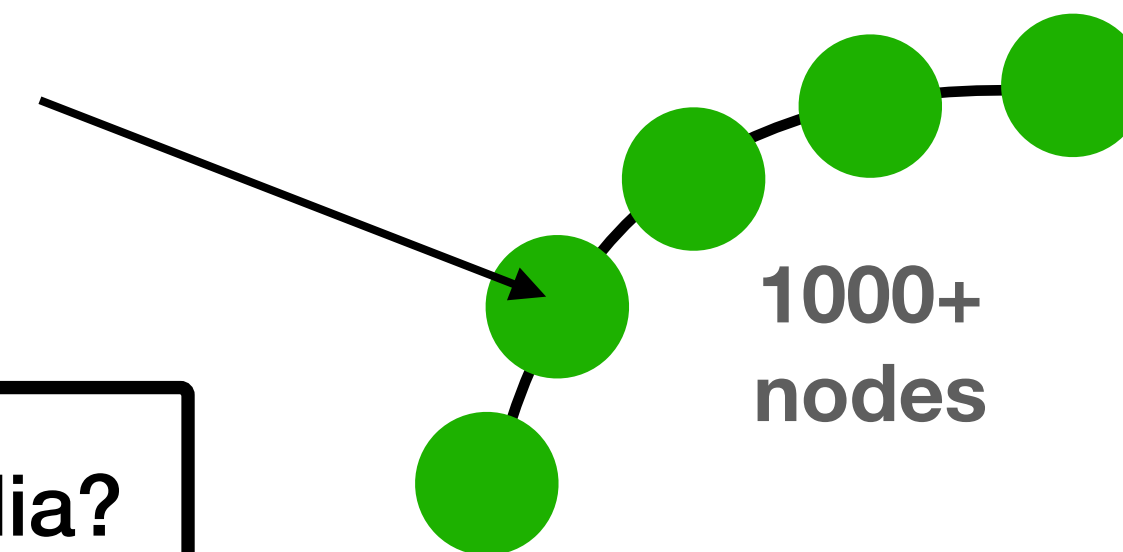
- Try a different model

```
SELECT * FROM users  
WHERE country = "israel"
```

users	
country	K
user_id	▼C
name	
birth_year	
...	

What happen if the country is India?

How can you solve this issue?



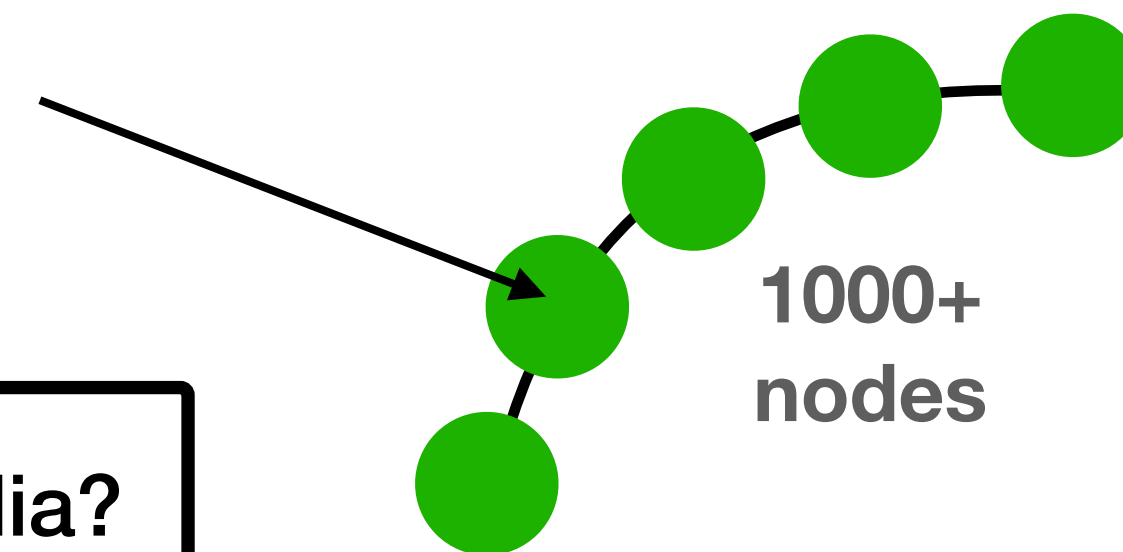
SELECT - partitions and keys

- Try a different model

```
SELECT * FROM users  
WHERE country = "israel"
```

users	
country	K
user_id	▼C
name	
birth_year	
...	

What happen if the country is India?



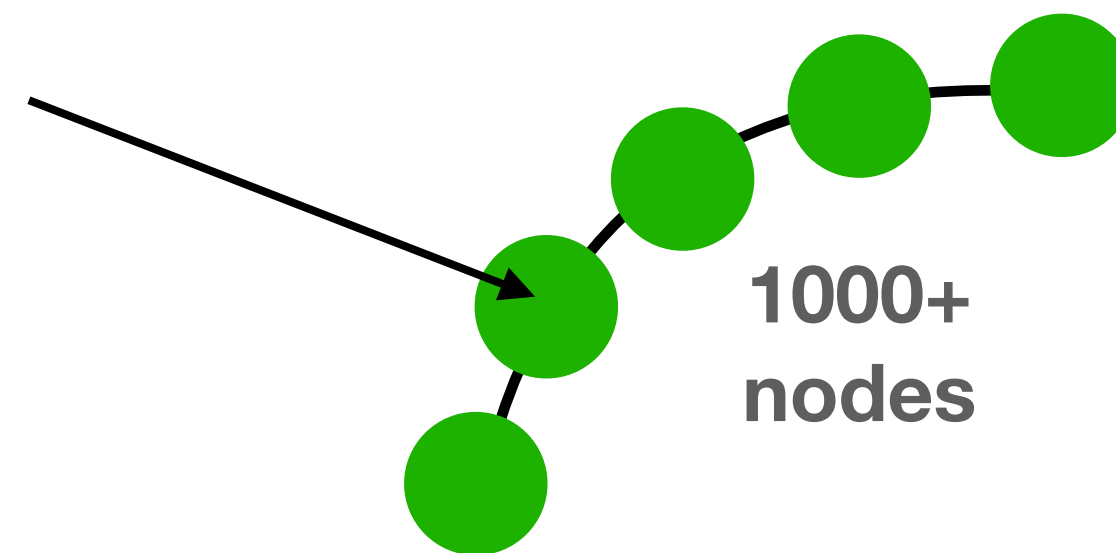
How can you solve this issue?

We can add "buckets" - more on this later

SELECT - partitions and keys

- What happens now?

```
SELECT * FROM users  
WHERE country = "israel"  
AND birth_year = 1982
```



users	
country	K
user_id	▼C
name	
birth_year	
...	

SELECT - partitions and keys

- What happens now?

```
SELECT * FROM users
WHERE country = "israel"
AND birth_year = 1982
```

users	
country	K
user_id	▼C
name	
birth_year	
...	



Error - why?

SELECT - partitions and keys

- What happens now?

```
SELECT * FROM users
WHERE country = "israel"
AND birth_year = 1982
```

users	
country	K
user_id	▼C
name	
birth_year	
...	



Error - why?

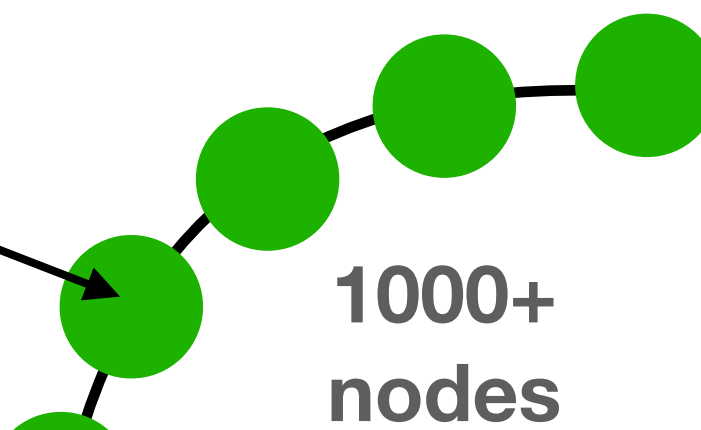
Cassandra will need to read the entire partition.
If there are 1m users, and only 10k were born in 1982,
there would be an unnecessary read/filter of 990k users

SELECT - partitions and keys

- What happens now?

```
SELECT * FROM users
WHERE country = "israel"
AND birth_year = 1982
ALLOW FILTERING
```

users	
country	K
user_id	▼C
name	
birth_year	



With “ALLOW FILTERING” Cassandra will approve the query
(ANTI PATTERN)

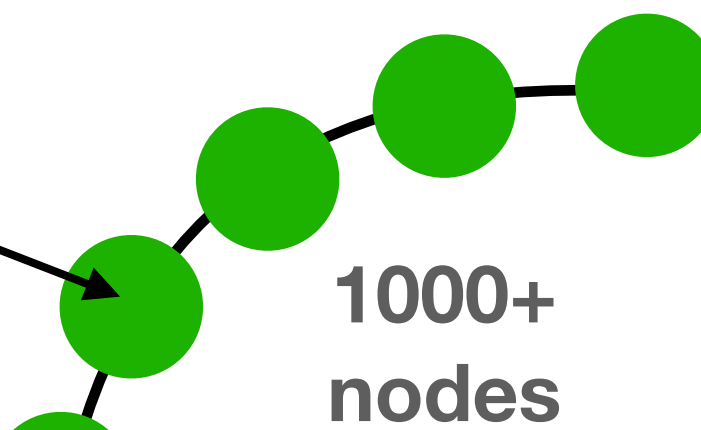
SELECT - partitions and keys

- What happens now?

How can you support the query without
“ALLOW FILTERING”?

```
SELECT * FROM users
WHERE country = "israel"
AND birth_year = 1982
ALLOW FILTERING
```

users	
country	K
user_id	▼C
name	
birth_year	

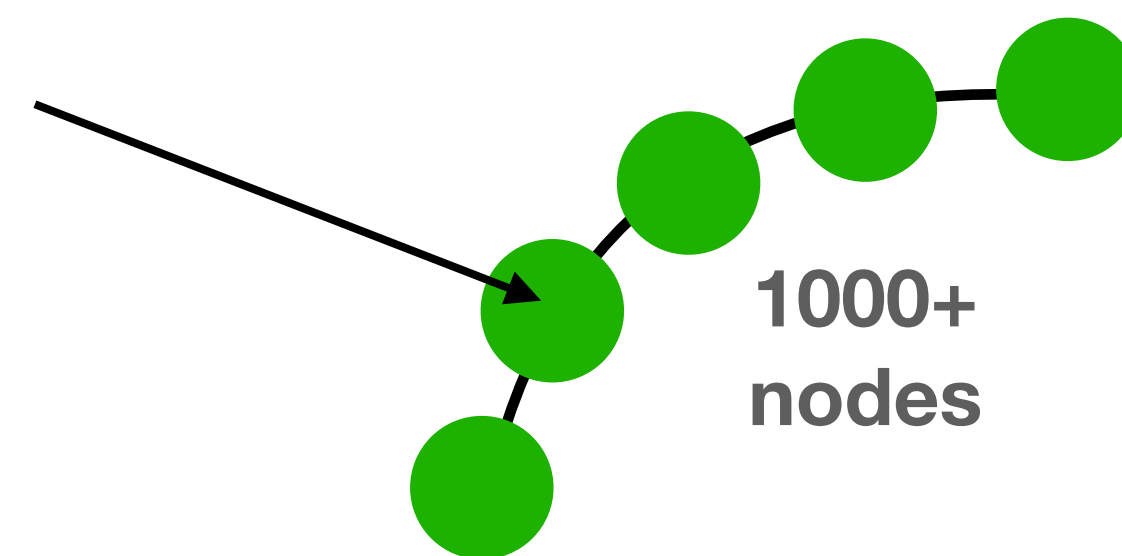


With “ALLOW FILTERING” Cassandra will approve the query
(ANTI PATTERN)

SELECT - partitions and keys

- Solved with denormalization

```
SELECT * FROM users_by_birth_year  
WHERE country = "israel"  
AND birth_year = 1982
```



- (we will talk about correct modeling later)

users	
country	K
user_id	▼C
name	
birth_year	
...	

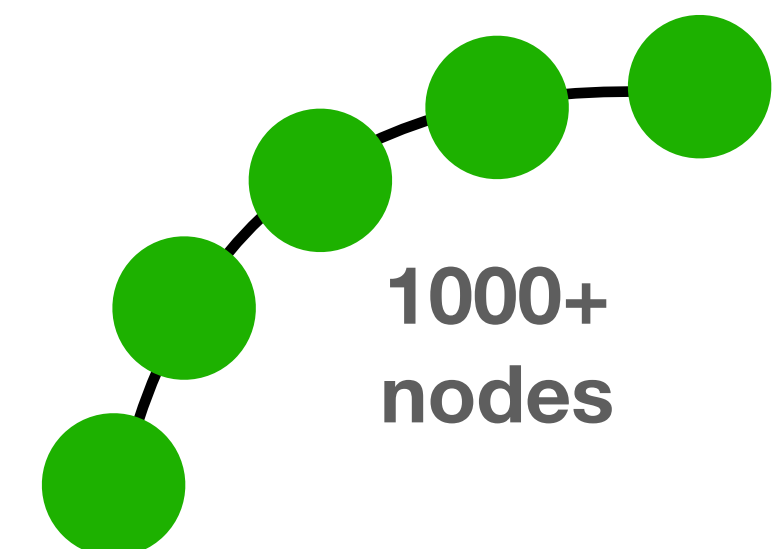
users_by_birth_year	
country	K
birth_year	▼C
user_id	▼C
name	
...	

SELECT - partitions and keys

- And what about this case?

```
SELECT * FROM users  
WHERE city = "tel aviv"
```

users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	



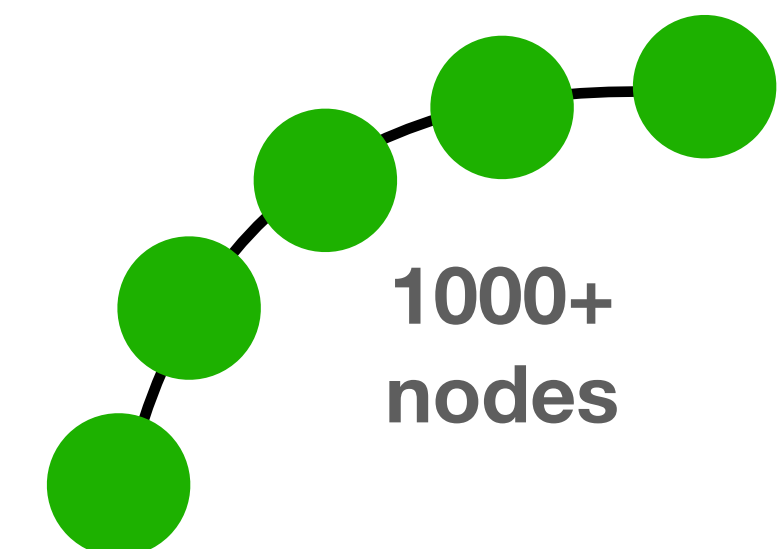
SELECT - partitions and keys

- And what about this case?

```
SELECT * FROM users  
WHERE city = "tel aviv"
```

users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	

Error - why?



SELECT - partitions and keys

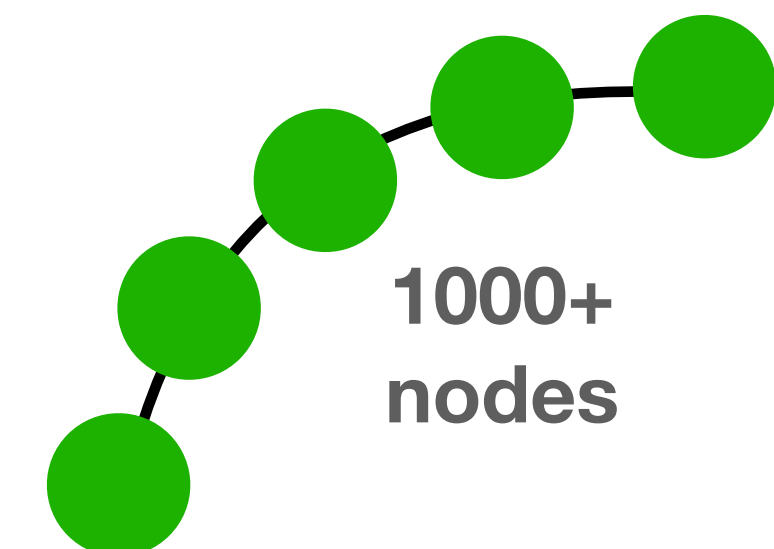
- And what about this case?

```
SELECT * FROM users  
WHERE city = "tel aviv"
```

users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	

Error - why?

Cassandra will need to contact all nodes and to check if such partition exists



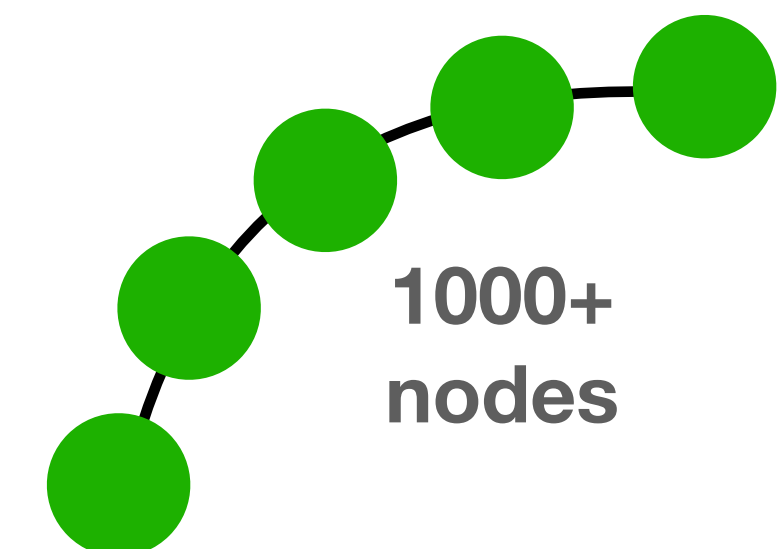
SELECT - partitions and keys

- And what about this case?

```
SELECT * FROM users
WHERE city = "tel aviv"
ALLOW FILTERING
```

With "ALLOW FILTERING" Cassandra will approve the query
(again - ANTI PATTERN)

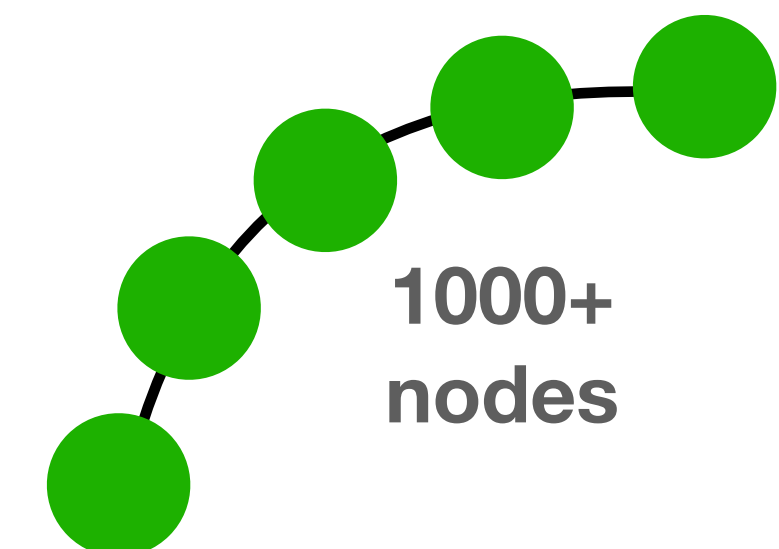
users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	



SELECT - ALLOW FILTERING

- Almost always ANTI PATTERN
- We saw these use cases
 - To “filter” columns in a single partition
 - To “filter” partitions across nodes

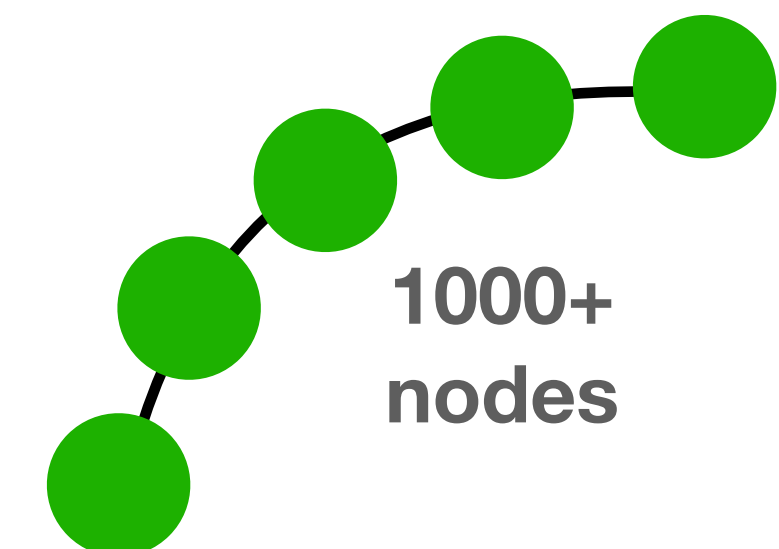
users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	



SELECT - ALLOW FILTERING

- Almost always ANTI PATTERN
- We saw these use cases
 - To “filter” columns in a single partition
 - To “filter” partitions across nodes
 - Can you think of another example?

users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	



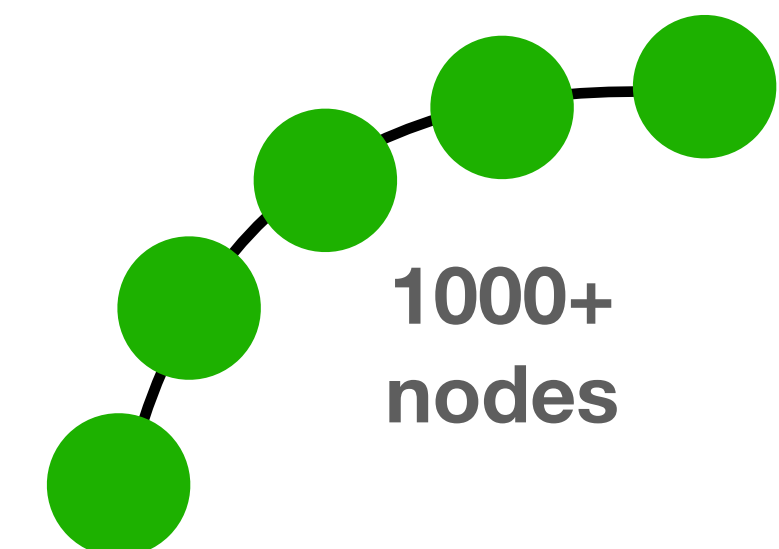
SELECT - ALLOW FILTERING

- Almost always ANTI PATTERN

```
SELECT * FROM users  
WHERE name = "rubi boim"  
ALLOW FILTERING
```

- We save
 - To “filter” columns in a single partition
 - To “filter” partitions across nodes
 - To “filter” columns across partitions

users	
country	K
city	K
neighborhood	K
user_id	▼C
name	
birth_year	



INSERT

- Primary key is obviously required

```
INSERT INTO BigDataCourse (column1 , column2)  
VALUES (123 , "name")
```

INSERT - IF NOT EXISTS

- Requires read before write!
- Use with caution

```
INSERT INTO BigDataCourse (column1 , column2)  
IF NOT EXSITS  
VALUES (123 , "name")
```

INSERT - IF NOT EXISTS

- Requires read before write!
- Use with caution

```
INSERT INTO BigDataCourse (column1 , column2)  
IF NOT EXSITS  
VALUES (123 , "name" )
```

Note - writes are cheaper than reads. If there are not too many writes, it is better to overwrite the same data instead of using "if not exists"

INSERT - USING TTL

- Time To Live - allows for automatic expiration (delete)
in seconds

```
INSERT INTO BigDataCourse (column1, column2)
VALUES (123, "name")
USING TTL 86400 // 24 hours
```

INSERT - USING TTL

- Time To Live - allows for automatic expiration (delete) in seconds

```
INSERT INTO BigDataCourse (column1 , column2)  
VALUES (123, "name")  
USING TTL 86400 // 24 hours
```



Creates tombstones
more on this later

UPDATE

- Primary key is obviously required

```
UPDATE BigDataCourse
SET column2 = "name", column3 = "abc"
WHERE column1 = 123
```

DELETE

- Warning:
DELETES in distributed databases is NOT TRIVIAL
- In Cassandra in particular
- Deleted data is not removed immediately
a tombstone is created
- More on this later

DELETE

- Delete data from a row

```
DELETE name FROM users
WHERE country = "israel"
AND user_id = "123"
```

- Delete an entire row

```
DELETE FROM users
WHERE country = "israel"
```

users	
country	K
user_id	▼C
name	
birth_year	
...	

Truncate

- Removes all SSTables holding data
- Use with care
- (Avoids tombstones)

TRUNCATE `users`

ALTER TABLE

- Add / drop / rename existing columns
- *change datatypes (with restrictions)
- Change table properties
- Can NOT alter PRIMARY KEY columns
- RTFM :)

```
ALTER TABLE [keyspace_name.] table_name  
[ALTER column_name TYPE cql_type]  
[ADD (column_definition_list) ]  
[DROP column_list | COMPACT STORAGE ]  
[RENAME column_name TO column_name]  
[WITH table_properties];
```