# I MapReduced a Neo store

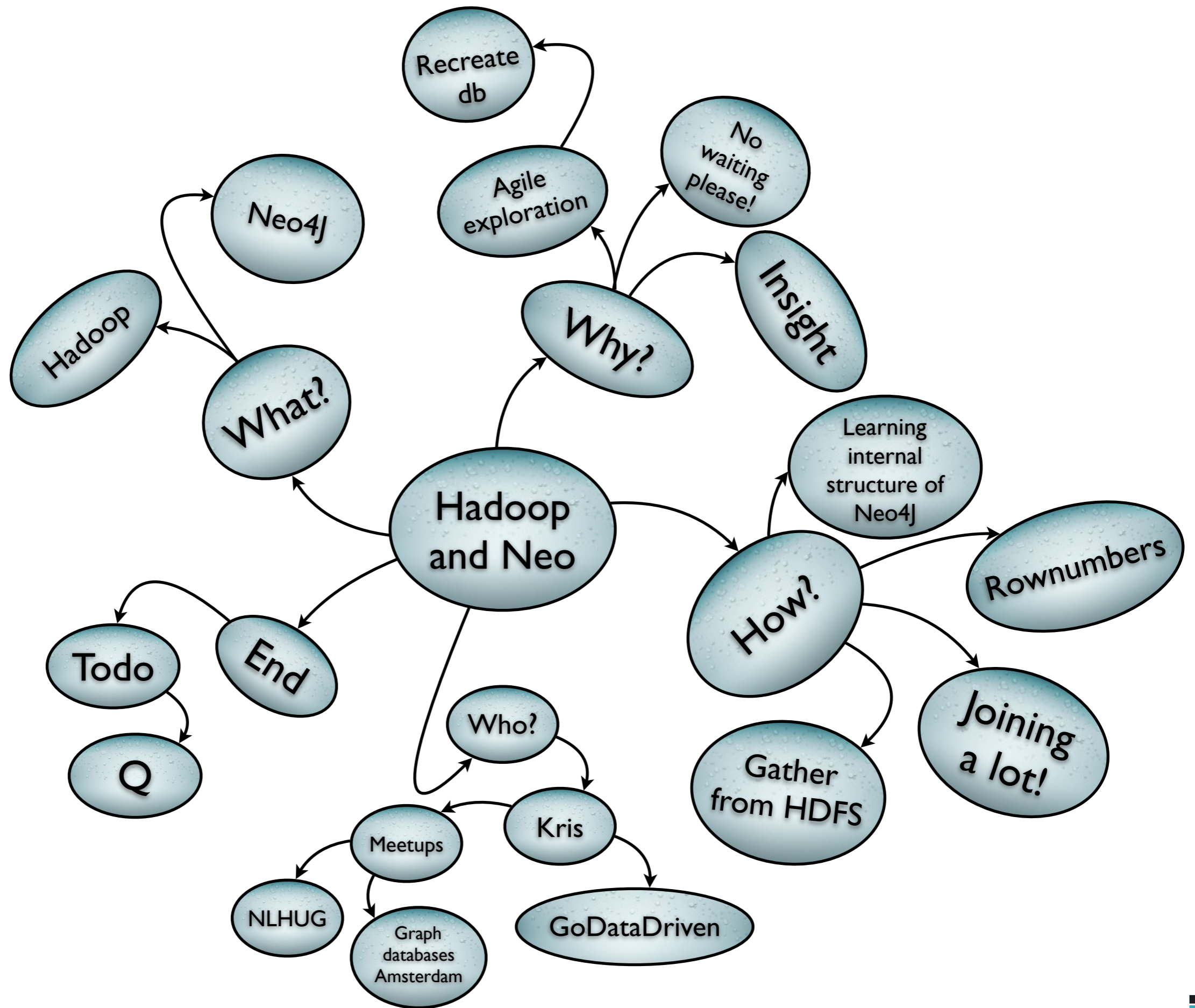## Creating large neo4j databases with hadoop
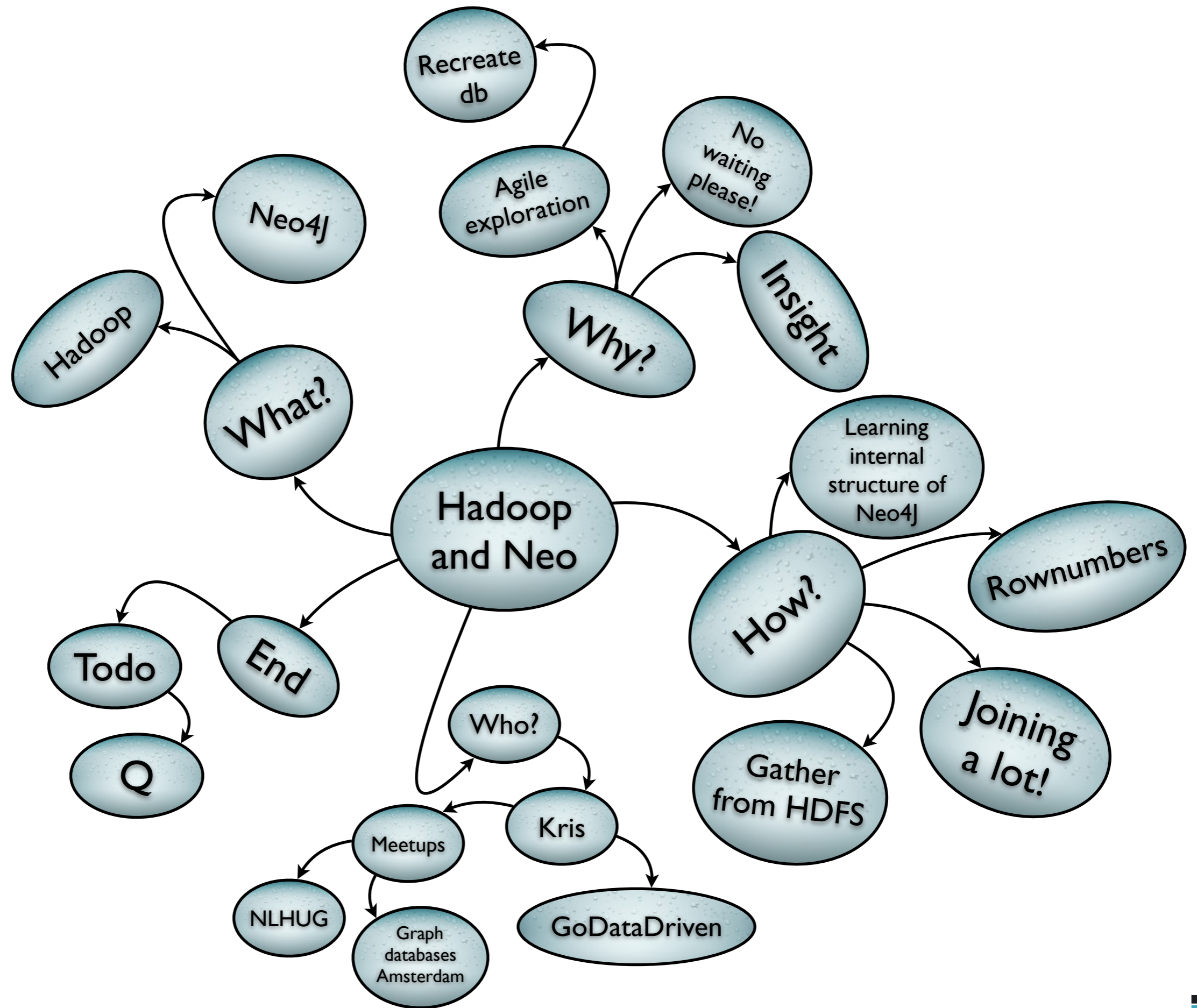
*Kris Geusebroek*
*Big Data Hacker*

*@krisgeus*
*krisgeusebroek@godatadriven.com*
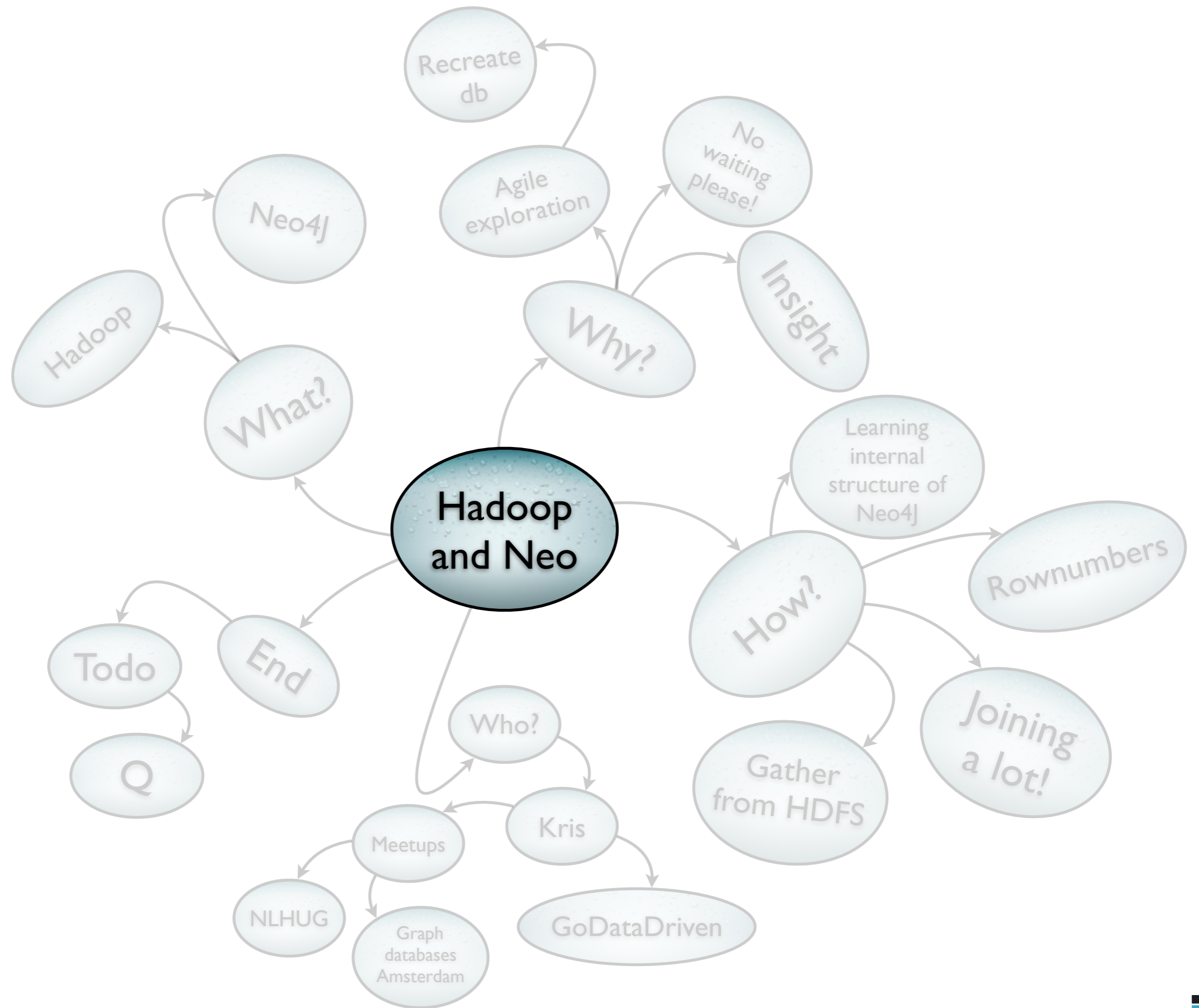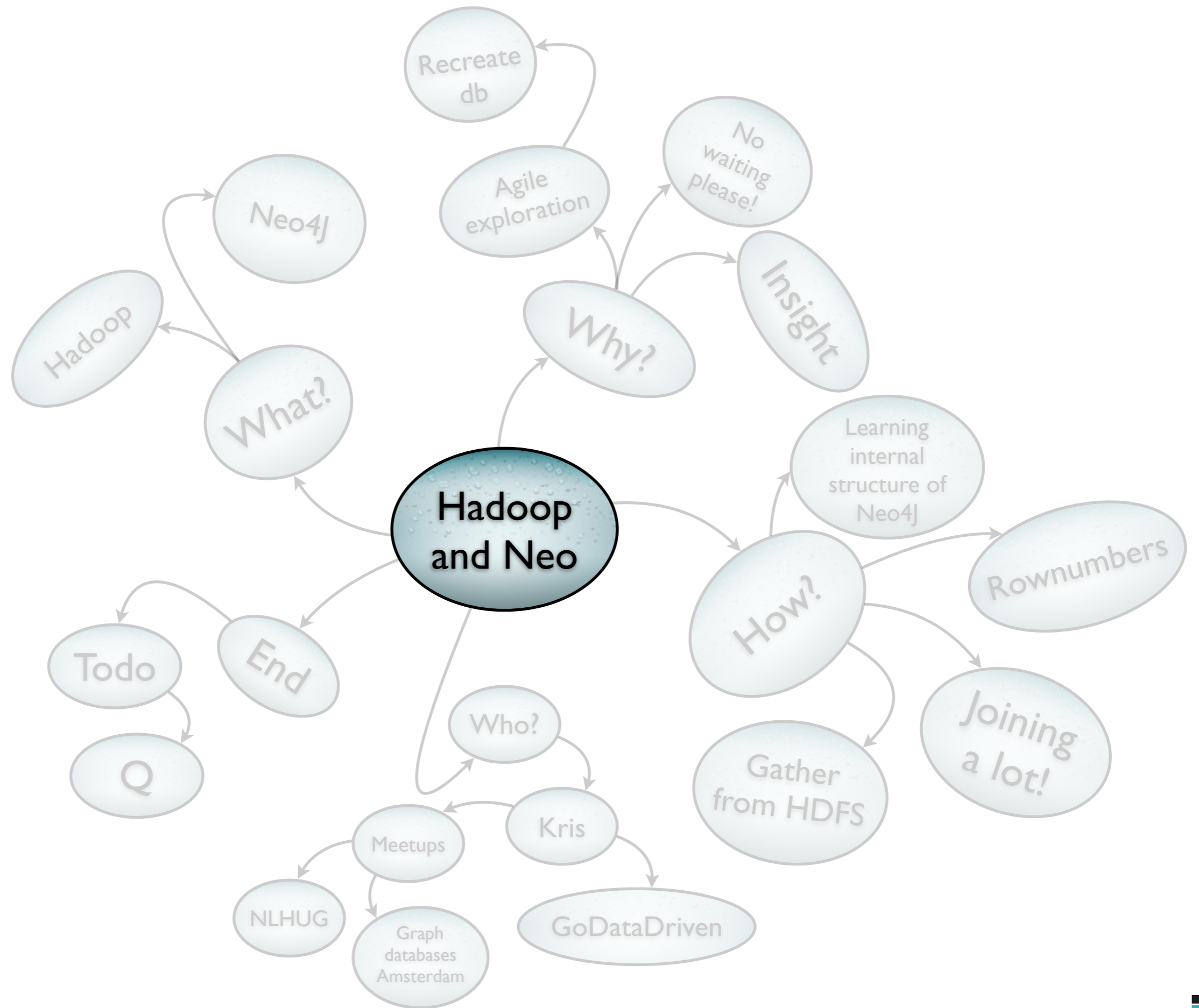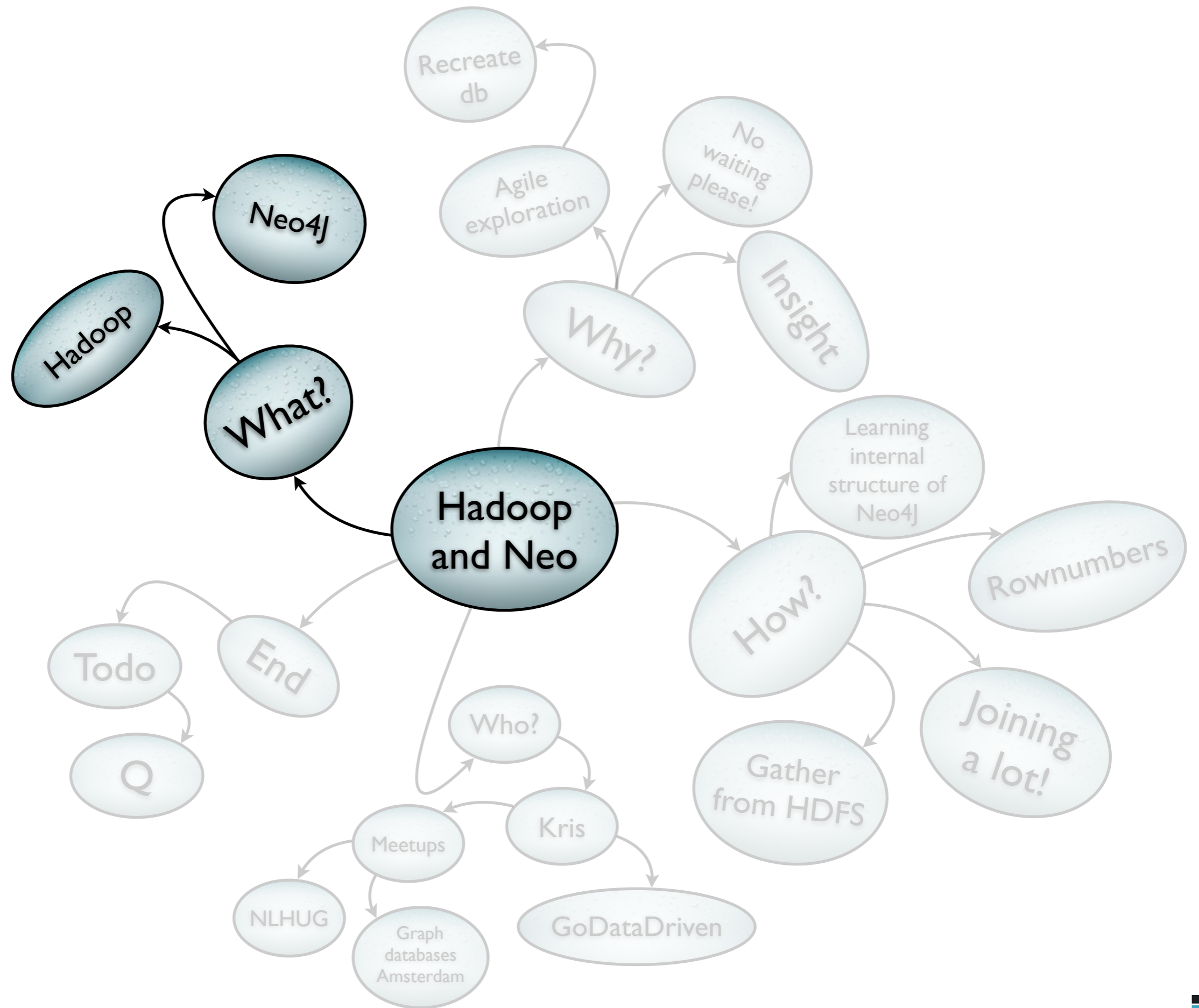
GoDataDriven

PROUDLY PART OF THE XEBIA GROUP
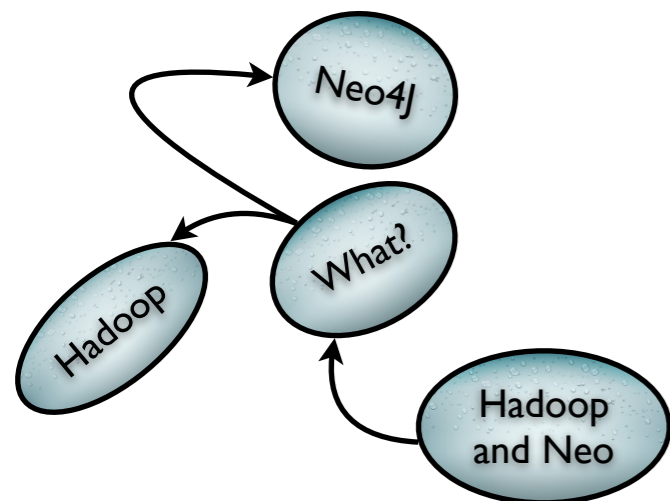
Hadoop

Apache Hadoop is an open-source framework that supports data-intensive distributed applications

Neo4J



WIKIPEDIA
The Free Encyclopedia



Neo4j is an open-source graph database, implemented in Java. The developers describe Neo4j as "embedded, disk-based, fully transactional Java persistence engine that stores data structured in graphs rather than in tables"

Large

Use case: Create large neo4j database

Multiple times!

about 30.000.000 nodes and 700.000.000 edges between them.

each node has 9 properties
each edge has 4 properties

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

## From

Insight

## To

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Agile exploration

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Re-create db

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Re-create db

Neo4j
the world's leading
graph database

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Re-create db

Neo4j
the world's leading
graph database

# Re-create db

Multiple ways to create a Neo4J database
- Just add nodes and edges with:
  - Java api
  - Rest api
  - Cypher via rest
- Batch import functionality
- My way ;-)

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Multiple ways to create a Neo4J database

Re-create db

Just adding nodes and edges has some problems.
Neo4J being transactional will give a lot of transaction overhead for large graphs.

Recreate
db

Agile
exploration

No
waiting
please!

Why?

Insight

Hadoop
and Neo

Re-create db

# Multiple ways to create a Neo4J database

Batch import is non transactional so thats a good way to start

# Some batch import code

Re-create db

```
String line;
while ((line = reader.readLine()) != null) {
    String[] parts = line.split(SPLIT_STRING);

    long fileNodeId = Long.parseLong(parts[0]);
    long nodeId;
    for (int c = 0; c < nodeFields.length; c++) {
        properties[1 + c * 2] = objectFromProperty(parts[c + 1].trim());
    }
    nodeId = db.createNode(map(properties));
    index.add(nodeId, map(properties));
}
```

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

Recreate db

Agile exploration

No waiting please!

Why?

Insight

Hadoop and Neo

No waiting please!

16 (+2) hours

code improved: 3 (+2) hours

Hadoop and Neo

Learning internal structure of Neo4J

How?

Rownumbers

Gather from HDFS

Joining a lot!

Using Neo4J

Learning internal structure of Neo4J

Hadoop and Neo

How?

Rownumbers

Joining a lot!

Gather from HDFS

Using Neo4J

# Learning internal structure of Neo4J

Hadoop and Neo

Learning internal structure of Neo4J

How?

Rownumbers

Joining a lot!

Gather from HDFS

| Name | Date Modified | Size |
| --- | --- | --- |
| messages.log | Friday, March 8, 2013 3:18 PM | 169 KB |
| neostore | Friday, March 8, 2013 3:26 PM | 54 bytes |
| neostore.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.nodestore.db | Friday, March 8, 2013 3:26 PM | 900 KB |
| neostore.nodestore.db.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.propertystore.db | Friday, March 8, 2013 3:26 PM | 16.8 MB |
| neostore.propertystore.db.arrays | Friday, March 8, 2013 3:26 PM | 12…ytes |
| neostore.propertystore.db.arrays.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.propertystore.db.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.propertystore.db.index | Friday, March 8, 2013 3:26 PM | 10 KB |
| neostore.propertystore.db.index.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.propertystore.db.index.keys | Friday, March 8, 2013 3:26 PM | 1 KB |
| neostore.propertystore.db.index.keys.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.propertystore.db.strings | Friday, March 8, 2013 3:26 PM | 1.4 MB |
| neostore.propertystore.db.strings.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.relationshipstore.db | Friday, March 8, 2013 3:26 PM | 20.6 MB |
| neostore.relationshipstore.db.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.relationshiptypestore.db | Friday, March 8, 2013 3:26 PM | 25 bytes |
| neostore.relationshiptypestore.db.id | Friday, March 8, 2013 3:18 PM | 9 bytes |
| neostore.relationshiptypestore.db.names | Friday, March 8, 2013 3:26 PM | 22…ytes |
| neostore.relationshiptypestore.db.names.id | Friday, March 8, 2013 3:18 PM | 9 bytes |

Learning internal structure of Neo4J

Hadoop and Neo

How?

Gather from HDFS

Joining a lot!

Rownumbers

## Learning internal structure of Neo4J

Lucky most work already done by

Chris Gioran

In essence it's all: doubly linked lists sequentially stored on disk

In essence it's all: doubly linked lists sequentially stored on disk

**Node**

| 1 byte | 4 byte | 4 byte |
|---|---|---|
| use flag | 1st rel | 1st prop |

**Property**

| 1 byte | 4 byte | 4 byte | 8 byte | 8 byte | 8 byte | 8 byte |
|---|---|---|---|---|---|---|
| use flag | prev | next | encoded value | name and type | | |

**Relationship**

| 1 byte | 4 byte | 4 byte | 4 byte | 4 byte | 4 byte | 4 byte | 4 byte | 4 byte |
|---|---|---|---|---|---|---|---|---|
| use flag | from | to | type | prev from | next from | prev to | next to | 1st prop |

Learning internal structure of Neo4J

1 {name : BerlinBuzzwords}

speaks_at

is_in {at : 03/04 june 2013}

3 {place : Berlin}
{venue : KulturBrauerei}

comes_to

2 {name : Kris}

Learning internal structure of Neo4J

In essence it's all: doubly linked lists sequentially stored on disk

Nodes

Properties

| 1 | - | - | name : BerlinBuzzwords |
| 2 | - | - | name : Kris |
| 3 | - | 4 | place : Berlin |
| 4 | 3 | - | venue : KulturBrauerei |
| 5 | - | - | at : 03/04 june 2013 |

Relationships

| 1 | 1 | 3 | T | - | 2 | - | 3 | |
| 2 | 2 | 1 | T | - | 3 | 1 | - | - |
| 3 | 2 | 3 | T | 2 | - | 1 | - | - |

{name : BerlinBuzzwords}

speaks_at    is_in    {at : 03/04 june 2013}

comes_to

{place : Berlin}
{venue : KulturBrauerei}

{name : Kris}

Rownumbers

01 ff ff ff ff ff ff ff ff 01 00 00 01 ec 00 00
00 00 01 00 00 00 67 00 00 00 01 01 00 00 00 7e
00 00 00 02 01 00 00 00 02 00 00 00 03 01 00 00
00 08 00 00 00 04 01 00 00 00 03 00 00 00 05 01
00 00 00 04 00 00 00 07 01 00 00 02 5c 00 00 00
08 01 00 00 01 f6 00 00 00 09 01 00 00 00 08 00
00 00 0a 01 00 00 00 0c 00 00 00 0b 01 00 00 00
52 00 00 00 0c 01 00 00 00 24 00 00 00 0d 01 00
00 00 0c 00 00 00 0e 01 00 00 02 00 00 00 00 0f
01 00 00 00 0e 00 00 00 10 01 00 00 00 10 00 00

Position in the file matters

Learning internal structure of Neo4J

Hadoop and Neo

How?

Rownumbers

Joining a lot!

Gather from HDFS

# Rownumbers

We need a rownumber generator

Hadoop and Neo

Learning internal structure of Neo4J

How?

Rownumbers

Joining a lot!

Gather from HDFS

# Rownumbers

*nix has one build into `cat`

Rownumbers

```
CAT(1)                      BSD General Commands Manual                        CAT(1)

NAME
     cat -- concatenate and print files

SYNOPSIS
     cat [-benstuv] [file ...]

DESCRIPTION
     The cat utility reads files sequentially, writing them to the standard output.  The file operands are pro-
     cessed in command-line order.  If file is a single dash (`-') or absent, cat reads from the standard
     input.  If file is a UNIX domain socket, cat connects to it and then reads it until EOF.  This complements
     the UNIX domain binding capability available in inetd(8).

     The options are as follows:

     -b        Number the non-blank output lines, starting at 1.

     -e        Display non-printing characters (see the -v option), and display a dollar sign (`$') at the end of
               each line.

     -n        Number the output lines, starting at 1.

     -s        Squeeze multiple adjacent empty lines, causing the output to be single spaced.

     -t        Display non-printing characters (see the -v option), and display tab characters as `^I'.

     -u        Disable output buffering.

     -v        Display non-printing characters so they are visible.  Control characters print as `^X' for con-
               trol-X; the delete character (octal 0177) prints as `^?'.  Non-ASCII characters (with the high bit
               set) are printed as `M-' (for meta) followed by the character for the low 7 bits.

EXIT STATUS
     The cat utility exits 0 on success, and >0 if an error occurs.
```
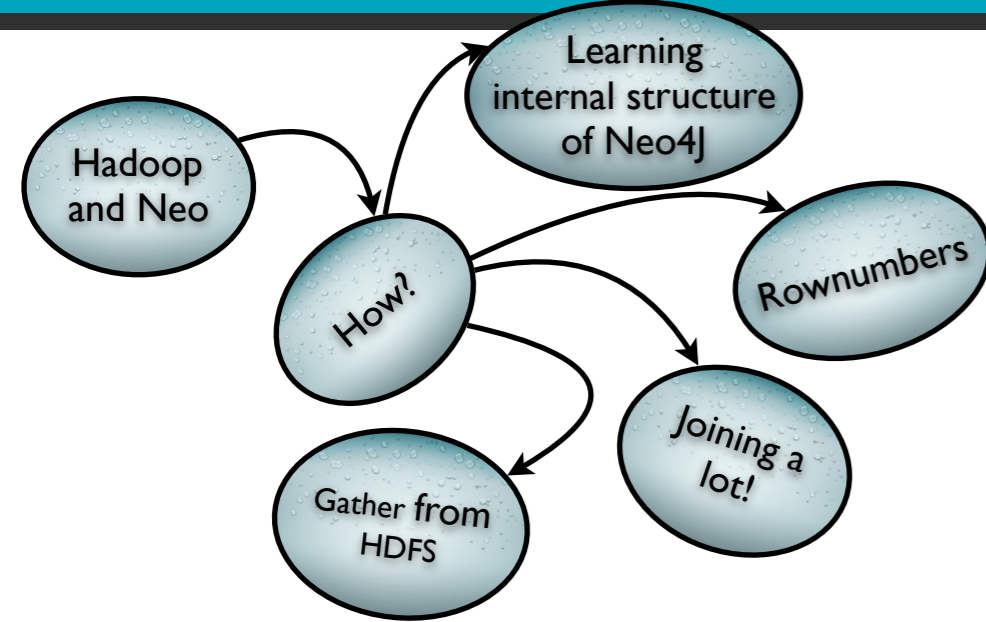
Rownumbers

We need a distributed 'cat -n'

# Rownumbers

We need a distributed 'cat -n'

| to convert: | | into: | | |
|---|---|---|---|---|
| ABC | | 0 | ABC |
| DEF | | 1 | DEF |
| GHI | | 2 | GHI |
| JKL | | 3 | JKL |
| MNO | | 4 | MNO |
| PQR | | 5 | PQR |
| STU | | 6 | STU |
| VWX | | 7 | VWX |
| YZ0 | | 8 | YZ0 |

Rownumbers

Easy way out is a 1 (one) reducer job doing the numbering. But that's no fun right!

# Rownumbers

1. Run the data through a mapper and have the mapper emit each record AS IS
2. Have each mapper keep track of how many records it sends to each reducer
3. Have each reducer process the count records and accumulate the counts from each mapper it receives
4. Have each reducer emit each record prepended by a row ID starting the ID sequence at the number calculated

# Rownumbers

1. Run the data through a mapper and have the mapper emit each record AS IS
2. Have each mapper keep track of how many records it sends to each reducer

**mapper:**

```
setup:
    initialize counters

map:
    read input ==> emit as is, increment the correct counter

cleanup:
    emit all counters
```

ABC  (hash=0)
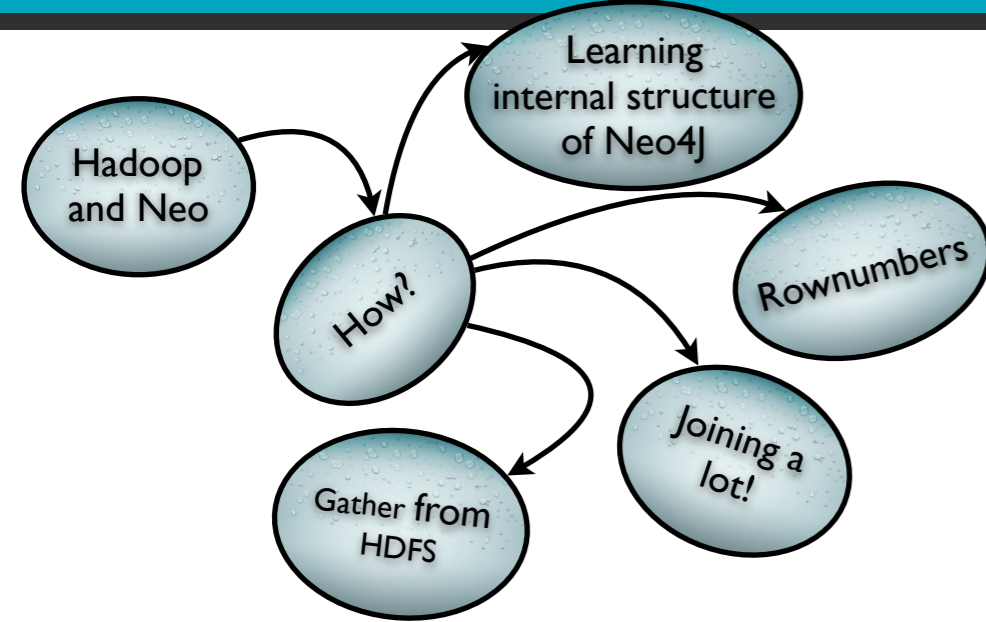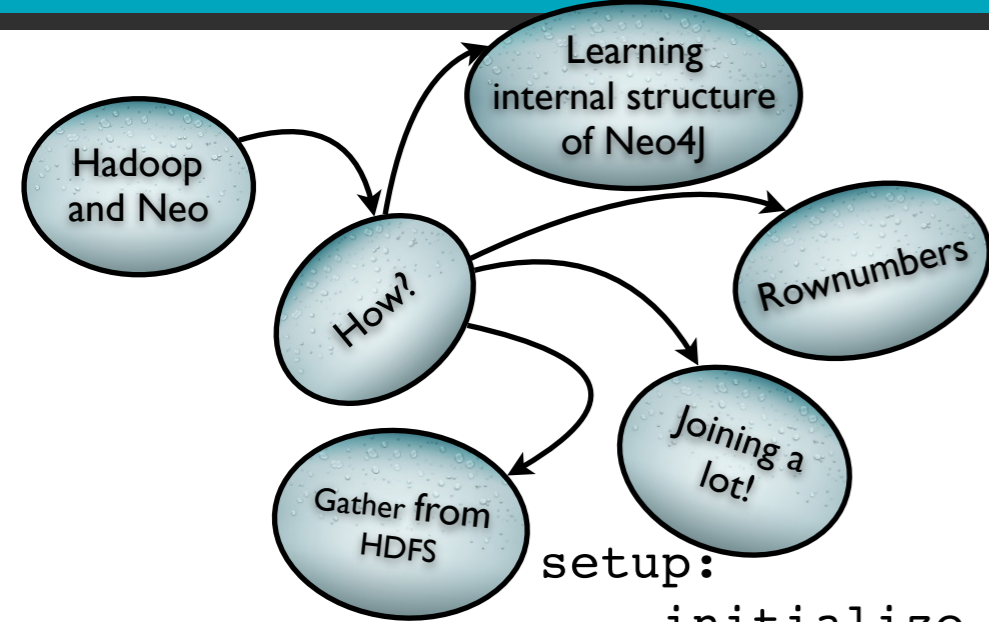DEF  (hash=1)
GHI  (hash=2)
JKL  (hash=0)
MNO (hash=1)
------------- split
PQR  (hash=2)
STU  (hash=0)
VWX (hash=1)
YZ0  (hash=1)

**mapper 0:**

```
setup:
    initialize counters = [0, 0, 0]
map:
    input "ABC" ==> emit "ABC", increment counters[0]
    input "DEF" ==> emit "DEF", increment counters[1]
    input "GHI" ==> emit "GHI", increment counters[2]
    input "JKL" ==> emit "JKL", increment counters[0]
    input "MNO" ==> emit "MNO", increment counters[1]    //now [2, 2, 1]
cleanup:
    emit counter for partition 0: none
    emit counter for partition 1: counters[0] = 2
    emit counter for partition 2: counters[0] + counters[1] = 4
```

**mapper 1:**

```
setup:
    initialize counters = [0, 0, 0]
map:
    input "PQR" ==> emit "PQR", increment counters[2]
    input "STU" ==> emit "STU", increment counters[0]
    input "VWX" ==> emit "VWX", increment counters[1]
    input "YZ0" ==> emit "YZ0", increment counters[1]    //now [1, 2, 1]
cleanup:
    emit counter for partition 0: none
    emit counter for partition 1: counters[0] = 1
    emit counter for partition 2: counters[0] + counters[1] = 3
```
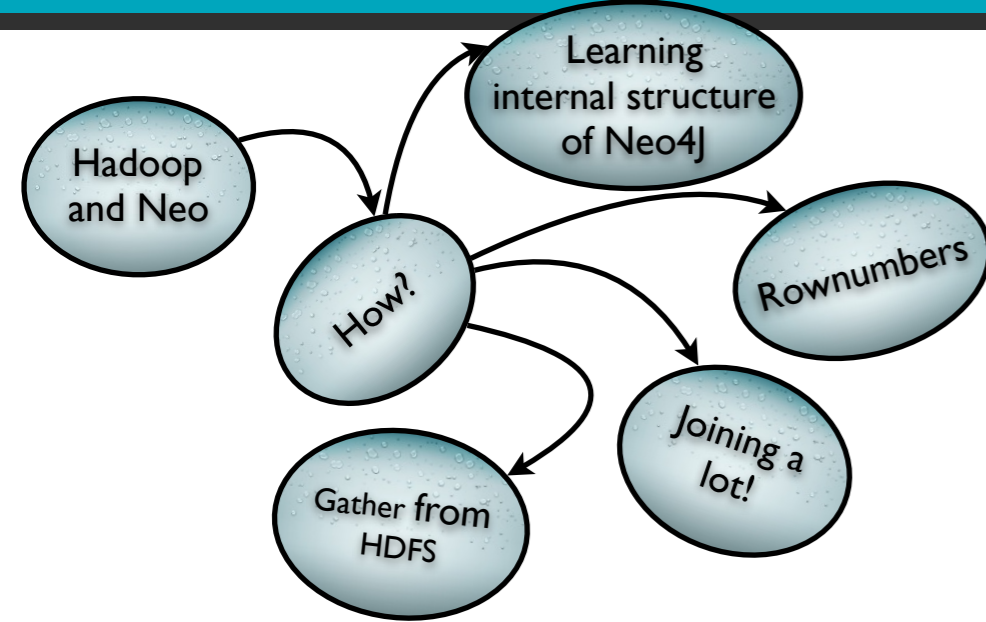
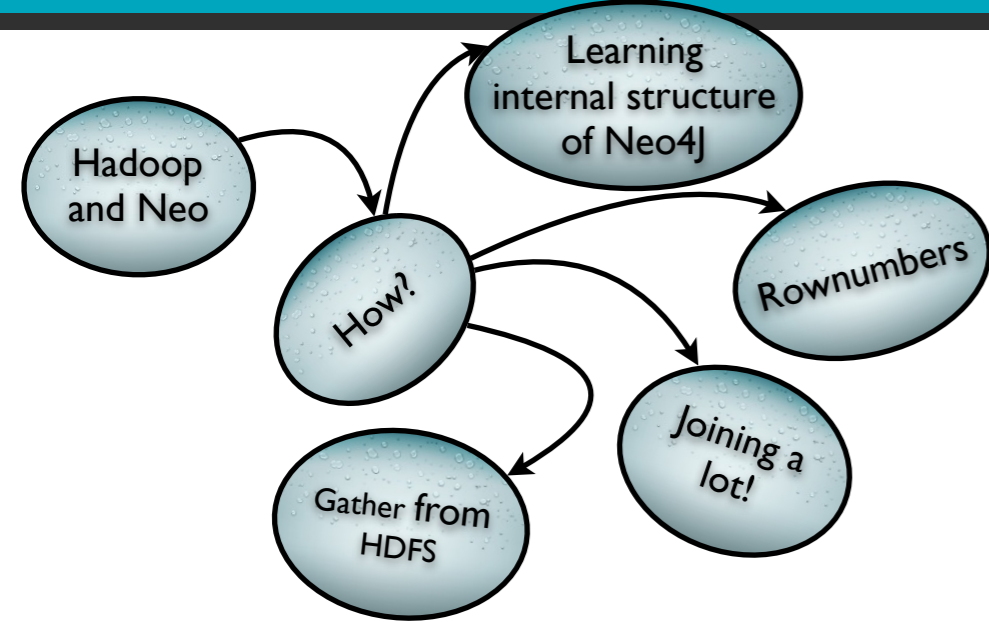1. Have each reducer process the count records and accumulate the counts from each mapper it receives
2. Have each reducer emit each record prepended by a row ID starting the ID sequence at the number calculated

**reducer:**

```
reduce:
    initialize offset = 0
    calculate offset to start numbering based on all counter records from the mappers
    read input ==> emit current offset + input, increment offset
```

# Rownumbers

**reducer 0:**

```
reduce:
    initialize offset = 0
    input "ABC" ==> emit "offset <tab> ABC", increment offset    //emits 0<tab>ABC
    input "JKL" ==> emit "offset <tab> JKL", increment offset    //emits 1<tab>ABC, and so on...
    input "STU" ==> emit "offset <tab> STU", increment offset    //last emitted offset is 2
```

**reducer 1:**

```
reduce:
    initialize offset = 0
    input counter with value 2 ==> offset = offset + 2
    input counter with value 1 ==> offset = offset + 1   //now offset == 3
    input "DEF" ==> emit "offset <tab> DEF", increment offset    //emits 3<tab>DEF
    input "MNO" ==> emit "offset <tab> MNO", increment offset    //emits 4<tab>MNO, and so on...
    input "VWX" ==> emit "offset <tab> VWX", increment offset
    input "XY0" ==> emit "offset <tab> XY0", increment offset    //last emitted offset is 6
```
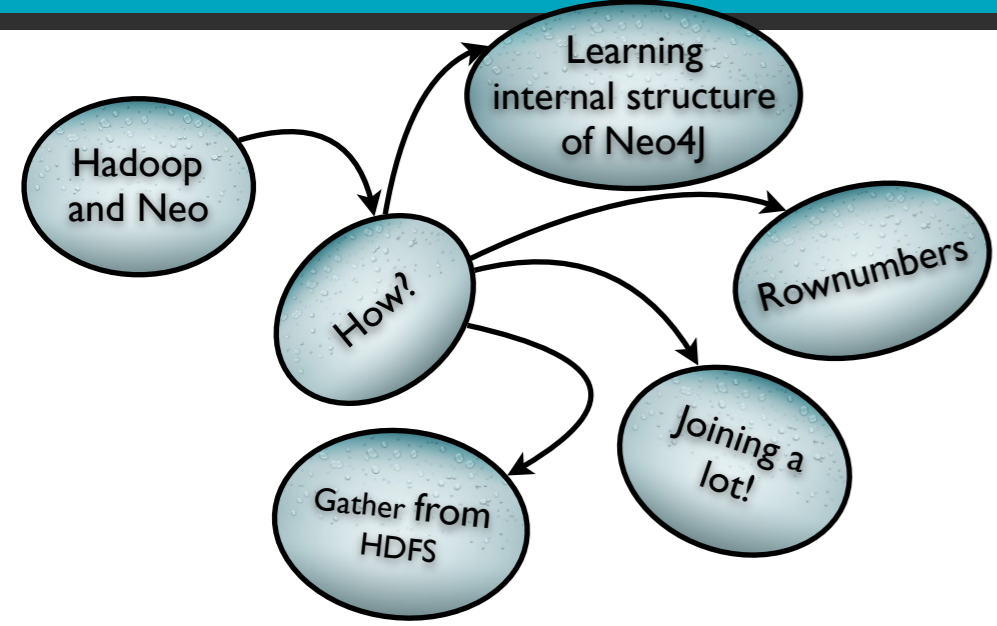
**reducer 2:**

```
reduce:
    initialize offset = 0
    input counter with value 4 ==> offset = offset + 4
    input counter with value 3 ==> offset = offset + 3   //now offset == 7
    input "GHI" ==> emit "offset <tab> GHI", increment offset    //emits 7<tab>GHI
    input "PQR" ==> emit "offset <tab> PQR", increment offset    //emits 8<tab>PQR
```
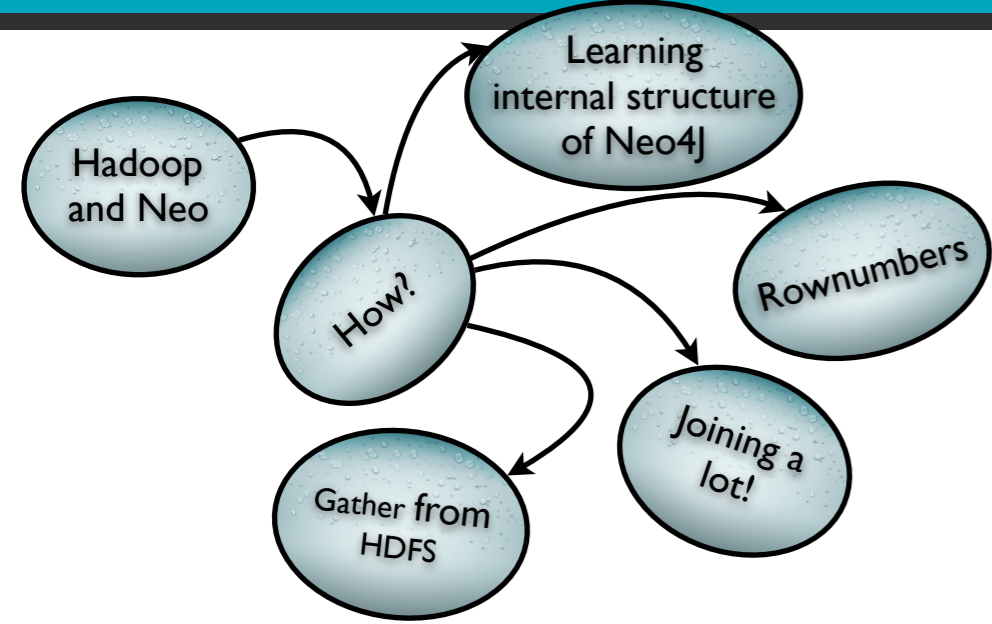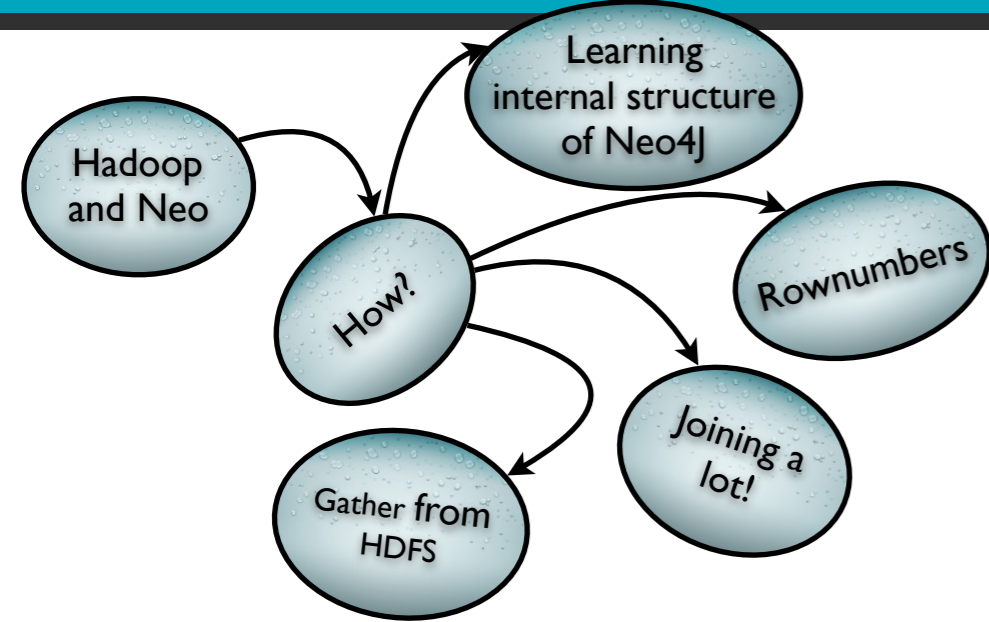
Rownumbers

Learning internal structure of Neo4J

Hadoop and Neo

How?

Rownumbers

Joining a lot!

Gather from HDFS

Is that all?

Rownumbers

# Is that all?

Rownumbers

## A custom partitioner:

No!

```java
public static class Partitioner extends Partitioner<ByteWritable, RowNumberWritable> {
    @Override
    public int getPartition(ByteWritable key, RowNumberWritable value, int numPartitions) {
        if (key.get() == (byte) RowNumberJob.COUNTER_MARKER) {
            return value.getPartition();
        } else {
            return Partitioner.partitionForValue(value, numPartitions);
        }
    }

    public static int partitionForValue(RowNumberWritable value, int numPartitions) {
        return (value.getValue().hashCode() & Integer.MAX_VALUE) % numPartitions;
    }
}
```
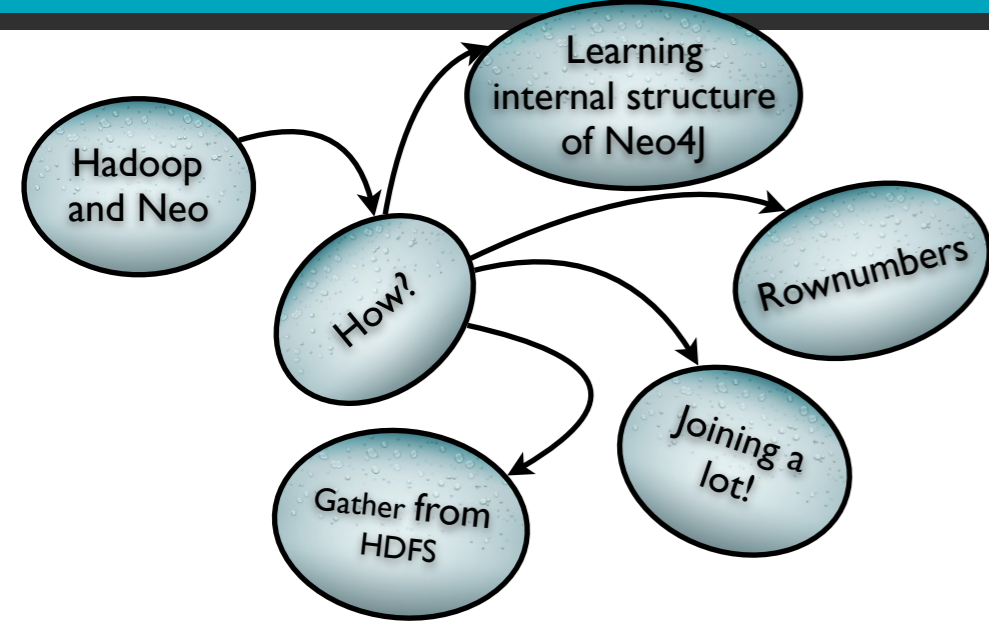
## A custom grouping comparator:

```java
public class IndifferentComparator implements RawComparator<ByteWritable> {
    @Override
    public int compare(ByteWritable left, ByteWritable right) {
        return 0;
    }
}
```
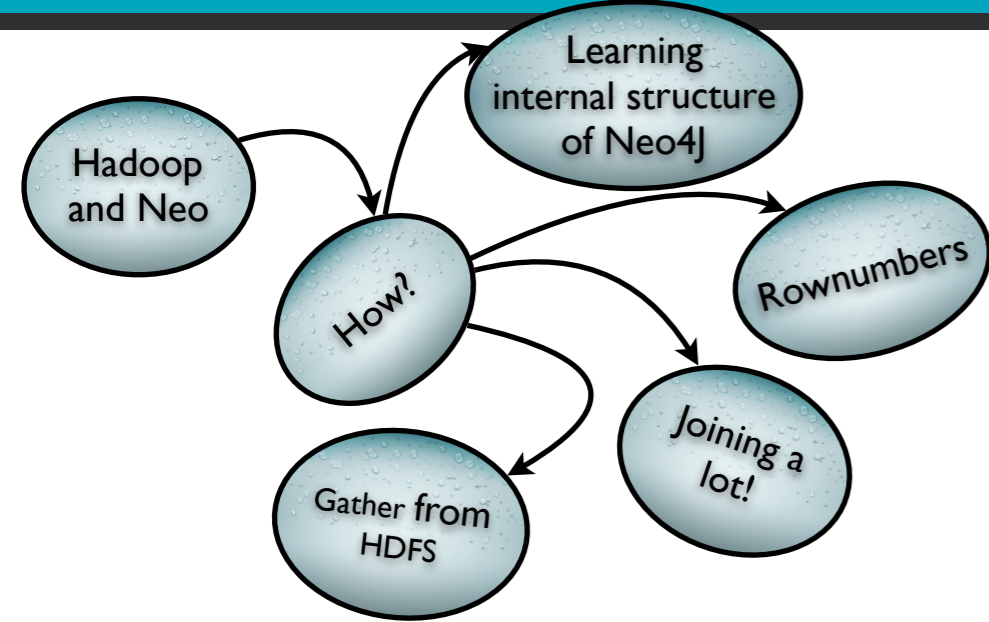
## A custom writable

Step 1:     Prepare properties
Step 2a:   Output properties
Step 2b:   Output nodes/edges and first property reference
Step 3:     Join edges and nodes to get from-id and to-id
Step 4:     Output nodes with first edge and first property reference
Step 5:     Output edges
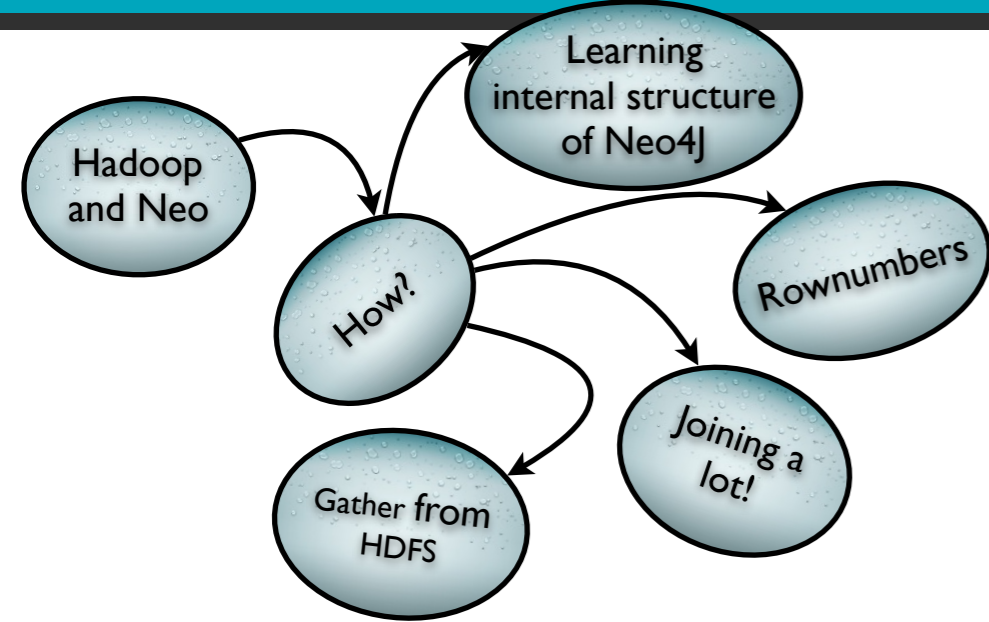
Joining a lot!

Step 1:    Prepare properties

Propertierecords can hold multiple properties
We need to find the first property reference for each node and edge,
so we need to know the total properties structure first

Joining a lot!
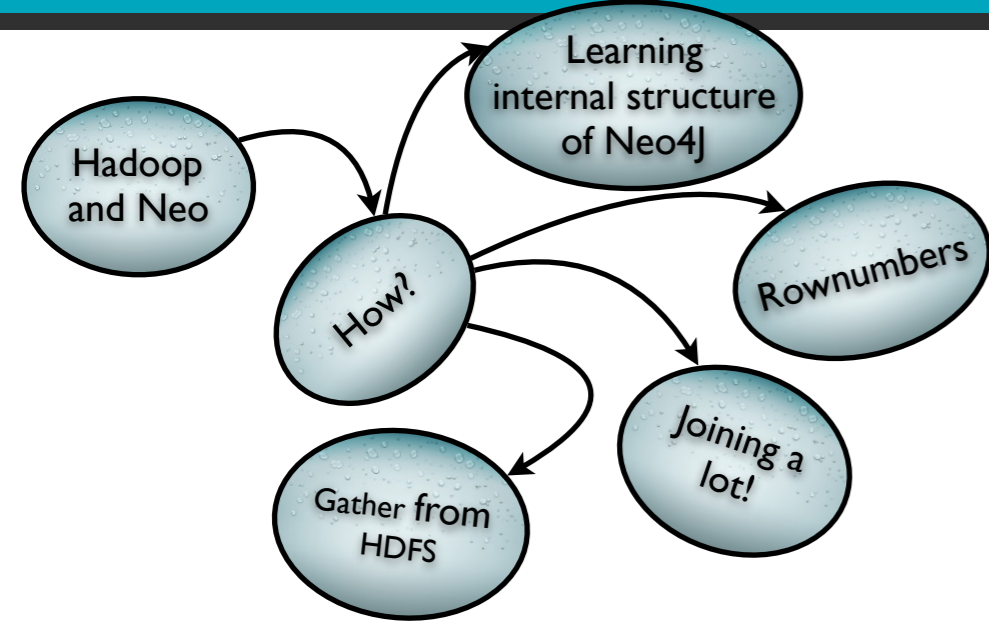
Step 2a:   Output properties
Step 2b:   Output nodes/edges and first property reference

We output the bytearray structure of the neo4j property files

And we can output the node id with functional id and the first property pointer

For edges we output the edgeid, fromnode, tonode and first property pointer
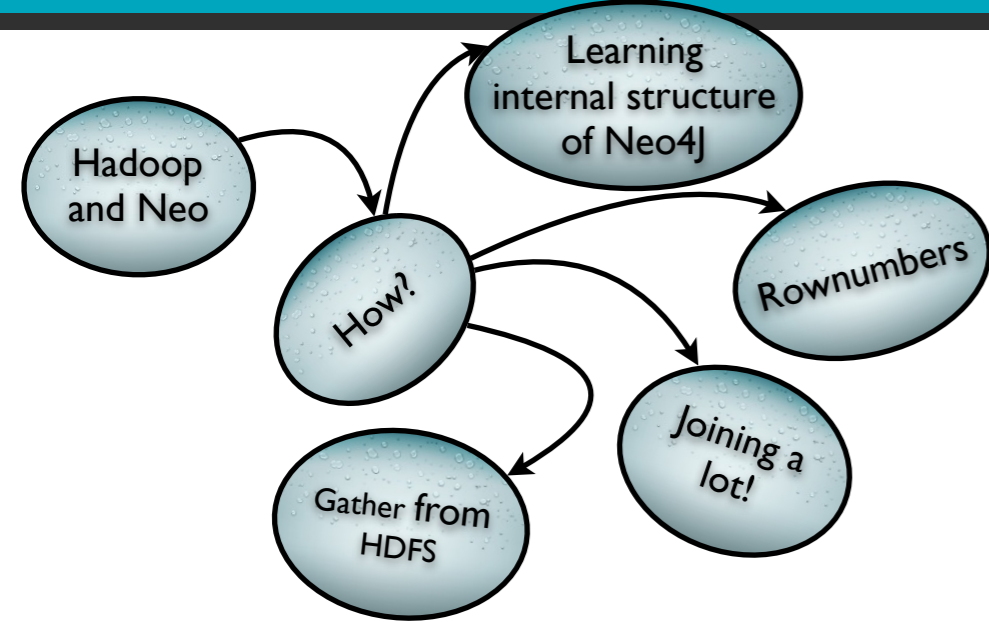
Joining a lot!

Step 3:    Join edges and nodes to get from-id and to-id

Remember we added rownumbers to the data and we need to use them as our pointers. Not the original functional id.

Learning internal structure of Neo4J

Hadoop and Neo

How?

Rownumbers

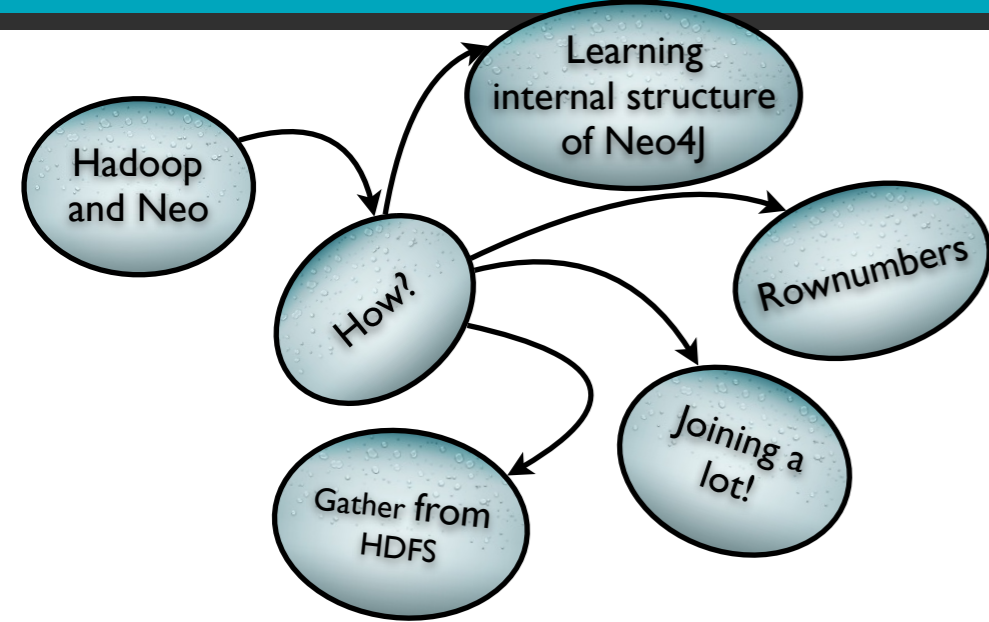Joining a lot!

Gather from HDFS

Joining a lot!

Step 4:     Output nodes with first edge and first property reference

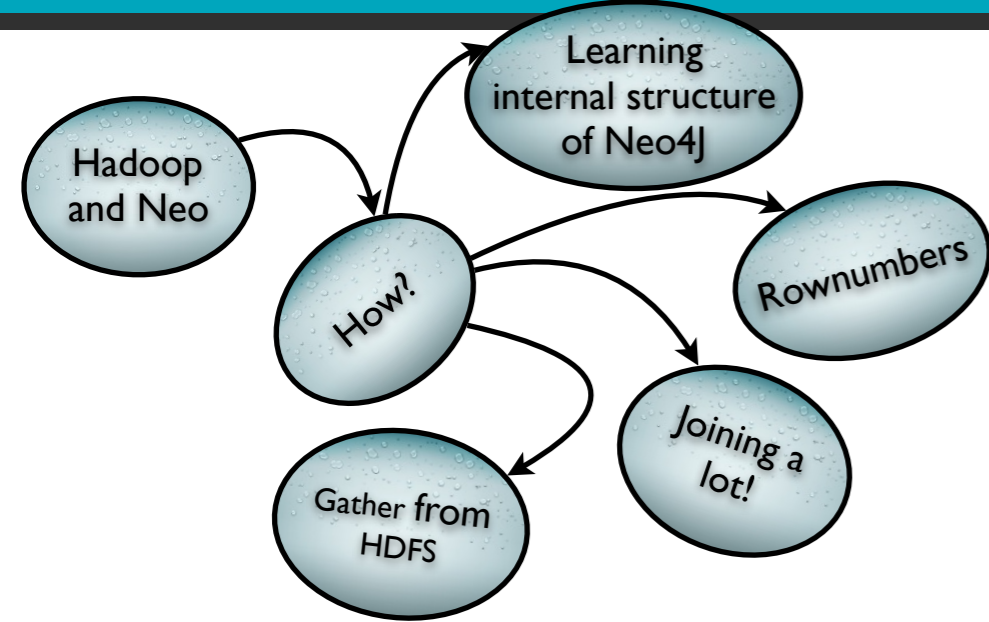It's actually more complicated. Need to selfjoin all edges of a node to determine the first.

Joining a lot!

Step 5:    Output edges

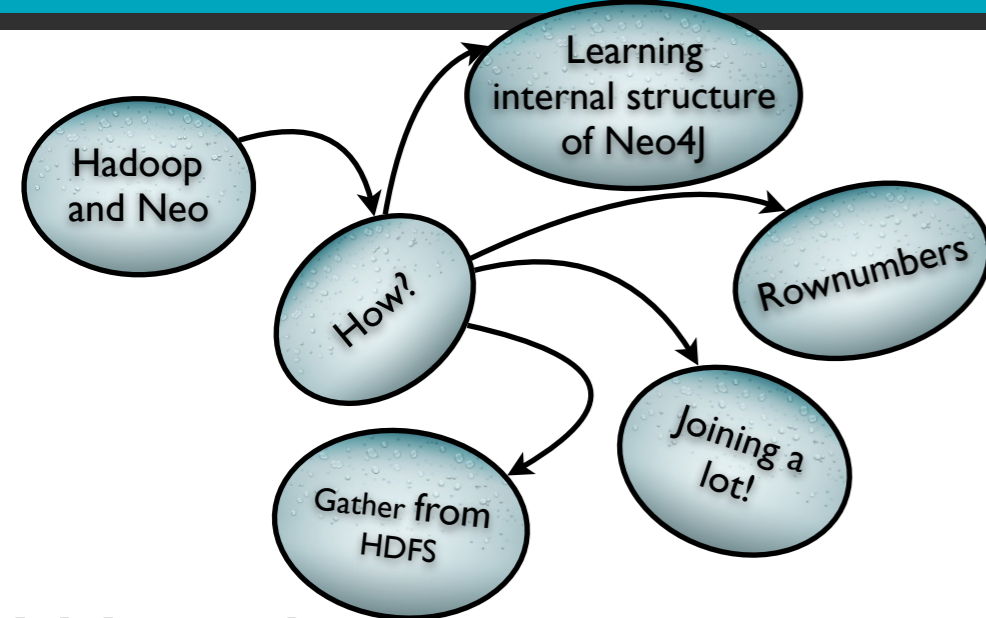Making sure you have them sorted to put in the next and previous
edge reference

Gather from HDFS

Just a simple

```
hadoop fs -cat <HDFS PATH>/neostore.nodestore.db/part-r-*
```
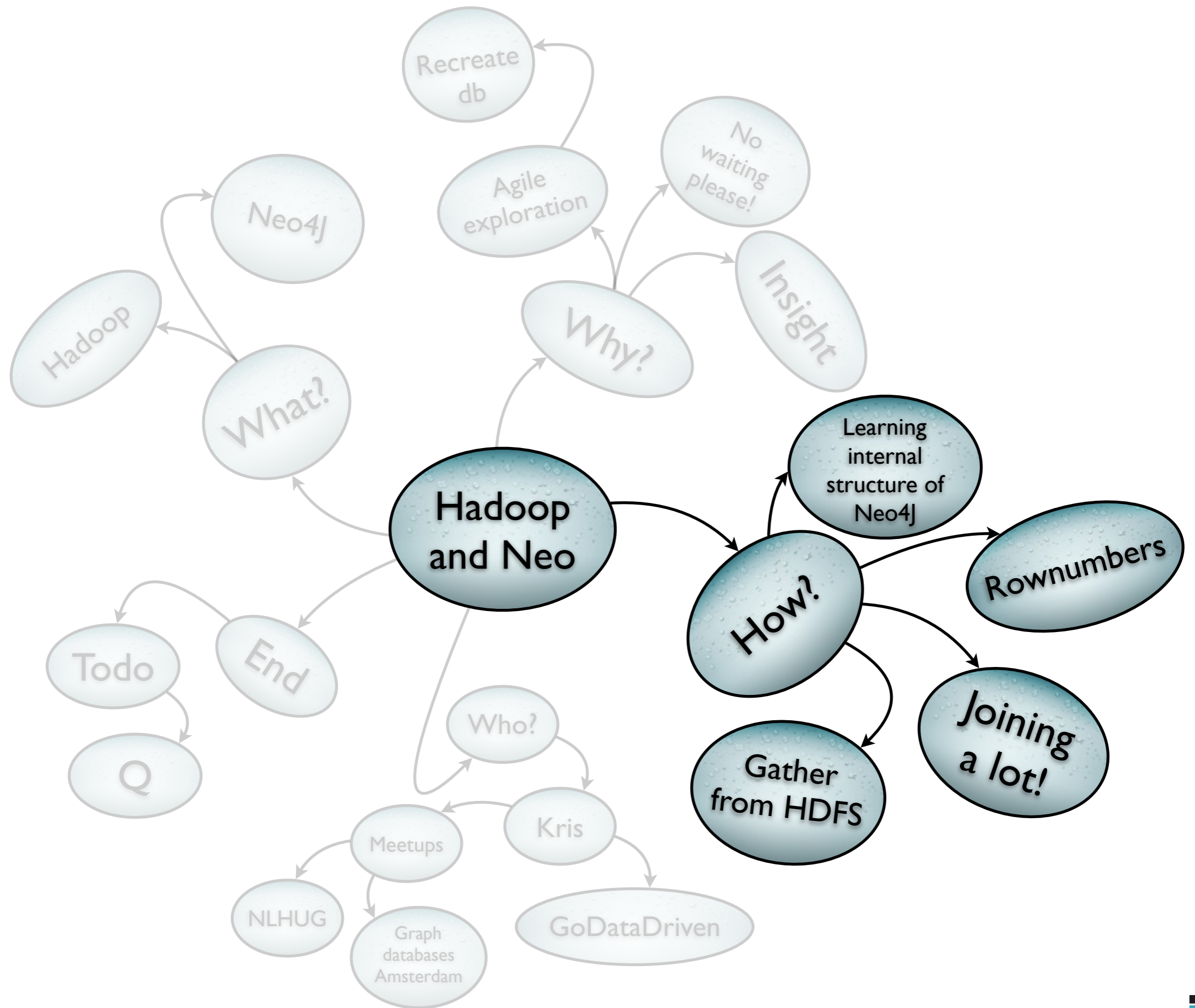
## Gather from HDFS
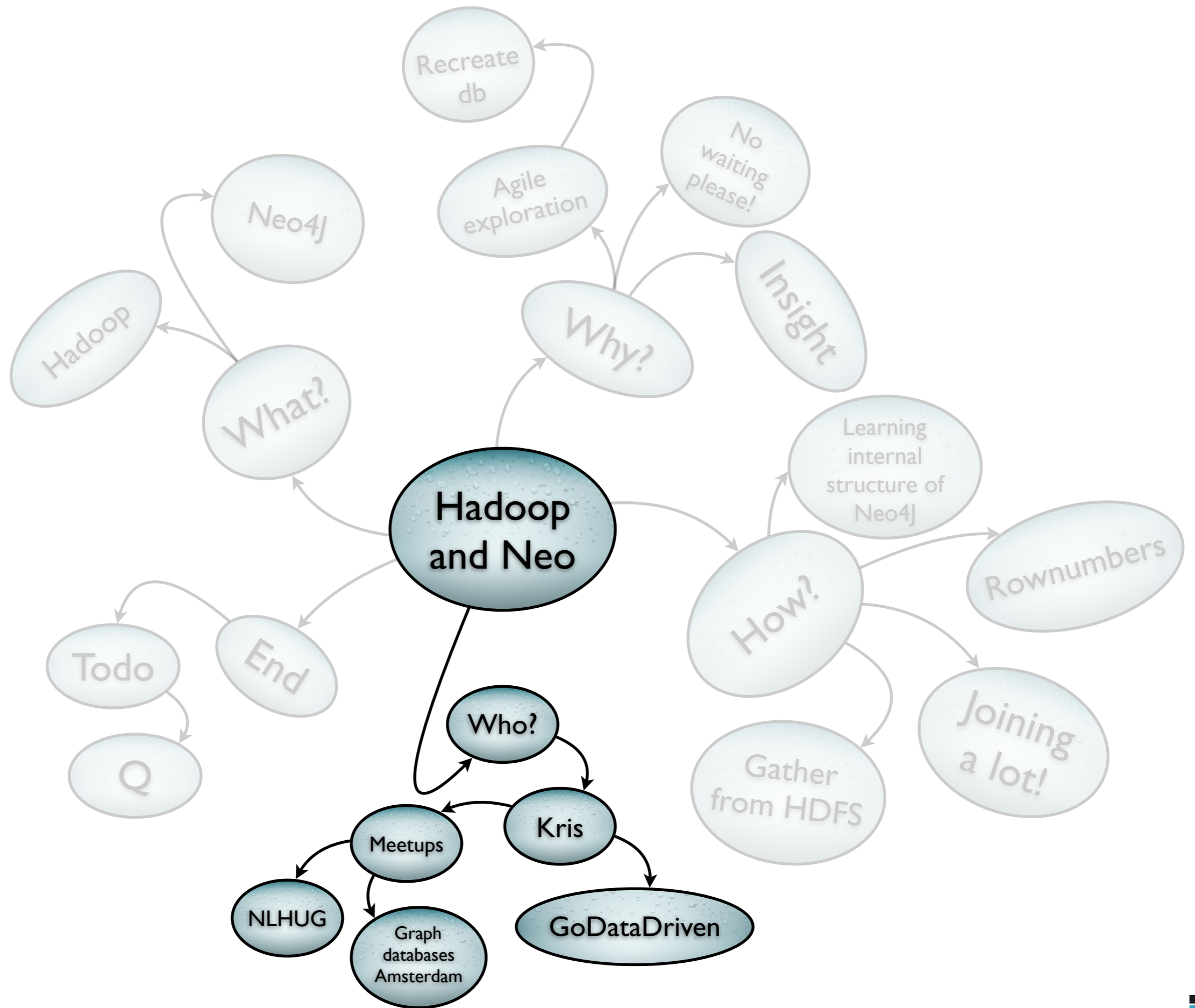
# Well a bit more to get all data

```bash
1   #!/bin/bash
2
3   rm -rf ./graph.db
4   mkdir graph.db/
5
6   TO=./graph.db/
7   FROM=${1}
8   hadoop fs -get ${FROM}/neostore ${TO}
9   hadoop fs -get ${FROM}/neostore.id ${TO}
10  hadoop fs -get ${FROM}/neostore.nodestore.db.id ${TO}
11  hadoop fs -get ${FROM}/neostore.relationshipstore.db.id ${TO}
12  hadoop fs -get ${FROM}/neostore.relationshiptypestore.db ${TO}
13  hadoop fs -get ${FROM}/neostore.relationshiptypestore.db.id ${TO}
14  hadoop fs -get ${FROM}/neostore.relationshiptypestore.db.names ${TO}
15  hadoop fs -get ${FROM}/neostore.relationshiptypestore.db.names.id ${TO}
16
17  hadoop fs -get ${FROM}/properties/neostore.propertystore.db.* ${TO}
18
19  hadoop fs -cat ${FROM}/neostore.nodestore.db/part-r-* > ${TO}/neostore.nodestore.db
20  hadoop fs -cat ${FROM}/neostore.relationshipstore.db/part-r-* > ${TO}/neostore.relationshipstore.db
21
22  hadoop fs -cat ${FROM}/nodeproperties/propertystore.db/props-r-* ${FROM}/edgeproperties/propertystore.db/props-r-* ${FROM}/p
23  hadoop fs -cat ${FROM}/properties/neostore.propertystore.db.strings.header ${FROM}/nodeproperties/propertystore.db/strings-
24
25  rm ${TO}/*.footer
26  rm ${TO}/*.header
27  exit
28
```
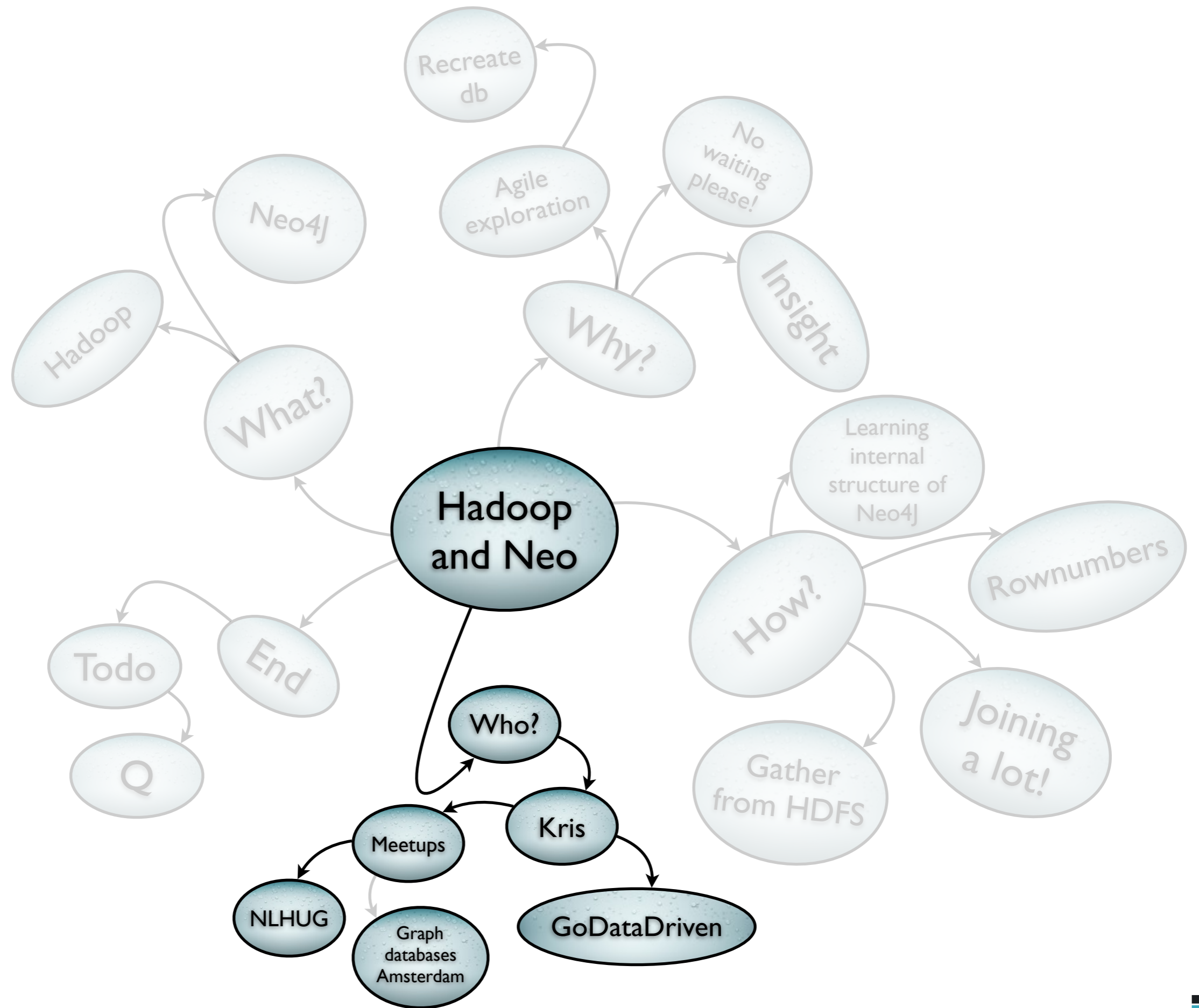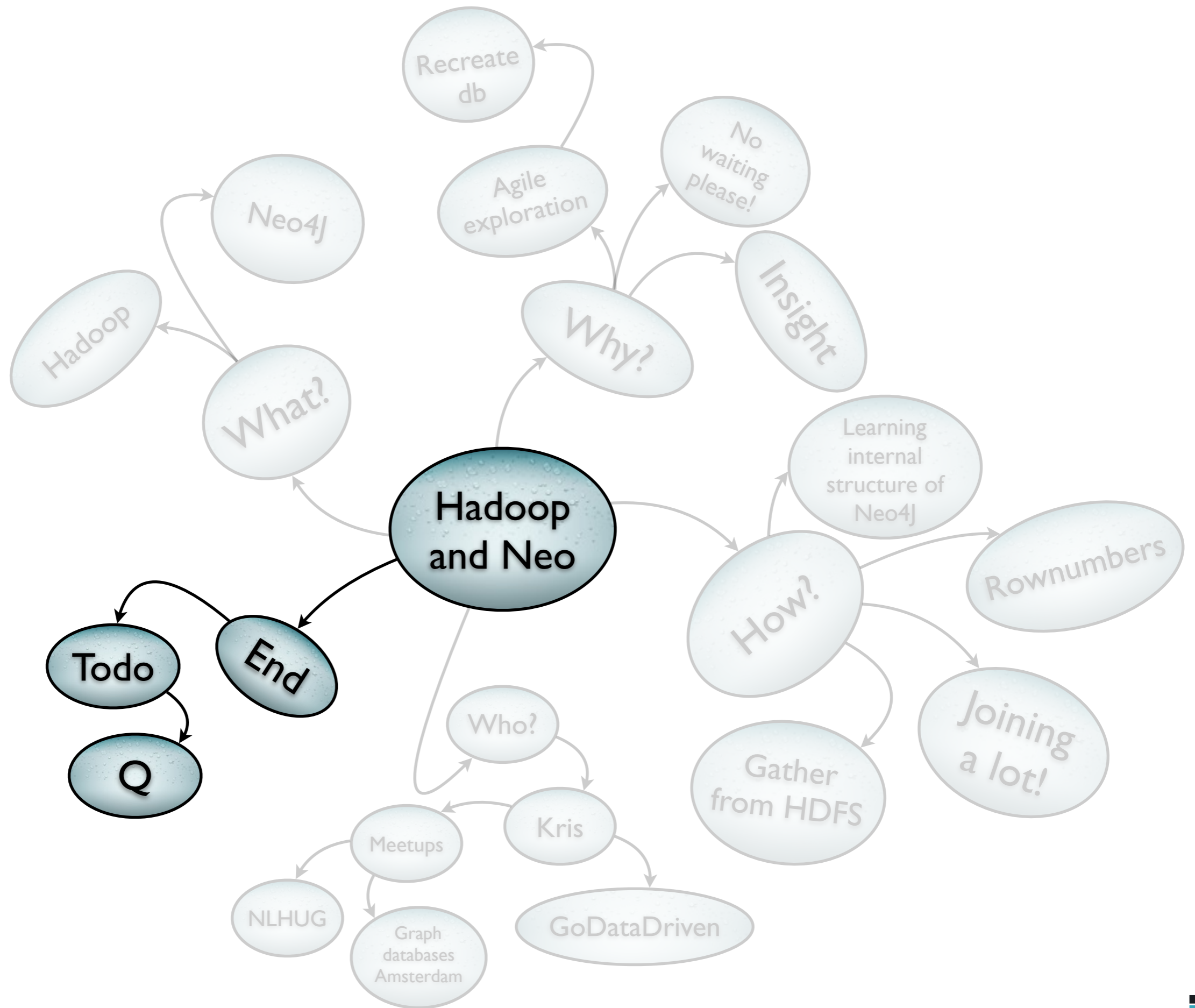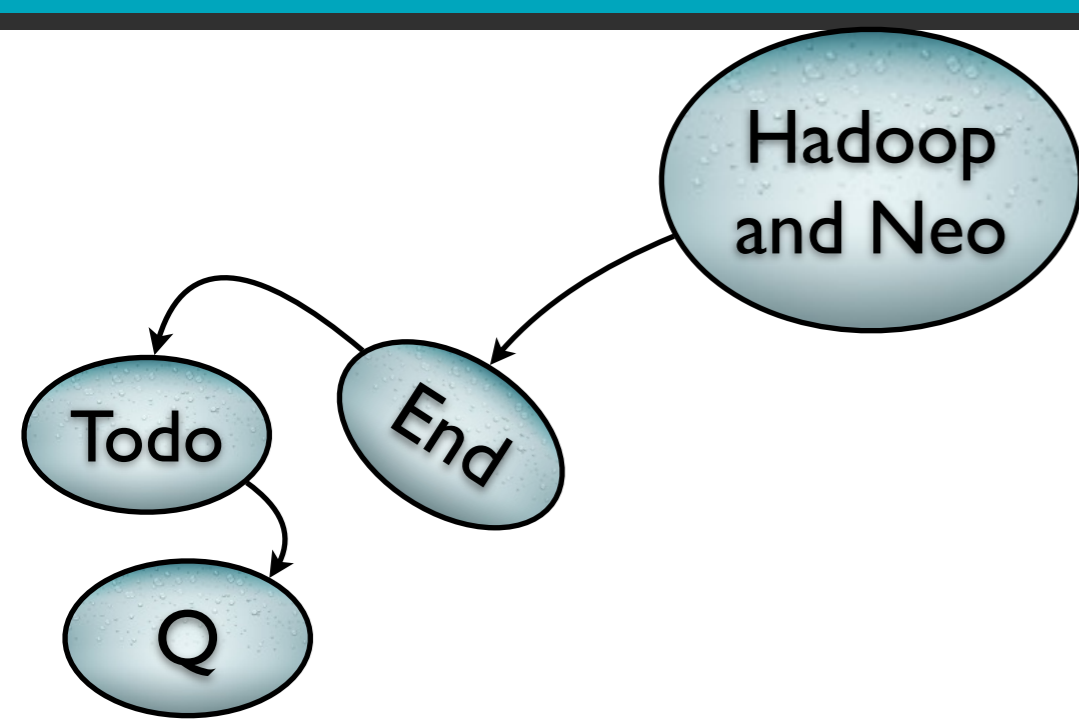
```
START me = node("Kris Geusebroek")
MATCH me-[:WORKS_AT]-("GoDataDriven")
    , me-[:MEMBER]-("meetup.com/NLHUG")
    , me-[:MEMBER]-("meetup.com/GraphdbAmsterdam")
WHERE me.twitter = "@krisgeus"
  AND me.github = "github.com/krisgeus"
  AND me.email = "krisgeusebroek@godatadriven.com"
RETURN "THANK YOU!"
```
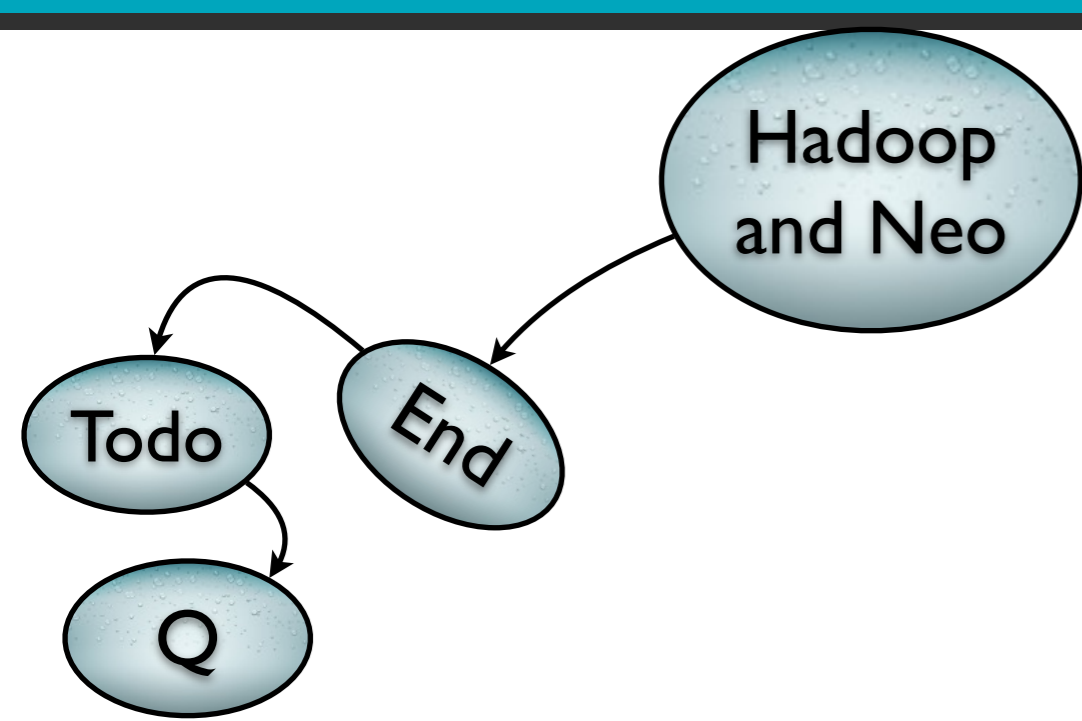
# Todo:
- Indexes
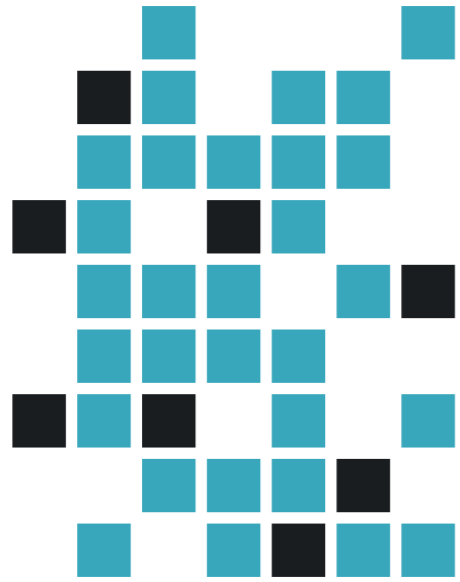- Array properties
- Neo4J 2.0 compatibility
- Inverse direction

Code is on

https://github.com/krisgeus/graphs

GoDataDriven

We're hiring / Questions? / Thank you!

*Kris Geusebroek*
*Big Data Hacker*

*@krisgeus*
*krisgeusebroek@godatadriven.com*