

# Neo4j - a graph DB

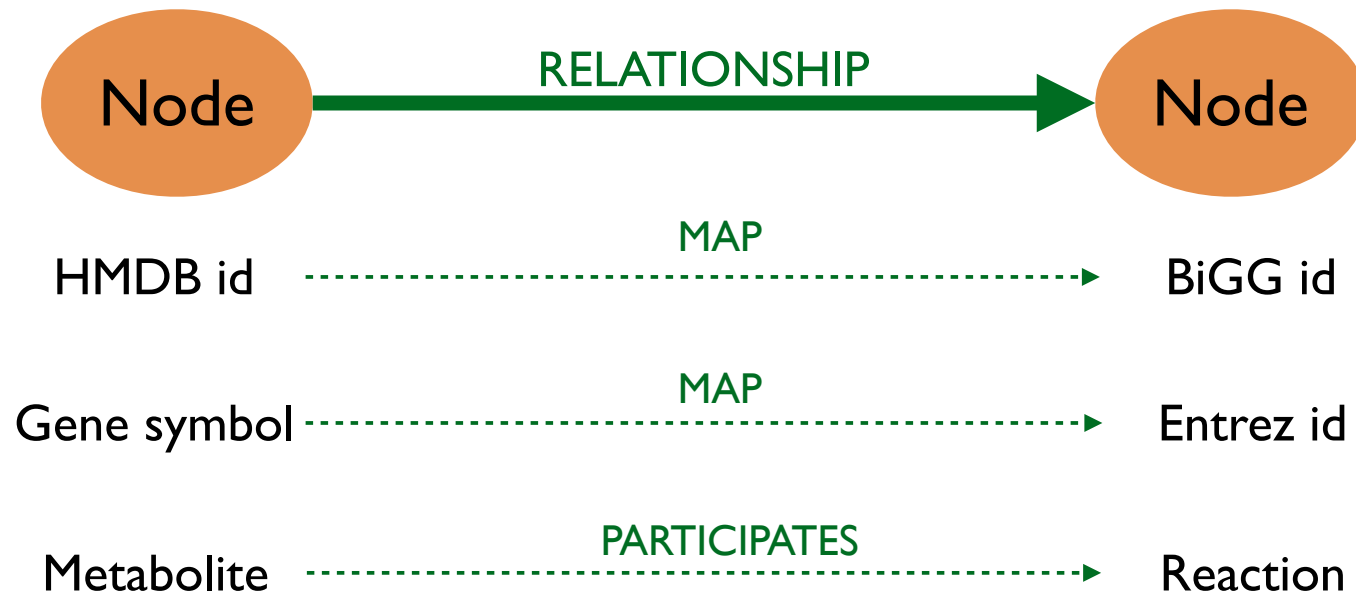
**Maria Wörheide     July 27, 2018**

# Overview

- Introduction to Neo4j and Cypher
- Data model
- Git repository
- Neo4j Servers
- Code

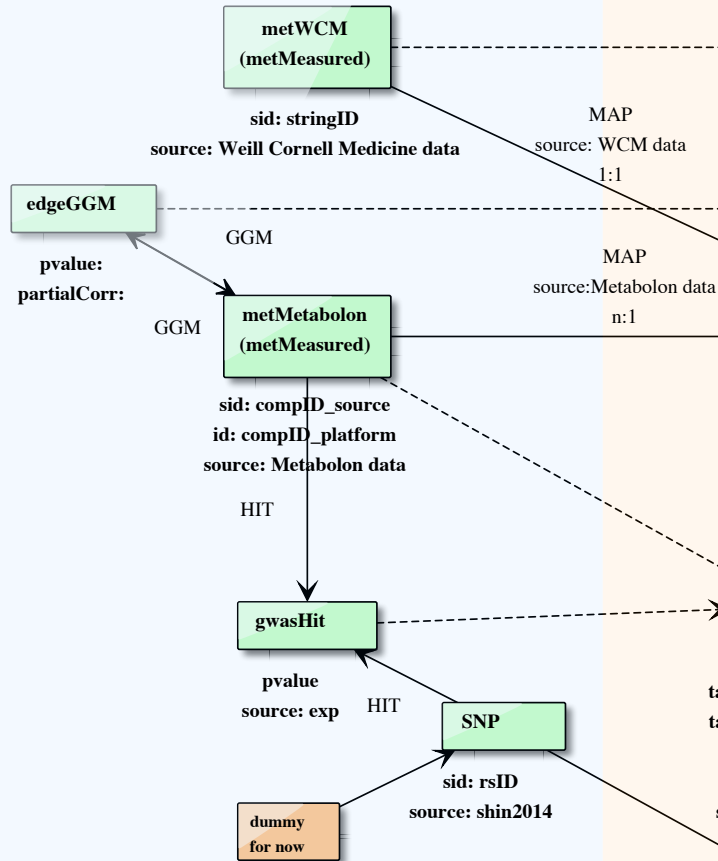
# neo4j - a native graph DB

“Today's world is no longer driven by data – it's driven by the connections between them<sup>1)</sup>. “

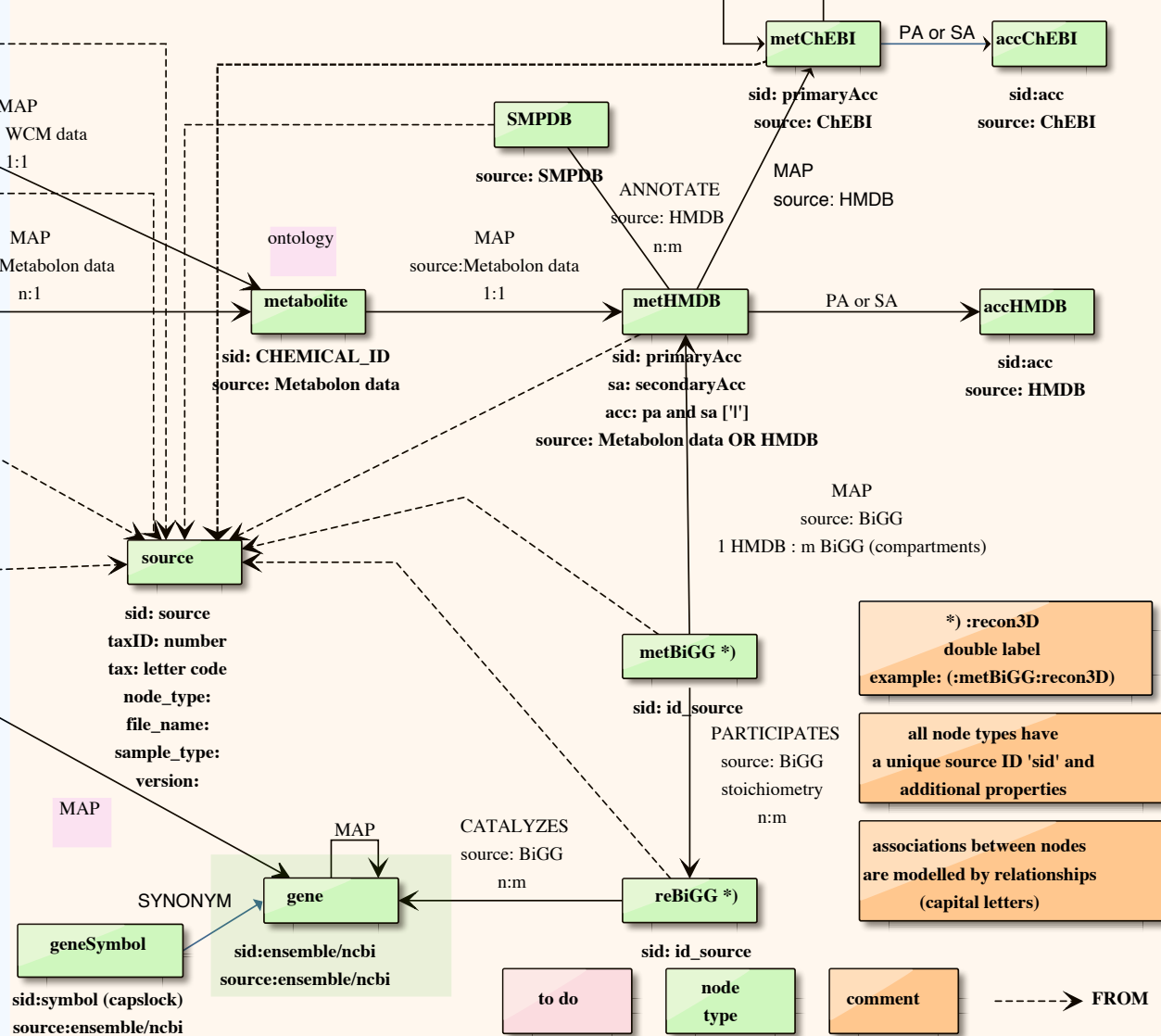


<sup>1)</sup><https://neo4j.com/product/>

## Experimental-World



## Knowledge-World



# Data Model

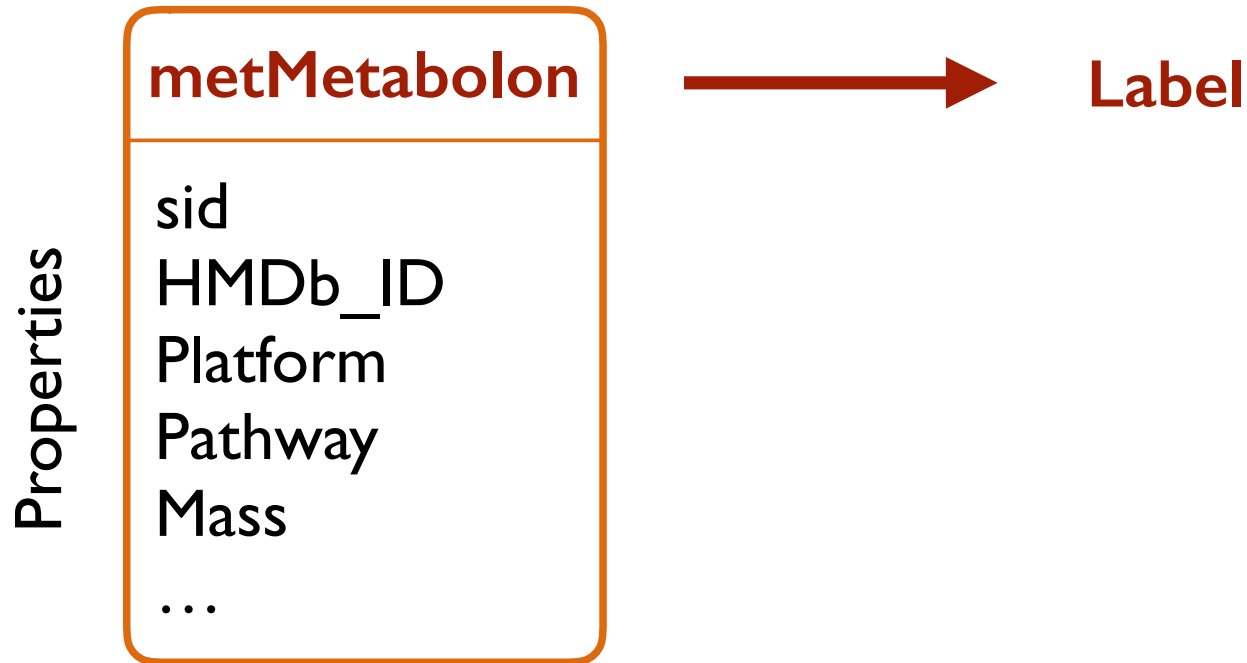
# Neo4j - Nodes

Properties

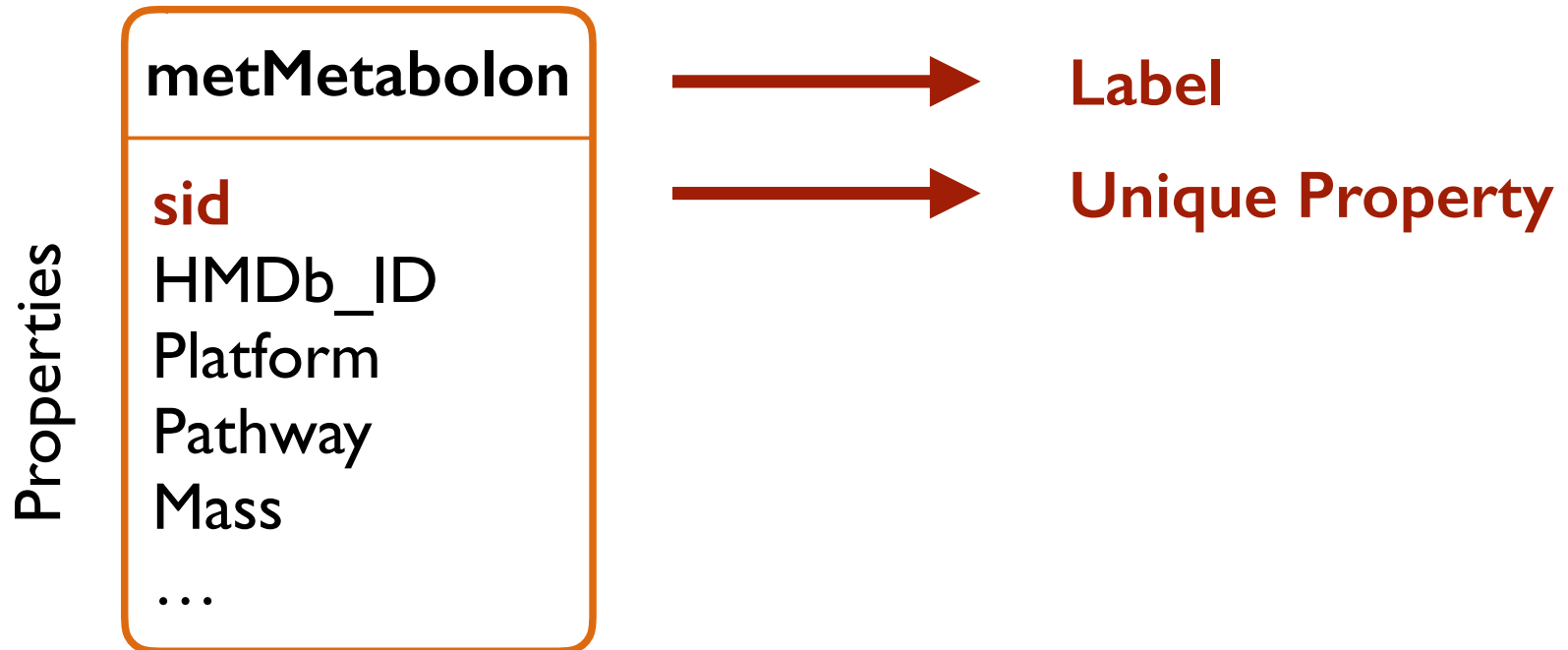
**metMetabolon**

sid  
HMDB\_ID  
Platform  
Pathway  
Mass  
...

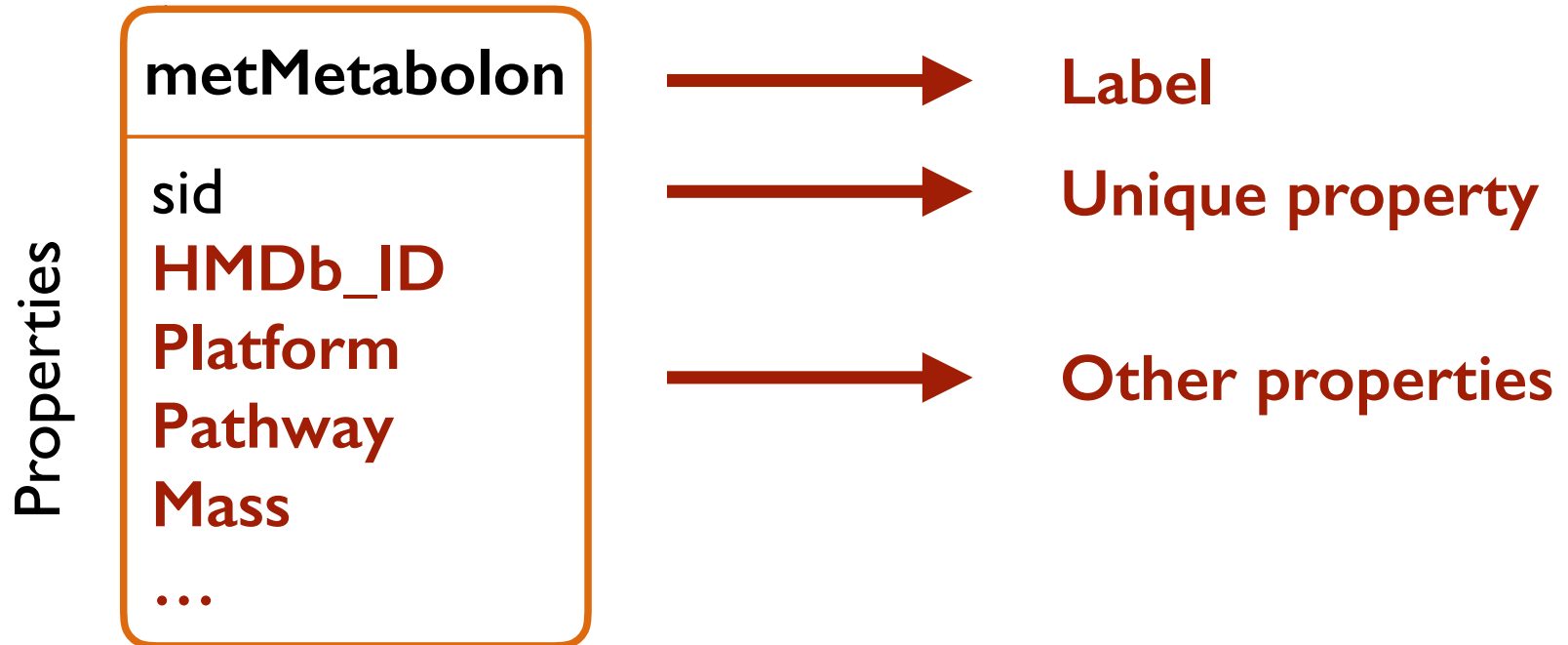
# Neo4j - Nodes



# Neo4j - Nodes

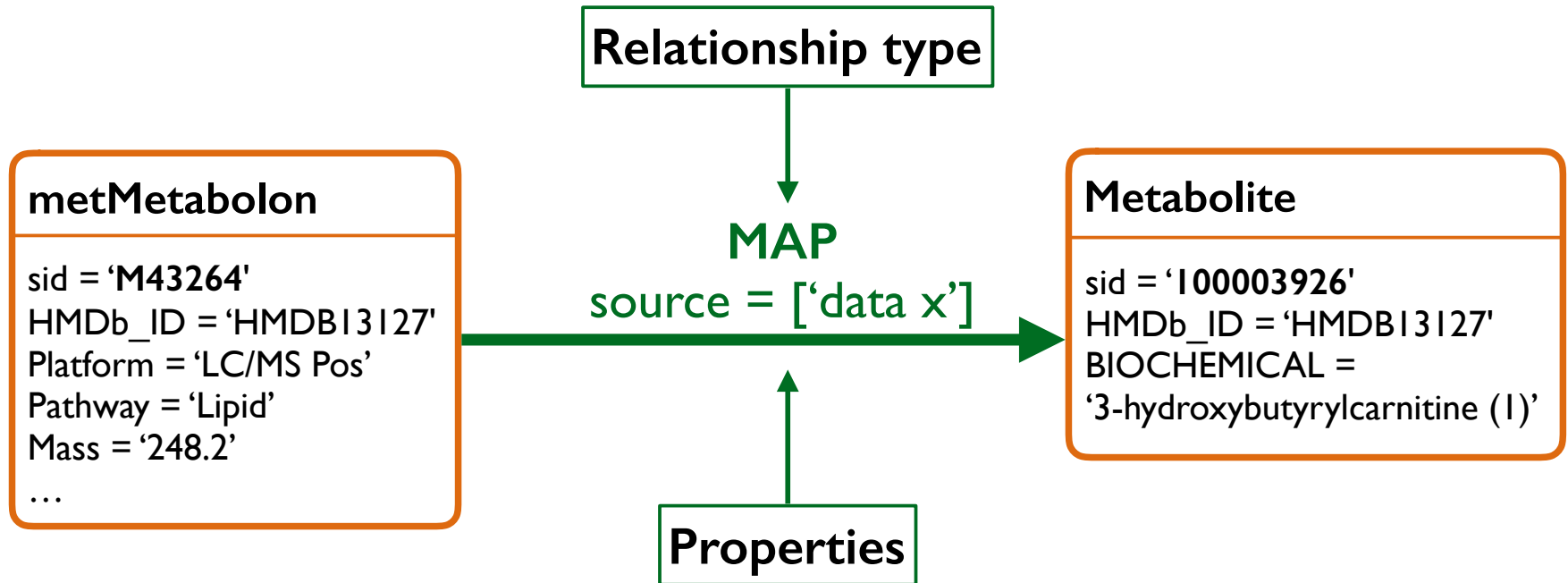


# Neo4j - Nodes





# Neo4j - Relationships



- ➡ relationships always have a direction
- ➡ direction can be ignored in queries

# Cypher

- declarative graph query language
  - uses patterns to describe graph data
  - can be used for querying and updating
  - queries built up using various clauses
- ➡ familiar, SQL-like
  - ➡ can be chained together
  - ➡ intermediate results will be context for next clause

# Queries - Simple

Match a **specific node** and return it:

**MATCH** (m : metabolite {sid:“1021”})  
**RETURN** m

variable                      label                      property

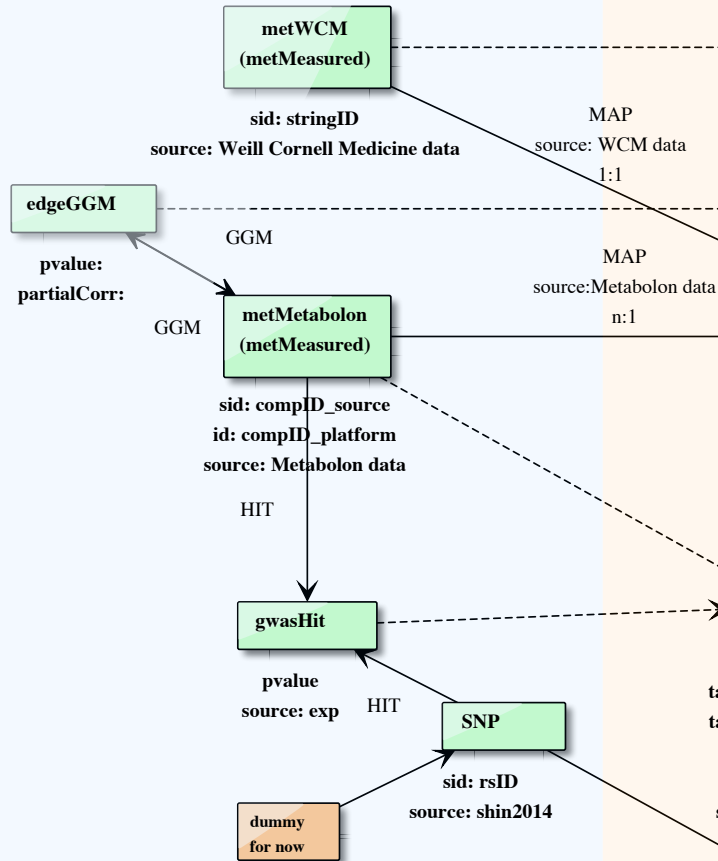
Table with **all reactions** that a **metabolite** participates in:

**MATCH** (m)-[:PARTICIPATES]-(r)  
**RETURN** m.sid **AS** metabo, collect(r.sid)

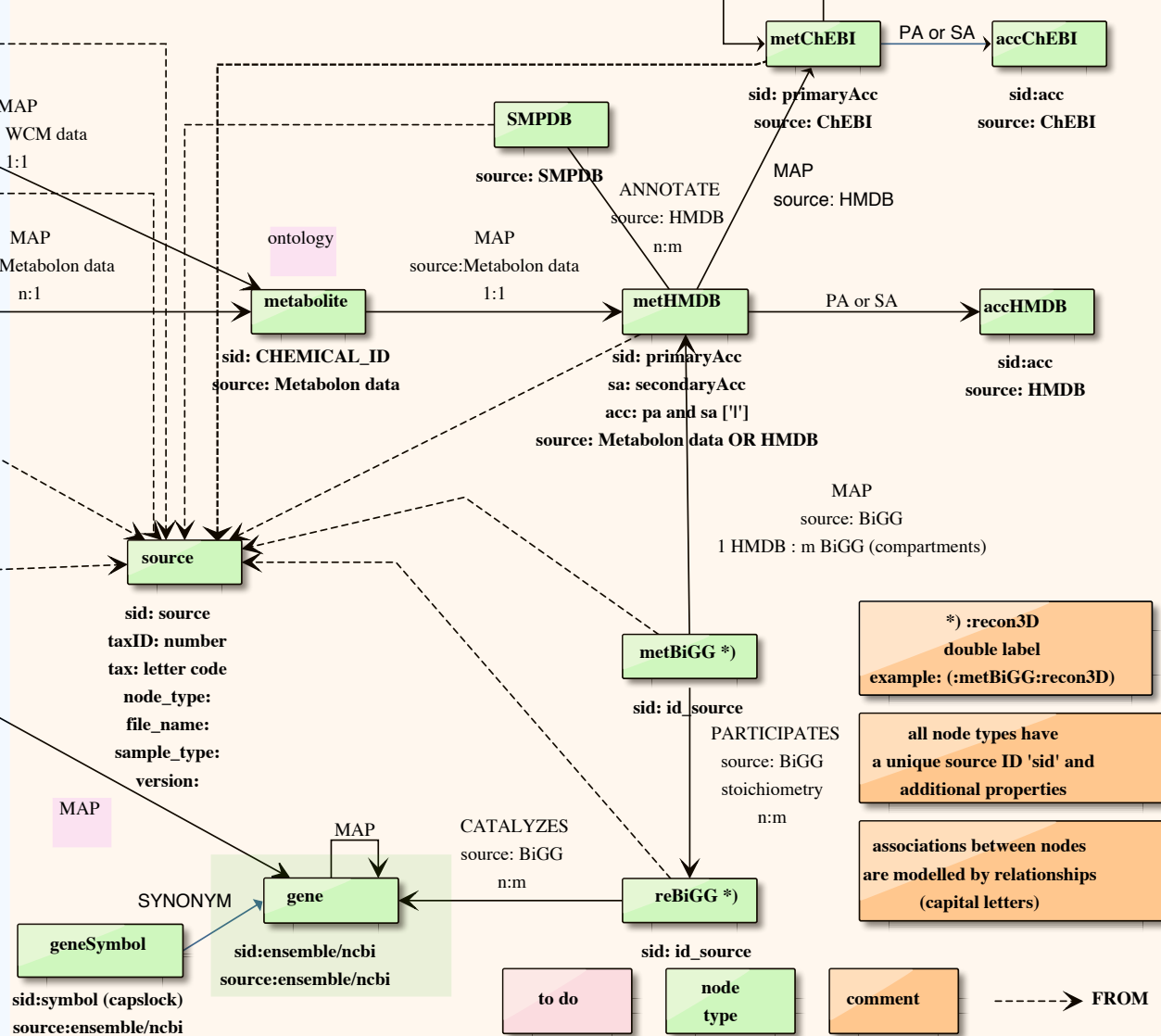
relationship type                      any node

only value                      aggregation function

## Experimental-World



## Knowledge-World



# Data Model

# Git Repository

```
.  
|-- code  
|   |-- scripts  
|   `-- setup  
|-- data  
|   |-- BiGGModel  
|   |-- GGM  
|   |-- SMPDB  
|   |-- SNP  
|   |-- genes  
|   |-- gwasHit  
|   |-- metChEBI  
|   |-- metHMDB  
|   |-- metMetabolon  
|   |-- metWCM  
|   `-- recon3d  
`-- misc
```

<https://gitlab.com/sysdiab/neo4j>

1. code


2. data

3. misc

# Git Repository

```
.  
|-- code  
|   |-- scripts  
|   |-- setup  
|-- data  
|   |-- BiGGModel  
|   |-- GGM  
|   |-- SMPDB  
|   |-- SNP  
|   |-- genes  
|   |-- gwasHit  
|   |-- metChEBI  
|   |-- methMDB  
|   |-- metMetabolon  
|   |-- metWCM  
|   |-- recon3d  
|-- misc
```

I. code



```
code/setup  
|-- README.md  
|-- setup.R  
|-- toLoad  
|   |-- parseMetabolonFile.R  
|   |-- setupFunctions.R
```

# Git Repository

```
.  
|-- code  
|   |-- scripts  
|   `-- setup  
|-- data  
|   |-- BiGGModel  
|   |-- GGM  
|   |-- SMPDB  
|   |-- SNP  
|   |-- genes  
|   |-- gwasHit  
|   |-- metChEBI  
|   |-- metHMDB  
|   |-- metMetabolon  
|   |-- metWCM  
|   `-- recon3d  
`-- misc
```

## 2. data

- File types:

- .xlsx
- .csv

- Info sheet:

- Properties for source node
- CSV: file with \_info.csv
- XLSX: additional sheet “info”

# Git Repository

```
.
|-- code
|   |-- scripts
|   `-- setup
|-- data
|   |-- BiGGModel
|   |-- GGM
|   |-- SMPDB
|   |-- SNP
|   |-- genes
|   |-- gwasHit
|   |-- metChEBI
|   |-- metHMDB
|   |-- metMetabolon
|   |-- metWCM
|   `-- recon3d
`-- misc
```

## 2. data

taxID	9606
tax	hsa
sid	HELM-16-16ML+ CDT (161123) liver
file_name	HELM-16-16ML+ CDT (161123)
node_type	metMetabolon
sample_type	liver
study	HELM-16-16ML+

- Info sheet:
  - Properties for source node
  - CSV: file with \_info.csv
  - XLSX: additional sheet “info”



# Git Repository

```
.  
|-- code  
|   |-- scripts  
|   `-- setup  
|-- data  
|   |-- BiGGModel  
|   |-- GGM  
|   |-- SMPDB  
|   |-- SNP  
|   |-- genes  
|   |-- gwasHit  
|   |-- metChEBI  
|   |-- metHMDB  
|   |-- metMetabolon  
|   |-- metWCM  
|   `-- recon3d  
`-- misc
```

## 2. data

```
data/gwasHit  
`-- 9606_hsa  
    |-- snp_metabolite_assocs_shin2014.csv  
    `-- snp_metabolite_assocs_shin2014_info.csv
```

# Git Repository

```
.  
|-- code  
|   |-- scripts  
|   `-- setup  
|-- data  
|   |-- BiGGModel  
|   |-- GGM  
|   |-- SMPDB  
|   |-- SNP  
|   |-- genes  
|   |-- gwasHit  
|   |-- metChEBI  
|   |-- metHMDB  
|   |-- metMetabolon  
|   |-- metWCM  
|   `-- recon3d  
`-- misc
```

## 2. data

```
data/gwasHit  
`-- 9606_hsa  
    |-- snp_metabolite_assocs_shin2014.csv  
    `-- snp_metabolite_assocs_shin2014_info.csv
```

```
metMetabolon  
|-- 10090_mmu  
|   |-- HELM-14-15ML+CDT\ \ (160202).xlsx  
|   |-- HELM-16-16ML+\ CDT\ (161123)\ WAT.xlsx  
|   |-- HELM-16-16ML+\ CDT\ (161123)\ liver.xlsx  
|   `-- NonTargeted_Test\ sample_results_HELM-13-1  
`-- 9606_hsa  
    |-- CORN-0301-09VWBL(110819).xlsx  
    |-- CORN-0302-09VWBL\ (120126).xlsx  
    |-- CORN-0402-11VWBL(120824).xlsx  
    |-- CORN-0802-13MLB1\ lung\ cancer\ copy_fixed  
    `-- Karpas-20422-20met_alt.xlsx
```

# Git Repository - Data - README

## FILE FORMAT CONVENTIONS:

DATA_TYPE	FORMAT	INFO	SPECIAL_REMARKS
metMetabolon	xlsx	yes	2 sheets (first named info)
metWCM (dummy)	xlsx	yes	must have an "id" and "hmdb" column
GGM	xlsx	yes	2 sheets: "info" and "GGM list": contains in order (metabolite a, metabolite b, partialCorr, pvalue)
gwasHit	csv	yes	column names: rsID, metabolite, pvalue
SNP (dummy)	csv	-	column names: sid, chr, position

## INFO SHEET:

DATA_TYPE	INFO	FORMAT	ENTRIES
metMetabolon	yes	xlsx	sid, file_name, node_type, tax, taxID
metWCM (dummy)	yes	xlsx	sid, file_name, node_type, tax, taxID
GGM	yes	xlsx	sid, tax, taxID, file_name, node_type, metMetabolon_id
gwasHit	yes	csv	sid, tax, taxID, file_name, node_type

# Neo4j Servers

- **Production Server:**

- stable, read-only
- no password
- <http://dzdcon1.helmholtz-muenchen.de:6464/browser/>

- **Development Server:**

- password protected
- <http://dzdcon1.helmholtz-muenchen.de:9494/browser/>

➡ both running Neo4j version 3.3.0

# Running setup.R

- clone git repository
  - set-url to git@gitlab.com:sysdiab/neo4j.git (not https)
  - not prompted for your username/password
- if Neo4j is running on different machine/server
  - ssh connection without password

# Code

## I) Paths

## II) Setup

- `setwd()`
- load scripts from `toLoad` directory
- create log files
- connect to db
- define unique constraints (“sid”)

## III) SCP

- `git pull local & remote`

## IV) Fill databank

# Code

## I) Paths

```
#####  
## PATHS  
#####  
git <- "~/Documents/ICB/neo4j/" # path to git repository  
log_path <- "~/Desktop/"  
db_addr <- "http://dzdcon1:9494/db/data/"  
## only required if Neo4j is running on different server  
neo4j_server <- "ssh alex@dzdcon1"  
wd_path_server <- "/home/alex/sysdiab/"
```

# Code

## II) Setup

- `library(RNeo4j)`

```
## CONNECT TO NEO4J  
db = startGraph(db_addr, username = "neo4j", password = "sysdiab")  
clear(db)
```

```
## define constraints of uniqueness for different node labels  
nodes = c(  
  "metHMDB",  
  "metWCM",  
  "metBiGG",  

```

```
add_constr(nodes, "sid", log, db)
```



# Code

## III) SCP

```
#####  
## SCP - only needed if neo4j on different server  
#####  
system("git pull")  
system(paste0(neo4j_server," '(cd ",wd_path_server,"; git pull)'"))
```

## IV) Fill databank

- SMPDB
- HMDB
- metWCM
- metMetabolon
- ...

## IV) Fill databank - csv

```
for (file in getFiles("SMPDB")) {  
  
  # read info on dataset  
  source <-getInfo(file, c("sid", "file_name", "node_type", "version"), err)  
  #skip file if no/incomplete info  
  if (!all(is.na(source))) {  
    tic(paste0("Added SMPDB file ",source["file_name"]))  
    #check/add source node for dataset  
    checkSource(source, db, log)  
    query=paste("USING PERIODIC COMMIT 1000  
                LOAD CSV WITH HEADERS FROM 'file:///",<file,"' AS props  
                MERGE (n:SMPDB {sid:props.sid})  
                SET n += props with n  
                MATCH (s:source {sid:\"<sid\"},source[\"sid\"],\"<sid\")  
                MERGE (n)-[:FROM]->(s)  
                ",sep="")  
    cypher(db, query) # send query  
    info(log, capture.output(toc())) # log  
  }  
}
```

## IV) Fill databank - xlsx

```
## MAP METABOLON TO METABOLITE
tx <- newTransaction(db) # new transaction
a_ply(metabo,1,function(x) { # process df row-by-row
  query <- paste0("
    UNWIND {properties} AS prop
    MATCH (m:metMetabolon {sid:prop.sid})
    where exists(prop.CHEMICAL_ID)
    MATCH (h:metabolite {sid:prop.CHEMICAL_ID})
    MERGE (m)-[:MAP]->(h)")
  appendCypher(tx, query, properties = x) # add query to transaction
})
logit(commit(tx),log) # send all queries
```