# SpideR

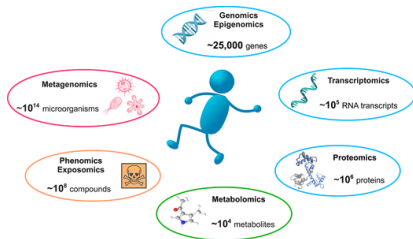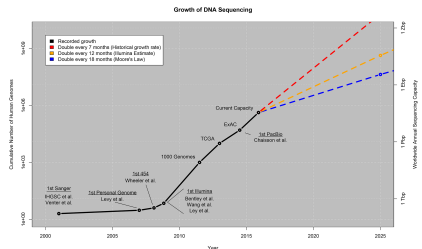## R Package for Search, Integration and Retrieval of Big Data

### Anna M. Sozanska

Supervisors: Dr Shamith Samarajiwa, Dr Dora Bihary
MRC Cancer Unit
University of Cambridge
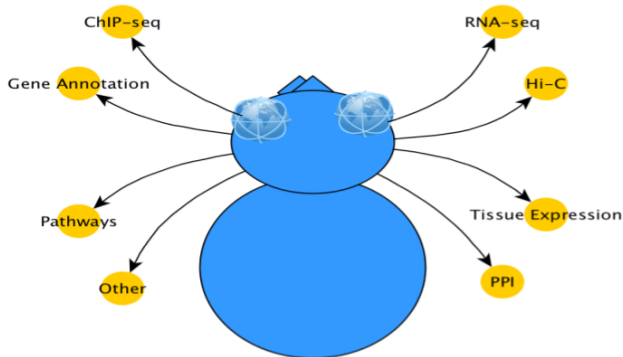
# Biological datasets are rapidly growing in number



(Gligorijević et al, 2016)



(Stephens et al, 2015)

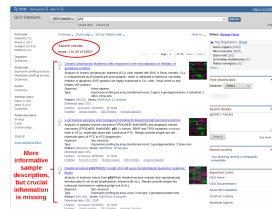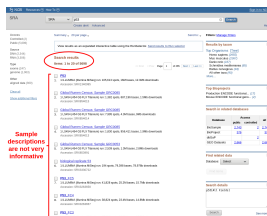# SpideR: Tool for integrated management of biological datasets

# SpideR takes **minutes** to output data that it would take **months** to collect manually

## Task: High-throughput bioinformatics analysis

Collect *all\** ChIP-seq and RNA-seq samples for **p53** from public databases

- Identify special category samples ('inputs' for ChIP-seq)
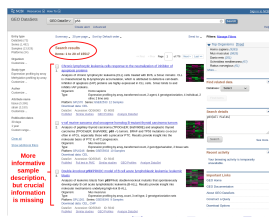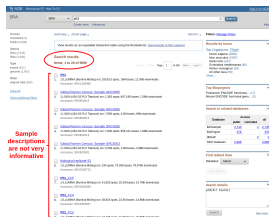- Integrate information from multiple databases (SRA and GEO)



Sample descriptions are not very informative

More informative sample description, but crucial information is missing

# SpideR takes **minutes** to output data that it would take **months** to collect manually

## Task: High-throughput bioinformatics analysis

Collect *all\** ChIP-seq and RNA-seq samples
for **A FEW HUNDRED FACTORS** from public databases

- Identify 'inputs' for ChIP-seq
- Integrate information from multiple databases (SRA and GEO)

# SRA Database: Search results



**Sample descriptions are not very informative**

# GEO Database: Search results

# GEO and SRA use different accession hierarchies

| SRA Name | SRA Acc. | SRA ↔ GEO | GEO Acc. | GEO Name |
|----------|----------|-----------|----------|----------|
| ■ Project | ■ SRP | ⟺ | ■ GSE | ■ Series |
| ■ **Sample!** | ■ SRS | | ■ | ■ |
| ■ Exp. | ■ SRX | ⟺ | ■ GSM | ■ **Sample!** |
| ■ Run | ■ SRR | | ■ | ■ |

# Database Accession Numbers

| SRA # | SRA | SRA → GEO # | GEO | GEO # |
|-------|-----|-------------|-----|-------|
| ■ 98 028 | ■ SRP | ■ 17 300 | ■ GSE | ■ 80 782 |
| ■ 2 974 554 | ■ SRS | ■ 396 513 | ■ | ■ |
| ■ 2 377 047 | ■ SRX | ■ 400 975 | ■ GSM | ■ 2 231 166 |
| ■ 3 366 463 | ■ SRR | ■ 510 890 | ■ | ■ |

# Main SpideR Functions

### searchForTerm

- IN:
    - library_strategy [R]
    - gene [*]
    - antibody [*]
    - cell_type [*]
    - treatment [*]
    - species [*]
    - platform
    - secondary_lib._strat.

- OUT:
  files for db & pipeline

For example:
Find human HiC data
Find ChIP-seq data
 with STAT1 antibody

### searchForAccession

- IN: accession list
- OUT: df or files for db &
  pipeline

For example:
Find runs within SRP052871
Find runs within GSE34715

Anna M. Sozanska

SpideR

# Other SpideR Functions

## superseriesFinder

- **IN: GSM list**
  List of samples of interest

- **OUT: GSM list**
  List of all other samples within the same superseries as the samples of interest

## convertAccession

- **IN: accession list**
  In **one** accession format
- **OUT: accession list**
  In **all** accession formats

### inputDetector
Labels inputs in a sample sheet

### controlDetector
Labels controls in a sample sheet

### dbExtractGenerator
Generates db extract from df

### sampleSheetGenerator
Generates sample sheet from df

Anna M. Sozanska

SpideR

# SpideR uses a **custom database** to solve the non-trivial problem of SRA/GEO mapping

Projects (SRP/GSE)



- No 1:1 mapping of samples
- Different accession hierarchy
- Different information in each database
- Mapping information scattered across different columns of the database

Samples (SRX/GSM)



All numbers in 1000s of entries.

# Exploring hidden sample hierarchy can help find related experiments

**Task**

Find related ChIP-seq and RNA-seq experiments

**Case 1**
Relationship easy to establish

| SRP#1 |
| --- |
| ■ A (ChIP) |
| ■ B (RNA) |
| ■ C (RNA) |

| GSE#1 |
| --- |
| ■ A (ChIP) |
| ■ B (RNA) |
| ■ C (RNA) |

**Case 2**
Relationship difficult to establish

| SRP#1 |
| --- |
| ■ A (ChIP) |

| GSE#1 |
| --- |
| ■ A (ChIP) |

*Superseries*?

| SRP#2 |
| --- |
| ■ B (RNA) |
| ■ C (RNA) |

| GSE#2 |
| --- |
| ■ B (RNA) |
| ■ C (RNA) |

# Overcoming superseries challenge can help find hidden hierarchies

Format in the database

| GSE | GSM |
|------|-----|
| GSE1, GSE3 | A |
| GSE1, GSE3 | B |
| GSE1, GSE3 | C |
| GSE2, GSE3 | **D** |
| GSE2, GSE3 | **E** |
| GSE2, GSE3 | **F** |
| GSE2, GSE3 | **G** |

### Approach

Given a list of samples (GSMs), get all the samples that belong to the same GSEs (excluding the original samples).

Anna M. Sozanska

SpideR

# Overcoming superseries challenge can help find hidden hierarchies

Format in the database

| GSE | GSM |
|-----------|-----|
| GSE1, GSE3 | A |
| GSE1, GSE3 | B |
| GSE1, GSE3 | C |
| GSE2, GSE3 | **D** |
| GSE2, GSE3 | **E** |
| GSE2, GSE3 | **F** |
| GSE2, GSE3 | **G** |

```
            GSE3
          /      \
      GSE1        GSE2
      / | \       / | \ \
     A  B  C     D  E  F  G
```

## Approach

Given a list of samples (GSMs), get all the samples that belong to the same GSEs (excluding the original samples).

# Overcoming superseries challenge can help find hidden hierarchies

Format in the database

| GSE | GSM |
|------|-----|
| GSE1, GSE3 | A |
| GSE1, GSE3 | B |
| GSE1, GSE3 | C |
| GSE2, GSE3 | **D** |
| GSE2, GSE3 | **E** |
| GSE2, GSE3 | **F** |
| GSE2, GSE3 | **G** |



### Approach

Given a list of samples (GSMs), get all the samples that belong to the same GSEs (excluding the original samples).

Anna M. Sozanska

SpideR

# Overcoming superseries challenge can help find hidden hierarchies

Format in the database

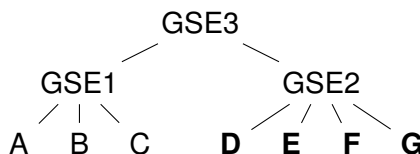| GSE | GSM |
|-----|-----|
| GSE1, GSE3 | A |
| GSE1, GSE3 | B |
| GSE1, GSE3 | C |
| GSE2, GSE3 | **D** |
| GSE2, GSE3 | **E** |
| GSE2, GSE3 | **F** |
| GSE2, GSE3 | **G** |



### Approach

Given a list of samples (GSMs), get all the samples that belong to the same GSEs (excluding the original samples).

Anna M. Sozanska

SpideR

# Overcoming superseries challenge can help find hidden hierarchies

Format in the database

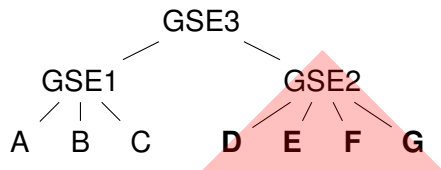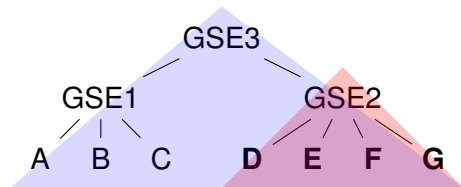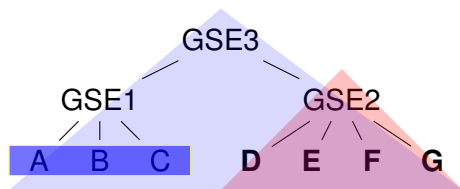| GSE | GSM |
|------------|-----|
| GSE1, GSE3 | A |
| GSE1, GSE3 | B |
| GSE1, GSE3 | C |
| GSE2, GSE3 | **D** |
| GSE2, GSE3 | **E** |
| GSE2, GSE3 | **F** |
| GSE2, GSE3 | **G** |



### Approach

Given a list of samples (GSMs), get all the samples that belong to the same GSEs (excluding the original samples).

## Other SpideR Features

### SpideR tackles inconsistencies in categories

For example, in sample_attribute field which contains comma-separated-categories.
E.g. tissue, source name, cell, cell_type, cell type - same/similar???

### SpideR provides easily manipulable output and is easily **reproducible**

.*Rda* outputs with **data frames**
.*tab* outputs with **data frames**
.*Rda* output with the **function call** .*tab* output with function **call parameters**

Anna M. Sozanska

SpideR

# SpideR addresses most challenges with database design and integration

☑ Misleading fields in SRAdb and GEOmetadb
☑ Inconsistent attributes or categories
☑ Superseries in GEO
☑ SRA-GEO conversion
☐ Erroneous or inconsistent entries

Anna M. Sozanska

SpideR

# Plans for future development

- A new function for searching everywhere in the database (to get a list of all *potentially matching* samples, including those rejected by the more specific **searchForTerm** function)
- Integration of *elasticsearchr*
- Linking the search functions to gene synonyms database
- SQL(ite) database for storing samples of interest

# Acknowledgements



**Samarajiwa Lab @ MRC Cancer Unit**

Shamith Samarajiwa
Dora Bihary
Charlie Fletcher

Kirschner & Samarajiwa et al, PLoS Genetics 2015