# EpiGraphDB

A database and data mining platform for health data science

**IEU Monthly Meeting 10 December 2019** 

http://docs.epigraphdb.org/slides/2019-12-ieu-meeting.pdf

Yi Liu, Benjamin Elsworth, Valeriia Haberland, Pau Erola, Jie Zheng, Matt Lyon, Tom R Gaunt



#### http://docs.epigraphdb.org/slides/2019-12-ieu-meeting.pdf

Introduction

Graph

EpiGraphDB project

Outline

- Use case 1: Pleiotropy
- Use case 2: Therapy response
- Use case 3: Literature
- EpiGraphDB platform
- Summary







- Emerging trends in bioinformatics and health data science:
  - Rise of systematic approaches using computational methods in mining epidemiological relationships
  - Increasing availability of complex, high-dimensional epidemiological data
- EpiGraphDB as a project seeks to develop innovative and scalable approaches to harness their potentials to address research questions of biomedical importance.



## EpiGraphDB project



- Integration of a range of data sources:
  - Systematic MR
  - Observational and genetic correlations
  - Literature-mined relationships
  - Molecular pathways
  - Protein-protein interactions
  - Drug-target relationships
- Data mining on the mechanisms of complex networks of association of risk factor / disease relationships



#### EpiGraphDB http://epigraphdb.org

- DB, API, web UI, R pkg, etc
- v0.2 (v0.3 in the works!)



#### Integrated Epidemiological Evidence



#### Systematic evidence from IEU studies

- IEU GWAS Database (Elsworth et al., forthcoming a);
- MR-EVE (Hemani et al., 2017)
- MELODI (Elsworth et al., 2017)
- pQTL MR (Zheng et al., 2019)
- p/eQTL MR (Zheng et al., forthcoming)
- PRS atlas (Richardson et al., 2019)
- Vectology (Elsworth et al., forthcoming b)
- Research studies by EpiGraphDB group members (<u>http://docs.epigraphdb.org/</u>)

#### External data sources

- EFO
- Gtex
- IntAct
- MeSH
- OpenTargets
- Reactome
- SemMedDB
- STRING-db
- ...









Integrated epidemiological evidence http://docs.epigraphdb.org

- Causal relationships
- Association relationships
- Molecular pathways
- Literature mined / derived evidence
- Others



# Use case 1 Pleiotropy (pQTL)





 Can we distinguish vertical and horizontal pleiotropic instruments using biological pathway data?

Vertical pleiotropy







For any instrument associated with multiple proteins, if

- these proteins are mapped to the same biological pathway
- exists a protein-protein interaction (PPI) between them

then, by definition, the instrument is more likely to act through vertical pleiotropy and it is more likely to be a valid instrument for MR.



Hypothesis





10 December 2019





- We checked the number of pathways and PPIs each protein is involved in for all the instruments associated with 2 to 5 proteins
- We used EpiGraphDB to extract high confidence PPIs from StringDB (confidence score >0.7)
  - How many PPIs they have
  - How many PPIs are shared between groups of proteins that are associated with the same SNP (or SNPs in strong LD)













- We checked the number of pathways and PPIs each protein is involved in for all the instruments associated with 2 to 5 proteins
- We used EpiGraphDB to extract pathway information from Reactome (lower level pathways)
  - Number of pathways each protein is involved in (either directly or as part of a complex)
  - How many pathways are shared between groups of proteins that are associated with the same SNP (SNPs in strong LD)



#### Pathways – examples









Jie Zheng *et al.*, Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases, *under revision* 

- 263 *tier 1* instruments associated with between two and five proteins
  - Test if mapped to the same pathway or PPI
- After the analysis, 68 instruments were considered valid instruments
- Limitation: some pathways and PPIs that may be not included in Reactome and STRING





### Summary

**EpiGraphDB** allows the users to evaluate the potential pleiotropic profile of genetic instruments.



# Use case 2 Therapy response

Resequencing of positional candidates identifies low frequency IL23R coding variants protecting against inflammatory bowel disease.

Momozawa et al. Nat. Genetics

2011

A genome-wide association study identifies IL23R as an inflammatory bowel disease gene.

Duerr et al. Science

2006

### IL23R and Inflammatory Bowel Disease (IBD)

EpiGraphDB IEU Monthly Meeting Talk

10 December 2019

IL23R polymorphisms influence phenotype and response to therapy in patients with ulcerative colitis.

Cravo et al. Eur J Gastroenterol Hepatol.

2014





Search for interacting\* druggable\*\* proteins:

Druggability Tier	Number of interacting proteins
Tier 1	25
Tier 2	8
Tier 3A and 3B	9

\* STRING (https://string-db.org/) and IntAct (https://www.ebi.ac.uk/intact/)

\*\* Finan et al, "The druggable genome and support for target identification and validation in drug development", *Sci. Transl. Med.* **9**, eaag1166 (2017)



### Interacting proteins (Tier 1)



Gene A	Druggability Tier A	Uniprot ID A	Gene B	Druggability Tier B	Uniprot ID B
IL23R	null	Q5VWK5	IL12B	Tier 1	P29460
IL23R	null	Q5VWK5	IL12RB1	Tier 1	P42701
IL23R	null	Q5VWK5	IL23A	Tier 1	Q9NPF7
IL23R	null	Q5VWK5	JAK1	Tier 1	P23458
IL23R	null	Q5VWK5	JAK2	Tier 1	O60674
IL23R	null	Q5VWK5	STAT3	Tier 1	P40763

#### Summaries for IL23R Gene

#### Entrez Gene Summary for IL23R Gene 🕑

The protein encoded by this gene is a subunit of the receptor for IL23A/IL23. This protein pairs with the receptor molecule IL12RB1/IL12Rbeta1, and both are required for IL23A signaling. This protein associates constitutively with Janus kinase 2 (JAK2), and also binds to transcription activator STAT3 in a ligand-dependent manner. [provided by RefSeq, Jul 2008]

#### GeneCards Summary for IL23R Gene

IL23R (Interleukin 23 Receptor) is a Protein Coding gene. Diseases associated with IL23R include Inflammatory Bowel Disease 17 and Psoriasis 7. Among its related pathways are Cytokine Signaling in Immune system and Th17 cell differentiation. Gene Ontology (GO) annotations related to this gene include *interleukin-12 receptor binding* and *interleukin-23 receptor activity*. An important paralog of this gene is LEPR.

?





Search for MR results\* for strongly related proteins (Tier 1) and their effect on IBD:

Gene	Effect size	SE	P-value	ID	Outcome
IL23R	1.5	0.0546	2.21E-166	294	Inflammatory bowel disease
IL12RB1	-0.0097	0.0142	0.49	294	Inflammatory bowel disease
IL12B	0.42	0.0345	9.59E-34	294	Inflammatory bowel disease

\* Zheng, Haberland, Baird, et al. "<u>Phenome-wide Mendelian randomization</u> <u>mapping the influence of the plasma proteome on complex diseases</u>", Submitted revised version to Nat. Genetics (2019)





### Summary

**EpiGraphDB** allows the users to search for therapy response related information either for the intended target or along its pathway.



# Use case 3 Literature



### Literature data v0.2



- Limited literature data
- Links to Publications from various places
- Links to SemMedDB via ontology matches





Literature data v0.3



v0.3 (some time next year)

All SemMedDB and PubMed
MELODI Lite enrichment for each GWAS



### MELODI



http://melodi.biocompute.org.uk/

# JELODI

#### SemMedDB

Database of triples extracted from MEDLINE titles and abstracts, e.g.

PCSK9 (subject) PREDISPOSES (predicate) Cardiovascular Diseases (object)





#### **MELODI** Lite





Restricted to certain types and predicates



100x quicker



Multiple exposures and outcomes



Application programming interface:

http://textbase.biocompute.org.uk/docs/



#### v0.3 Literature Metrics



Genes represented in the literature	15,489 / 57,736
Pathways represented by genes in the literature	2,249 / 2,259
GWAS with literature evidence	4,226 / 11,016
Trait-MR->Trait (p<1e-20) pairs with no literature connection	1,839 / 8,830



1. Find potential risk factors (no ncase value - continuous) match (g1:Gwas)-[mr:MR]->(g2:Gwas) where g2.id = 'ieu-a-30'
and not exists(g1.ncase) and mr.pval<1e-5 with order by
mr.pval asc, mr.moescore desc limit 5 with g1,g2,mr
collect(distinct(g1.id))+collect(distinct(g2.id))as g\_list</pre>



- 1. Find potential risk factors (no ncase value - continuous)
- 2. Get SNP-gene data for all GWAS

match (g1:Gwas)-[mr:MR]->(g2:Gwas) where g2.id = 'ieu-a-30'
and not exists(g1.ncase) and mr.pval<1e-5 with order by
mr.pval asc, mr.moescore desc limit 5 with g1,g2,mr
collect(distinct(g1.id))+collect(distinct(g2.id))as g\_list</pre>

match (gene1:Gene) <- [vg:VARIANT\_TO\_GENE] - (v:Variant) [gv:GWAS\_TO\_VARIANT] - (gwas:Gwas) where gwas.id in g\_list and
gv.pval<1e-20 with gene1,gwas,v,vg</pre>



- 1. Find potential risk factors (no ncase value - continuous)
- 2. Get SNP-gene data for all GWAS
- 3. Get literature data for all GWAS

match (g1:Gwas)-[mr:MR]->(g2:Gwas) where g2.id = 'ieu-a-30'
and not exists(g1.ncase) and mr.pval<1e-5 with order by
mr.pval asc, mr.moescore desc limit 5 with g1,g2,mr
collect(distinct(g1.id))+collect(distinct(g2.id))as g list</pre>

match (gene1:Gene) <- [vg:VARIANT\_TO\_GENE] - (v:Variant) [gv:GWAS\_TO\_VARIANT] - (gwas:Gwas) where gwas.id in g\_list and
gv.pval<1e-20 with gene1,gwas,v,vg</pre>

optional match (gwas)-[gs:GWAS\_SEM]-(s:SemmedTriple)-[:SEM\_SUB|:SEM\_OBJ]-(st:SemmedTerm)-[sg:SEM\_GENE]-(gene2:Gene) where gs.pval<1e-20 with gwas,gene1,gene2,v,s,st



- 1. Find potential risk factors (no ncase value - continuous)
- 2. Get SNP-gene data for all GWAS
- 3. Get literature data for all GWAS
- 4. Find shared pathways between SNP-genes and literature genes

match (g1:Gwas)-[mr:MR]->(g2:Gwas) where g2.id = 'ieu-a-30'
and not exists(g1.ncase) and mr.pval<1e-5 with order by
mr.pval asc, mr.moescore desc limit 5 with g1,g2,mr
collect(distinct(g1.id))+collect(distinct(g2.id))as g\_list</pre>

match (gene1:Gene) <- [vg:VARIANT\_TO\_GENE] - (v:Variant) [gv:GWAS\_TO\_VARIANT] - (gwas:Gwas) where gwas.id in g\_list and
gv.pval<1e-20 with gene1,gwas,v,vg</pre>

optional match (gwas)-[gs:GWAS\_SEM]-(s:SemmedTriple)-[:SEM\_SUB|:SEM\_OBJ]-(st:SemmedTerm)-[sg:SEM\_GENE]-(gene2:Gene) where gs.pval<1e-20 with gwas,gene1,gene2,v,s,st

optional match (gene1)-[:GENE\_TO\_PATHWAY]->(p:Pathway)<[:GENE\_TO\_PATHWAY]-(gene2)</pre>



- 1. Find potential risk factors (no ncase value continuous)
- 2. Get SNP-gene data for all GWAS
- 3. Get literature data for all GWAS
- 4. Find shared pathways between SNP-genes and literature genes

## 5. Map SNP-genes to literature genes

match (g1:Gwas)-[mr:MR]->(g2:Gwas) where g2.id = 'ieu-a-30'
and not exists(g1.ncase) and mr.pval<1e-5 with order by
mr.pval asc, mr.moescore desc limit 5 with g1,g2,mr
collect(distinct(g1.id))+collect(distinct(g2.id))as g list</pre>

match (gene1:Gene) <- [vg:VARIANT\_TO\_GENE] - (v:Variant) [gv:GWAS\_TO\_VARIANT] - (gwas:Gwas) where gwas.id in g\_list and
gv.pval <1e-20 with gene1,gwas,v,vg</pre>

optional match (gwas)-[gs:GWAS\_SEM]-(s:SemmedTriple)-[:SEM\_SUB|:SEM\_OBJ]-(st:SemmedTerm)-[sg:SEM\_GENE]-(gene2:Gene) where gs.pval<1e-20 with gwas,gene1,gene2,v,s,st

optional match (gene1)-[:GENE\_TO\_PATHWAY]->(p:Pathway)<[:GENE\_TO\_PATHWAY]-(gene2)</pre>

optional match (gene1)-[:SEM\_GENE]-(st2:SemmedTerm) return
gwas,gene1,gene2,v,s,st,p,st2;



University of BRISTOL MRC Unit

> 9 genes 6 semantic

6 SNPs

- 292 papers
- 232 journals





### Summary

# **EpiGraphDB** allows the users to search for literature evidence.



# Platform





### EpiGraphDB platform





10 December 2019



Web UI



#### http://epigraphdb.org

Query

#### Network plot

Confounder MR causal estimate

Network plot

DO VR

== Table

</>> Query

Coronary heart disease

Exposure Body mass index

Outcome

Search

B Documentation

CFullscreen 3D

MR evidence on confounding traits between exposure and outcome



10 December 2019



Web UI



#### Explorer

#### http://epigraphdb.org/explorer

# Gallery http://epigraphdb.org/gallery



Diagnoses - main ICD10: C91.0 Acute lymphoblastic leukaemia





### epigraphdb R package



#### https://mrcieu.github.io/epigraphdb-r

- API client package
- tidyverse compliant
- pkgdown doc site
- Need user feedbacks!

epiaraphdb 0.1 CHANGELOG Home Reference Articles -EpiGraphDB R package epigraphdb License Full license GPL-3 Graph DB University of BRISTOL Developers Yi Liu EpiGraphDB is an analytical platform and database to support data mining in epidemiology. The platform incorporates a graph of Author, maintainer causal estimates generated by systematically applying Mendelian randomization to a wide array of phenotypes, and augments this with Valeriia Haberland a wealth of additional data from other bioinformatic sources. EpiGraphDB aims to support appropriate application and interpretation of Author causal inference in systematic automated analyses of many phenotypes Pau Erola epigraphdb is an R package to provide ease of access to EpiGraphDB services. We will refer to epigraphdb as the name of the R Author package whereas "EpiGraphDB" as the overall platform. Benjamin Elsworth Author Installation Tom Gaunt Author devtools is required to install from github: Dev status # install.packages("devtools") CRAN not publishe devtools::install\_github("MRCIEU/epigraphdb-r") build passing NOTE: while the package repository is "epigraphdb-r", the R package name is "epigraphdb"

#### Using epigraphdb

lib	rary("epigraphdb")
#>	EpiGraphDB v0.2
#>	
#>	Web API: http://api.epigraphdb.org
#>	
#>	To turn off this message, use
#>	<pre>suppressPackageStartupMessages({library("epigraphdb")})</pre>
mr(	outcome = "Body mass index")
#>	# A tibble: 370 x 12
#>	exposure_id exposure_name outcome_id outcome_name estimate se

install.packages("devtools")
devtools::install\_github(
 "MRCIEU/epigraphdb-r"

library("epigraphdb")

codecov 88%

0



## Query EpiGraphDB: colliders



API

epigraphdb::confounder

R package

```
confounder(
    exposure = "Body mass index",
    outcome = "Coronary heart disease",
    type = "collider"
```

#### epigraphdb::query\_epigraphdb

```
query_epigraphdb(
  endpoint = "/confounder",
  params = list(
    exposure = "Body mass index",
    outcome = "Coronary heart disease",
    type = "collider"
  )
```

#### Python requests

#### import requests

```
api_url = "http://api.epigraphdb.org
endpoint = "/confounder"
response = requests.get(
    url=api_url + endpoint,
    params={
        "exposure": "Body mass index",
        "outcome": "Coronary heart
disease",
        "type": "collider"
    }
```

#### print(response.json())

Documentation (pkgdown):
https://mrcieu.github.io/epigraphdb-r

Documentation (swagger interface):
http://api.epigraphdb.org

10 December 2019

# Summary





http://docs.epigraphdb.org/slides/2019-12-ieu-meeting.pdf

- EpiGraphDB as a database and platform powers the the data mining of high dimensional, complex epidemiological relationships.
- We are actively developing EpiGraphDB and working on associated research studies.
- Please use EpiGraphDB and get in touch!
- Email: <u>feedback@epigraphdb.org</u>

#### Acknowledgements

#### **EpiGraphDB**

Yi Liu Benjamin Elsworth Valeriia Haberland Pau Erola Jie Zheng Matt Lyon Tom R Gaunt

pQTL project

Jie Zheng Valeriia Haberland Benjamin Elsworth Denis Baird Venexia Walker Tom Richardson Kurt Taylor James Staley George Davey Smith Philip Haycock Gibran Hemani Robert Scott Biogen & GSK collaborators







- Elsworth, B., et al., 2018. MELODI: mining enriched literature objects to derive intermediates.
- Elsworth, B., et al., forthcoming a. IEU GWAS database <u>https://gwas.mrcieu.ac.uk/</u>
- Elsworth, B., et al., forthcoming b. Vectology: exploring biomedical variable relationships using sentence embedding and vectors.
- Hemani, G., et al., 2017. Automating Mendelian randomization through machine learning to construct a putative causal map of the human phenome. *BioRxiv*, p.173682.
- Richardson, T.G., et al., 2019. An atlas of polygenic risk score associations to highlight putative causal relationships across the human phenome. *Elife*, *8*, p.e43657.
- Zheng, J. et al., 2019. Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases. *BioRxiv*, doi: https://doi.org/10.1101/627398
- Zheng, J., et al., forthcoming. Systematic Mendelian randomization and colocalization analyses of the plasma proteome and blood transcriptome to prioritize drug targets for complex disease.