# EpiGraphDB

## Case studies and version 0.3 features

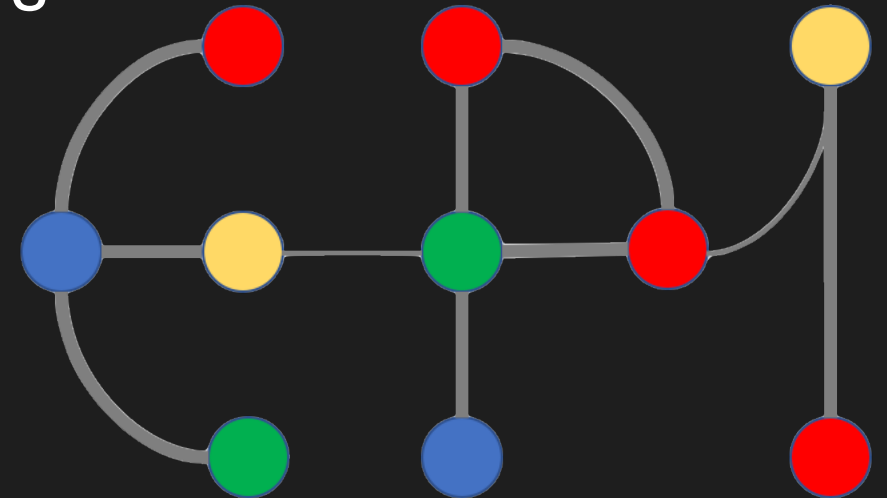**IEU P4 Meeting 09 March 2020**

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

- Introduction

- EpiGraphDB version 0.3

- Use case 1: Pleiotropy

- Use case 2: Alternative drug targets

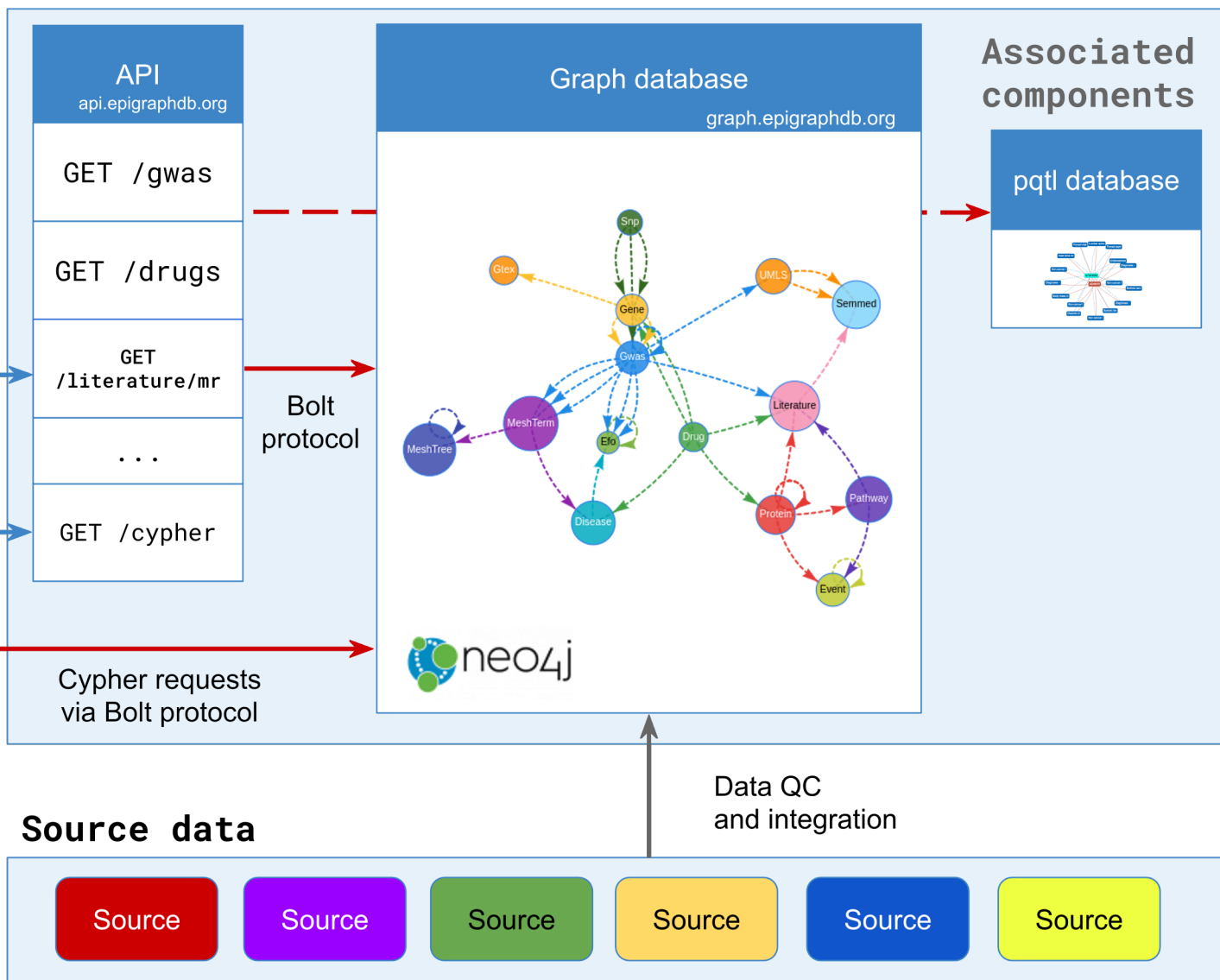- Use case 3: Literature

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

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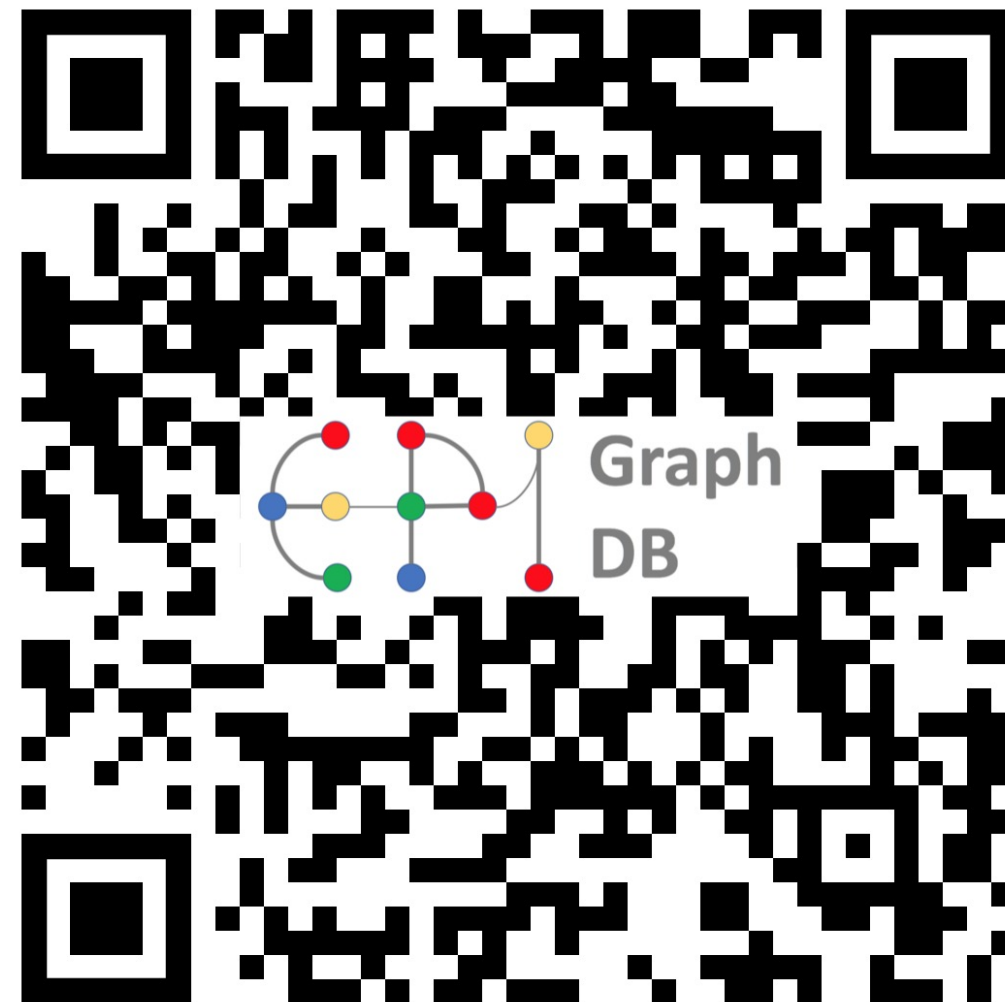- Literature mined / derived evidence
- Others

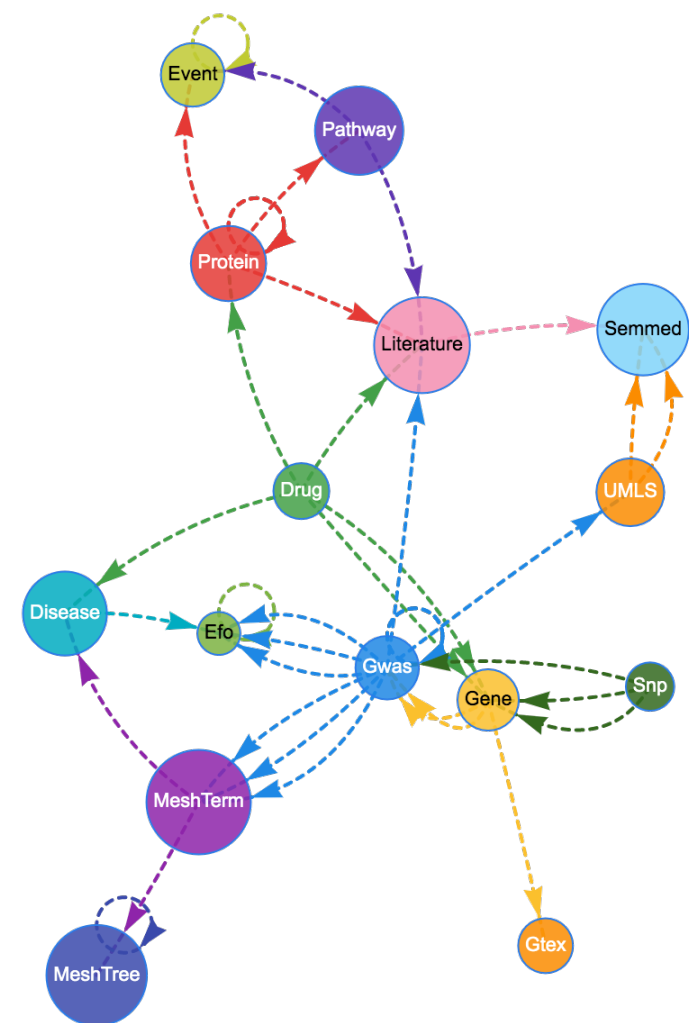**http://docs.epigraphdb.org/slides/2019-12-ieu-meeting.pdf**

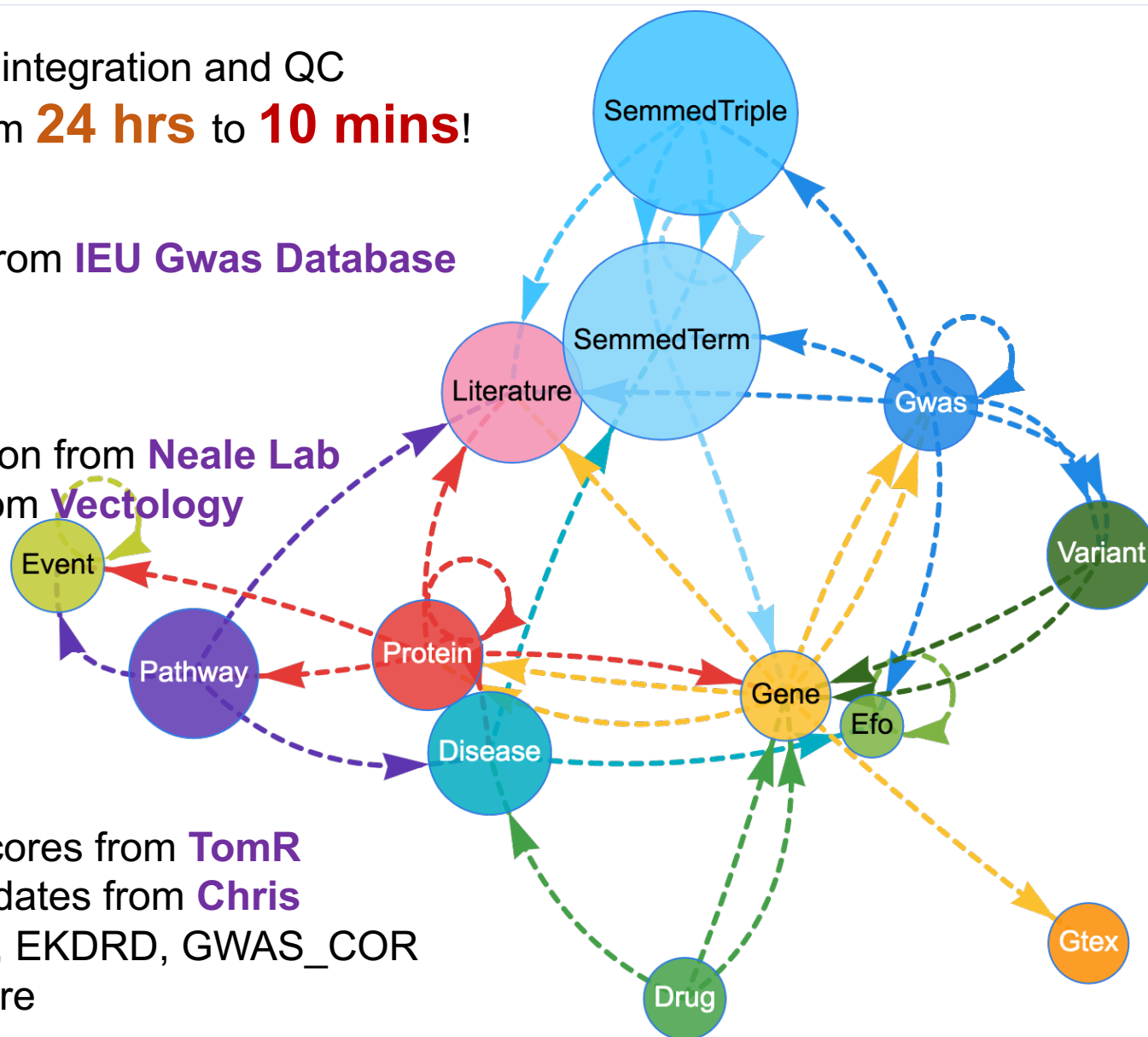- Research project

- Case studies

- Platform

- Version 0.2

EpiGraphDB v0.3

# A new graph database



- Refactored data integration and QC
- Graph builds from **24 hrs** to **10 mins**!
- **10x** bigger!
- Updated Gwas from **IEU Gwas Database**
- SemMedDB
- MELODI
- EFO
- Genetic correlation from **Neale Lab**
- NLP mapping from **Vectology**

- Polygenic risk scores from **TomR**
- eQTL / pQTL updates from **Chris**
- Dropped: MeSH, EKDRD, GWAS_COR
- … and many more

# A new Web UI

- v0.2: Flask + Jinja + Bootstrap3
- v0.3: FastAPI + Vue.js + Bootstrap 4
- Improved event handling
- Improved mobile support
- Improved caching

`/mr`

Mendelian randomisation results

**1. Query for exposure trait**

**Script**

EpiGraphDB API endpoints - EpiGraphDB

```
1  import requests
2
3
4  url = f'{EPIGRAPHDB_URL}/mr'
5  params = {'exposure_trait': 'Body mass index'}
6  r = requests.get(url, params=params)
7  r.raise_for_status()
8  r.json()
```

**Results**

```
1   {'metadata': {'query': 'MATCH (exposure:Gwas)-[mr:MR]->(outcome:Gwas) WHERE '
2                          'exposure.trait = "Body mass index" AND mr.pval < 1e-05 '
3                          'RETURN exposure {.id, .trait}, outcome {.id, .trait}, '
4                          'mr {.b, .se, .pval, .method, .selection, .moescore} '
5                          'ORDER BY mr.pval ;'},
6    'results': [{'exposure': {'id': 'ieu-a-2', 'trait': 'Body mass index'},
7                 'mr': {'b': 0.034558869898319244,
8                        'method': 'FE IVW',
9                        'moescore': 0.9300000071525574,
10                       'pval': 0.0,
11                       'se': 0.002418252406641841,
12                       'selection': 'DF'},
13                'outcome': {'id': 'ukb-a-74',
14                            'trait': 'Non-cancer illness code  self-reported: '
15                                     'diabetes'}},
16               {'exposure': {'id': 'ieu-a-2', 'trait': 'Body mass index'},
17                'mr': {'b': 0.7241045236587524
```
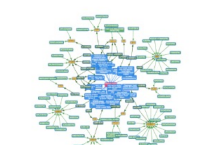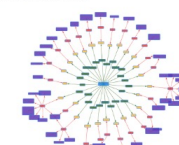
## EpiGraphDB meta nodes

### Disease

**Schema**:

```
1   {'additionalProperties': False,
2    'properties': {'definition': {'title': 'Definition', 'type': 'string'},
3                   'doid': {'items': {'type': 'string'},
4                            'title': 'Doid',
5                            'type': 'array'},
6                   'efo': {'items': {'type': 'string'},
7                           'title': 'Efo',
8                           'type': 'array'},
9                   'icd10': {'items': {'type': 'string'},
10                            'title': 'Icd10',
11                            'type': 'array'},
12                  'icd9': {'items': {'type': 'string'},
13                           'title': 'Icd9',
14                           'type': 'array'},
15                  'id': {'title': 'Id', 'type': 'string'},
16                  'label': {'title': 'Label', 'type': 'string'},
17                  'mesh': {'items': {'type': 'string'},
18                           'title': 'Mesh',
19                           'type': 'array'},
20                  'umls': {'items': {'type': 'string'},
21                           'title': 'Umls',
22                           'type': 'array'}},
23   'required': ['id', 'label', 'definition'],
24   'title': 'Disease',
25   'type': 'object'}
```

# • Still working in progress

Case 1

**Vertical pleiotropy**

SNP → Protein 1 → Protein 2

SNP affect proteins on the same pathway

Valid instrument for MR

**Horizontal pleiotropy**

SNP → Protein 1
SNP → Protein 2
X

SNP associated with proteins from different pathways

Invalid instrument for MR

violates the "exclusion restriction criterion"

**Reality**

SNP → Protein 1
SNP → Protein 2
?

SNP associated with two proteins but relationship between the two proteins are missing

Valid instrument for MR?

Case 2

# IL23R and IBD

Search for interacting druggable proteins

**Using PPI networks for alternative drug targets search**

```python
In [4]: def get_drug_targets_ppi(gene_name):
            endpoint = "/gene/druggability/ppi"
            url = f"{API_URL}{endpoint}"
            params = {
                "gene_name": gene_name
            }
            r = requests.get(url, params=params)
            r.raise_for_status()
            df = pd.json_normalize(r.json()["results"])
            return df

ppi_df = get_drug_targets_ppi(gene_name=GENE_NAME)
ppi_df
```

Out[4]:

| | g1.name | p1.uniprot_id | p2.uniprot_id | g2.name | g2.druggability_tier |
|---|---|---|---|---|---|
| 0 | IL23R | Q5VWK5 | P04141 | CSF2 | Tier 1 |
| 1 | IL23R | Q5VWK5 | P01562 | IFNA1 | Tier 1 |
| 2 | IL23R | Q5VWK5 | P01579 | IFNG | Tier 1 |
| 3 | IL23R | Q5VWK5 | P22301 | IL10 | Tier 1 |
| 4 | IL23R | Q5VWK5 | P29460 | IL12B | Tier 1 |
| 5 | IL23R | Q5VWK5 | P42701 | IL12RB1 | Tier 1 |
| 6 | IL23R | Q5VWK5 | P35225 | IL13 | Tier 1 |

# IL23R and IBD

Search for MR results for the Tier 1 interacting proteins

# IL23R and IBD

Search for literature evidence of the interacting proteins

# Systematic analysis

Case 3

Explore the literature evidence connecting two (or more) traits

1. Given an exposure, find disease traits with causal evidence

2. Select one (or more) exposure -> disease pair and extract literature evidence

3. Select subgraph of literature and extract publication information

# 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

# Reference