# Linking & Exploring Open Government Data: Business Investment and Grants

## COMP3003 Report

William Vigolo da Silva

BSc Hons Computer Science with Year in Industry

psywv@nottingham.ac.uk

14279225

May 4, 2020

# Contents

# 1  Introduction

Many governments and their institutions regularly publish information due to legal requirements of freedom of information and policies regarding open data. The purposes of such policies are often to improve and maintain transparency in government.

Transparency is part of the principle of open government (Lathrop and Ruma 2010), and is considered to strengthen democracy, regulate government behaviour, and promote government efficiency (Schauer 2011).
Open government aims to encourage citizen participation in overseeing the actions of government and their officials - a tradition that can be traced back to ancient Greece (Dornum 1997).
This idea that data should be pro-actively published and freely available is known as the principle of *open data* ("Open Definition 2.1" 2019). This transfers some of the responsibility of oversight; governments must make available facts that could reveal misbehaviour or areas that could be improved.

Open government data often includes topics such as election statistics, department budgets & expenditures, and information about significant individuals within organisations. These topics are useful in identifying undesirable behaviour like corruption and nepotism.
Other areas such as statistics in crime, justice, and healthcare could be used to check that a government is acting in the best interests of its citizens.

There are some unfortunate limitations in the way that governments publish open data. They do not properly follow the principles of open data, failing to use open formats and formats that are not machine-readable.
This makes the data harder to use by both citizens interested in finding out facts about their government, and by researchers looking to perform quantitative research on the data available.

Consequentially, analysing open government data requires a few extra steps, including:

- **Collecting** datasets from multiple sources.
- **Cleaning** data that is invalid or incomplete
- **Linking** data from the different sources, i.e. combining the available information.

# 2  Motivation

This project's goal is to explore the value of open data and government transparency.
It discovers how existing open data and existing government data publishing platforms can used to perform specific research and produce useful findings.

My intention is to apply the full process needed to conduct in-depth data analysis and address the challenges that one faces while working with open data. Outcomes of this effort will include insights into ways to improve government publishing of open data in order to enable a range of use scenarios and potential applications.
The focus is on open data published by the government of the United Kingdom.

# 3  Related work and resources

## 3.1  Papers

Shadbolt et al. (2012): Open data from the UK government was collected, migrated, and linked into a 'semantic web'. Use of standard technologies like the Resource Description Framework data

format and Sparql query language enabled a variety of interesting analyses and visualisations, including exploring data point variations over regions and time. A user interface was created that allowed less technical users to explore data using only their spreadsheet skills.

This paper is a valuable reference into the process needed for transforming a collection of datasets into a database that can be analysed using data analysis methods and to which machine learning methods can be applied.

Winkler (2006): An overview of existing methods and research in the area of record linkage. Records from different datasets need to be linked together in order for the relationships between entities to be explored, and to provide extra information about individual entities, such as businesses and grants.

The methods described in the paper were helpful through the process of linking records in the dataset, particularly in describing algorithms and heuristics for record comparison.

Elmagarmid, Ipeirotis, and Verykios (2006): Techniques for detecting duplicate records in a database are explored and explained in detail.

Among other things, duplicate records in the datasets used in this project need to be merged in order for the data analysis to produce high quality and accurate results. This paper provides a valuable reference for methods used during this process, including Q-gram matching.

## 3.2  Data sources

UK Research and Innovation[1] (UKRI): A non-governmental organisation that funds research and innovation in the UK through grants. It is an umbrella organisation for several 'research councils' in several industries and areas. Data on funding from each council is made publicly available through the Gateway to Research[2] database.

This provides details on each project, organisations & individuals involved, and catalogues research outcomes as summarised by the researchers. Relevant research papers produced during the project are also listed.

This is the main dataset that will be used throughout this project, as it includes most information relevant to the intended research.

Data.gov.uk[3], European Data Portal[4], UK Transparency and FOI Releases[5], Data.gov[6]: catalogues datasets released by government institutions of the United Kingdom, European Union, and United States of America respectively.

These could be used to provide additional context to the UKRI dataset, if needed. Should time allow, research data from other regions such as the United States may be explored, in a similar manner to that of the United Kingdom.

## 3.3  Software

PostgreSQL[7]: A database management system. As I have prior experience with it, I intend to use it throughout this project to manage any databases created throughout the project.

Stanford Natural Language Processor (Manning et al. 2014): A toolkit for natural language processing. This software could be used to extract data points from textual documents, or

---

[1]https://www.ukri.org
[2]https://gtr.ukri.org
[3]https://data.gov.uk
[4]https://europeandataportal.eu
[5]https://gov.uk/search/transparency-and-freedom-of-information-releases
[6]https://data.gov
[7]https://www.postgresql.org

during analysis of grant descriptions, for example.

WorldMap[8]: An open-source platform for overlaying data on a map, which could be used to visualise data points over different regions.

NodeXL[9] and Gephi[10] (Bastian, Heymann, and Jacomy 2009): graph visualisation and analysis tools. These tools could be used to explore the relationship between entities in collected datasets. Gephi supports attaching to a database management system, such as PostgreSQL. This makes it easier to use with datasets stored in a single database, as is the case in this project.

---

[8]https://worldmap.harvard.edu
[9]https://nodexl.com
[10]https://gephi.org

# 4   Description of work

The aim is to leverage computational techniques and open government data about investment in research & innovation and businesses in different development stages, to determine the relationship between levels of investment and the impact on businesses within various sectors. Such analysis is of interest to companies that analyse markets, e.g. to support business decision making. This includes The Decision Project[11], who expressed an interest in being involved in the project.

## 4.1   Research questions

The following questions were identified as potential interesting analysis from preliminary research and initial insights gained by exploring the available data.

- What is the structure of the ecosystem of publicly funded research - how does it change over time?
- What are the significant factors that influence collaboration between organisations?

## 4.2   Technical approach

This project explores open government data from the United Kingdom relating to grants and investment in research and development. Specifically, it uses the Gateway to Research database from UK Research and Innovation.

This dataset is aggregated, normalised, and linked so that all available information can be used during the project's research. This involves translating the schema of this dataset into a new schema which encapsulates both existing and new data. This schema is then implemented using the PostgreSQL database management system.

Focus is placed on applying existing data analysis and machine learning techniques on these datasets which contain both structured and unstructured data.
Network analysis methods are used to identify the relationships between entities in the datasets.
Correlation analysis is performed to answer questions based on historical data.
Visualisations are created to show the identified relationships and other findings.
Machine learning methods such as decision trees and artificial neural networks are applied to perform predictive analyses required by some research questions. This requires identifying the most significant attributes within the dataset that contribute best to answering the questions. Each method applied is evaluated to identify each performs with respect to the accuracy or usefulness of findings or predictions. For example, a predictive model for the popularity of a subject over time could be tested against historical data.

---

[11]https://www.decisionplatform.io

# 5 Methodology

(1) **Data aggregation**
  (i) Datasets exported from the UKRI Gateway to Research database
  (ii) Data points extracted and normalised where necessary
  (iii) Design a database schema to accommodate storage of data attributes and metadata
  (iv) Import all records into a single relational database

  In order to link and ultimately analyse the data, each dataset was parsed and exported from their original format into a single database. Individual data points are be extracted to normalise the data, such that it can be stored in a relational database.

  A database schema is needed so that records from all datasets are available in a single relational database. This enables relationship analysis and other analyses.

  Focus was placed on data that is well-formatted and machine-readable. Besides numeric and categorical data, text was also processed.

(2) **Linking datasets**
  (i) Discard low-quality or invalid records
  (ii) Research appropriate linking & de-duplication methods
  (iii) Merge duplicate records
  Some records refer to non-existent entities, such as unnamed businesses listed as participants in research projects. These records are unusable for the purposes of this project, so they need to be identified and removed. These records are identifiable by their similar names and lack of other details.

  Datasets can refer to the same business or grant by different names, causing duplicate entries in the database. In such cases similarity comparison algorithms is be applied to automatically identify the single entity which is being referred to. Some manual work was performed to clean up remaining duplicates.

  The data was sampled randomly to test the efficacy of linking & de-duplication methods.

(3) **Analysing data**
  (i) Research and test useful tools
  (ii) Perform network analysis to identify relationships
  (iii) Analyse data to will to answer the desired questions
  (iv) Develop software to perform analysis, visualisation, and learning
  (v) Create visualisations to show relationships and other results

  Research was performed to identify tools that could be useful in analysis or visualisation. Software was developed to apply existing analysis and visualisation tools to the collected data.

  An initial network analysis was performed on the linked data to identify and visualise some interesting relationships between the entities.

  Several algorithms were then applied to explore these relationships. For network analysis, this included degree distributions and clustering coefficients. Some statistical analysis of relationships were also applied.

(4) **Evaluating results**
  (i) Evaluate outcomes of analysis using metrics appropriate for each method

(ii) Make conclusions to answer the desired questions
Relationships identified through network analysis were evaluated to explore how well they answer the research questions.

# 6  Data acquisition and preparation

## 6.1  Data aggregation

The UKRI Gateway to Research (GtR) service is available online[12]. The service also provides APIs that allow programmatic access to the database, which are documented on the Gateway to Research website[13].

A shell script was written that uses the GtR-2 API. As the API restricts downloading to 100 records per query, the script downloads all available pages into individual XML files. In order to avoid overtaxing the service, the script pauses for some time after each page is downloaded. Initially developed for downloading only organisation records, once fully functioning the script was extended to support downloading all types of records available through the API.

For example, the first page of organisations stored in the GtR database can be accessed at the URL https://gtr.ukri.org/gtr/api/organisations. Though one might have to 'View page source' in an internet browser, one can see that a series of records have been returned in XML format, for example:

```xml
<?xml version="1.0" encoding="UTF-8"?>
<ns2:organisations ...>
    <ns2:organisation ns1:created="2020-01-28T17:07:16Z" ...>
        <ns1:links>
            <ns1:link ns1:href="https://gtr.ukri.org:443/gtr/api/projects/..."
                ns1:rel="PROJECT"/>
            <!-- More project links -->
            <ns1:link ns1:href="https://gtr.ukri.org:443/gtr/api/persons/..."
                ns1:rel="EMPLOYEE"/>
            <!-- More employee links -->
        </ns1:links>
        <ns2:name>Tsinghua University</ns2:name>
        <ns2:addresses>
            <ns1:address
                ns1:created="2020-01-28T17:07:15Z"
                ns1:id="41B25E49-A52C-46D5-A83B-8CF1723184C7">
                <ns1:line1>Tsinghua</ns1:line1>
                <ns1:line2>Hai Dian District</ns1:line2>
                <ns1:postCode>100084</ns1:postCode>
                <ns1:region>Unknown</ns1:region>
                <ns1:type>MAIN_ADDRESS</ns1:type>
            </ns1:address>
        </ns2:addresses>
    </ns2:organisation>
</ns2:organisations>
```

Downloading all records involves parsing this file to identify how many pages of records are available, then performing a query to download each one in turn. This process needs to be repeated for each type of record: organisations, individuals, projects, and project outcomes.

The shell script was written in Bourne Again Shell (Bash)[14] for the GNU/Linux operating

---

[12]https://gtr.ukri.org
[13]https://gtr.ukri.org./resources/api.html
[14]https://www.gnu.org/software/bash

system. Additional command-line tools were used, including: xmllint from libxml2[15] for XML manipulation, GNU Awk[16] for text processing, and curl[17] for downloading web pages.

Since these software packages are free, open source, and commonly used, the research can be relatively easily reproduced and results verified. The tools support automation of the main process downloading a series of files from the internet.

In order to merge the downloaded XML files an auxiliary program was created.

Each page of downloaded records is an individual XML file, with records contained in a single wrapping XML element. As an example, pages of organisation records are in the following structure:

```xml
<?xml version="1.0" encoding="utf8"?>
<ns0:organisations xmlns:ns0="...">
    <ns0:organisation ns1:created="...">
        <!-- Record information -->
    </ns0:organisation>
    <!-- repeat ns0:organisation -->
</ns0:organisations>
```

In order to combine the pages, the program collects the inner records, e.g. `ns0:organisation`s, into the outer `ns0:organisations` of the first record.

This could not be performed by the main script itself, as the XML tool being used, xmllint, did not support this combination of records. Python was used instead, as I am similarly familiar with it, it is cross-platform, free and open-source, and it comes with a fully-featured library for manipulating XML data in that manner that was needed.

The source code of the main download program is provided in the appendix section ukriDownload, and the auxiliary combing script in mergeXML.

## 6.2 Data importing

Aggregated data is imported into a single relational database management system, which will enable analysis through queries composed in the Structured Query Language (SQL).

The schema of the database is designed to be *normalised*, such that each atomic piece of information is stored in an individual field. This minimises data duplication - particularly helpful when handling large databases such as in this project - as well as enabling analysis of individual components of records.

All aggregated data is kept in the database so that it can be used throughout the project. As a result some information is duplicated, as records exist in both the original dataset and the cleaned dataset. In order to differentiate between the sources of data, tables containing data from the Gateway to Research database are prefixed with 'gtr'. All other tables contain data formed from this project.

The design for data from the Gateway to Research database is based partially on the contents of the GtR API manual (*Gateway to Research Api 2* 2018), which explains the contents of the records returned by the API.

In occasions where the manual provided inaccurate or incomplete information, the XML schema of records and some exported records were visually inspected to identify how best to normalise

---

each type of record. For example, the XML schema for individuals is available through the REST API[18].

The final database design is visualised figure 16, with each table and field annotated in the appendix section Database schema manual. It was applied to the PostgreSQL database management system, and is implemented in the SQL file provided in appendix section setup.sql.

PostgreSQL has good support for importing structured XML data, though the structure of the XML records must be mapped to the relational and normalised structure of the database schema.

An SQL procedure was produced for the importing of each record type: organisations, individuals, projects, and project outcomes. These are provided in the appendix: importGtrOrgs.sql, importGtrPersons.sql, importGtrProjects.sql, and importGtrOutcomes.sql.

The SQL procedures extract individual data points from records through XPath queries against the XML data ("XML Path Language (Xpath) 3.1" 2017). Similarly to the design of the database schema, these were developed through a combination of inspecting the XML schema as well as some extracted records.

Once again, visual inspection during this process resulted in additional adjustments being made to the database design, including adding new fields to existing tables and creating new ones where necessary, e.g. the `gtrTopics` table.

Some difficulties encountered during the development of these procedures included:

- Very high memory usage during importing.
  The exported XML files for project records total to over 700 megabytes at the time they were last exported. While creating the import procedures, I found an that PostgreSQL would use large amount of the system's memory, dramatically slowing down both the system and the import process.
  I noticed that memory usage increased dramatically each time I implemented an additional sub-record to be imported. For example, the list of organisations associated with a project is specified in the XML records for projects. Initially, these were extracted using an XPath query, then iterated over to insert each sub-record, in this case into `gtrProjectOrgs`. This iteration is performed using the FOR-IN syntax[19], which has the following structure:

  ```
  FOR identifier IN
    -- QUERY
  LOOP
    -- STATEMENTS
  END LOOP
  ```

  I found that each unique identifier used within a for loop resulted in greater memory usage.
  My solution was to simply use a single variable for these sub-loops, when necessary, sacrificing some clarity in the meaning of the loop variables for more efficient use of memory. Most of the queries were instead adjusted to use a SELECT query within the INSERT statement[20], eliminating the need for temporary variables.
  The project importing procedure remains the most memory-intensive, but the PostgreSQL process peaks at just under 6 gigabytes during its execution, which is feasible for most computers to handle.

---

[18]https://gtr.ukri.org/gtr/api/person
[19]https://www.postgresql.org/docs/12/plpgsql-control-structures.html#PLPGSQL-RECORDS-ITERATING
[20]https://www.postgresql.org/docs/current/sql-insert.html

Memory usage could be further reduced by using temporary tables for sub-records, and additional research into minimising memory usage in PostgreSQL procedures.

- XML entities are not replaced automatically.
  An XPath query in PostgreSQL that targets the text of an element does not automatically decode any entities that reside within the text.
  XML data will normally encode specific characters as *entities* to avoid conflicting with the syntax. Among these are ampersands, `&`: in XML data these are encoded as `&amp;`.
  Once this was noticed, manual text replacement was added for fields where it was noticed.

- Money values being imported incorrectly.
  Initially monetary values, such as the funding organisations provided to projects, were being imported by casting the text of the fields directly to the PostgreSQL type `money`, which can store currency values.
  For example, an exported project record will contain a list of participating organisations and the value of their contribution as decimal number:

```
<ns2:project>
  <!-- ... -->
  <ns2:participantValues>
      <ns2:participant>
      <ns2:organisationId>
          AAC0E6BF-AC85-4577-AA5A-2E8725FFAF7B
      </ns2:organisationId>
      <ns2:organisationName>Wsp UK Limited</ns2:organisationName>
      <ns2:role>LEAD_PARTICIPANT</ns2:role>
      <ns2:projectCost>100000.0</ns2:projectCost>
      <ns2:grantOffer>100000.0</ns2:grantOffer>
      </ns2:participant>
  </ns2:participantValues>
  <!-- ... -->
</ns2:project>
```

However, once these values become sufficiently large, the exported records will store these decimal values in scientific form, which PostgreSQL cannot cast to a monetary type directly. Attempting to do so results in errors such as:
`ERROR:  invalid input syntax for type money: "3.4442927E7"`
The solution was to cast to an intermediary type, `numeric`, before then casting to `money`:
`(xpath('//pro:projectCost/text()', x, nss))[1]::text::numeric::money AS cost`
instead of simply
`(xpath('//pro:projectCost/text()', x, nss))[1]::text::money AS cost`.

- Conflicts occurring upon repeating the import process
  Testing the import procedures involved running them repeatedly to test new adjustments. An expected result of the primary key constraints implemented in the database design is that no two records may share the same primary key. As a result, these repeated executions would result in unique constraint violations.
  In order to work around these, I used PostgreSQL's `ON CONFLICT` clause[21] to handle these conflicts. In cases where new fields were added, the `DO UPDATE SET field = value` action was used to set the value of this new field for already existing records. These

---

[21]https://www.postgresql.org/docs/12/sql-insert.html#SQL-ON-CONFLICT

were eventually replaced with the `DO NOTHING` action to simply skip over already existing records.

## 6.3   Data cleaning

Data cleaning is required to ensure that the data used during analysis is of high quality, i.e. free of errors and sufficient for the planned analysis. The use of low quality data negatively impacts analysis, resulting in inaccuracies and less useful outcomes (Rahm and Do (2000)).

The Gateway to Research data is entered by researchers and individuals involved in the grant funding process. As a result, human error can result in erroneous data and duplicate records. These problems are classified as 'single-source' problems, as they occur even a single database is being considered. The introduction of multiple databases then introduce 'multi-source problems', such as different encodings for the roles of organisations, or different monetary systems being used for currency values.

In order to make the development and further steps of this project simpler, the processes undertaken to clean data avoided removing or modifying any data imported from the Gateway to Research database. Instead, new database tables and/or fields were created to store information during the cleaning process, as well as storing the outcomes of cleaning: merged records.

### 6.3.1   Eliminating invalid records

The first step taken during data cleaning was to identify records that contain no useful information and/or do not refer to a real entity. For example, an organisation record with the name 'Unknown'[22] doesn't represent a real organisation, limiting how its relationship with other entities can be explored. Such issues exist at the *record* scope of problems (as categorised by Rahm), due to this problem applying to the record as a whole.

Initial exploration of the data was undertaken, revealing a number of other organisation and person records with the name 'Unknown' (with varying capitalisation). By querying the database for organisation records with the name specified as 'Unknown', I found that a dozen other records used this name, each providing no information in the fields other than the name.
By querying the database for most popular names used by organisation records (`SELECT name, COUNT(name) FROM gtrOrgs GROUP BY name ORDER BY count desc;`), I found a similar pattern with the name 'Unlisted': over 100 records used the name (in various forms of capitalisation), each also lacking any other information.

In order to indicate these records were not to be used during analysis, I created a table `junkGtrOrgs`, which contains the UUIDs of such organisation records. All records with the name 'Unknown' or 'Unlisted' were inserted into this table.

From the name frequency analysis, I also noticed a significant number of organisation records shared what appeared to be names of organisation departments - 'Research', 'Economics', 'Psychology', etc.. Similarly, these records contained other information but the name.
An example of such a record is the Economics Department of Queen Mary University of London (QMUL)[23]. On the web portal, the associated projects for this record all indicate that these projects are/were led by QMUL - the same is indicated by my records in `gtrProjectOrgs`.
I inferred that organisation entries are automatically created for organisation departments when the 'lead organisation department' is specified for a project, for example in an associated project:

```
<ns2:project ... ns1:id="86616057-4857-4DAC-A1AE-3D354EB9C34A">
    <ns1:links>
        <ns1:link
            ns1:href=".../organisations/D5337A10-AC8A-402A-8164-C5F9CC6B0140"
```

---

[22]https://gtr.ukri.org/organisation/39BBE949-0333-428F-864F-C3B196D3D92D
[23]https://gtr.ukri.org/organisation/3A5E126D-C175-4730-9B7B-E6D8CF447F83

```
        ns1:rel="LEAD_ORG"
    />
    <!-- ... -->
  </ns1:links>
  <ns2:title>The marriage premium revisited: the case of baseball</ns2:title>
  <ns2:leadOrganisationDepartment>Economics</ns2:leadOrganisationDepartment>
  <!-- ... -->
</ns2:project>
```

The only linked organisation in this project record is QMUL. However, if we inspect organisation records for the Economics department of QMUL, we see this project is linked.

During the database importing process explained in [Importing data], only the links specified in project records are imported. Hence, organisation records which are actually departments will have no entries for projects that they were associated with.

These records do not refer to a real organisation entity, and are isolated from other records due to their links not being imported, hence relationships between these entities cannot be analysed. For these reasons, I decided to add organisation records with no linked projects to the invalid records collection.

As a result, 5213 of 47822 (10.9%) organisation records were marked as invalid records.

### 6.3.2 Merging organisation roles enumerations

Through visual inspection, the enumerations that specify an organisation's role in a project (in the `gtrProjectOrgs` table) appeared to some duplicates. This is likely a result of the merging of several databases without considering the meanings of these enumerations.

In order to solve this, I created a new database type, `gtrOrgRole`, and added a utility function that mapped the enumerations within the GtR database to this new type:

| GtR enumeration | `gtrOrgRole` mapping |
| --- | --- |
| LEAD_PARTICIPANT | Lead |
| PARTICIPANT | Participant |
| LEAD_ORG | Lead |
| COLLAB_ORG | Collaborating |
| FELLOW_ORG | Fellow |
| PP_ORG | Project Partner |
| FUNDER | Funder |
| COFUND_ORG | Co-Funder |
| PARTICIPANT_ORG | Participant |
| STUDENT_PP_ORG | Student Project Partner |

These decisions were made by investigating the GtR API documentation, manually gauging the semantics of each enumeration, and determining which refer to the same role.

For example, the first two enumerations, LEAD_PARTICIPANT and PARTICIPANT are specified to only be used by the subsidiary Innovate UK ((*GtR Data Dictionary*, n.d.)). Both LEAD_PARTICIPANT and LEAD_ORG both refer to "[an] organisation receiving project funding which is accountable for ensuring that the planned outcomes for the project are achieved...". Hence, these enumerations can be merged into a single role. The PARTICIPANT enumeration is undocumented, but since it is named almost identically to PARTICIPANT_ORG, I decided that this pair were semantically identical and thus mapped them to the same value.

Mapping is performed by `ParseOrgRole` in [importGtrProjects.sql](importGtrProjects.sql).

## 6.4 Merging duplicate records

As explored by Winkler ([2006](#)) and Elmagarmid, Ipeirotis, and Verykios ([2006](#)) there are many existing methods for identifying duplicate records, most of which involve comparing the similarity of text within records.

In order to track the de-duplication process, additional database tables were created for each entity being de-duplicated.
One set of tables stores the similarity of records: e.g., the `similarGtrOrgs` table tracks pairs of similar organisations within the Gateway to Research database, including metrics measuring the level of similarity, and the result of any manual checking.
Another set of tables stores primary keys of pairs of records that have been determined to be duplicates of each other.
Finally, when appropriate, another table stores records with the merged information of all duplicates. This is the purpose of the `orgs` table, which has all the fields that the main `gtrOrgs` table does, associated with the primary organisation UUID. New records will merge each field from pairs of its duplicates using the `COALESCE` function, which will take the first non-null value. Thus, records in the `orgs` table will contain as much data as the Gateway to Research database contains for each individual entity.

Data normalisation involves splitting data into its minimal atomic components, and enforcing well-structured relations between entities stored in the database. Doing so provides confidence that the stored data is free of anomalies and always correct. One part of this process involves eliminating fields whose values can be calculated from the values of other fields - if the record is updated these values can become outdated, resulting in anomalous information being stored (Lee ([1995](#))).
Storing similarity metrics for pairs of records breaches this principle. However, I believe this decision is justified for several reasons:

1. Original data needs to be maintained so that the changes made during data cleaning can be identified and statistics can be reported within this report.

2. It takes a significant amount of time to calculate these metrics:
   As of the date of data aggregation, the Gateway to Research database contained just under 48,000 records of organisations, and just over 82,000 records of people. Performing unique pairwise comparisons for each type of record is thus rather computationally intensive: in total, the triangle number of $n - 1$, (where $n$ is the number of records) comparisons must be made. Given 48,000 organisations, that results in $\frac{n(n-1)}{2} = 1151976000$ comparisons.
   Re-calculating each time these values are required would be a great inconvenience during the development of the project, as these values are required regularly during de-duplication and analysis.
   The similarity tables instead act as a cache that can be referred back to when needed to speed up processes throughout the project.

3. The database is not regularly being updated:
   Data from Gateway to Research was exported only twice during the project, and the project does not aim to keep up-to-date with the current state of the Gateway to Research database, so it is appropriate to calculate metrics once so that they can be reused.

Database tables were created to store pairs of similar records and potentially useful values for each pair. `similarGtrOrgs` stores similar pairs from `gtrOrgs`, and `similarGtrPeople` stores

similar pairs from `gtrPeople`.

For each pair in `similarGtrOrgs`, the calculated trigram similarity of the names and addresses are stored, as well as a boolean field for the result of any manual checking: `NULL` indicates no manual check has been performed, `TRUE` that the pair were confirmed to be duplicates, and `FALSE` that the pair were confirmed to be unique entities.

### 6.4.1   Entity name analysis

Manually exploring records revealed small variations with how organisation names were entered into the database, as expected of data entered by untrained users.

For example, private companies incorporated as 'limited' often contain a variation of 'Limited', 'Ltd', or 'Ltd.' in their names. The names of duplicate records will vary in how this is specified in the name. The same is true of 'Corporation' ('Corp.'), 'Company' ('Co.'), ampersands ('and'), and 'University' ('Uni.').

Names sometimes included 'The', other times it was omitted - e.g. 'The Institute of . . . ' would sometimes be entered as 'Institute of . . .'.

Another notable feature was the significant number organisation records with postcodes or any form of address omitted: 42.5% provided no postcode, and 32.2% provided neither postcode nor any other part of an address.

As addresses enable two organisations to be clearly distinguished from each other, this means it is not possible to rely on the address information to de-duplicate organisations.

As for records on people, another common pattern is users entering nicknames in the place of first names, such as 'Tom' instead of 'Thomas'. This complicates the de-duplication process, as if one record uses a nickname and the other the full name, these will have a lower similarity metric despite both referring to the same name.

To compensate for this, the similarity comparison metrics were adjusted to ensure that variations of the same words are replaced with a single variation. This resulted in comparisons of records that used these variations having a higher trigram similarity than they did before these replacements.

### 6.4.2   Q-Gram matching

One approach for detecting duplicates is the q-gram matching method, which compares the number of shared substrings of length $q$ within two bodies of text.

PostgreSQL has built-in support for performing this comparison through the pg_trgm module[24]. This module performs q-gram comparison with $q$ as three, hence comparing 3-character sequences. These are known as *trigrams.*

As an example, the word 'trigram' is broken down into the following trigrams:
`tri`, `rig`, `igr`, `gra`, and `ram`.

In PostgreSQL's implementation, a "string is considered to have two spaces prefixed and one space suffixed", hence these additional trigrams are generated (the symbol '␣' indicates a space):
`␣␣t`, `␣tr`, and `am␣`.

These additional trigrams add additional weight to the beginnings and ends of individual words.

I chose to use this method due to existing support in the PostgreSQL database management system, as well as its particular affinity to catching typographical errors in data entered by users.

Compared to some other similarity comparison methods, q-grams are not as sensitive to the

---

[24]https://www.postgresql.org/docs/12/pgtrgm.html

order of words within a strings. This is particularly important for data entry by users, as names or addresses of organisations may be entered in slightly different variations.

For example the Levenshtein distance (or edit distance) "between two strings [. . . ] is the minimum number of edit operations of single characters needed to transform" one string to another. The trigram similarity of the strings "The University of Nottingham" and "Nottingham University" is 75.9% (3sf), whereas the Levenshtein distance is 23. A significant number of characters must be changed to transform one string to the other, despite both strings sharing two entire words.

Similarity was calculated and recorded through SQL procedures, similarGtrOrgs.sql and similarGtrPeople.sql.

The calculated values for trigram similarity were stored for names of organisations (`similarGtrOrgs.simTrigramName`) and people (`similarGtrPeople.simTrigramName`).

An initial attempt at running these scripts involved calculating the similarity of every pair of records, but this failed after a significant amount of time, once it had completely filled the remaining 50GB or so available on the storage device the database was being stored.

In order to avoid this, I filtered pairs by requiring that they share at least 50% of trigrams in their names. This allows the amount of data generated from this procedure to be significantly reduced, down to 243363 from 1151976000 rows for organisations.

An initial attempt to compute the pairwise similarity of records proved to be computational infeasible. It required over 30 minutes and a storage space of more than 50GB before the procedure failed due to running out of disk space. In order to address this issue, record pairs were filtered to include those that share at least 50% of trigrams in their names. This reduced the amount of data generated down to 243363 from the original 1151976000 rows for organisations.



Figure 1: Similarity distribution of similarly-named organisation records. Generated with `OrgSimDist` from stats.sql

**6.4.2.1 Organisations**  Once all these similarity records were created for organisations, I applied an hypothetical heuristic: assume two records are duplicates if their names share at least 90% of trigrams. This results in 2710 or 5.6% organisations being duplicates.

I then chose 100 random records and manually verified whether the pair of organisations were indeed duplicates, which was determined if one of the following conditions held:

1. Both organisation records listed addresses, and the addresses matched
2. Both organisation records listed addresses, and there existed public records of a single organisation being registered at both addresses throughout its history (Using the Companies House[25] service)
3. Both organisations records worked on an identical project

If one of the records did not specify any form of address, then the pair was skipped.

A script was created to present the user with each pair, who can then specify the result of manual checking and update the database records accordingly.

My reasons for choosing Python for this utility are similar to those explained in [Aggregating Data] - Python additionally has support for performing PostgreSQL queries through the psycopg2 module[26].

This necessitated adding an additional field to `similarGtrOrgs` named `manualResult`. The field is a boolean flag indicating the result of manual checking, where a true value indicates the records refer to the same entity. Should this field be `NULL`, this indicates no manual checking has been performed.

The results are shown in the following plot:



Figure 2: Results of manual verification of duplicate organisations whose names have a trigram similarity of 90% or above

A significant portion number of pairs contained an organisation that had no address specified and lacked sufficient information to determine whether the pair were indeed duplicates. This was expected, as a significant portion of organisation records lack address information.

However, of those pairs that both contained addresses, 99 were confirmed to refer to the same entity, normally through sharing an identical address. A smaller number of those were confirmed by checking the historical addresses for a real organisation, and finding that it has previously been registered in both. Less than 10 were determined to be duplicates through sharing similar partners on projects, or similar topics for projects.

The single pair that was confirmed to not be a duplicate was the University of Los Lagos (in Chile) and the University of Lagos (in Nigeria), with a trigram name similarity of 91%.

---

[25]https://beta.companieshouse.gov.uk
[26]https://psycopg.org

This accuracy is acceptable enough for the purposes of this project, and so I decided to apply the heuristic of 90% trigram similarity during the de-duplication process, merging records where this property applies.

Additional confidence could be provided to the de-duplication of organisations by automatically querying an API for records of these organisations, such as the Companies House API[27]. The addresses of records in the Gateway to Research database could be matched against historical names and addresses automatically, similarly to the manual process I performed.
Data from the Companies House database could also be exported and linked with that of the Gateway to Research database.

**6.4.2.2   People**   Trigram similarity was calculated for pairs of records on people, taking the average trigram similarity of both the first name and surname. Where both of the pair specified the 'other names' field, the similarity was calculated as the average trigram similarity of the first name, surname, and other names.

First names were filtered to replace nicknames with possible full names before comparison, to minimise the negative impact of users providing nicknames in the place of full names.

### 6.4.3   Manual inspection

**6.4.3.1   Organisations**   Through manually inspecting organisation records with high name similarity, I identified that the specified postcodes can be used to reliably identify whether similarly named records are indeed duplicates.
Hence, I decided to merge organisation records if the following conditions hold:

1. Neither of the pair is listed as an invalid record
2. The pair share at least 90% of trigrams in their names, or
3. The pair share at least 50% of trigrams in the names, both have specified postcodes, and the postcodes match

As explored in Q-Gram matching, using a threshold of 90% in name similarity is a rather effective heuristic, and hence I have applied it when deciding to merge organisation records.
Matching postcodes provides enough additional confidence that a pair of organisation records are duplicates, hence I apply a much lower trigram similarity threshold for merging them.
The accuracy of the de-duplication is not paramount to the purposes of this project, so 100% accuracy is not necessary. The number of records and proportion of those that lack much information also makes it infeasible to complete this process manually.

---

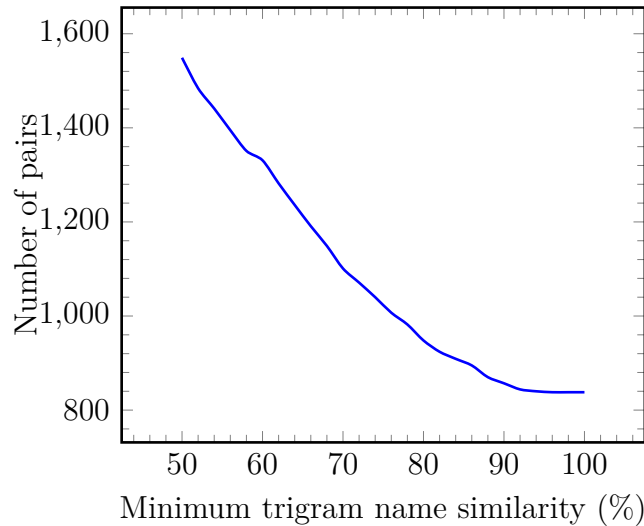[27]https://developer.companieshouse.gov.uk/api/docs

Figure 3: Similarity distribution of similarly-named organisation records, who also share the same postcode. Generated with `OrgSimDistPostcode` from stats.sql

The mergeGtrOrgs.sql procedure performs the record merging, storing the resulting records in the `orgs` table. The `COALESCE` function on each duplicate pair's fields so that the resulting record will take all information available on a single entity.

As a result, 3091 of 47822 (6.5%) organisation records were marked as duplicates.

**6.4.3.2 People** I identified that 617 of 14456 (4.3%) of similar person records shared the employer, which could be used as an additional heuristic in determining whether the pair are duplicates. By linking the employers against the duplicate organisation list in `duplicateGtrOrgs` as generated in the previous section, this number rises to 625.

Using this employer information in conjunction with name similarity results in much more confidence when determining if two records indeed refer to the same individual.

From inspecting the most similar records, I noted that despite having similar names, most had different employers. While it's possible for the same individual to be employed by different organisations over time, the records lack sufficient personally-identifiable information to reliably link them to a single real individual.

As individuals' names are significantly less unique than that of organisations, relying on name similarity alone in de-duplication appears infeasible. Hence, I decided that person records should only be merged if the following conditions hold:

1. Both records list the same employer
2. The name of both records share at least 90% of trigrams

The mergeGtrPeople.sql procedure performs the record merging, storing the resulting records in the `people` table. The `COALESCE` function on each duplicate pair's fields so that the resulting record will take all information available on a single entity.

As a result, 499 of 82015 (0.61%) of records on people were marked as duplicates.

### 6.4.4 Summary

Name trigram similarity is an effective heuristic for detecting duplicate organisations, as these tend to use more unique names to distinguish themselves. The same is not true of people: due

to the wide re-use of names within populations, names alone cannot be used to uniquely identify a real individual. In these cases, additional contextual information must be used to ascertain duplication, such as an individual's employer.

Using name comparison is more effective when common variations of the same words or names are substituted for a single variation before comparison. However, this approach is limited by the requirement of domain knowledge being applied or a significant amount visual inspection being done.

Combining name similarity with the comparison of other information available - such as the addresses of organisations - results in more records being linked. 4871 of 47822 (10.2%) organisation records are detected as duplicates when the heuristic is that they share 90% or more of trigrams in their name. By checking whether postcodes match, this number rises to 5935 (12.4%).
Manual testing of this method found 99 of 100 linked organisation records were indeed duplicates (though an additional 36 could not be determined). This suggests combining such contextual information is more effective than simply using name similarity.

Limitations to these heuristics include:

- Incomplete records: Many omitted address information, meaning solely name comparison could be performed.
- User error: Postcodes being specified in the incorrect field or differing by a couple of characters.
- Limited amount of information: Analysing sets of organisations collaborated with is less reliable for records with few projects associated with them.

Public government records could be queried to search previous names and addresses of organisations. These combinations of names and addresses could be compared between pairs of records to provide additional information to the de-duplication process.
The dates that these records were created is also provided and could be used to identify what an organisation was named and what address they were using at a certain point in time.
This information could provide more confidence in the decisions made during the de-duplication process.

Another possible heuristic is to compare the set of organisations that a pair of organisations have previously collaborated with, or to compare the research topics of projects they have been involved in.

# 7 Data analysis and visualisations

## 7.1 Grouping organisations by type

A variety of organisations exist in the database, from universities, local councils, medical institutions, to private companies.
For the purposes of analysis, they have been grouped into four categories:

- **Academic**: For academic institutions like universities, colleges, and other schools
- **Medical**: Medical institutions like hospitals, medical research institutions, and other healthcare providers
- **Private**: Privately owned corporations
- **Public**: Government bodies and publicly funded institutions

Organizations are placed in each category based on pattern matching of the organization names. More precisely, organisation records (and records that are listed as duplicates) are searched for keywords that indicate their type:

| Type | Exemplary keywords |
|------|--------------------|
| Academic | University, College, Academy, Academic |
| Medical | Hospital, NHS, Medical, Health |
| Private | Limited/LTD, Corporation, Company, Incorporated |
| Public | Council, Government, Governorate |

The organisation type is stored as an additional field within the organisations database table (`orgs`), which contains records which have been de-duplicated and cleansed of invalid records. Grouping is performed in the procedure [classifyOrgs.sql].

After applying this procedure, 21694 of 39578 organisation records (54.8%) were given a type.

The procedure was manually enhanced by refining the keyword list and controlling the scope of pattern matching. For example, abbreviations, such as LTD, LLP, and PLC are also used for identifying private limited companies alongside the 'Limited' keyword. The procedure also considers the optional period at the end of these abbreviations, such as "Uni." for "University". Another issue is that of keyword overlap: e.g. is Albany Medical College a medical institution or an academic one? Through manual research we can discover that it is indeed an academic one, but as the procedure sets the group of medical organisations after academic ones, these are categorised as medical.

An improvement to this procedure could be to consider combinations of words that appear in the outliers, ordered to prioritise these patterns over single keywords.
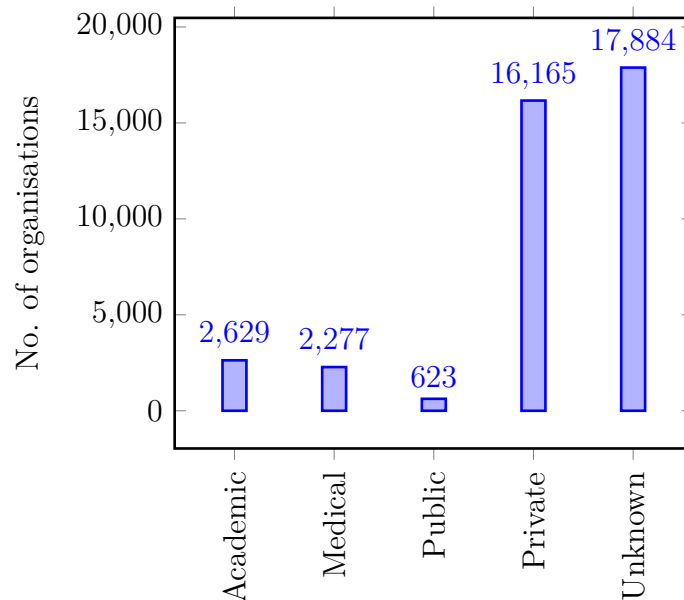


Figure 4: Distribution of organisations within the UKRI dataset, grouped by type. The type is determined using keywords contained in their names, such as 'Limited', or 'University'. 54.8% of records were grouped, with the remainder given the 'Unknown' type.

Due to the limitation of time, I have not endeavour to collect more information about individual

organizations. For that purpose, a public database such as the UK's Companies House[28] service could be used to determine both the legal classification and the nature of business of an organisation. The nature of business is reported by the company themselves, and so may not be entirely reliable: e.g., the Nottingham University Hospitals Trust Charity[29] lists their nature of business as "hospital activities" despite also being an educational institution.

The categories used by Companies House are provided on their website: ("Nature of Business: Standard Industrial Classification (Sic) Codes," n.d.).

## 7.2 East Midlands network analysis and visualisation with Gephi

In order to explore the structure of research ecosystems, network analysis was employed. A subset of the network was chosen so that patterns could be explored at a smaller scale. Organisations were selected if they had 'East Midlands' specified as their region and had an active projected within the years 2002 and 2008.

Gephi[30] is a Java-based graph analysis tool, and was used to plot the network of organisations based in the East Midlands region or Nottingham.

This software was chosen as it is free and open source, and so it costs nothing to use. Documentation also suggested it was easy to use, which would enable me to quickly get to grips with the application in order to use it in this project. Through research I also found that it includes several features that enable network analysis and visualisation, the possibilities of which I was keen to apply to explore.

Limiting the dataset used resulted in analysis being much faster to perform, as only 238 of 39578 organisations match these criteria. Fewer records mean less computational power is needed to manipulate the dataset in Gephi.

While Gephi is able to load the entire dataset if the Java virtual machine is allowed to use more system memory (e.g. by setting `_JAVA_OPTIONS="-Xms1024m -Xmx10000m"`), adjusting the layout, filtering, or calculating statistics still took a significant amount of time.

I also found that limiting the dataset significantly increased the legibility of the resulting graph, as having fewer nodes and edges means less overlap. With all nodes and edges visible, it was impossible to distinguish individual edges due to the massive number of them - 283197 in total. In the filtered dataset this number is reduced to a much more manageable 940.

While Gephi supports importing data directly from a database through queries, I encountered difficulties when attempting to use this feature with long queries and large amounts of data.

I instead opted to create procedures that exported the needed data to Comma Separated Values (CSV) files, a common plaintext format, and then imported these files manually into Gephi.

Data was exported from the database using [scripts/eastMidlandsGraphGephi.sql], then imported into Gephi through the import spreadsheet[31] wizard. This includes organisation records (used as the nodes) and projects that pairs of those organisations were both involved in (used as the edges). Organisation records were inspected for duplicates by sorting by name, and some were manually merged within Gephi.

### 7.2.1 Analysis

Gephi allows some graph properties to be analysed, including the degree distribution of nodes and (weighted) clustering coefficient.

---

[28]https://beta.companieshouse.gov.uk

[29]https://beta.companieshouse.gov.uk/company/09978675

[30]https://gephi.org

[31]https://github.com/gephi/gephi/wiki/Import-CSV-Data

Figure 5: Degree distribution of the East Midlands funding network (2002 to 2008), where nodes are organisations and edges are projects those organisations are collaborating on.

The degree distribution shown in figure 5 is similar to that of the network as a whole (see Amount of research), with the vast majority of nodes (i.e. organisations) having a low number of connections, indicating that they are involved in only a few collaborations.

Frequencies drop dramatically past degree 10, and very few organisations have the highest degrees, and these are particularly sparse within the distribution. This supports the expectation that there are a few organisations that make up the most of research and collaboration, and that most are only involved in a few projects.

As the distribution is similar to that of the database as a whole, this suggests that publicly funded research in the East Midlands region does not stand out significantly from other regions in terms of amount of research and its distribution.



Figure 6: Weighted degree distribution of the East Midlands funding network (2002-2008), where nodes are organisations, edges are projects those organisations are collaborating on, and the weight of edges is the total spent on the project between both organisations.

The weighted degree distribution (where each edge is multiplied by its weight) for the graph is shown in figure 6, where the weight is the cost of the project between the two organisations. A similar pattern is shown to the degree distribution; lower values have significantly greater frequencies, and higher costs becoming rarer as indicated by the sparsity at higher values.



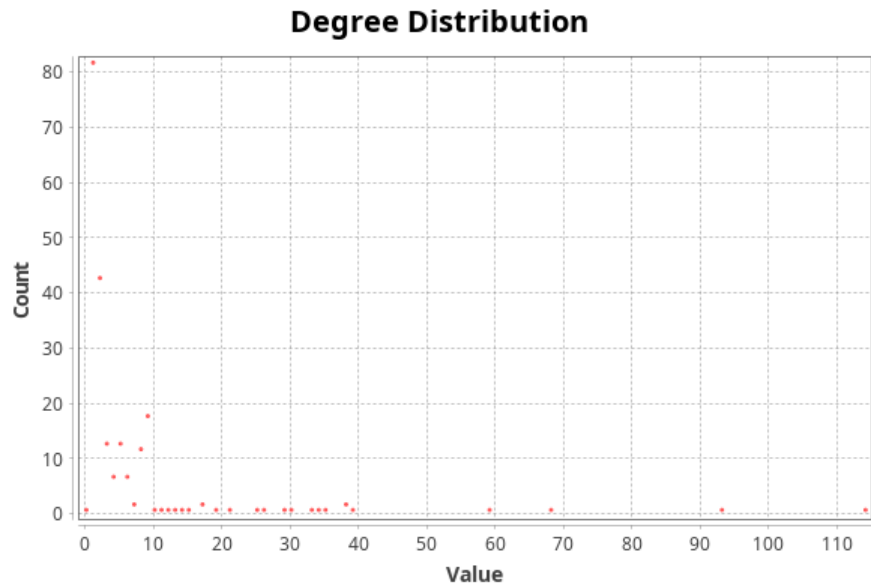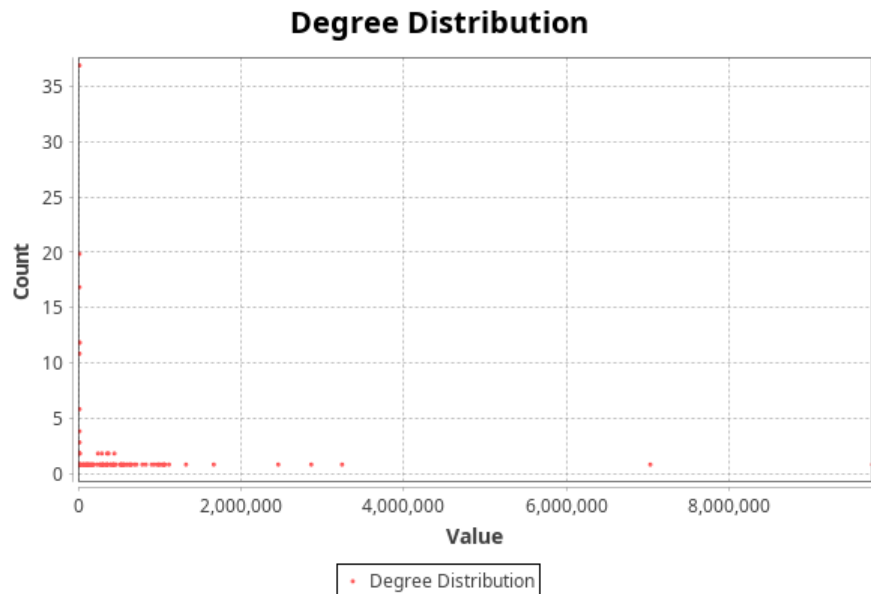Figure 7: Clustering coefficient distribution of the East Midlands funding network (2002 to 2008), where nodes are organisations and edges are projects those organisations are collaborating on. It shows that organisations are much more likely to only be involved in a few collaborations, and organisations with high amounts of collaboration are few and far between.

The clustering coefficient is a measure of "the degree to which nodes tend to cluster together" (Opsahl and Panzarasa (2009)), originally attempted by Luce and Perry (1949).

When considering a single node, its coefficient indicates how close that node's neighbours are to being a fully-connected graph: a lower value indicates less connection, higher indicates more. Gephi calculates this for triplets (referred to as 'triangles') of nodes in the graph.

Gephi was used to calculate the distribution of clustering coefficients for all nodes in the network, with the results shown in figure 7. We can see that a large number of triangles have a clustering coefficient of 0, and slightly fewer have a coefficient of 1. This indicates that there are a similar number of triplets of organisations that have collaborated with each other, as there are those that have never collaborated with each other. The remaining node triangles are distributed between these extremes 0 and 1, with coefficients higher than 0.5 being slightly more frequent. These results show that most of the network is highly connected, suggesting there are many organisations in the East Midlands that engage in research with many other organisations in this region. There is also a significant portion that does not: these organisations are likely involved in few projects with few collaborators.

### 7.2.2 Visualisation

Nodes were filtered to those with at least one connection, then Yifan Yu's proportional graph layout algorithm (Hu (n.d.)) was applied. This layout algorithm clusters related nodes while minimising the amount of edge crossing. Clustering is a result of edges causing connected vertices to attract one another while all vertices repel each other.

Edges were scaled according to the amount spent on projects between those organisations.

Nodes were scaled by their (unfiltered) total number of connections and coloured according to their type:

- Purple: Academic
- Blue: Private
- Green: Medical
- Gray: Unknown

Nodes with no connections were removed, minor repositioning was performed to make the graph visualisation more compact, and some nodes were adjusted to fix overlapping labels.

The resulting visualisation can be found in figure 18.

The network shows academic organisations being involved in the majority of projects in this period & area, having many connections to other organisations through projects. Among these are the University of Nottingham, Loughborough University, and Loughborough College.
Medical organisations are also heavily involved in research, including the University Hospitals of Nottingham and Leicester. These have formed large hubs of connected organisations, showing the amount of research and variety of organisations they are involved in.
A few private and other classes of organisations stand out including PERA Innovation[32], a research association for the manufacturing sector, and Experian[33], a credit reporting company. Both have a significant number of connections and spend a lot on research.

Almost all organisations are connected to a single contiguous graph, but there are a few outliers that reside in their own disjoint networks, shown in the bottom left of the graph. These include Sun Chemical[34] who produce "printing inks, coatings and supplies".
This suggests these organisations are in niche industries or involved in research that does not benefit most organisations.

## 7.3 East Midlands network analysis and visualisation with NodeXL

NodeXL[35] is an add-on for Microsoft Office Excel by the Social Media Research Foundation, originally developed to explore and analyse social media networks. It can be used to explore any network, however. NodeXL was recommended by my supervisor, who had previously used it in her research, and kindly provided a training session on using it.
As the University provides complimentary access to the Microsoft Office suite, and a free version is provided, I decided to try applying it to this project.

Similarly to the Gephi network, I filtered organisations to the East Midlands region, for the same reasons as explained in [East Midlands network with Gephi].
Microsoft Excel performed much worse compared to Gephi when I attempted to handle the entire dataset. The application regularly became unresponsive when manipulating the dataset or when NodeXL was generating graphs.
I decided to explore a different year range for projects than was explored in the Gephi visualisation, from 2005 to 2010, so the differences in the networks can be explored.
This resulted in 327 nodes and 806 edges being selected for this network.

The procedure for exporting this data is provided in eastMidlandsGraphNodeXL.sql.

---

[32]https://www.perainternational.com/about

[33]https://www.experian.com

[34]https://www.sunchemical.com

[35]https://www.smrfoundation.org/nodexl

### 7.3.1 Analysis

In the free version of NodeXL, the available analyses are greatly limited. While it supports several algorithms Gephi also supports (such as clustering coefficients and eigenvector centrality), these are limited to the paid version. One available analysis available is degree distribution, which is shown in figure 8.

Just as in the 2005 to 2008 network, the vast majority of organisations are involved in only a few collaborations, with higher amounts of collaboration being increasingly rare. Degrees above 20 seem to have become rarer in the 2005 to 2010 time range, but the top degree is 172 compared to 2002 to 2008's 114. This suggests a lower amount of collaboration from most organisations in this later time period, but the top organisations are involved in more.



Figure 8: Degree distribution of the East Midlands funding network (2005-2010), where nodes are organisations and edges are projects those organisations are collaborating on. It shows that organisations are much more likely to only be involved in a few collaborations, and organisations with high amounts of collaboration are few and far between.

### 7.3.2 Visualisation

Nodes were scaled by the number of project connections they had within that time period, and edges scaled by how much was spent by either organisation on that project.

I applied the Fruchterman-Reingold layout algorithm (Fruchterman and Reingold (1991)), which considers edges in positioning. By adjusting the configuration of the repulsive force to 25.0 and number of iterations to 25, this resulted in nodes with greater degrees being positioned towards the centre of the graph, and other nodes being placed towards the periphery.

The layout algorithm was also configured to position small networks in the bottom-left of the graph.

Labels were added for nodes with the highest degrees (greater than 13), and some other ones close to the centre of the graph. Minor repositioning was performed to reduce overlap of these labels.

Nodes were coloured according to their type:

- Dark blue: Academic
- Red: Private
- Light green: Medical
- Dark green: Unknown
- Light blue: Public

The resulting visualisation can be found in figure 19.

Academic and Medical institutions are shown towards the centre of the graph and with high degrees, indicating they are heavily involved in collaborative research, as one might expect. This pattern is also visible in the 2002 to 2008 network. This graph also shows that most academic organisations have higher degrees than most, hence are involved in more collaboration.
With the exception of PERA Innovation, all the organisations mentioned in the analysis of the 2002 to 2008 network also appear as nodes with high degree in this new graph. This suggests the remaining organisations have maintained the large amount of research they perform, whereas PERA Innovation is less involved in publicly funded research within the East Midlands region.

Some particularly notable edges appear near the centre of the graph and travel downwards to the bottom.
These both represent a project that a significant amount of money was invested into to develop hydrogen fuel cell systems[36], led by Intelligent Energy Limited[37] and partnered with Frost Electronics Limited.

## 7.4 Important factors in collaboration

### 7.4.1 Amount of research

We can take the number of projects an organisation has been involved in to be an indicator of the amount of research activity they are involved in. If we explore organisations by the number of projects they are involved in, we can see that the top organisations are involved in significantly more projects than lower organisations. Below the 90th percentile, these organisations are involved in less than 10 distinct projects each, whereas this number becomes exponential past the 90th percentile:



Figure 9: Percentiles for the number of projects an organisation is involved in, across the entire UKRI database. Generated with `OrgProjectPercentiles` in stats.sql

Visualising the distribution of organisations within these percentiles shows a similar pattern: the vast majority are involved with only one or two projects:

---

[36]https://gtr.ukri.org/projects?ref=113057
[37]https://www.intelligent-energy.com

Figure 10: Cumulative distribution of all organisations in the UKRI database by percentile, ranked by the total number of projects they are involved in. The organisation counts are broken down by their type. E.g. the red vertical line indicates just under 20,000 organisations are ranked at the 40th percentile or below. Generated with `OrgProjectPercentileDist` in stats.sql

By exploring the relationship between project involvement percentiles and collaboration, we can see that the majority of organisations have collaborated with organisations of lower percentile ranks:



Figure 11: Cumulative distribution of all organisations in the UKRI database, filtered by ones that have collaborated with at least one other organisation at a certain percentile. Organisations are ranked by the total number of projects they are involved in. E.g. the red vertical line indicates that there are just under 30,000 organisations that have collaborated with at least one other organisation at or below the 60th percentile, when ranked by the number of projects this organisation has been involved in. Generated with `OrgProjectPercentileCollab` in stats.sql

While there are many organisations that have collaborated with lower ranked organisations, we still see a sharp increase past the 90th percentile. This indicates a similar proportion of organisations have only collaborated with the top rank researchers, suggesting that the number

of projects one organisation has taken part in is a factor in whether another organisation decides to collaborate with them, though not a very important one.

### 7.4.2 Amount of funding

I expect that the amount of funding an organisation has previously received is a strong indicator of the likelihood that another organisations will collaborate with them.
Hypothetically, the UKRI would provide more funding to organisations who are the leading experts in their area compared to those who are less reputable. Organisations would be more interested in collaborating with experts and those who would bring more funding to projects, as this could result in the outcomes of the project being more valuable.

Of projects that reported any offered funding, the mean funded amount was £471,995, with a large (population) standard deviation of £5,535,964. This indicates that there is a great variety in the amount of funding received by projects.



Figure 12: Percentiles of public funding that organisations have received for the projects the are involved in, as well as the total funding received by projects. E.g. the red line indicates that when ranked by the amount of funding received, the top 20least £300,000 in total funding. Generated with `ProjectFundingPercentiles` and `OrgFundingPercentiles` in stats.sql

The percentile distribution is very similar for the total funding received by both projects and organisations. Both curves show that the top 20% of organisations and projects account for the majority of funding received from UKRI.
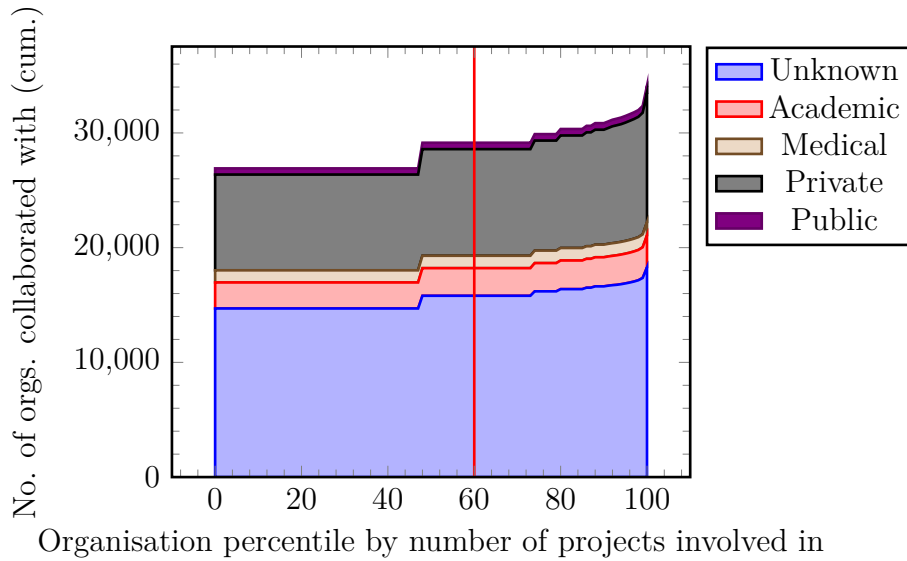
Figure 13: Cumulative distribution of all organisations in the UKRI database, filtered by ones that have collaborated with at least one other organisation at a certain percentile. Organisations are ranked by the total amount of funding they have received. E.g. the red vertical line indicates that there are just over 10,000 organisations that have collaborated with at least one other organisation at or below the 40th percentile, when ranked by the total amount of funding this organisation has received. Generated with `OrgFundingPercentileCollab` in stats.sql

By exploring the relationship between funding percentile and the number of organisations that have been collaborated with, one can see a rather positively linear relationship for the majority of the percentile distribution. This indicates that for organisations within the 20th to 90th percentile, the funding percentile does not have much influence in the likelihood that another organisation will collaborate with them.

However, from the 90th percentile and above, the number of organisations collaborated with dramatically increases in a more exponential manner. This suggests that a significant portion of organisations in the Gateway to Research database have only collaborated with the top 10% most funded organisations.

Figure 14: Cumulative distribution of projects by percentile, across the entire UKRI database. Projects are ranked by how much total funding they received. Generated with `OrgFundingPercentileProjects` in stats.sql

If we instead explore the total number of projects that involve at least one organisation at certain funding percentiles, it is trivial to identify a strong positive correlation between the funding percentile and number of projects these organisations are involved in.

A significant majority of projects involve organisations above the 90th percentile, suggesting that these organisations are also the most prolific in their involvement in publicly funded research.

These analyses show that organisations are more likely to have collaborated with only the most funded organisations, and that there are many projects involving only the most funded organisations.

While it is impossible to prove causality from this, the results support the theory that, when an organisation is deciding whether to collaborate with another, the total research funding which that other organisation has received is a factor.

# 8 Summary and Reflections

## 8.1 Conclusions

- What is the structure of the ecosystem of publicly funded research - how does it change over time?

Through aggregate analysis and specific analysis of the East Midlands network, a common pattern appears within the ecosystem wherein a few organisations (4252, or 9.92%) and projects (9723, or 10.0%) receive the vast majority of funding from UKRI and its subsidiaries. This pattern applies to all types of organisation, though academic and privately-owned organisations are likelier to be among these top researchers.

Private organisations represent the overwhelming majority of all those taking part in this research.

Analysis of the East Midlands network reveals that the top researching organisations form hubs of research, wherein many organisations collaborate with solely this top organisation. Within the 2002-2008 period, such hubs include:

- The University of Nottingham: 114 collaborations

- Loughborough University: 93
- University Hospitals of Leicester: 68
- Nottingham University Hospitals: 59
- PERA Innovation: 38
- Experian: 17

Comparing two different time periods (2002-2008, and 2005-2010) in this network there reveals not much change within the ecosystem in terms of distribution of collaboration among the involved organisations. However, while the top academic and medical organisations remain at the top, the private organisations involved vary over time. This suggests private organisations do not always maintain the amount of research they are involved with as consistently as these other types of organisations.

- What are the significant factors that influence collaboration between organisations?

Further aggregate analysis revealed that organisations are more likely to collaborate with the top organisation when ranked by either amount of funding received or the number of projects involved.

This suggests a relationship between these attributes and one organisation's decision to collaborate with these top organisations, perhaps due to this indicating the quality or value of the research they engage in.

## 8.2 Contributions

The project succeeded in going through the entire process of data warehousing and analysis for a set of open data, as originally intended.

Research did not bring up any other works that focused entirely on UKRI's Gateway to Research database and explored the dataset in its entirety. Previous research involved creating systems around such open data, including EnAKTing (Shadbolt et al. (2012)) and Dbpedia (Auer et al. (2007)). While this project did not go as far as creating an entire platform for users to perform their own analysis, I believe some useful analyses were performed and that it provides a valuable starting point from which to create such a platform.

By applying existing computational methods to an entire dataset in this project, I believe it provides a useful case study that supports future research that applies the chosen methods on open datasets. This includes difficulties encountered in applying them, methods for circumventing or alleviating these difficulties, and a data point from which to gauge their efficacy or usefulness.

I also believe the projects provides additional insights into possible improvements to open data platforms that would make similar future research easier. Details provided about de-duplication and schema merging performed during data cleaning could be applied by these open data platforms.

E.g. by showing similar existing records in data entry forms in order to avoid the creation of duplicate records. The project found overlapping semantics in the schema used by the Gateway to Research database, which may be a result of it encapsulating data from several subsidiary organisations. UKRI could apply these findings by reviewing the existing schema and applying changes to reduce this overlap.

Lack of documentation about the Gateway to Research's schema required manual work in exploring the structure of the dataset in order to convert it into a relational schema that could be implemented by standard relational databases. Similarly, documentation about the contents of the database tended to be outdated or lacking in detail. In order to support use of their database by researchers, improving this documentation is a step that UKRI could take.

## 8.3 Reflections

What the project achieved differs from the original plan, mostly due to the time constrictions of this project.

While the data warehousing tasks were achieved in some capacity, the originally planned analyses had to be cut down for the project to be completed on time. This includes using correlation analysis and machine learning for predicting possible collaboration, identifying technologies being researched, and estimating the market readiness of those technologies.

Time spent on individual tasks was typically longer than planned (with the plan laid out in the Gantt chart below). This was both due to the allocated time being insufficient, and due to less focus being put on the project than it required during the earlier periods of the project. This is main reason for reducing the amount of analysis within the project.

Initially The Decision Project (TDP), a business consultancy, intended to have a more active part in the project by providing additional contextual data about organisations within their own datasets. This would have provided more possibilities for analysing the impact of research on businesses, but required more work be done to incorporate these datasets into the database formed in this project.

While I planned to perform tasks in the order specified in the Methodology, I found myself continually going back to adjust previously done work, such as adjusting the schema or de-duplication procedures. This made keeping track of the pace of the project to be difficult, resulting in tasks and priorities continually shifting as the project progressed.

I ceased updating the work plan shown in the Gantt chart, as well not tracking tasks with Taskwarrior[38] as originally planned - I found having to do this regularly as plans shifted both unproductive and demoralising. Instead of setting fixed dates by which to accomplish certain things, I would have preferred to use a more agile project management methodology such as Kanban Ahmad, Markkula, and Oivo (2013). In such a method, the tasks that compose the project could be laid out individually, prioritised, and their dependencies linked. This would allow the project plan to be flexibly adjusted when tasks take a different amount of time to complete than predicted.

I also regret not using version control software from the beginning of the project. I believe it would've provided additional help in managing the pace of work and could have been paired nicely with a Kanban system - where logged changes work towards completing individual tasks. A version control system was applied to the project very late, where most of the benefits were in the ability to tracking changes in documentation.

---

[38]https://taskwarrior.org

Figure 15: Gannt chart visualising the work plan for the project.

# 9 Bibliography

Ahmad, Muhammad Ovais, Jouni Markkula, and Markku Oivo. 2013. "Kanban in Software Development: A Systematic Literature Review." In *2013 39th Euromicro Conference on Software Engineering and Advanced Applications*, 9–16. IEEE.

Auer, Sören, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. "Dbpedia: A Nucleus for a Web of Open Data." In *The Semantic Web*, 722–35. Springer.

Bastian, Mathieu, Sebastien Heymann, and Mathieu Jacomy. 2009. "Gephi: An Open Source Software for Exploring and Manipulating Networks." http://www.aaai.org/ocs/index.php/ICWSM/09/paper/view/154.

Dornum, Deirdre Dionysia von. 1997. "The Straight and the Crooked: Legal Accountability in Ancient Greece." *Columbia Law Review*, 1483–1518.

Elmagarmid, Ahmed K, Panagiotis G Ipeirotis, and Vassilios S Verykios. 2006. "Duplicate Record Detection: A Survey." *IEEE Transactions on Knowledge and Data Engineering* 19 (1): 1–16.

Fruchterman, Thomas MJ, and Edward M Reingold. 1991. "Graph Drawing by Force-Directed Placement." *Software: Practice and Experience* 21 (11): 1129–64.

*Gateway to Research Api 2.* 2018. Version 1.7.4. UK Research and Innovation. https://gtr.ukri.org/resources/GtR-2-API-v1.7.4.pdf.

*GtR Data Dictionary.* n.d. Version 8.0. UK Research and Innovation. https://gtr.ukri.org/resources/GtRDataDictionary.pdf.

Hu, Yifan. n.d. "Efficient and High Quality Force-Directed Graph Drawing." Wolfram Research Inc. http://yifanhu.net/PUB/graph_draw_small.pdf.

Lathrop, D., and L. Ruma. 2010. Theory in Practice. O'Reilly Media.

Lee, Heeseok. 1995. "Justifying Database Normalization: A Cost/Benefit Model." *Information Processing & Management* 31 (1): 59–67.

Luce, R Duncan, and Albert D Perry. 1949. "A Method of Matrix Analysis of Group Structure." *Psychometrika* 14 (2): 95–116.

Manning, Christopher, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. "The Stanford Corenlp Natural Language Processing Toolkit." In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 55–60.

"Nature of Business: Standard Industrial Classification (Sic) Codes." n.d. Companies House. https://resources.companieshouse.gov.uk/sic.

"Open Definition 2.1." 2019. Open Knowledge Foundation. 2019. https://opendefinition.org/od/2.1/en/.

Opsahl, Tore, and Pietro Panzarasa. 2009. "Clustering in Weighted Networks." *Social Networks* 31 (2): 155–63.

Rahm, Erhard, and Hong Hai Do. 2000. "Data Cleaning: Problems and Current Approaches." *IEEE Data Eng. Bull.* 23 (4): 3–13.

Schauer, Frederick. 2011. "Transparency in Three Dimensions." *University of Illinois Law Review*, 1339–58.

Shadbolt, Nigel, Kieron O'Hara, Tim Berners-Lee, Nicholas Gibbins, Hugh Glaser, Wendy Hall, and others. 2012. "Linked Open Government Data: Lessons from Data.gov.uk." *IEEE Intelligent Systems* 27 (3): 16–24.

Winkler, William E. 2006. "Overview of Record Linkage and Current Research Directions." In *Bureau of the Census*. Citeseer.

"XML Path Language (Xpath) 3.1." 2017. Standard. World Wide Web Consortium. https://www.w3.org/TR/2017/REC-xpath-31-20170321.

# 10 Appendix

## 10.1 ukriDownload

```bash
#!/usr/bin/env bash
URL='https://gtr.ukri.org/gtr/api/'
PAUSE=4
NUM_PAGES_XPATH='string(/*/@*[local-name()="totalPages"])'
OUTPUT='.'

# Parse options: https://stackoverflow.com/a/14203146
TYPES=()
while [[ $# -gt 0 ]]; do
    key="$1"

    case $key in
        -o|--output)      OUTPUT="$2";          shift 2;;
        -u|--url)         URL="$2";             shift 2;;
        -p|--pause)       PAUSE="$2";           shift 2;;
        -x|--pages-xpath) NUM_PAGES_XPATH="$2"; shift 2;;
        *)                TYPES+=("$1");        shift 1;;
    esac
done

set -- "${TYPES[@]}"

if [ "${#TYPES[@]}" == 0 ]; then
    echo "No record types specified"
    exit 1
fi

for typ in "$TYPES"; do
    case "$typ" in
        or*) recordType='organisations';;
        ou*) recordType='outcomes';;
        pe*) recordType='persons';;
        pr*) recordType='projects';;
        *) (>&2 echo "Unknown record type '$typ' - skipping"); continue;;
    esac

    url="$URL$recordType"
    mkdir -p "$OUTPUT/$recordType"

    function getPageURL {
        echo "$url?p=$1&s=100"
    }

    function getPageFile {
        echo "$OUTPUT/$recordType/$(printf '%04d' "$1").xml"
    }

    function xmlCount {
        awk '{s+=$1} END {print s}' <(xmllint --xpath 'count(/*/*)' $*)
    }

    file1="$(getPageFile 1)"
    curl "$(getPageURL 1)" -so "$file1"
    numPages=$(xmllint --xpath "$NUM_PAGES_XPATH" "$file1")

    echo "Total pages: $numPages"
```

```
57
58      for i in $(seq 2 $numPages); do
59          file="$(getPageFile $i)"
60          [ -f "$file" ] && continue
61          sleep "$PAUSE"
62          echo "Downloading page $i of $numPages"
63          curl "$(getPageURL $i)" -so "$file"
64          [ $? ] || break
65      done
66
67      merged="$OUTPUT/$recordType.xml"
68      pages="$OUTPUT/$recordType/*.xml"
69      if ! [ -f "$merged" ] || \
70          ! [ "$(xmlCount "$merged")" = "$(xmlCount "$pages")" ]; then
71          echo "Merging into a single file..."
72          "$(basename $0)"/mergeXML "$pages" >"$merged"
73      else
74          echo "Merged file contains all records: $merged"
75      fi
76  done
```

## 10.2   mergeXML

```python
#!/usr/bin/env python
# Adapted from https://stackoverflow.com/a/11315257
import sys

from xml.etree import ElementTree

def run(files):
    first = None

    for filename in files:
        data = ElementTree.parse(filename).getroot()

        if first is None:
            first = data
        else:
            first.extend(data)

    if first is not None:
        print(ElementTree.tostring(first, encoding='utf8').decode('utf-8'))

if __name__ == "__main__":
    run(sys.argv[1:])
```

## 10.3   ukriXmlToCsv

```python
#!/usr/bin/env python3
'''Convert a UKRI XML dataset to CSV.
Expected format is:
<records>
    <record></record>
    ...
</records>

Children of records will be flattened, with only the first sub-child being kept.
```

```python
'''

import argparse, csv, re, io, sys, os, xml.etree.ElementTree as ET
from typing import List, Optional
import os.path as path

usage_example = '''example: \
    ukriXmlToCsv \
        --fields name line1 line2 line3 line4 line5 postCode country \
        --output . \
        organisations.xml'''

parser = argparse.ArgumentParser(
    description='Convert UKRI XML record collections to CSV',
    epilog=usage_example)
parser.add_argument('--fields', '-f',
    nargs='+', metavar='pattern', action='extend', type=str,
    help='patterns of fields to include in the output')
parser.add_argument('--output', '-o',
    default='.',
    help='output directory for the converted files')
parser.add_argument('files',
    metavar='file', nargs='+',
    help='XML files to convert')

args = parser.parse_args()

os.makedirs(args.output, exist_ok=True)

def strip_namespace(field_tag: str) -> str:
    '''Remove the namespace from a field's tag'''
    i = field_tag.find('}')
    if i >= 0:
        field_tag = field_tag[i + 1:]
    return field_tag

def match_field_name(field: ET.Element) -> Optional[str]:
    '''Return whether a field with a specific tag should be filtered, according
    to the ignored field filters'''
    field_name = strip_namespace(field.tag)
    for pattern in args.fields:
        if re.match(pattern, field_name):
            return pattern
    return None

def filter_fields(fields: List[ET.Element]):
    return list(filter(lambda f: match_field_name(f) != None, fields))

def xml_to_csv(fname: str) -> str:
    output = io.StringIO()
    tree = ET.parse(fname, parser=ET.XMLParser(encoding='utf-8'))

    writer = csv.DictWriter(output, fieldnames=args.fields, delimiter=',')
    writer.writeheader()

    for record in tree.getroot():
        record_dict: dict[(str, str)] = dict()
        fields = filter_fields(list(record))

        # Flatten fields
```

```
70          for field in record:
71              children = list(field)
72              if children:
73                  if len(children) > 1:
74                      print('Warning: more than one child in flattened field ' \
75                          f"{field.tag}", file=sys.stderr)
76                  fields.extend(filter_fields(field[0]))
77
78          for field in fields:
79              field_name = match_field_name(field)
80              if field_name:
81                  record_dict[field_name] = field.text.strip()
82
83          writer.writerow(record_dict)
84
85      return output.getvalue()
86
87  for f in args.files:
88      csv = xml_to_csv(f)
89      csv_name = path.join(args.output,
90          path.splitext(path.basename(f))[0] + '.csv')
91
92      with open(csv_name, 'w', encoding='utf-8') as out:
93          out.write(csv)
```

## 10.4   setup.sql

```
1   -- Set up the database schema
2   -- Execute using:
3   -- $ psql -U <username> -d <database> -f setup.sql
4
5   -- Provides UUID data types
6   CREATE EXTENSION IF NOT EXISTS "uuid-ossp";
7
8   -- Provides similarity comparisons
9   CREATE EXTENSION IF NOT EXISTS "pg_trgm";
10
11  -- Provides crosstab()
12  CREATE EXTENSION IF NOT EXISTS "tablefunc";
13
14  -- GtR organisations
15  -- based on https://gtr.ukri.org/gtr/api
16  -- and https://gtr.ukri.org/gtr/api/organisation
17
18  DO $$ BEGIN
19      CREATE TYPE gtrRegion as ENUM(
20          'Channel Islands/Isle of Man',
21          'East Midlands',
22          'East of England',
23          'London',
24          'North East',
25          'Northern Ireland',
26          'North West',
27          'Outside UK',
28          'Scotland',
29          'South East',
30          'South West',
31          'Unknown',
32          'Wales',
```

```sql
33          'West Midlands',
34          'Yorkshire and The Humber'
35      );
36  EXCEPTION
37      WHEN duplicate_object THEN null;
38  END $$; -- BEGIN
39
40  DO $$ BEGIN
41      CREATE TYPE gtrSector as ENUM(
42          'Academic/University',
43          'Charity/Non-profit',
44          'Public',
45          'Private'
46      );
47  EXCEPTION
48      WHEN duplicate_object THEN null;
49  END $$; -- BEGIN
50
51  CREATE TABLE IF NOT EXISTS gtrOrgs(
52      orgUuid uuid PRIMARY KEY NOT NULL,
53      name text not null,
54      address1 text,
55      address2 text,
56      address3 text,
57      address4 text,
58      address5 text,
59      postCode text,
60      city text,
61      region gtrRegion not null default 'Unknown',
62      country text,
63      recorded date
64  );
65
66  -- GtR projects
67  -- based on https://gtr.ukri.org/gtr/api/project
68
69  DO $$ BEGIN
70      CREATE TYPE gtrProjectStatus as ENUM(
71          'Active',
72          'Closed'
73      );
74  EXCEPTION
75      WHEN duplicate_object THEN null;
76  END $$; -- BEGIN
77
78  DO $$ BEGIN
79      CREATE TYPE gtrGrantCategory as ENUM(
80          'BIS-Funded Programmes',
81          'Centres',
82          'Collaborative R&D',
83          'CR&D Bilateral',
84          'EU-Funded',
85          'European Enterprise Network',
86          'Fast Track',
87          'Feasibility Studies',
88          'Fellowship',
89          'GRD Development of Prototype',
90          'GRD Proof of Concept',
91          'GRD Proof of Market',
92          'Intramural',
```

```
 93            'Knowledge Transfer Network',
 94            'Knowledge Transfer Partnership',
 95            'Large Project',
 96            'Launchpad',
 97            'Legacy Department of Trade & Industry',
 98            'Legacy RDA Collaborative R&D',
 99            'Legacy RDA Grant for R&D',
100            'Missions',
101            'Other Grant',
102            'Procurement',
103            'Research Grant',
104            'Small Business Research Initiative',
105            'SME Support',
106            'Special Interest Group',
107            'Studentship',
108            'Study',
109            'Third Party Grant',
110            'Training Grant',
111            'Unknown',
112            'Vouchers'
113        );
114 EXCEPTION
115        WHEN duplicate_object THEN null;
116 END $$; -- BEGIN
117
118 DO $$ BEGIN
119        CREATE TYPE gtrFunder as ENUM(
120            'AHRC',
121            'BBSRC',
122            'EPSRC',
123            'ESRC',
124            'Innovate UK',
125            'MRC',
126            'NC3Rs',
127            'NERC',
128            'STFC',
129            'UKRI'
130        );
131 EXCEPTION
132        WHEN duplicate_object THEN null;
133 END $$; -- BEGIN
134
135 CREATE TABLE IF NOT EXISTS gtrProjects(
136        projectUuid uuid PRIMARY KEY NOT NULL,
137        title text not null,
138        status gtrProjectStatus,
139        category gtrGrantCategory,
140        leadFunder gtrFunder,
141        abstract text,
142        techAbstract text,
143        potentialImpact text,
144        startDate date,
145        endDate date,
146        recorded date
147 );
148
149 CREATE TABLE IF NOT EXISTS gtrSubjects(
150        subjectUuid uuid PRIMARY KEY NOT NULL,
151        name text
152 );
```

```sql
CREATE TABLE IF NOT EXISTS gtrProjectSubjects(
    subjectUuid uuid REFERENCES gtrSubjects,
    projectUuid uuid REFERENCES gtrProjects,
    percent numeric,
    PRIMARY KEY (subjectUuid, projectUuid)
);

CREATE TABLE IF NOT EXISTS gtrTopics(
    topicUuid uuid PRIMARY KEY NOT NULL,
    name text
);

CREATE TABLE IF NOT EXISTS gtrProjectTopics(
    topicUuid uuid REFERENCES gtrTopics,
    projectUuid uuid REFERENCES gtrProjects,
    percent numeric,
    PRIMARY KEY (topicUuid, projectUuid)
);

-- GtR outcomes
-- based on https://gtr.ukri.org/gtr/api/outcome

CREATE TABLE IF NOT EXISTS gtrDisseminations(
    disUuid uuid PRIMARY KEY NOT NULL,
    title text,
    description text,
    form text,
    primaryAudience text,
    yearsOfDissemination text,
    results text,
    impact text,
    typeOfPresentation text,
    geographicReach text,
    partOfOfficialScheme boolean,
    supportingUrl text
);

CREATE TABLE IF NOT EXISTS gtrCollaborations(
    collabUuid uuid PRIMARY KEY NOT NULL,
    description text,
    parentOrganisation text,
    childOrganisation text,
    principalInvestigatorContribution text,
    partnerContribution text,
    startDate date,
    endDate date,
    sector gtrSector,
    country text,
    impact text,
    supportingUrl text
);

CREATE TABLE IF NOT EXISTS gtrKeyFindings(
    keyFindingUuid uuid PRIMARY KEY NOT NULL,
    description text,
    nonAcademicUses text,
    exploitationPathways text,
    sectors text,
    supportingUrl text
```

```sql
213  );
214
215  CREATE TABLE IF NOT EXISTS gtrFurtherFundings(
216      furtherFundingUuid uuid PRIMARY KEY NOT NULL,
217      title text,
218      description text,
219      narrative text,
220      amount money,
221      organisation text,
222      department text,
223      fundingId text,
224      startDate date,
225      endDate date,
226      sector gtrSector,
227      country text
228  );
229
230  CREATE TABLE IF NOT EXISTS gtrImpactSummaries(
231      impactSummaryUuid uuid PRIMARY KEY NOT NULL,
232      title text,
233      description text,
234      impactTypes text,
235      summary text,
236      beneficiaries text,
237      contributionMethod text,
238      sector gtrSector,
239      firstYearOfImpact int
240  );
241
242  CREATE TABLE IF NOT EXISTS gtrPolicyInfluences(
243      policyInfluenceUuid uuid PRIMARY KEY NOT NULL,
244      influence text,
245      type text,
246      guidelineTitle text,
247      impact text,
248      methods text,
249      areas text,
250      geographicReach text,
251      supportingUrl text
252  );
253
254  CREATE TABLE IF NOT EXISTS gtrResearchMaterials(
255      researchMatUuid uuid PRIMARY KEY NOT NULL,
256      title text,
257      description text,
258      type text,
259      impact text,
260      softwareDeveloped boolean,
261      softwareOpenSourced boolean,
262      providedToOthers boolean,
263      yearFirstProvided int,
264      supportingUrl text
265  );
266
267  -- GtR persons
268  -- based on https://gtr.ukri.org/gtr/api/person
269
270  CREATE TABLE IF NOT EXISTS gtrPeople(
271      personUuid uuid PRIMARY KEY NOT NULL,
272      firstName text,
```

```
273      otherNames text,
274      surname text,
275      email text,
276      orcId text
277  );
278
279  -- Similarity indices
280
281  CREATE INDEX IF NOT EXISTS gtrOrgs_name_gist ON gtrOrgs USING gist(
282      name gist_trgm_ops
283  );
284
285  CREATE INDEX IF NOT EXISTS gtrPeople_name_gist ON gtrPeople USING gist(
286      firstName gist_trgm_ops,
287      surname gist_trgm_ops
288  );
289
290  -- GtR links
291
292  DO $$ BEGIN
293      CREATE TYPE gtrPersonRole as ENUM(
294          'Principal Investigator',
295          'Co-Investigator',
296          'Project Manager',
297          'Fellow',
298          'Training Grant Holder',
299          'Primary Supervisor',
300          -- These do not appear in the API reference, so may be deprecated
301          'Researcher Co-Investigator',
302          'Researcher'
303      );
304  EXCEPTION
305      WHEN duplicate_object THEN null;
306  END $$; -- BEGIN
307
308  DO $$ BEGIN
309      CREATE TYPE gtrOrgRole as ENUM(
310          'Lead',
311          'Collaborating',
312          'Fellow',
313          'Project Partner',
314          'Funder',
315          'Co-Funder',
316          'Participant',
317          'Student Project Partner'
318      );
319  EXCEPTION
320      WHEN duplicate_object THEN null;
321  END $$; -- BEGIN
322
323  DO $$ BEGIN
324      CREATE TYPE gtrProjectRelation as ENUM(
325          'TRANSFER',
326          'STUDENTSHIP_FROM',
327          'TRANSFER_FROM',
328          'STUDENTSHIP'
329      );
330  EXCEPTION
331      WHEN duplicate_object THEN null;
332  END $$; -- BEGIN
```

```
333
334  CREATE TABLE IF NOT EXISTS gtrOrgPeople(
335      personUuid uuid REFERENCES gtrPeople,
336      orgUuid uuid REFERENCES gtrOrgs,
337      PRIMARY KEY (personUuid, orgUuid)
338  );
339
340  CREATE TABLE IF NOT EXISTS gtrProjectPeople(
341      projectUuid uuid REFERENCES gtrProjects,
342      personUuid uuid REFERENCES gtrPeople,
343      role gtrPersonRole,
344      PRIMARY KEY (projectUuid, personUuid, role)
345  );
346
347  CREATE TABLE IF NOT EXISTS gtrProjectOrgs(
348      projectUuid uuid REFERENCES gtrProjects,
349      orgUuid uuid REFERENCES gtrOrgs,
350      role gtrOrgRole NOT NULL,
351      startDate date,
352      endDate date,
353      cost money,
354      offer money,
355      PRIMARY KEY (projectUuid, orgUuid, role)
356  );
357
358  CREATE TABLE IF NOT EXISTS gtrRelatedProjects(
359      projectUuid1 uuid REFERENCES gtrProjects,
360      projectUuid2 uuid REFERENCES gtrProjects,
361      relation gtrProjectRelation NOT NULL,
362      startDate date,
363      endDate date,
364      PRIMARY KEY (projectUuid1, projectUuid2, relation)
365  );
366
367  -- GtR cleaning
368
369  CREATE TABLE IF NOT EXISTS junkGtrOrgs(
370      orgUuid uuid REFERENCES gtrOrgs PRIMARY KEY
371  );
372
373  CREATE TABLE IF NOT EXISTS similarGtrOrgs(
374      orgUuid1 uuid REFERENCES gtrOrgs,
375      orgUuid2 uuid REFERENCES gtrOrgs,
376      simTrigramName float,
377      simTrigramAddress float,
378      manualResult bool,
379      PRIMARY KEY (orgUuid1, orgUuid2)
380  );
381
382  CREATE TABLE IF NOT EXISTS similarGtrPeople(
383      personUuid1 uuid REFERENCES gtrPeople,
384      personUuid2 uuid REFERENCES gtrPeople,
385      simTrigram float,
386      manualResult bool,
387      PRIMARY KEY (personUuid1, personUuid2)
388  );
389
390  -- Merged organisations
391
392  CREATE TABLE IF NOT EXISTS duplicateGtrOrgs(
```

```sql
393        orgUuid uuid REFERENCES gtrOrgs,
394        duplicateUuid uuid UNIQUE REFERENCES gtrOrgs,
395        PRIMARY KEY (orgUuid, duplicateUuid)
396    );
397
398    DO $$ BEGIN
399        CREATE TYPE orgType as ENUM(
400            'Academic',
401            'Medical',
402            'Private',
403            'Public',
404            'Unknown'
405        );
406    EXCEPTION
407        WHEN duplicate_object THEN null;
408    END $$; -- BEGIN
409
410    CREATE TABLE IF NOT EXISTS orgs(
411        gtrOrgUuid uuid PRIMARY KEY NOT NULL,
412        name text not null,
413        address1 text,
414        address2 text,
415        address3 text,
416        address4 text,
417        address5 text,
418        postCode text,
419        city text,
420        region gtrRegion NOT NULL DEFAULT 'Unknown',
421        country text,
422        type orgType DEFAULT 'Unknown'
423    );
424
425    -- Merged people
426
427    CREATE TABLE IF NOT EXISTS duplicateGtrPeople(
428        personUuid uuid REFERENCES gtrPeople,
429        duplicateUuid uuid UNIQUE REFERENCES gtrPeople,
430        PRIMARY KEY (personUuid, duplicateUuid)
431    );
432
433    CREATE TABLE IF NOT EXISTS people(
434        gtrPersonUuid uuid PRIMARY KEY NOT NULL,
435        firstName text,
436        surname text,
437        otherNames text
438    );
439
440    -- Utility functions
441
442    -- Remove all whitespace in some text
443    CREATE OR REPLACE
444    FUNCTION STRIP(txt text)
445    RETURNS text AS $$
446    BEGIN
447        RETURN regexp_replace(txt, '\s', '');
448    END; $$ -- FUNCTION
449    LANGUAGE plpgsql IMMUTABLE;
450
451    -- Generate a list of percentile fractions in the specified range
452    CREATE OR REPLACE
```

```
453   FUNCTION PERCENTILE_FRACTIONS(
454       lowerBound int,
455       upperBound int,
456       step int DEFAULT 1
457   )
458   RETURNS numeric[] AS $$
459   BEGIN
460       RETURN (SELECT ARRAY(
461           SELECT (a.n::numeric) / 100
462           FROM GENERATE_SERIES(lowerBound, upperBound, step) AS a(n)
463       ));
464   END; $$ -- FUNCTION
465   LANGUAGE plpgsql IMMUTABLE;
466
467   CREATE OR REPLACE FUNCTION CoalesceAll(state anyelement, new anyelement)
468   RETURNS anyelement
469   AS $$
470   BEGIN
471       RETURN COALESCE(state, new);
472   END; $$ --FUNCTION
473   LANGUAGE plpgsql IMMUTABLE;
474
475   CREATE OR REPLACE AGGREGATE Merge(anyelement) (
476       sfunc = CoalesceAll,
477       stype = anyelement
478   );
```

## 10.5   importGtrOrgs.sql

```
1    -- Run with
2    -- $ psql -U <username> -d <database> -f importGtrOrgs.sql
3    -- Ensure the organisations.xml file exists in the PostgreSQL data directory
4    DO $$
5    DECLARE
6        file text = 'organisations.xml';
7        xmlRecords xml;
8        org record;
9        nss text[][2] := array[
10           array['api', 'http://gtr.rcuk.ac.uk/gtr/api'],
11           array['org', 'http://gtr.rcuk.ac.uk/gtr/api/organisation']
12       ];
13   BEGIN
14
15   xmlRecords := XMLPARSE(
16       DOCUMENT convert_from(pg_read_binary_file(file), 'UTF8'));
17
18   DROP TABLE IF EXISTS xmlImport;
19   CREATE TEMP TABLE xmlImport AS
20   SELECT
21       (xpath('//@api:id', x, nss))[1]::text AS uuid,
22       (xpath('//@api:created', x, nss))[1]::text::date AS recorded,
23       (xpath('//org:name/text()', x, nss))[1]::text AS name,
24       (xpath('//org:addresses[1]/api:address/api:line1/text()',   x, nss))[1]::text AS line1,
25       (xpath('//org:addresses[1]/api:address/api:line2/text()',   x, nss))[1]::text AS line2,
26       (xpath('//org:addresses[1]/api:address/api:line3/text()',   x, nss))[1]::text AS line3,
27       (xpath('//org:addresses[1]/api:address/api:line4/text()',   x, nss))[1]::text AS line4,
28       (xpath('//org:addresses[1]/api:address/api:line5/text()',   x, nss))[1]::text AS line5,
29       (xpath('//org:addresses[1]/api:address/api:postCode/text()', x, nss))[1]::text AS
         ↪ postCode,
```

```sql
      (xpath('//org:addresses[1]/api:address/api:city/text()',    x, nss))[1]::text AS city,
      (xpath('//org:addresses[1]/api:address/api:region/text()',  x, nss))[1]::text AS
          ↪  region,
      (xpath('//org:addresses[1]/api:address/api:country/text()', x, nss))[1]::text AS
          ↪  country
FROM unnest(xpath('//org:organisations/org:organisation', xmlRecords, nss)) x;

FOR org IN
    SELECT * FROM xmlImport
LOOP
    INSERT INTO gtrOrgs VALUES(
        uuid(org.uuid),
        replace(org.name, '&amp;', '&'),
        org.line1,
        org.line2,
        org.line3,
        org.line4,
        org.line5,
        org.postCode,
        org.city,
        COALESCE(gtrRegion(org.region), gtrRegion('Unknown')),
        org.country,
        org.recorded
    ) ON CONFLICT DO UPDATE SET
        recorded = org.recorded;
END LOOP;

END$$;

DROP TABLE IF EXISTS xmlImport;
```

## 10.6   importGtrPersons.sql

```sql
-- Run with
-- $ psql -U <username> -d <database> -f importGtrPersons.sql
-- Ensure the persons.xml file exists in the PostgreSQL data directory
DO $$
DECLARE
    file text := 'persons.xml';
    xmlRecords xml;
    per record;
    link record;
    nss text[][2] := array[
        array['api', 'http://gtr.rcuk.ac.uk/gtr/api'],
        array['per', 'http://gtr.rcuk.ac.uk/gtr/api/person']
    ];
BEGIN

xmlRecords := XMLPARSE(
    DOCUMENT convert_from(pg_read_binary_file(file), 'UTF8'));

DROP TABLE IF EXISTS xmlImport;
CREATE TEMP TABLE xmlImport AS
SELECT
    (xpath('//@api:id', x, nss))[1]::text::uuid AS uuid,
    (xpath('//per:firstName/text()', x, nss))[1]::text AS firstName,
    (xpath('//per:otherNames/text()', x, nss))[1]::text AS otherNames,
    (xpath('//per:surname/text()', x, nss))[1]::text AS surname,
    (xpath('//per:email/text()', x, nss))[1]::text AS email,
```

```
27      (xpath('//per:orcidId/text()', x, nss))[1]::text AS orcId,
28      (xpath('//api:links/api:link[@api:rel="EMPLOYED"]/@api:href', x, nss))[1]::text AS
        ↪ employerLinks,
29      (xpath('//api:links', x, nss))[1] AS links
30  FROM unnest(xpath('//per:persons/per:person', xmlRecords, nss)) x;

31

32  FOR per IN
33      SELECT * FROM xmlImport
34  LOOP
35      INSERT INTO gtrPeople VALUES(
36          per.uuid,
37          per.firstName,
38          per.otherNames,
39          per.surname,
40          per.email,
41          per.orcId
42      ) ON CONFLICT (personUuid) DO UPDATE SET
43          firstName  = per.firstName,
44          otherNames = per.otherNames,
45          surname    = per.surname,
46          email      = per.email,
47          orcId      = per.orcId;

48

49      IF per.employerLinks IS NOT NULL THEN
50          INSERT INTO gtrOrgPeople VALUES(
51              per.uuid,
52              -- Extract organisation UUID from HREF
53              substring(per.employerLinks, '/organisations/(.*)$')::uuid
54          ) ON CONFLICT DO NOTHING;
55      END IF;

56

57      FOR link IN
58          SELECT
59              (xpath('//@api:rel', x, nss))[1]::text AS rel,
60              (xpath('//@api:href', x, nss))[1]::text AS href
61          FROM unnest(xpath('//api:link', per.links, nss)) x
62      LOOP
63          IF EXISTS (SELECT projectUuid FROM gtrProjects WHERE
64              projectUuid = substring(link.href, '/projects/(.*)$')::uuid)
65          THEN
66              INSERT INTO gtrProjectPeople VALUES(
67                  substring(link.href, '/projects/(.*)$')::uuid,
68                  per.uuid,
69                  (CASE
70                      WHEN link.rel = 'PI_PER'            THEN 'Principal Investigator'
71                      WHEN link.rel = 'COI_PER'           THEN 'Co-Investigator'
72                      WHEN link.rel = 'PM_PER'            THEN 'Project Manager'
73                      WHEN link.rel = 'FELLOW_PER'        THEN 'Fellow'
74                      WHEN link.rel = 'TGH_PER'           THEN 'Training Grant Holder'
75                      WHEN link.rel = 'SUPER_PER'         THEN 'Primary Supervisor'
76                      WHEN link.rel = 'RESERACH_COI_PER' THEN 'Researcher Co-Investigator'
77                      WHEN link.rel = 'RESERACH_PER'     THEN 'Researcher'
78                      ELSE link.rel::gtrPersonRole
79                  END)::gtrPersonRole
80              ) ON CONFLICT DO NOTHING;
81          END IF;
82      END LOOP;
83  END LOOP;

84

85  END$$;
```

## 10.7 importGtrProjects.sql

```
1   -- Run with(xmlImport,(xmlImport,
2   -- $ psql -U <username> -d <database> -f importGtrProjects.sql
3   -- Ensure the projects.xml file exists in the PostgreSQL data directory
4
5   CREATE OR REPLACE
6   FUNCTION ParseOrgRole(txt text)
7   RETURNS gtrOrgRole AS $$
8   BEGIN
9       RETURN (CASE
10          -- Project participants
11          WHEN txt = 'LEAD_PARTICIPANT' THEN 'Lead'
12          WHEN txt = 'PARTICIPANT'      THEN 'Participant'
13          -- Project links
14          WHEN txt = 'LEAD_ORG'         THEN 'Lead'
15          WHEN txt = 'COLLAB_ORG'       THEN 'Collaborating'
16          WHEN txt = 'FELLOW_ORG'       THEN 'Fellow'
17          WHEN txt = 'PP_ORG'           THEN 'Project Partner'
18          WHEN txt = 'FUNDER'           THEN 'Funder'
19          WHEN txt = 'COFUND_ORG'       THEN 'Co-Funder'
20          WHEN txt = 'PARTICIPANT_ORG'  THEN 'Participant'
21          WHEN txt = 'STUDENT_PP_ORG'   THEN 'Student Project Partner'
22          ELSE txt
23      END)::gtrOrgRole;
24  END; $$ -- FUNCTION
25  LANGUAGE plpgsql IMMUTABLE;
26
27  DO $$
28  DECLARE
29      file text := 'projects.xml';
30      xmlRecords xml;
31      rec record;
32      nss text[][2] := array[
33          array['api', 'http://gtr.rcuk.ac.uk/gtr/api'],
34          array['pro', 'http://gtr.rcuk.ac.uk/gtr/api/project']
35      ];
36  BEGIN
37
38  xmlRecords := XMLPARSE(
39      DOCUMENT convert_from(pg_read_binary_file(file), 'UTF8'));
40
41  CREATE TEMP TABLE IF NOT EXISTS xmlImport AS
42  SELECT
43      (xpath('//@api:id', x, nss))[1]::text::uuid AS uuid,
44      (xpath('//@api:created', x, nss))[1]::text::date AS recorded,
45      (xpath('//pro:title/text()', x, nss))[1]::text AS title,
46      (xpath('//pro:status/text()', x, nss))[1]::text AS status,
47      (xpath('//pro:grantCategory/text()', x, nss))[1]::text AS category,
48      (xpath('//pro:leadFunder/text()', x, nss))[1]::text AS funder,
49      (xpath('//pro:abstract/text()', x, nss))[1]::text AS abstract,
50      (xpath('//pro:techAbstract/text()', x, nss))[1]::text AS techAbstract,
51      (xpath('//pro:potentialImpact/text()', x, nss))[1]::text AS potentialImpact,
52      (xpath('//api:links/api:link[@api:rel="FUND"]/@api:start', x, nss))[1]::text::date AS
        ↪  startDate,
53      (xpath('//api:links/api:link[@api:rel="FUND"]/@api:end', x, nss))[1]::text::date AS
        ↪  endDate,
54      (xpath('//api:links', x, nss))[1] AS links,
55      (xpath('//pro:participantValues', x, nss))[1] AS participants,
56      (xpath('//pro:researchSubjects', x, nss))[1] AS subjects,
```

```sql
57          (xpath('//pro:researchTopics', x, nss))[1] AS topics
58      FROM unnest(xpath('//pro:projects/pro:project', xmlRecords, nss)) x;
59
60      -- Projects
61
62      INSERT INTO gtrProjects
63      SELECT
64          uuid,
65          REPLACE(title, '&amp;', '&'),
66          status::gtrProjectStatus,
67          -- XML entities aren't parsed, so manually convert ampersands
68          REPLACE(category, '&amp;', '&')::gtrGrantCategory,
69          funder::gtrFunder,
70          abstract,
71          techAbstract,
72          potentialImpact,
73          startDate,
74          endDate,
75          recorded
76      FROM xmlImport
77      ON CONFLICT (projectUuid) DO NOTHING;
78
79      -- Project organisations
80
81      INSERT INTO gtrProjectOrgs(
82          projectUuid,
83          orgUuid,
84          role,
85          cost,
86          offer
87      )
88      SELECT
89          proUuid,
90          (xpath('//pro:organisationId/text()', x, nss))[1]::text::uuid AS orgUuid,
91          ParseOrgRole((xpath('//pro:role/text()', x, nss))[1]::text) AS role,
92          (xpath('//pro:projectCost/text()', x, nss))[1]::text::numeric::money AS cost,
93          (xpath('//pro:grantOffer/text()', x, nss))[1]::text::numeric::money AS offer
94      FROM (
95          SELECT
96              uuid AS proUuid,
97              unnest(xpath('//pro:participant', participants, nss)) AS x
98          FROM xmlImport
99      ) q
100     ON CONFLICT DO NOTHING;
101
102     -- Subjects
103
104     FOR rec IN
105         SELECT
106             proUuid,
107             (xpath('//pro:id/text()', x, nss))[1]::text::uuid AS subUuid,
108             (xpath('//pro:text/text()', x, nss))[1]::text AS name,
109             (xpath('//pro:percentage/text()', x, nss))[1]::text::numeric AS percent
110         FROM (
111             SELECT
112                 uuid AS proUuid,
113                 unnest(xpath('//pro:researchSubject', subjects, nss)) x
114             FROM xmlImport
115         ) q
116     LOOP
```

```sql
117        INSERT INTO gtrSubjects VALUES(
118            rec.subUuid,
119            rec.name
120        ) ON CONFLICT DO NOTHING;
121
122        INSERT INTO gtrProjectSubjects VALUES(
123            rec.subUuid,
124            rec.proUuid,
125            rec.percent
126        ) ON CONFLICT DO NOTHING;
127    END LOOP;
128
129    -- Topics
130
131    FOR rec IN
132        SELECT
133            proUuid,
134            (xpath('//pro:id/text()', x, nss))[1]::text::uuid AS topUuid,
135            (xpath('//pro:text/text()', x, nss))[1]::text AS name,
136            (xpath('//pro:percentage/text()', x, nss))[1]::text::numeric AS percent
137        FROM (
138            SELECT
139                uuid AS proUuid,
140                unnest(xpath('//pro:researchTopic', topics, nss)) x
141            FROM xmlImport
142        ) q
143    LOOP
144        INSERT INTO gtrTopics VALUES(
145            rec.topUuid,
146            rec.name
147        ) ON CONFLICT DO NOTHING;
148
149        INSERT INTO gtrProjectTopics VALUES(
150            rec.topUuid,
151            rec.proUuid,
152            rec.percent
153        ) ON CONFLICT DO NOTHING;
154    END LOOP;
155
156    -- Links
157
158    FOR rec IN
159        SELECT
160            proUuid,
161            (xpath('//@api:rel', x, nss))[1]::text AS rel,
162            href,
163            (xpath('//@api:start', x, nss))[1]::text::date AS startDate,
164            (xpath('//@api:end', x, nss))[1]::text::date AS endDate,
165            (SELECT orgUuid FROM gtrOrgs
166                WHERE orgUuid = substring(href, '/organisations/(.*)$')::uuid
167                LIMIT 1) AS orgUuid,
168            (SELECT proUuid FROM gtrProjects
169                WHERE projectUuid = substring(href, '/projects/(.*)$')::uuid
170                LIMIT 1) AS proUuid2
171        FROM (
172            SELECT
173                *,
174                (xpath('//@api:href', x, nss))[1]::text AS href
175            FROM (
176                SELECT
```

```
177                    uuid AS proUuid,
178                    unnest(xpath('//api:link', links, nss)) x
179                FROM xmlImport
180           ) q2
181       ) q
182   LOOP
183       -- Linked organisations
184       IF rec.orgUuid IS NOT NULL THEN
185           INSERT INTO gtrProjectOrgs VALUES(
186               rec.proUuid,
187               rec.orgUuid,
188               ParseOrgRole(rec.rel),
189               rec.startDate,
190               rec.endDate,
191               NULL,
192               NULL
193           ) ON CONFLICT DO NOTHING;
194
195       -- Related projects
196       ELSIF rec.proUuid2 IS NOT NULL THEN
197           INSERT INTO gtrRelatedProjects VALUES(
198               rec.proUuid,
199               rec.proUuid2,
200               rec.rel::gtrProjectRelation,
201               rec.startDate,
202               rec.endDate
203           ) ON CONFLICT DO NOTHING;
204       END IF;
205   END LOOP;
206
207   END$$;
```

## 10.8   importGtrOutcomes.sql

```
1    -- Run with
2    -- $ psql -U <username> -d <database> -f importGtrOutcomes.sql
3    -- Ensure the outcomes.xml file exists in the PostgreSQL data directory
4    DO $$
5    DECLARE
6        file text := 'outcomes.xml';
7        xmlRecords xml;
8        outcome record;
9        nss text[][2] := array[
10           array['api', 'http://gtr.rcuk.ac.uk/gtr/api'],
11           array['out', 'http://gtr.rcuk.ac.uk/gtr/api/project/outcome']
12       ];
13   BEGIN
14
15   xmlRecords := XMLPARSE(
16       DOCUMENT convert_from(pg_read_binary_file(file), 'UTF8'));
17
18   -- Disseminations
19
20   DROP TABLE IF EXISTS xmlImport;
21   CREATE TEMP TABLE xmlImport AS
22   SELECT
23       (xpath('//@api:id', x, nss))[1]::text AS uuid,
24       (xpath('//out:title/text()', x, nss))[1]::text AS title,
25       (xpath('//out:description/text()', x, nss))[1]::text AS description,
```

```
26      (xpath('//out:form/text()', x, nss))[1]::text AS form,
27      (xpath('//out:primaryAudience/text()', x, nss))[1]::text AS primaryAudience,
28      (xpath('//out:yearsOfDissemination/text()', x, nss))[1]::text AS yearsOfDissemination,
29      (xpath('//out:results/text()', x, nss))[1]::text AS results,
30      (xpath('//out:impact/text()', x, nss))[1]::text AS impact,
31      (xpath('//out:typeOfPresentation/text()', x, nss))[1]::text AS typeOfPresentation,
32      (xpath('//out:geographicReach/text()', x, nss))[1]::text AS geographicReach,
33      (xpath('//out:partOfOfficialScheme/text()', x, nss))[1]::text AS partOfOfficialScheme,
34      (xpath('//out:supportingUrl/text()', x, nss))[1]::text AS supportingUrl
35  FROM unnest(xpath('//out:outcomes/out:dissemination', xmlRecords, nss)) x;
36
37  FOR outcome IN
38      SELECT * FROM xmlImport
39  LOOP
40      INSERT INTO gtrDisseminations VALUES(
41          outcome.uuid::uuid,
42          outcome.title,
43          outcome.description,
44          outcome.form,
45          outcome.primaryAudience,
46          outcome.yearsOfDissemination,
47          outcome.results,
48          outcome.impact,
49          outcome.typeOfPresentation,
50          outcome.geographicReach,
51          bool(outcome.partOfOfficialScheme),
52          outcome.supportingUrl
53      ) ON CONFLICT DO NOTHING;
54  END LOOP;
55
56  -- Collaborations
57
58  DROP TABLE IF EXISTS xmlImport;
59  CREATE TEMP TABLE xmlImport AS
60  SELECT
61      (xpath('//@api:id', x, nss))[1]::text AS uuid,
62      (xpath('//out:description/text()', x, nss))[1]::text AS description,
63      (xpath('//out:parentOrganisation/text()', x, nss))[1]::text AS parentOrganisation,
64      (xpath('//out:childOrganisation/text()', x, nss))[1]::text AS childOrganisation,
65      (xpath('//out:principalInvestigatorContribution/text()', x, nss))[1]::text AS
        ↪ principalInvestigatorContribution,
66      (xpath('//out:partnerContribution/text()', x, nss))[1]::text AS partnerContribution,
67      (xpath('//out:start/text()', x, nss))[1]::text AS startDate,
68      (xpath('//out:end/text()', x, nss))[1]::text AS endDate,
69      (xpath('//out:sector/text()', x, nss))[1]::text AS sector,
70      (xpath('//out:country/text()', x, nss))[1]::text AS country,
71      (xpath('//out:impact/text()', x, nss))[1]::text AS impact,
72      (xpath('//out:supportingUrl/text()', x, nss))[1]::text AS supportingUrl
73  FROM unnest(xpath('//out:outcomes/out:collaboration', xmlRecords, nss)) x;
74
75  FOR outcome IN
76      SELECT * FROM xmlImport
77  LOOP
78      INSERT INTO gtrCollaborations VALUES(
79          outcome.uuid::uuid,
80          outcome.description,
81          outcome.parentOrganisation,
82          outcome.childOrganisation,
83          outcome.principalInvestigatorContribution,
84          outcome.partnerContribution,
```

```
85          date(outcome.startDate),
86          date(outcome.endDate),
87          outcome.sector,
88          outcome.country,
89          outcome.impact,
90          outcome.supportingUrl
91      ) ON CONFLICT DO NOTHING;
92  END LOOP;
93
94  -- Key findings
95
96  DROP TABLE IF EXISTS xmlImport;
97  CREATE TEMP TABLE xmlImport AS
98  SELECT
99      (xpath('//@api:id', x, nss))[1]::text AS uuid,
100     (xpath('//out:description/text()', x, nss))[1]::text AS description,
101     (xpath('//out:nonAcademicUses/text()', x, nss))[1]::text AS nonAcademicUses,
102     (xpath('//out:exploitationPathways/text()', x, nss))[1]::text AS exploitationPathways,
103     (xpath('//out:sectors/text()', x, nss))[1]::text AS sectors,
104     (xpath('//out:supportingUrl/text()', x, nss))[1]::text AS supportingUrl
105 FROM unnest(xpath('//out:outcomes/out:keyFinding', xmlRecords, nss)) x;
106
107 FOR outcome IN
108     SELECT * FROM xmlImport
109 LOOP
110     INSERT INTO gtrKeyFindings VALUES(
111         uuid(outcome.uuid),
112         outcome.description,
113         outcome.nonAcademicUses,
114         outcome.exploitationPathways,
115         outcome.sectors,
116         outcome.supportingUrl
117     ) ON CONFLICT DO NOTHING;
118 END LOOP;
119
120 -- Further fundings
121
122 DROP TABLE IF EXISTS xmlImport;
123 CREATE TEMP TABLE xmlImport AS
124 SELECT
125     (xpath('//@api:id', x, nss))[1]::text AS uuid,
126     (xpath('//out:title/text()', x, nss))[1]::text AS title,
127     (xpath('//out:description/text()', x, nss))[1]::text AS description,
128     (xpath('//out:narrative/text()', x, nss))[1]::text AS narrative,
129     (xpath('//out:amount/@api:amount', x, nss))[1]::text AS amount,
130     (xpath('//out:organisation/text()', x, nss))[1]::text AS organisation,
131     (xpath('//out:department/text()', x, nss))[1]::text AS department,
132     (xpath('//out:fundingId/text()', x, nss))[1]::text AS fundingId,
133     (xpath('//out:start/text()', x, nss))[1]::text AS startDate,
134     (xpath('//out:end/text()', x, nss))[1]::text AS endDate,
135     (xpath('//out:sector/text()', x, nss))[1]::text AS sector,
136     (xpath('//out:country/text()', x, nss))[1]::text AS country
137 FROM unnest(xpath('//out:outcomes/out:furtherFunding', xmlRecords, nss)) x;
138
139 FOR outcome IN
140     SELECT * FROM xmlImport
141 LOOP
142     INSERT INTO gtrFurtherFundings VALUES(
143         uuid(outcome.uuid),
144         outcome.title,
```

```
145          outcome.description,
146          outcome.narrative,
147          money(outcome.amount),
148          outcome.organisation,
149          outcome.department,
150          outcome.fundingId,
151          date(outcome.startDate),
152          date(outcome.endDate),
153          outcome.sector,
154          outcome.country
155      ) ON CONFLICT DO NOTHING;
156  END LOOP;
157
158  -- Impact summaries
159
160  DROP TABLE IF EXISTS xmlImport;
161  CREATE TEMP TABLE xmlImport AS
162  SELECT
163      (xpath('//@api:id', x, nss))[1]::text AS uuid,
164      (xpath('//out:title/text()', x, nss))[1]::text AS title,
165      (xpath('//out:description/text()', x, nss))[1]::text AS description,
166      (xpath('//out:impactTypes/text()', x, nss))[1]::text AS impactTypes,
167      (xpath('//out:summary/text()', x, nss))[1]::text AS summary,
168      (xpath('//out:beneficiaries/text()', x, nss))[1]::text AS beneficiaries,
169      (xpath('//out:contributionMethod/text()', x, nss))[1]::text AS contributionMethod,
170      (xpath('//out:sector/text()', x, nss))[1]::text AS sector,
171      (xpath('//out:firstYearOfImpact/text()', x, nss))[1]::text AS firstYearOfImpact
172  FROM unnest(xpath('//out:outcomes/out:impactSummary', xmlRecords, nss)) x;
173
174  FOR outcome IN
175      SELECT * FROM xmlImport
176  LOOP
177      INSERT INTO gtrImpactSummaries VALUES(
178          outcome.uuid::uuid,
179          outcome.title,
180          outcome.description,
181          outcome.impactTypes,
182          outcome.summary,
183          outcome.beneficiaries,
184          outcome.contributionMethod,
185          outcome.sector,
186          outcome.firstYearOfImpact::int
187      ) ON CONFLICT DO NOTHING;
188  END LOOP;
189
190  -- Policy influences
191
192  DROP TABLE IF EXISTS xmlImport;
193  CREATE TEMP TABLE xmlImport AS
194  SELECT
195      (xpath('//@api:id', x, nss))[1]::text AS uuid,
196      (xpath('//out:influence/text()', x, nss))[1]::text AS influence,
197      (xpath('//out:type/text()', x, nss))[1]::text AS type,
198      (xpath('//out:guidelineTitle/text()', x, nss))[1]::text AS guidelineTitle,
199      (xpath('//out:impact/text()', x, nss))[1]::text AS impact,
200      (xpath('//out:methods/text()', x, nss))[1]::text AS methods,
201      (xpath('//out:areas/text()', x, nss))[1]::text AS areas,
202      (xpath('//out:geographicReach/text()', x, nss))[1]::text AS geographicReach,
203      (xpath('//out:supportingUrl/text()', x, nss))[1]::text AS supportingUrl
204  FROM unnest(xpath('//out:outcomes/out:policyInfluence', xmlRecords, nss)) x;
```

```
205
206    FOR outcome IN
207        SELECT * FROM xmlImport
208    LOOP
209        INSERT INTO gtrPolicyInfluences VALUES(
210            outcome.uuid::uuid,
211            outcome.influence,
212            outcome.type,
213            outcome.guidelineTitle,
214            outcome.impact,
215            outcome.methods,
216            outcome.areas,
217            outcome.geographicReach,
218            outcome.supportingUrl
219        ) ON CONFLICT DO NOTHING;
220    END LOOP;
221
222    -- Research materials
223
224    DROP TABLE IF EXISTS xmlImport;
225    CREATE TEMP TABLE xmlImport AS
226    SELECT
227        (xpath('//@api:id', x, nss))[1]::text AS uuid,
228        (xpath('//out:title/text()', x, nss))[1]::text AS title,
229        (xpath('//out:description/text()', x, nss))[1]::text AS description,
230        (xpath('//out:type/text()', x, nss))[1]::text AS type,
231        (xpath('//out:impact/text()', x, nss))[1]::text AS impact,
232        (xpath('//out:softwareDeveloped/text()', x, nss))[1]::text AS softwareDeveloped,
233        (xpath('//out:softwareOpenSourced/text()', x, nss))[1]::text AS softwareOpenSourced,
234        (xpath('//out:providedToOthers/text()', x, nss))[1]::text AS providedToOthers,
235        (xpath('//out:yearFirstProvided/text()', x, nss))[1]::text AS yearFirstProvided,
236        (xpath('//out:supportingUrl/text()', x, nss))[1]::text AS supportingUrl
237    FROM unnest(xpath('//out:outcomes/out:researchMaterials', xmlRecords, nss)) x;
238
239    FOR outcome IN
240        SELECT * FROM xmlImport
241    LOOP
242        INSERT INTO gtrResearchMaterials VALUES(
243            outcome.uuid::uuid,
244            outcome.title,
245            outcome.description,
246            outcome.type,
247            outcome.impact,
248            bool(outcome.softwareDeveloped),
249            bool(outcome.softwareOpenSourced),
250            bool(outcome.providedToOthers),
251            outcome.yearFirstProvided::int,
252            outcome.supportingUrl
253        ) ON CONFLICT DO NOTHING;
254    END LOOP;
255    END$$; -- BEGIN
256
257    DROP TABLE IF EXISTS xmlImport;
```

## 10.9  similarGtrOrgs.sql

```
1    CREATE OR REPLACE FUNCTION REPLACE_MANY(txt text, replaces text[][2])
2    RETURNS text AS $$
3    DECLARE i int;
```

```
4    BEGIN
5        FOR i IN 1 .. array_upper(replaces, 1) LOOP
6            txt := REGEXP_REPLACE(txt, replaces[i][1], replaces[i][2], 'gi');
7        END LOOP;
8        RETURN txt;
9    END;
10   $$ LANGUAGE plpgsql IMMUTABLE;
11
12   DO $$
13   DECLARE
14       name_replaces text[][] := array[
15           array['\yltd(\.|\y)',  'Limited'],
16           array['\ycorp(\.|\y)', 'Corporation'],
17           array['\yco(\.|\y)',   'Company'],
18           array['\yuni(\.|\y)',  'University'],
19           array['&',             'and'],
20           array['The',           '']
21       ];
22   BEGIN
23       SET pg_trgm.similarity_threshold = 0.9;
24
25       INSERT INTO similarGtrOrgs(
26           orgUuid1,
27           orgUuid2,
28           simTrigramName,
29           simTrigramAddress
30       )
31       SELECT
32           uuid1,
33           uuid2,
34           similarity(name1, name2) AS simName,
35           similarity(addr1, addr2) AS simAddress
36       FROM
37           (SELECT
38               o1.orgUuid AS uuid1,
39               o2.orgUuid AS uuid2,
40               REPLACE_MANY(o1.name, name_replaces) AS name1,
41               REPLACE_MANY(o2.name, name_replaces) AS name2,
42               CONCAT(o1.address1, '\n',
43                       o1.address2, '\n',
44                       o1.address3, '\n',
45                       o1.address3, '\n',
46                       o1.address5, '\n',
47                       o1.postCode, '\n',
48                       o1.city) AS addr1,
49               CONCAT(o2.address1, '\n',
50                       o2.address2, '\n',
51                       o2.address3, '\n',
52                       o2.address3, '\n',
53                       o2.address5, '\n',
54                       o2.postCode, '\n',
55                       o2.city) AS addr2
56           FROM gtrOrgs o1
57           JOIN gtrOrgs o2
58               ON  o1.orgUuid < o2.orgUuid
59               AND REPLACE_MANY(o1.name, name_replaces)
60                 % REPLACE_MANY(o2.name, name_replaces)
61           WHERE NOT EXISTS(
62               SELECT
63                   *
```

```
64              FROM
65                  junkGtrOrgs
66              WHERE
67                  orgUuid IN (o1.orgUuid, o2.orgUuid)
68          )
69      ) q
70  ON CONFLICT (orgUuid1, orgUuid2) DO UPDATE SET (
71      simTrigramName,
72      simTrigramAddress
73  ) = (
74      excluded.simTrigramName,
75      excluded.simTrigramAddress
76  );
77  END $$; -- BEGIN
```

## 10.10  similarGtrPeople.sql

```
1   DO $$
2   DECLARE
3       rec record;
4   BEGIN
5       SET pg_trgm.similarity_threshold = 0.5;
6
7       INSERT INTO similarGtrPeople(
8           personUuid1,
9           personUuid2,
10          simTrigram)
11      SELECT
12          uuid1,
13          uuid2,
14          GREATEST((simFirst + simSur) / 2, (simFirst + simSur + simOther) / 3)
15      FROM (SELECT
16              p1.personUuid                        AS uuid1,
17              p2.personUuid                        AS uuid2,
18              similarity(p1.firstName,  p2.firstName)  AS simFirst,
19              similarity(p1.surname,    p2.surname)   AS simSur,
20              similarity(p1.otherNames, p2.otherNames) AS simOther
21          FROM gtrPeople p1
22          INNER JOIN gtrPeople p2
23              ON  p1.personUuid <  p2.personUuid
24              AND p1.firstName  %  p2.firstName
25              AND p1.surname    %  p2.surname
26      ) q
27      ON CONFLICT (personUuid1, personUuid2) DO UPDATE SET
28          simTrigram = excluded.simTrigram;
29  END $$; -- BEGIN
```

## 10.11  stats.sql

```
1   -- Perform and export statistical analyses on the dataset
2   -- Execute using:
3   -- $ psql -U <username> -d <database> -f setup.sql
4
5   -- Utility types
6   DO $$ BEGIN
7     CREATE TYPE PercentileInt AS (
8         percentile int,
```

```
 9          academic bigint,
10          medical bigint,
11          public bigint,
12          private bigint,
13          unknown bigint
14      );
15  EXCEPTION
16      WHEN duplicate_object THEN null;
17  END $$; -- BEGIN
18
19  CREATE TEMP TABLE IF NOT EXISTS orgTypeCounts AS
20  SELECT
21      type AS type,
22      COUNT(*) AS count
23  FROM
24      orgs
25  GROUP BY type;
26
27  \copy (SELECT * FROM orgTypeCounts) TO data/orgTypeCounts.csv (FORMAT CSV, HEADER)
28
29  -- Similarity distributions of organisations
30  CREATE OR REPLACE
31  FUNCTION OrgSimDist(
32      lowerBound int default 50,
33      step int default 2
34  )
35  RETURNS TABLE("lowerBound" int, "count" bigint) AS $$
36  DECLARE
37  BEGIN
38      RETURN QUERY
39      SELECT
40          *,
41          (SELECT
42              COUNT(*)
43          FROM similarGtrOrgs
44          WHERE simTrigramName >= (q.bound::float/100))
45      FROM (
46          SELECT * FROM GENERATE_SERIES(lowerBound, 100, step) AS bound
47      ) q;
48  END; $$ -- FUNCTION
49  LANGUAGE plpgsql;
50
51  \copy (SELECT * FROM OrgSimDist()) TO data/orgSimDist.csv (FORMAT CSV, HEADER)
52
53  -- Similarity distributions of organisations with the same postcode
54  CREATE OR REPLACE
55  FUNCTION OrgSimDistPostcode(
56      lowerBound int default 50,
57      step int default 2
58  )
59  RETURNS TABLE("lowerBound" int, "count" bigint) AS $$
60  DECLARE
61  BEGIN
62      RETURN QUERY
63      SELECT
64          *,
65          (SELECT
66              COUNT(*)
67          FROM similarGtrOrgs s
68              INNER JOIN gtrOrgs o1
```

```
 69                    ON s.orgUuid1 = o1.orgUuid
 70                INNER JOIN gtrOrgs o2
 71                    ON s.orgUuid2 = o2.orgUuid
 72            WHERE
 73                    simTrigramName >= (q.bound::float/100)
 74                AND o1.postcode IS NOT NULL
 75                AND STRIP(o1.postcode) = STRIP(o2.postcode))
 76        FROM (
 77            SELECT * FROM generate_series(lowerBound, 100, step) AS bound
 78        ) q;
 79    END; $$ -- FUNCTION
 80    LANGUAGE plpgsql;
 81
 82    \copy (SELECT * FROM OrgSimDistPostcode()) TO data/orgSimDistPostcode.csv (FORMAT CSV,
    ↪  HEADER)
 83
 84    CREATE OR REPLACE
 85    FUNCTION ProjectFundingPercentiles(
 86        lowerBound int DEFAULT 1,
 87        upperBound int DEFAULT 99
 88    )
 89    RETURNS TABLE("percentile" int, "value" numeric) AS $$
 90    DECLARE
 91        fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
 92    BEGIN
 93        RETURN QUERY
 94        SELECT
 95            ROUND(fraction * 100)::int,
 96            ROUND(result::numeric, 2)
 97        FROM (
 98            SELECT
 99                unnest(fractions) AS fraction,
100                unnest(
101                    (SELECT PERCENTILE_CONT(fractions) WITHIN GROUP (ORDER BY sum)
102                    FROM (
103                        SELECT
104                            SUM(COALESCE(offer::numeric, 0))
105                        FROM gtrProjectOrgs po
106                        WHERE offer > 0::money
107                        GROUP BY po.projectUuid
108                    ) q)
109                ) AS result
110        ) q;
111    END; $$ -- FUNCTION
112    LANGUAGE plpgsql;
113
114    \copy (SELECT * FROM ProjectFundingPercentiles()) TO data/projectFundingPercentiles.csv
    ↪  (FORMAT CSV, HEADER)
115
116    -- Percentiles for total number of projects organisations are involved in
117    CREATE OR REPLACE
118    FUNCTION OrgProjectPercentiles(
119        lowerBound int DEFAULT 1,
120        upperBound int DEFAULT 99
121    )
122    RETURNS TABLE("percentile" int, "value" double precision) AS $$
123    DECLARE
124        fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
125    BEGIN
126        RETURN QUERY
```

```sql
127        SELECT
128            ROUND(fraction * 100)::int,
129            result
130        FROM (
131            SELECT
132                unnest(fractions) AS fraction,
133                unnest(
134                    (SELECT PERCENTILE_CONT(fractions) WITHIN GROUP (ORDER BY count)
135                    FROM (
136                        SELECT
137                            COUNT(*)
138                        FROM gtrProjectOrgs po
139                        LEFT OUTER JOIN duplicateGtrOrgs d
140                            ON d.duplicateUuid = po.orgUuid
141                        GROUP BY COALESCE(d.orgUuid, po.orgUuid)
142                    ) q)
143                ) AS result
144        ) q;
145    END; $$ -- FUNCTION
146    LANGUAGE plpgsql;
147
148    \copy (SELECT * FROM OrgProjectPercentiles()) TO data/orgProjectPercentiles.csv (FORMAT
       ↪   CSV, HEADER)
149
150    -- Organisations' percentiles in total projects involved in
151    CREATE TEMP TABLE IF NOT EXISTS orgProjectPercentiles AS
152    SELECT
153        q.orgUuid,
154        q.duplicateUuid,
155        PERCENT_RANK() OVER (ORDER BY count) AS rank
156    FROM (
157        SELECT
158            COALESCE(d.orgUuid, po.orgUuid) AS orgUuid,
159            po.orgUuid AS duplicateUuid,
160            COUNT(*)
161        FROM
162            gtrProjectOrgs po
163            LEFT OUTER JOIN duplicateGtrOrgs d
164                ON d.duplicateUuid = po.orgUuid
165        GROUP BY COALESCE(d.orgUuid, po.orgUuid), po.orgUuid
166    ) q;
167
168    -- Percentile distributions for organisations, by number of projects involved
169    -- with
170    CREATE OR REPLACE
171    FUNCTION OrgProjectPercentileDist(
172        lowerBound int DEFAULT 0,
173        upperBound int DEFAULT 100
174    )
175    RETURNS SETOF PercentileInt AS $$
176    DECLARE
177        fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
178        fraction numeric;
179    BEGIN
180        DROP TABLE IF EXISTS res;
181        CREATE TEMP TABLE IF NOT EXISTS res OF PercentileInt;
182
183        FOREACH fraction in ARRAY fractions LOOP
184            INSERT INTO res
185            SELECT
```

```
186             COALESCE(ct.fraction, 0) AS fraction,
187             COALESCE(ct.academic, 0) AS academic,
188             COALESCE(ct.medical,  0) AS medical,
189             COALESCE(ct.public,   0) AS public,
190             COALESCE(ct.private,  0) AS private,
191             COALESCE(ct.unknown,  0) AS unknown
192         FROM crosstab('
193             SELECT
194                 ROUND('||fraction||' * 100)::int,
195                 type,
196                 COUNT(*)
197             FROM (
198                 SELECT
199                     COALESCE(type, ''Unknown'') AS type,
200                     COUNT(*)
201                 FROM
202                     orgProjectPercentiles opp
203                     INNER JOIN orgs o
204                         ON o.gtrOrgUuid = opp.orgUuid
205                 WHERE
206                     rank <= '||fraction||'
207                 GROUP BY orgUuid, type
208                 ORDER BY 1, 2
209             ) q
210             WHERE count > 0
211             GROUP BY type
212             ',
213             'VALUES
214                 (''Academic''::orgType),
215                 (''Medical''),
216                 (''Public''),
217                 (''Private''),
218                 (''Unknown'')'
219         ) AS ct (
220             fraction numeric,
221             academic bigint,
222             medical bigint,
223             public bigint,
224             private bigint,
225             unknown bigint
226         );
227     END LOOP;
228
229     RETURN QUERY SELECT * FROM res;
230 END; $$ -- FUNCTION
231 LANGUAGE plpgsql;
232
233 \copy (SELECT * FROM OrgProjectPercentileDist()) TO data/orgProjectPercentileDist.csv
    ↪  (FORMAT CSV, HEADER)
234
235 -- Percentage of organisations who collaborated with at least one other
236 -- organisation within project count percentile ranges
237 CREATE OR REPLACE
238 FUNCTION OrgProjectPercentileCollab(
239     lowerBound int DEFAULT 0,
240     upperBound int DEFAULT 100
241 )
242 RETURNS SETOF PercentileInt AS $$
243 DECLARE
244     fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
```

```
245       fraction numeric;
246   BEGIN
247       CREATE TEMP TABLE IF NOT EXISTS pairedOrgProjects AS
248       SELECT
249           COALESCE(d1.orgUuid, o1.orgUuid) AS orgUuid1,
250           type,
251           o2.orgUuid AS orgUuid2,
252           rank
253       FROM
254           -- First organisation
255           gtrOrgs o1
256           LEFT OUTER JOIN duplicateGtrOrgs d1
257               ON d1.duplicateUuid = o1.orgUuid
258           INNER JOIN orgs oo1
259               ON oo1.gtrOrgUuid = COALESCE(d1.orgUuid, o1.orgUuid)
260           -- Second organisation
261           INNER JOIN orgProjectPercentiles o2
262               ON o2.orgUuid <> o1.orgUuid
263           -- Projects both organisations are involved in
264           INNER JOIN gtrProjectOrgs po1
265               ON po1.orgUuid IN (d1.orgUuid, o1.orgUuid)
266           INNER JOIN gtrProjectOrgs po2
267               ON  po2.projectUuid = po1.projectUuid
268               AND po2.orgUuid = o2.orgUuid;
269
270       DROP TABLE IF EXISTS res;
271       CREATE TEMP TABLE IF NOT EXISTS res OF PercentileInt;
272
273       FOREACH fraction in ARRAY fractions LOOP
274           INSERT INTO res
275           SELECT
276               COALESCE(ct.fraction, 0) AS fraction,
277               COALESCE(ct.academic, 0) AS academic,
278               COALESCE(ct.medical,  0) AS medical,
279               COALESCE(ct.public,   0) AS public,
280               COALESCE(ct.private,  0) AS private,
281               COALESCE(ct.unknown,  0) AS unknown
282           FROM crosstab('
283               SELECT
284                   ROUND('||fraction||' * 100)::int,
285                   type,
286                   COUNT(*)
287               FROM (
288                   SELECT
289                       COALESCE(type, ''Unknown'') AS type,
290                       COUNT(*)
291                   FROM
292                       pairedOrgProjects
293                   WHERE
294                       rank <= '||fraction||'
295                   GROUP BY orgUuid1, type
296                   ORDER BY 1, 2
297               ) q
298               WHERE count > 0
299               GROUP BY type
300               ',
301               'VALUES
302                   (''Academic''::orgType),
303                   (''Medical''),
304                   (''Public''),
```

```
                    (''Private''),
                    (''Unknown'')'
        ) AS ct (
                fraction numeric,
                academic bigint,
                medical bigint,
                public bigint,
                private bigint,
                unknown bigint
        );
    END LOOP;

    RETURN QUERY SELECT * FROM res;
END; $$ -- FUNCTION
LANGUAGE plpgsql;

\copy (SELECT * FROM OrgProjectPercentileCollab()) TO data/orgProjectPercentileCollab.csv
↪  (FORMAT CSV, HEADER)

-- Percentiles for funding received by organisations
CREATE OR REPLACE
FUNCTION OrgFundingPercentiles(
    lowerBound int DEFAULT 1,
    upperBound int DEFAULT 99
)
RETURNS TABLE("percentile" int, "value" numeric) AS $$
DECLARE
    fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
BEGIN
    RETURN QUERY
    SELECT
        ROUND(fraction * 100)::int,
        ROUND(result::numeric, 2)
    FROM (
        SELECT
            unnest(fractions) AS fraction,
            unnest(
                (SELECT PERCENTILE_CONT(fractions) WITHIN GROUP (ORDER BY sum)
                FROM (
                    SELECT
                        SUM(COALESCE(offer::numeric, 0))
                    FROM gtrProjectOrgs po
                    LEFT OUTER JOIN duplicateGtrOrgs d
                        ON d.duplicateUuid = po.orgUuid
                    WHERE offer > 0::money
                    GROUP BY COALESCE(d.orgUuid, po.orgUuid)
                ) q)
            ) AS result
    ) q;
END; $$ -- FUNCTION
LANGUAGE plpgsql;

\copy (SELECT * FROM OrgFundingPercentiles()) TO data/orgFundingPercentiles.csv (FORMAT
↪  CSV, HEADER)

-- Organisations' percentiles in total funding received
CREATE TEMP TABLE IF NOT EXISTS orgFundingPercentiles AS
SELECT
    q.orgUuid,
    q.duplicateUuid,
```

```
363        PERCENT_RANK() OVER (ORDER BY sum) AS rank
364    FROM (
365        SELECT
366            COALESCE(d.orgUuid, po.orgUuid) AS orgUuid,
367            po.orgUuid AS duplicateUuid,
368            SUM(COALESCE(offer::numeric, 0))
369        FROM
370            gtrProjectOrgs po
371            LEFT OUTER JOIN duplicateGtrOrgs d
372                ON d.duplicateUuid = po.orgUuid
373        WHERE offer > 0::money
374        GROUP BY COALESCE(d.orgUuid, po.orgUuid), po.orgUuid
375    ) q;
376
377    -- Number of organisations who collaborated with at least one other
378    -- organisation within funding percentile ranges
379    CREATE OR REPLACE
380    FUNCTION OrgFundingPercentileCollab(
381        lowerBound int DEFAULT 0,
382        upperBound int DEFAULT 100
383    )
384    RETURNS SETOF PercentileInt AS $$
385    DECLARE
386        fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
387        fraction numeric;
388    BEGIN
389        CREATE TEMP TABLE IF NOT EXISTS pairedOrgFunding AS
390        SELECT
391            COALESCE(d1.orgUuid, o1.orgUuid) AS orgUuid1,
392            oo1.type AS type,
393            o2.orgUuid AS orgUuid2,
394            rank
395        FROM
396            -- First organisation
397            gtrOrgs o1
398            LEFT OUTER JOIN duplicateGtrOrgs d1
399                ON d1.duplicateUuid = o1.orgUuid
400            LEFT OUTER JOIN orgs oo1
401                ON oo1.gtrOrgUuid IN (d1.orgUuid, o1.orgUuid)
402            -- Second organisation
403            INNER JOIN orgFundingPercentiles o2
404                ON o2.orgUuid <> o1.orgUuid
405            -- Projects both organisations are involved in
406            INNER JOIN gtrProjectOrgs po1
407                ON po1.orgUuid IN (d1.orgUuid, o1.orgUuid)
408            INNER JOIN gtrProjectOrgs po2
409                ON  po2.projectUuid = po1.projectUuid
410                AND po2.orgUuid = o2.orgUuid;
411
412        DROP TABLE IF EXISTS res;
413        CREATE TEMP TABLE IF NOT EXISTS res OF PercentileInt;
414
415        FOREACH fraction in ARRAY fractions LOOP
416            INSERT INTO res
417            SELECT
418                COALESCE(ct.fraction, 0) AS fraction,
419                COALESCE(ct.academic, 0) AS academic,
420                COALESCE(ct.medical,  0) AS medical,
421                COALESCE(ct.public,   0) AS public,
422                COALESCE(ct.private,  0) AS private,
```

```
423              COALESCE(ct.unknown,  0) AS unknown
424          FROM crosstab('
425              SELECT
426                  ROUND('||fraction||' * 100)::int,
427                  type,
428                  COUNT(*)
429              FROM (
430                  SELECT
431                      COALESCE(type, ''Unknown'') AS type,
432                      COUNT(*)
433                  FROM
434                      pairedOrgFunding
435                  WHERE
436                      rank <= '||fraction||'
437                  GROUP BY orgUuid1, type
438                  ORDER BY 1, 2
439              ) q
440              WHERE count > 0
441              GROUP BY type
442              ',
443              'VALUES
444                  (''Academic''::orgType),
445                  (''Medical''),
446                  (''Public''),
447                  (''Private''),
448                  (''Unknown'')'
449          ) AS ct (
450              fraction numeric,
451              academic bigint,
452              medical bigint,
453              public bigint,
454              private bigint,
455              unknown bigint
456          );
457      END LOOP;
458
459      RETURN QUERY SELECT * FROM res;
460  END; $$ -- FUNCTION
461  LANGUAGE plpgsql;
462
463  \copy (SELECT * FROM OrgFundingPercentileCollab()) TO data/orgFundingPercentileCollab.csv
     ↪  (FORMAT CSV, HEADER)
464
465  -- Percentile distributions for organisations, by total amount of funding
466  -- received
467  CREATE OR REPLACE
468  FUNCTION OrgFundingPercentileProjects(
469      lowerBound int DEFAULT 0,
470      upperBound int DEFAULT 100
471  )
472  RETURNS TABLE("percentile" int, "value" bigint) AS $$
473  DECLARE
474      fractions numeric[] := PERCENTILE_FRACTIONS(lowerBound, upperBound);
475  BEGIN
476      CREATE TEMP TABLE IF NOT EXISTS projectPairedRanks AS
477      SELECT
478          projectUuid,
479          rank
480      FROM
481          orgFundingPercentiles o
```

```
482          INNER JOIN gtrProjectOrgs po
483              ON po.orgUuid IN (o.orgUuid, o.duplicateUuid);
484
485      RETURN QUERY
486      SELECT
487          ROUND(fractions.upperBound * 100)::int,
488          (SELECT
489              COUNT(*)
490          FROM (
491              SELECT
492                  COUNT(*)
493              FROM
494                  projectPairedRanks
495              WHERE
496                  rank <= fractions.upperBound
497              GROUP BY projectUuid
498          ) q
499          WHERE count > 0
500          ) AS result
501      FROM unnest(fractions) AS fractions(upperBound);
502 END; $$ -- FUNCTION
503 LANGUAGE plpgsql;
504
505 \copy (SELECT * FROM OrgFundingPercentileProjects()) TO
  ↪  data/orgFundingPercentileProjects.csv (FORMAT CSV, HEADER)
```

## 10.12   mergeGtrOrgs.sql

```
1  BEGIN;
2
3  INSERT INTO junkGtrOrgs(
4      orgUuid
5  )
6  SELECT
7      orgUuid
8  FROM
9      gtrOrgs
10 WHERE name ~* '^unknown*|^unlisted*'
11 ON CONFLICT DO NOTHING;
12
13 INSERT INTO duplicateGtrOrgs(
14     orgUuid,
15     duplicateUuid
16 )
17 SELECT
18     LEAST(o1.orgUuid, o2.orgUuid) AS orgUuid,
19     GREATEST(o1.orgUuid, o2.orgUuid) AS duplicateUuid
20 FROM
21     similarGtrOrgs s
22     INNER JOIN gtrOrgs o1
23         ON s.orgUuid1 = o1.orgUuid
24     INNER JOIN gtrOrgs o2
25         ON s.orgUuid2 = o2.orgUuid
26 WHERE
27         simTrigramName >= 0.9
28     AND COALESCE(manualResult, TRUE)
29     OR  (simTrigramName >= 0.5
30         AND o1.postCode IS NOT NULL
```

```
31          AND strip(o1.postcode) = strip(o2.postcode))
32      -- Check that the first rec isn't already marked as a duplicate
33      AND NOT EXISTS(SELECT *
34          FROM duplicateGtrOrgs
35          WHERE duplicateUuid IN (o1.orgUuid, o2.orgUuid))
36      AND NOT EXISTS(SELECT *
37          FROM junkGtrOrgs
38          WHERE orgUuid IN (o1.orgUuid, o2.orgUuid))
39  GROUP BY
40      LEAST(   o1.orgUuid, o2.orgUuid),
41      GREATEST(o1.orgUuid, o2.orgUuid)
42  ON CONFLICT (duplicateUuid) DO NOTHING;
43
44  INSERT INTO orgs(
45      gtrOrgUuid,
46      name,
47      address1,
48      address2,
49      address3,
50      address4,
51      address5,
52      postCode,
53      city,
54      region,
55      country
56  )
57  SELECT
58  DISTINCT ON (orgUuid)
59      LEAST(o.orgUuid, d.orgUuid, d.duplicateUuid) AS orgUuid,
60      MERGE(name)                                  AS name,
61      MERGE(address1)                              AS address1,
62      MERGE(address2)                              AS address2,
63      MERGE(address3)                              AS address3,
64      MERGE(address4)                              AS address4,
65      MERGE(address5)                              AS address5,
66      MERGE(postCode)                              AS postCode,
67      MERGE(city)                                  AS city,
68      COALESCE(MERGE(NULLIF(region, 'Unknown')),
69          'Unknown')                               AS region,
70      MERGE(country)                               AS country
71  FROM
72      gtrOrgs o
73      LEFT OUTER JOIN duplicateGtrOrgs d
74          ON o.orgUuid IN (d.orgUuid, d.duplicateUuid)
75  WHERE
76      NOT EXISTS(SELECT *
77          FROM junkGtrOrgs
78          WHERE orgUuid IN (o.orgUuid, d.orgUuid, d.duplicateUuid))
79  GROUP BY
80      LEAST(o.orgUuid, d.orgUuid, d.duplicateUuid)
81  ON CONFLICT (gtrOrgUuid) DO UPDATE SET
82      name     = excluded.name,
83      address1 = excluded.address1,
84      address2 = excluded.address2,
85      address3 = excluded.address3,
86      address4 = excluded.address4,
87      address5 = excluded.address5,
88      postCode = excluded.postCode,
89      city     = excluded.city,
90      region   = excluded.region,
```

```
91    country  = excluded.country;
92
93  COMMIT;
```

## 10.13   mergeGtrPeople.sql

```
1   BEGIN;
2
3   INSERT INTO duplicateGtrPeople(
4       personUuid,
5       duplicateUuid
6   )
7   SELECT
8       LEAST(s.personUuid1, s.personUuid2) AS personUuid,
9       GREATEST(s.personUuid1, s.personUuid2) AS duplicateUuid
10  FROM
11      similarGtrPeople s
12      INNER JOIN gtrOrgPeople op1
13          ON s.personUuid1 = op1.personUuid
14      INNER JOIN gtrOrgPeople op2
15          ON  s.personUuid2 = op2.personUuid
16      LEFT OUTER JOIN duplicateGtrOrgs d1
17          ON d1.duplicateUuid = op1.orgUuid
18      LEFT OUTER JOIN duplicateGtrOrgs d2
19          ON d2.duplicateUuid = op2.orgUuid
20  WHERE
21          COALESCE(d1.orgUuid, op1.orgUuid) = COALESCE(d2.orgUuid, op2.orgUuid)
22      AND simTrigram >= 0.9
23      AND COALESCE(s.manualResult, TRUE)
24      -- Check that the first rec isn't already marked as a duplicate
25      AND NOT EXISTS(SELECT *
26          FROM duplicateGtrPeople
27          WHERE duplicateUuid IN (s.personUuid1, s.personUuid2))
28  GROUP BY
29      LEAST(   s.personUuid1, s.personUuid2),
30      GREATEST(s.personUuid1, s.personUuid2)
31  ON CONFLICT (duplicateUuid) DO NOTHING;
32
33  INSERT INTO people(
34      gtrPersonUuid,
35      firstName,
36      surname,
37      otherNames
38  )
39  SELECT
40      LEAST(p1.personUuid, d.personUuid, d.duplicateUuid) AS personUuid,
41      MERGE(firstName)                                     AS firstName,
42      MERGE(surname)                                       AS surname,
43      MERGE(otherNames)                                    AS otherNames
44  FROM
45      gtrPeople p1
46      LEFT OUTER JOIN duplicateGtrPeople d
47          ON p1.personUuid IN (d.personUuid, d.duplicateUuid)
48  GROUP BY
49      LEAST(p1.personUuid, d.personUuid, d.duplicateUuid)
50  ON CONFLICT (gtrPersonUuid) DO UPDATE SET
51      firstName  = excluded.firstName,
52      surname    = excluded.surname,
53      otherNames = excluded.otherNames;
```

```
54
55  COMMIT;
```

## 10.14   eastMidlandsGraphGephi.sql

```
1   DROP TABLE IF EXISTS nodes;
2   CREATE TEMP TABLE IF NOT EXISTS nodes AS
3   SELECT
4       gtrOrgUuid AS id,
5       name AS label,
6       *
7   FROM
8       orgs o
9   WHERE
10      region = 'East Midlands';
11
12  \copy (SELECT * FROM nodes) TO data/eastMidlandsGraphGephiNodes.csv (FORMAT CSV, HEADER)
13
14  DROP TABLE IF EXISTS edges;
15  CREATE TEMP TABLE IF NOT EXISTS edges AS
16  SELECT
17      po1.projectUuid AS projectUuid,
18      o1.gtrOrgUuid AS source,
19      GREATEST(po1.offer::numeric::int, po1.cost::numeric::int, 1) AS weight,
20
21      o2.gtrOrgUuid AS target,
22      p.title AS label,
23      p.startDate AS startDate,
24      p.endDate AS endDate
25  FROM
26      -- First organisation
27      nodes o1
28      LEFT OUTER JOIN duplicateGtrOrgs d1
29          ON d1.orgUuid = o1.gtrOrgUuid
30      INNER JOIN gtrProjectOrgs po1
31          ON po1.orgUuid IN (o1.gtrOrgUuid, d1.duplicateUuid)
32
33      -- Second organisation
34      INNER JOIN gtrProjectOrgs po2
35          ON po2.orgUuid NOT IN (o1.gtrOrgUuid, d1.duplicateUuid)
36          AND po2.projectUuid = po1.projectUuid
37      LEFT OUTER JOIN duplicateGtrOrgs d2
38          ON po2.orgUuid IN (d2.orgUuid, d2.duplicateUuid)
39      INNER JOIN nodes o2
40          ON o2.gtrOrgUuid IN (d2.orgUuid, po2.orgUuid)
41
42      -- Shared project
43      INNER JOIN gtrProjects p
44          ON p.projectUuid = po1.projectUuid;
45
46  \copy (SELECT * FROM edges) TO data/eastMidlandsGraphGephiEdges.csv (FORMAT CSV, HEADER)
```

## 10.15   eastMidlandsGraphNodeXL.sql

```
1   DROP TABLE IF EXISTS edges;
2   CREATE TEMP TABLE edges AS
3   SELECT
```

```
4        o1.gtrOrgUuid AS source,
5        o1.name AS sourceName,
6        o1.region AS sourceRegion,
7        o1.city AS sourceCity,
8        o1.country AS sourceCountry,
9        o1.type AS sourceType,
10       po1.role AS sourceRole,
11       po1.offer::numeric::int AS offer,
12       po1.cost::numeric::int AS cost,
13
14       o2.gtrOrgUuid AS target,
15       o2.name AS targetName,
16       o2.region AS targetRegion,
17       o2.city AS targetCity,
18       o2.country AS targetCountry,
19       o2.type AS targetType,
20
21       p.projectUuid AS projectUuid,
22       p.startDate AS startDate,
23       p.endDate AS endDate
24   FROM
25       -- First organisation
26       orgs o1
27       LEFT OUTER JOIN duplicateGtrOrgs d1
28           ON d1.orgUuid = o1.gtrOrgUuid
29       INNER JOIN gtrProjectOrgs po1
30           ON po1.orgUuid IN (o1.gtrOrgUuid, d1.duplicateUuid)
31
32       -- Second organisation
33       INNER JOIN gtrProjectOrgs po2
34           ON po2.orgUuid NOT IN (o1.gtrOrgUuid, d1.duplicateUuid)
35           AND po2.projectUuid = po1.projectUuid
36       LEFT OUTER JOIN duplicateGtrOrgs d2
37           ON po2.orgUuid IN (d2.orgUuid, d2.duplicateUuid)
38       INNER JOIN orgs o2
39           ON o2.gtrOrgUuid IN (d2.orgUuid, po2.orgUuid)
40
41       -- Shared project
42       INNER JOIN gtrProjects p
43           ON p.projectUuid = po1.projectUuid
44   WHERE
45       p.startDate >= '1/1/2010'
46       AND (p.endDate IS NULL
47       OR p.endDate > '1/1/2015')
48       AND o1.region = 'East Midlands'
49       AND o2.region = 'East Midlands';
50
51   \copy (SELECT * FROM edges) TO data/eastMidlandsGraphNodeXL.csv (FORMAT CSV, HEADER)
```

## 10.16  Database schema manual

An technical description of the contents of the database schema, visualised in figure 16.

### 10.16.1  Organisations

**10.16.1.1  gtrOrgs**  Organisations stored in the Gateway to Research (GtR) database. 47822 records

| Field | Type | Description |
| --- | --- | --- |
| orgUuid | UUID | Unique identifier |
| name | text | Name of organisation |
| address1 | text | Address line 1 |
| address2 | text | Address line 2 |
| address3 | text | Address line 3 |
| address4 | text | Address line 4 |
| address5 | text | Address line 5 |
| postCode | text | Post code |
| city | text | City organisation is based in |
| region | gtrRegion | Region organisation is based in |
| country | text | Country organisation is based in |
| recorded | date | Date this record was created (in the GtR database) |

**10.16.1.2  gtrRegion: The regions an organisation can be based in:**

- Channel Islands/Isle of Man
- East Midlands
- East of England
- London
- North East
- Northern Ireland
- North West
- Scotland
- South East
- South West
- Wales
- West Midlands
- Yorkshire and The Humber
- Outside UK
- Unknown

### 10.16.2  Projects

**10.16.2.1  gtrProjects**  Projects stored in the Gateway to Research (GtR) database. 97277 records

| Field | Type | Description |
| --- | --- | --- |
| projectUuid | uuid | Unique identifier |
| title | text | Title of project |
| status | gtrProjectStatus | Current status of project: 'Active' or 'Closed' |
| category | gtrGrantCategory | Category of project |

| Field | Type | Description |
|---|---|---|
| leadFunder | gtrFunder | The main UKRI (sub)organisation providing the majority of the funding for this project |
| abstract | text | Abstract text summarising the project's background and goals |
| techAbstract | text | Technical summary of research being undertaken |
| potentialImpact | text | Planned impact that will result from the project |
| startDate | date | Date of the beginning of project funding |
| endDate | date | Date of the end of project funding |
| recorded | date | Date this record was created (in the GtR database) |

**10.16.2.2 gtrGrantCategory** Categories that projects can be contained within

- BIS-Funded Programmes
- Centres
- Collaborative R&D
- CR&D Bilateral
- EU-Funded
- European Enterprise Network
- Fast Track
- Feasibility Studies
- Fellowship
- GRD Development of Prototype
- GRD Proof of Concept
- GRD Proof of Market
- Intramural
- Knowledge Transfer Network
- Knowledge Transfer Partnership
- Large Project
- Launchpad
- Legacy Department of Trade & Industry
- Legacy RDA Collaborative R&D
- Legacy RDA Grant for R&D
- Missions
- Other Grant
- Procurement
- Research Grant
- Small Business Research Initiative
- SME Support
- Special Interest Group
- Studentship
- Study
- Third Party Grant
- Training Grant
- Unknown
- Vouchers

**10.16.2.3 gtrFunder** UKRI, or sub-organisations of UKRI, that provide funding

- UKRI
- AHRC

- BBSRC
- EPSRC
- ESRC
- Innovate UK
- MRC
- NC3Rs
- NERC
- STFC

**10.16.2.4  gtrSubjects**   Research subjects associated with projects.
83 records

| Field | Type | Description |
|-------|------|-------------|
| subjectUuid | uuid | Unique identifier |
| name | text | Name of subject, as specified by users |

**10.16.2.5  gtrProjectSubjects**   Subjects that particular projects are associated with.
Percentages are specified by users, each project can be associated with multiple research subjects.
78270 records

| Field | Type | Description |
|-------|------|-------------|
| subjectUuid | uuid | Unique identifier of subject in gtrSubjects |
| projectUuid | uuid | Unique identifier of project in gtrProjects |
| percent | numeric | Percentage of project that is related to this subject |

**10.16.2.6  gtrTopics**   Research topics associated with projects.
610 records

| Field | Type | Description |
|-------|------|-------------|
| topicUuid | uuid | Unique identifier |
| name | text | Name of topic, as specified by users |

**10.16.2.7  gtrProjectTopics**   Topics that particular projects are associated with.
Percentages are specified by users, each project can be associated with multiple research Topics.
150016 records

| Field | Type | Description |
|-------|------|-------------|
| topicUuid | uuid | Unique identifier of topic in gtrTopics |
| projectUuid | uuid | Unique identifier of project in gtrProjects |
| percent | numeric | Percentage of project that is related to this topic |

### 10.16.3  Project outcomes

The outcomess of a project as logged by users

**10.16.3.1  gtrDisseminations**  Disseminations published as a result of a project. 58 records.

| Field | Type | Description |
|---|---|---|
| disUuid | uuid | Unique identifier |
| title | text | Title |
| description | text | Description of the contents |
| form | text | The form of publishing, e.g. presentations, conferences, newsletter |
| primaryAudience | text | Audience this dissemination is targeted at, e.g. policymakers, students |
| yearsOfDissemination | text | Comma-separated list of years it was active |
| results | text | Description of the results of publishing (unused) |
| impact | text | Descriptions of the outcomes of publishing |
| typeOfPresentation | text | Type of presentation (mostly unused) |
| geographicReach | text | Local, Regional, National, or International |
| partOfOfficialScheme | boolean | Whether engagement activity is part of an official scheme |
| supportingUrl | text | URLs for accessing supporting material |

**10.16.3.2  gtrCollaborations**  Further collaborations that formed as a result of a project. 8 records

| Field | Type | Description |
|---|---|---|
| collabUuid | uuid | Unique identifier |
| description | text | Description of collaboration |
| parentOrganisation | text | Parent organisation of collaboration partner |
| childOrganisation | text | Child organisation of collaboration partner |
| principalInvestigatorContribution | text | Benefits provided to new collaborator by principal investigator |
| partnerContribution | text | Benefits provided to principal investigator by new collaborator |
| startDate | date | Date collaboration began |
| endDate | date | Date collaboration ended |
| sector | gtrSector | Sector of collaborating organisation |
| country | text | Country collaborator is based in |
| impact | text | Resulting impact of the collaboration |
| supportingUrl | text | URLs for accessing supporting material |

**10.16.3.3  gtrKeyFindings**  Key findings from a project. 1 record

| Field | Type | Description |
|---|---|---|
| keyFindingUuid | uuid | Unique identifier |
| description | text | Description of findings |
| nonAcademicUses | text | Non-academic applications of findings |
| exploitationPathways | text | Possible further research and development |
| sectors | text | Sectors that findings could be applied to |
| supportingUrl | text | URLs for accessing supporting material |

**10.16.3.4  gtrFurtherFundings**  Further project funding after the initial funding period expired.
9 records

| Field | Type | Description |
|---|---|---|
| furtherFundingUuid | uuid | Unique identifier |
| title | text | Tile of grant/funding received (unused) |
| description | text | Description of further funding received |
| narrative | text | Reasoning for further funding (unused) |
| amount | money | Amount of funding received |
| organisation | text | Organisation providing funding |
| department | text | Department of organisation providing funding |
| fundingId | text | External reference identifier for the funding received |
| startDate | date | Start of funding period |
| endDate | date | End of funding period |
| sector | gtrSector | Sector of organisation providing funding |
| country | text | Country of organisation providing funding |

**10.16.3.5  gtrImpactSummaries**  Further project funding after the initial funding period expired.
5 records

| Field | Type | Description |
|---|---|---|
| impactSummaryUuid | uuid | Unique identifier |
| title | text | Unused |
| description | text | Description of impact |
| impactTypes | text | Unused |
| summary | text | Unused |
| beneficiaries | text | Unused |
| contributionMethod | text | Unused |
| sector | text | Comma-delimited list of areas affected, e.g. 'Healthcare', 'Manufacturing' |
| firstYearOfImpact | int | Year that the impact took effect |

**10.16.3.6  gtrPolicyInfluences**  Any influence on an organisation/government's policy as a result of a project.
5 records

| Field | Type | Description |
|---|---|---|
| policyInfluenceUuid | uuid | Unique identifier |
| influence | text | Description of influence, e.g. committee influenced |
| type | text | Form of incluence, e.g. membership of a guideline committee |
| guidelineTitle | text | Unused |
| impact | text | Unused |
| methods | text | Unused |
| areas | text | Unused |
| geographicReach | text | Region influenced, e.g. 'National', 'Europe' |
| supportingUrl | text | URLs for accessing supporting material |

**10.16.3.7 gtrResearchMaterials** Research material published as a result of a project.
0 records

| Field | Type | Description |
| --- | --- | --- |
| researchMatUuid | uuid | Unique identifier |
| title | text | Title |
| description | text | Description |
| type | text | Type |
| impact | text | Impact from releasing |
| softwareDeveloped | boolean | Whether software was developed |
| softwareOpenSourced | boolean | Whether software was open-sourced |
| providedToOthers | boolean | Whether material was provided to other entities |
| yearFirstProvided | int | First year research material was provided |
| supportingUrl | text | URLs for accessing supporting material |

### 10.16.4 Links

**10.16.4.1 gtrOrgPeople** Organisations that people in `gtrPeople` are linked to.
82015 records

| Field | Type | Description |
| --- | --- | --- |
| personUuid | uuid | UUID of person in `gtrPeople` |
| orgUuid | uuid | UUID of organisation in `gtrOrgs` |

**10.16.4.2 gtrProjectPeople** Individuals working on a particular project.
201596 records

| Field | Type | Description |
| --- | --- | --- |
| projectUuid | uuid | UUID of project in `gtrProjects` |
| personUuid | uuid | UUID of person in `gtrPersonRole` |
| role | gtrPersonRole | Role of individual in the project |

**10.16.4.3 gtrPersonRole** Possible roles an individual can have in a project:

- `Principal Investigator`
- `Co-Investigator`
- `Project Manager`
- `Fellow`
- `Training Grant Holder`
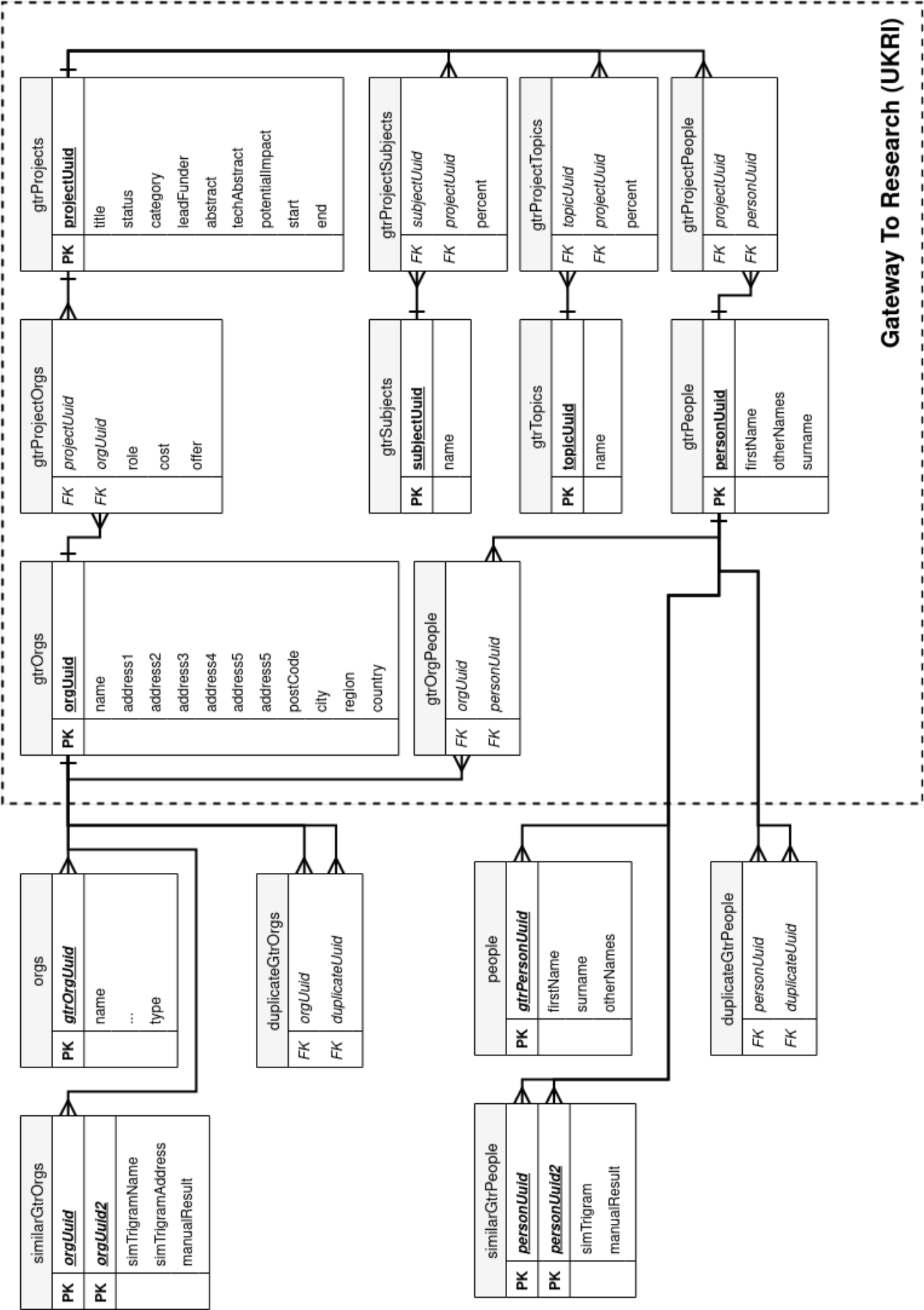- `Primary Supervisor`
- `Researcher Co-Investigator`
- `Researcher`

**10.16.4.4 gtrProjectOrgs** Organisations working on a particular project.
283527 records

| Field | Type | Description |
| --- | --- | --- |
| projectUuid | uuid | UUID of project in `gtrProjects` |

| Field | Type | Description |
| --- | --- | --- |
| orgUuid | uuid | UUID of project in `gtrOrgs` |
| role | gtrOrgRole | Role of organisation in the project |
| startDate | date | Start of collaboration period |
| endDate | date | End of collaboration period |
| cost | money | Total funds spent by this organisation on the project |
| offer | money | Total funding offered by the funder to this organisations |

**10.16.4.5  gtrRelatedProjects**  Projects that are related to one another in some way.
25392 records

| Field | Type | Description |
| --- | --- | --- |
| projectUuid1 | uuid | UUID of project in `gtrProjects` |
| projectUuid2 | uuid | UUID of project in `gtrProjects` |
| relation | gtrProjectRelation | Relationship between the projects |
| startDate | date | Start of related (2nd) project |
| endDate | date | End of related (2nd) project |

## 10.17  Database schema



Figure 16: Entity Relationship Diagram describing the database schema
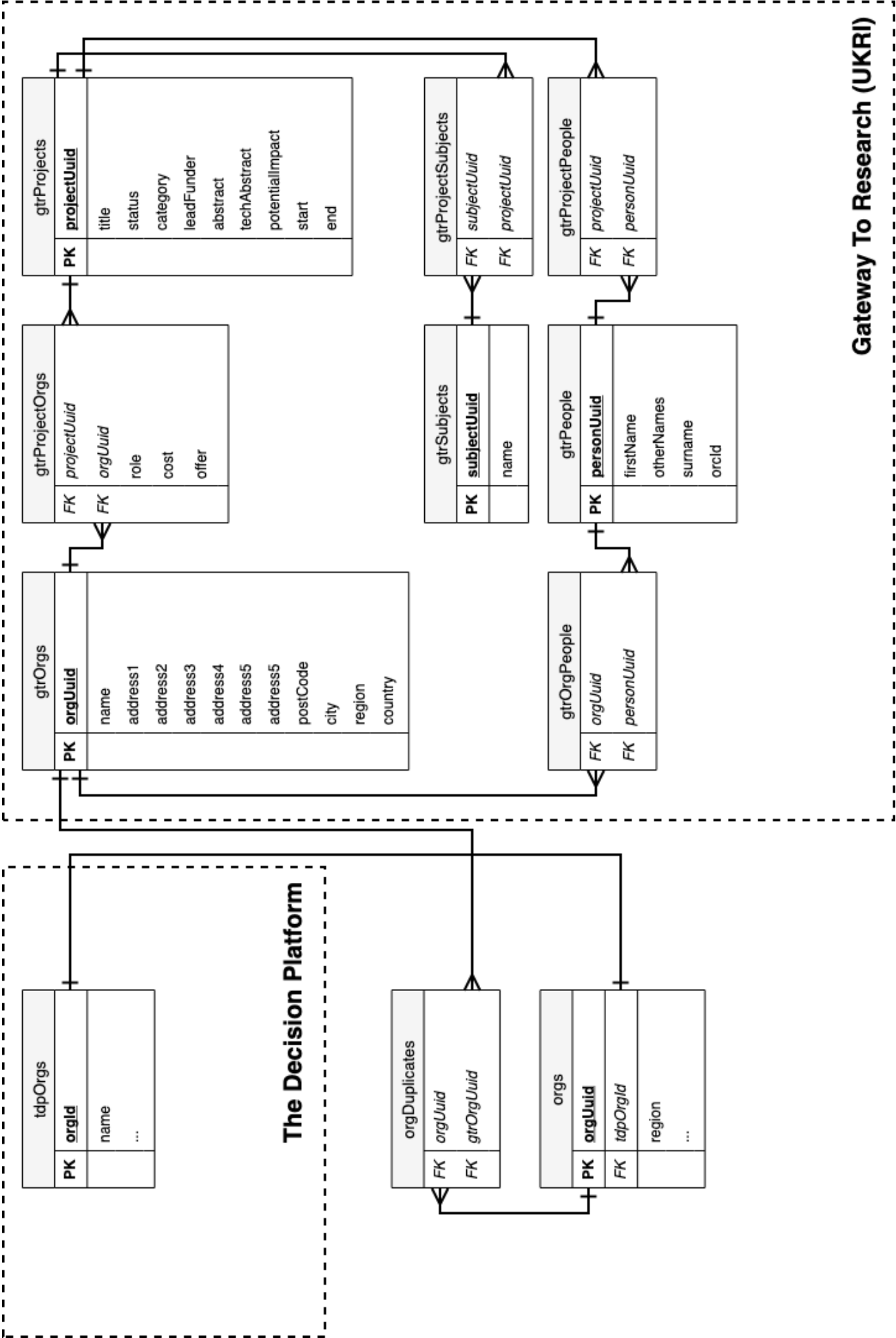
## 10.18 Database schema (Initial design)



Figure 17: Entity Relationship Diagram describing the planned database schema, before changes were made throughout the implementation of the project
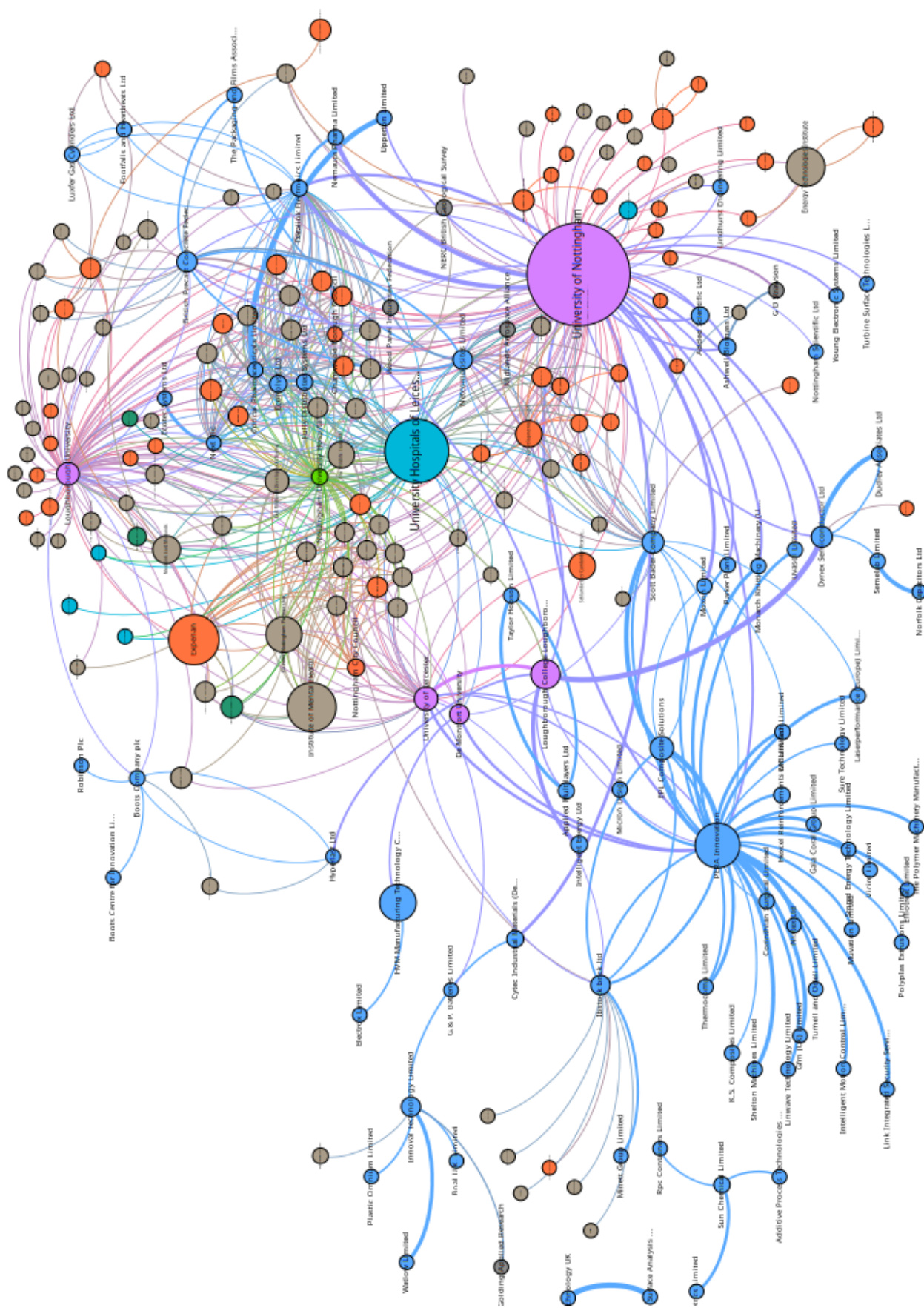
## 10.19 East Midlands network (2002-2008)



Figure 18: Funding network for organisations in the East Midlands within the years 2002 to 2008

## 10.20 East Midlands network (2005-2010)



Figure 19: Funding network for organisations in the East Midlands within the years 2005 to 2010