# Large Scale Semantic Annotation, Indexing, and Search at The National Archives

**Diana Maynard**

**Mark Greenwood**

University of Sheffield, UK

# Burning questions you may have....

- In the last 3 years, which female politicians have talked about hospitals spending more than 1 million pounds?

- Which cabinet ministers currently in power were born in Sheffield?

- Which schools spent between 1 and 10 million pounds in the period July 2010-July 2011?

- Which speeches did Tony Blair make about foot and mouth disease while Prime Minister, and when did he make them?`

# Government Web Archive Project

---

- Project aims to help open up TNA's records of .gov.uk websites (going back to 1997 and comprising some 340 million pages)

- Government funding has been allocated to publishing more and more material on data.gov.uk in open and accessible forms

- But it's still pretty hard to find the information you're looking for

- Aim is basically to improve access to this enormous volume of data - both for the general public and for specialist researchers

# What are the National Archives?

- UK government's official archive, containing over 1,000 years of history and making this information publicly available

- Work with 250 government and public sector bodies, helping them to manage and use information more effectively.

- Over 11 million historical government and public records - one of the largest in the world.

- In general, government records that have been selected for permanent preservation are sent to The National Archives when they are 30 years old, but many are released earlier under the Freedom of Information Act

# What kinds of information do they hold?

- **Government, diplomacy and the armed forces**: e.g. documents from all government departments

- **Court records**

- Approx 6 million **historical maps**

- **DocumentsOnline**: digitised public records, e.g. famous historical wills, selected records from MI5 and MI6, a range of UFO-related files from the MoD, WW1 and WW2 selected records

- Census records, alien arrivals, birth, marriage and death certificates etc.

# Basic National Archives search tool



european archive

About | Contact
Terms, Privacy & Copyright

## Advanced search (beta)

Find **all** these words:

**None** of these words:

This **exact** phrase: Harriet Harman

Results in this format: all formats

Search ▸

Quick search | Search help

**Search within categories**

☐ Business, industry, economics and finance
☑ Central government
☐ Culture and leisure
☐ Environment
☐ Health, well-being and care
☐ Home affairs, public order, justice and rights
☐ Honours, awards, appointments and titles

☐ International affairs and defence
☐ People, community and housing
☐ Public inquiries and Royal Commissions
☐ Regional Government
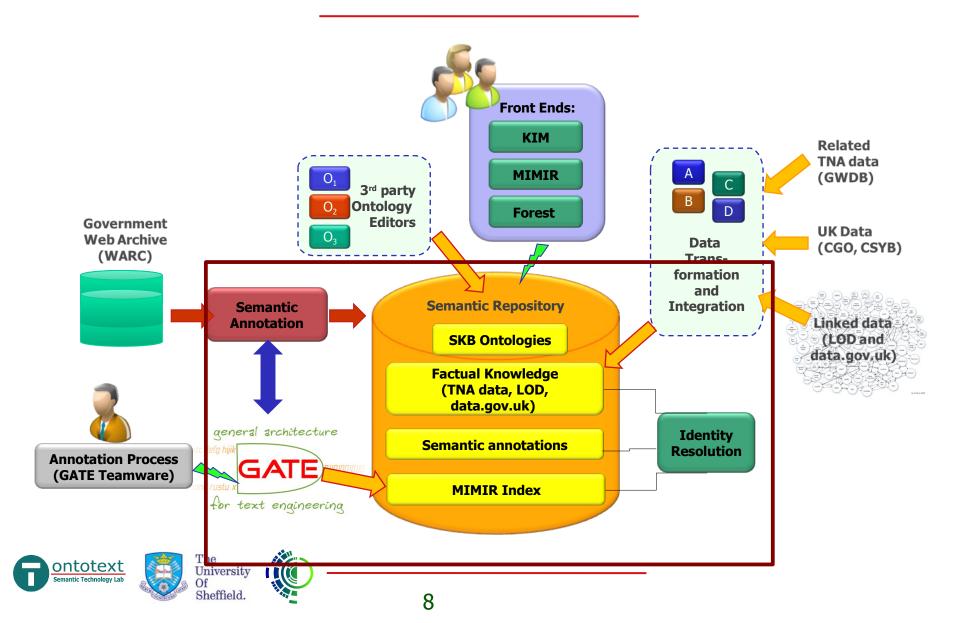☐ Transport, communication and technology
☐ Work, education and skills

Done

# We want to do better than this!

- Enable semantic-based search for categories of things, e.g. all Cabinet Ministers, all cities in the UK

- Search results include morphological variants of words and synonyms

- Search for specific phrases with some unknowns, e.g. a Person and a monetary amount in the same sentence

- Search for ranges, e.g. all monetary amounts greater than a million pounds

- Restrict search to certain date periods, domains etc.

- and so on...

# System Architecture: GATE/MIMIR



**Front Ends:**
- KIM
- MIMIR
- Forest

$O_1$ $O_2$ $O_3$ **3rd party Ontology Editors**

A B C D

**Related TNA data (GWDB)**

**UK Data (CGO, CSYB)**

**Data Transformation and Integration**

**Government Web Archive (WARC)**

**Semantic Annotation**

**Semantic Repository**
- SKB Ontologies
- Factual Knowledge (TNA data, LOD, data.gov.uk)
- Semantic annotations
- MIMIR Index

**Identity Resolution**

**Linked data (LOD and data.gov.uk)**

**Annotation Process (GATE Teamware)**

general architecture

GATE

for text engineering

ontotext
Semantic Technology Lab

The University Of Sheffield.

# How does it work?

- Step 1: Annotate the documents using GATE
  - Ontology-based information extraction (OBIE) tools to annotate entities in text and relates them to ontology where appropriate

- Step 2: Index the documents using GATE/MIMIR
  - Create a MIMIR index based on the annotations produced in Step 1

- Step 3: Search the documents using MIMIR
  - Query the MIMIR index using full text or annotation-based queries

- Step 4: Browse the results
  - Search results point back to the texts in the original archive

# Annotation Types

**General NEs:**

- standard ANNIE NEs (with some additional features)

**Measurements:**

- measurement (dimension, type, unit, value, normalised, scalar, interval)
- ratio (value)

**Posts:**

- cabinet, civil service, military, medical, other (e.g. MP, CEO)

**Official Documents**

- legislation, other (e.g. white paper)

**Projects/Initiatives/Campaigns**

**Wars/Military Conflicts**

# Measurements

- Measurement plugin based on the GNU units application

- Recognise numbers in documents, and then use the parser to determine if the following words are valid units

- Number and units are combined to give a measurement

- Each measurement is then normalised to SI units, e.g. all time measurements are normalised to seconds, distances to meters etc.

- This enables searching for measurements no matter how they are expressed in the text, e.g. searching for 0.03m would also find references to 3cm, 30mm and 1.18 inches.

# Measurement annotation

# Relative date normalisation

- Absolute dates have an annotation feature showing the normalised form

- Special plugin that converts the fully specified date into standard numerical format

- If we know the date of the article, we can also calculate the actual value of relative dates (e.g. "last year", "in the next fortnight", etc.)

- We can often find the date of the article from the title or other information in the body of the article: this is added as a document feature

- The default option is to use the date the document was crawled

- Once we have a reference date for the document, we can normalise the values against it

# Relative Date annotation

# Semantic Knowledge Base (SKB)

- Semantic Knowledge Base developed (by Ontotext) on the basis of information extracted from the archive, along with existing KBs and ontologies

- Pre-processed, transformed and integrated to form a consistent KB

- Continuously enriched and extended

- Contains linked data from LOD cloud, data from data.gov.uk, data from related TNA projects, geographical data from Ontotext

- We don't try to extract, for example, a consistent description of London from the archive

  - Instead we integrate various knowledgebases and ontologies where extensive and well-maintained geographical features are already present

  - So a reference to London from a web page in the archive will be annotated with reference to the profile of London in the background knowledge bases
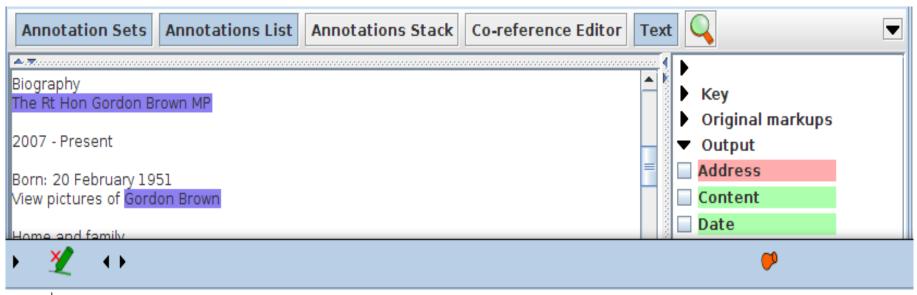
# Using SKB to Annotate the Archive

- Standard IE uses small hand-crafted gazetteers to annotate relevant government entities.

- Semantic Knowledge Base developed (by Ontotext) on the basis of information extracted from the archive, along with existing KBs and ontologies

- The SKB contains a lot more relevant information useful for annotation

- Use the LKB gazetteer to annotate certain concepts directly from the ontology

    - his leads to greater coverage and (theoretically) better precision

    - Where appropriate, annotated entities are linked to specific instances in the ontology, for example Post

- Identity resolution module also links co-referring mentions with a known URI

# Co-referring mentions linked to instance

# GATE Cloud Paralleliser

- One of the main challenges in the task of annotating, indexind and searching millions of documents is the sheer size of the data

- At the time of processing, the archive contained 42TB of data (around 700 million documents, of which around 150 million were unique)

- We processed the data using GATE Cloud Parallelizer, installed on the Amazon Cloud https://gatecloud.net/

- GCP is a platform for parallel semantic annotation of documents, designed as a parallel version of the execution engine installed in GATE

- It takes a language processing pipeline created using GATE Developer, and executes it using a series of parallel threads

- The job control is performed by executing batches, which are XML files describing outstanding tasks

# GATE Mímir: Answering Questions Google Can't

# Mímir

- Mímir is an IR engine that can index and search over:
  - text
  - semantic annotations
  - ontologies and KBs
- Allows queries that arbitrarily mix full-text, structural, semantic and linguistic annotations
- Scales to millions of documents

# What can GATE Mímir do that Google can't?

Show me:

- all documents mentioning a temperature between 30 and 90 degrees F (expressed in any unit)

- all abstracts written in French on Patent Documents from the last 6 months which mention any form of the word "transistor" in the English abstract

- the names of the patent inventors of those abstracts

- all documents mentioning steel industries in the UK, along with their location

# Search News Articles for Politicians born in Sheffield

http://demos.gate.ac.uk/mimir/gpd/search/gus

# Summary

- Use of simple mechanisms can add value and increase usage in the short and medium terms to a huge archive of information.
- Our GATE-based text-mining solution provides search paradigms over semantic annotation that relates archival content to Linked Data and other structured sources.
- Using a Semantic KB makes it possible to integrate external knowledge with facts extracted from texts
- Initial evaluations by TNA (comparing it with their existing search tools) were very promising
- Tools will continue to be developed and enhanced by in-house TNA team, and new interfaces developed