



**British
Geological Survey**
Expert | Impartial | Innovative



Gateway to the Earth

Geosemantics: Text Mining and Linked Data at the British Geological Survey

Jo Walsh – Senior Software Engineer, British Geological Survey

jowalsh@bgs.ac.uk

British Geological Survey (BGS)

The British Survey can claim to be the oldest such undertaking to have functioned continuously since its inception in May 1832 and formal establishment on 11 July 1835.



Sir Henry Thomas De la Beche, by William Brockedon, 1842 (© National Portrait Gallery, London).

Geosemantics

- Natural Language Processing with open source tools
- Named Entity Recognition model from labelled text
- Linked Open Data vocabularies for BGS resources
- Classifying relations between Named Entities

Geological Memoirs

The memoirs,
compiled by BGS
geologists, date from
the late 1890s to
present day and
provide a
comprehensive and
detailed account of all
aspects of the geology
of an area.

South-East Caldera and Concentric Folds.

175

Ghleannain, which lies in a hollow excavated along soft 'schists' (slates and limestone), and is flanked by ridges of Old Red Sandstone lavas, in their turn separated by hollows, worn in Mesozoic sediments, from uplands of Tertiary lavas on either side. The shores of the Firth of Lorne and Loch Spelve help to complete the structural presentation.

The best district to study the relationship of the Loch Spelve Anticline to the South-East Caldera lies about Sgùrr Dearg, and is illustrated in Figs. 30, 35, pp. 204, 237.

In reading the following descriptions, it is important to bear in mind that, despite the time-intervals involved, there is no appreciable difference of dip between the Tertiary lavas of Mull and the underlying Mesozoic sediments. This fact often greatly simplifies the study of the Tertiary tectonic features of the district.

Marginal Tilt. The Marginal Tilt, as it is styled in Fig. 25, furnishes the outer limb of the Duart Bay Syncline, and must be held responsible for the preservation of basalt-lavas in the Java Point peninsula as compared with gneiss in Sgùrr nan Gobhar, Glas Eileanan, and Eilean Rudha an Ridire, in the Sound of Mull.

From Duart Point southwards, the Marginal Tilt is clearly discernable in the lie of the lavas as shown in Plate V. The arcuate trend of the structure is also particularly well-shown in the course of the coast-line between Duart Point and the entrance to Loch Buie; for this gently curved coast marks approximately the line along which the lavas affected by the Marginal Tilt pass through sea-level.

West of Loch Buie, there is some difficulty in deciding what to refer to as Marginal Tilt. Near Carsaig Bay, the tilt clearly must form part of the eastern limb of a faint anticline which, assisted by faulting, leads to an important exposure of Mesozoic sediments; but, whereas this anticline is vaguely traceable east-north-eastwards (p. 181) into country unaffected by arcuate folding, the marginal tilt seems to continue as an independent and often well-defined flexure northwards across Loch Beg at the head of Loch Scridain.

Duart Bay Syncline. The continuity of the Duart Bay Syncline from Craignure Bay to near Corra-bheinn as shown on Plate V. is easily verified in the field, except where, for a space south of the entrance to Loch Spelve, the fold is very shallow. Sometimes the lavas are so clearly seen dipping in from either side towards the axis of the syncline that the structure becomes quite a scenic feature, as for instance north of Loch Spelve and in the coastal cliffs either side of Loch Buie.

Comparison of Plates III. and V. shows that an additional indication of the syncline is afforded by the preservation of relatively late rocks in its embrace; thus in the north-east part of its course, one finds Tertiary lavas flanked by Pre-Tertiary rocks, while north of the Port Donain Fault, and again, both east and west of Loch Buie, there are outcrops of big-felspar basalt, or mugearite, as the case may be, contrasting with the normal unbroken sequence of Plateau Basalts met with either side.

Geological memoirs

Scanned texts
processed with optical
character recognition
software have been
corrected by hand and
annotated with
identifiers for other BGS
data (borehole logs,
core samples)

MUDSTONE FORMATIONS 75

turbidites are up to 5 mm thick, but are commonly discontinuous and lenticular, occurring as a series of low-angle, sigmoidal lenses.

The proportion of hemipelagite in the formation varies from 30 to 70 per cent. Individual beds range up to 5 cm thick, although predominantly the laminated variety, rare burrowed hemipelagites are locally present. They are most common in the top few metres of the formation, beneath the Derwenslas Formation.

Individual sandstones are thin, massive (type C1), occurring either scattered or as thin lamellae, or largely restricted to regions of intercalation with the Caban Conglomerate and Yatrad Meurig Grits formations. Individual sandstones, in type C1 and C1t turbidites, are up to 20 and 5 cm thick, respectively.

BIOSTRATIGRAPHY
Graptoites from the Rhondda sheet area indicate that the Cwmre Formation ranges in age from the *percivalensis* Biozone to the *magnus* Biozone. The Mottled Mudstone Member is entirely of *percivalensis* Biozone age. The distribution and approximate thicknesses of individual bio-zones are depicted in Figure 19. Faunal lists from localities in the Cwmre Formation, beneath the Cwmre Formation and the intercalations of Caban Conglomerate are given in Tables 9 and 12.

On the north-western limb of the Ystwyth Anticline, north-east of the Afon Elan, the mudstones overlying facies of the Caban Conglomerate Formation range from possibly as low as upper *acuminatus* Biozone through to the *percarinatus* Biozone. There, in the Wye valley (Figure 19; Table 9), *acuminatus* Biozone mudstones are overlain 25 m above the base of these mudstones in the lowest part of the A470 road section [SN 9799 6738 to SN 9749 6749]. Graptolites from the remainder of the A470 section and from a stratigraphically higher section in the River Wye at Ddole Farm [SN 9776 6737 to SN 9762 6745], have demonstrated that the *acuminatus*, *cylindrus* and *truncatus* bio-zones are 47 m and 20 m thick respectively (Table 9). The top 14 m of the mudstones contain *magnus* Biozone graptolites.

In the Llannau sheet area, *engyulus* Biozone graptolites also occur in the uppermost part of the Cwmre Formation, while upper *acuminatus* Biozone faunas have been preserved in the basal 10 m of the mudstones of the tongue of Yatrad Meurig Grits Formation (Edward Goch facies) (Table 13). Although graptolite faunas have been found in the remainder of the formation (Table 13) none of the intervening bio-zones has been definitely proved.

The only shelly taxon from the Cwmre Formation in the district is a probable novakid terebratulid, *Novakia genensis*, from the Wye valley [SN 9799 6738] east of Rhayader. This may be associated here to the *acuminatus* Biozone rather than the *truncatus* Biozone reported by Tunnicliffe (1989).

The mudstones in Nant Paradesy [SN 8918 6088 to SN 8933 6002] have yielded sphaeromorph acritarchs and indeterminate acanthomorphs of no biostratigraphical value (Appendix 2bii). More diverse assemblages including *Dicyctidium dictyonum*, *Lophophoridium parvulum*, *Mulicosphaeridium fischeri* and *Tychotella costatum*

were obtained from *magnus* Biozone mudstones at Ddole Farm [SN 9761 6747] (Appendix 2biii). They are referred to as *O. opalinoides* M. Johnson or to the *O. opalinoides* Biozone. Recycled acritarchs with similar morphologies from both localities include *Anomolodiscus* spp. and *Valecania sanguinostrotum*, *Stellerites* sp. cf. *bryocystoides* and *Valecania sanguinostrotum* from the mid- to late Ordovician.

DETAILS

MOTTLED MUDSTONE MEMBER

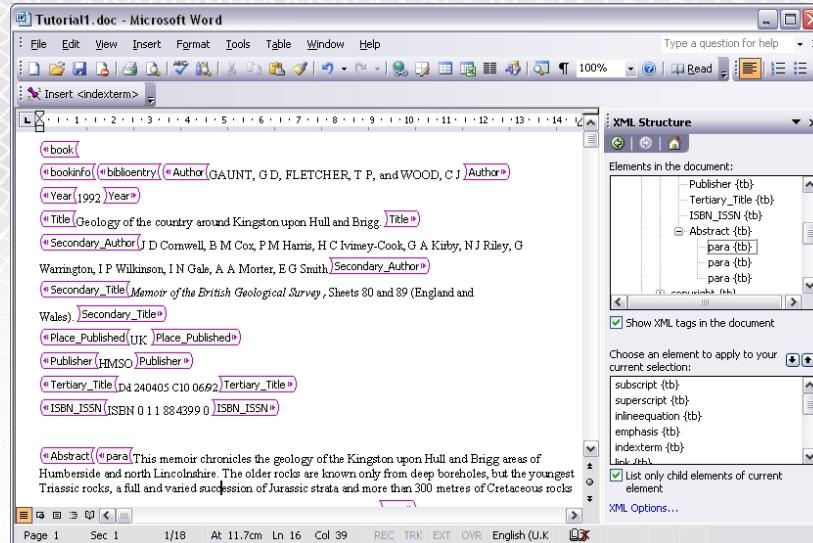
On the north-western limb of the Rhondda Anticline part of the Mottled Mudstone is exposed in Flos y Rhest east of Llyn Garw (Figure 21b). Scattered sections in the member are common around the closure and on the south-eastern limb of the Ystwyth Anticline, between Craig y Bodach and Craig y Ddol, on the north side of the Claerwen valley, the sharp basal contact with the Yr Alt Formation and the gradational contact with the overlying Ddyllyn Flaser are well exposed. Rare acanthite sandstone/mudstone complex (type D) are seen in the lower mudstones, up to 2 cm thick, occurring in the top 1.5 m. A kilometre to the southeast on Craig y Bodach [SN 8991 6205] (Figure 18a, locality 1), *Normalgraptols* spp. were obtained 60 cm above the base of the member in an incomplete section (Table 12). On the south-eastern limb of the Rhondda Anticline, between Nant-y-Groes and Nant Paradesy [SN 8917 6097] (Figure 21c) and 6 m thick in Nant y Groes [SN 8779 5993], both sections exhibit a gradational passage up into Ddyllyn Flaser.

Scattered thin sections of the Mottled Mudstone Member on the north-western limb of the Ystwyth Anticline, where it is overlain by Cerig Gwynion Grits, include a stream section [SN 9533 6450] north of Cwm and Cerig Gwynion quarry [SN 9710 6559] (Figure 20; Table 9; Plate 12a), both south-west of Rhayader. At the latter locality the following section is seen:

CABAN CONGLOMERATE FORMATION, CERIG GWYNION GRITS FACIES	Thickness m
SANDSTONE, medium-grained with sharp base	
CWMRE FORMATION, MOTTLED MUDSTONE MEMBER	
Hemipelagite, pale and medium grey, colour banded, diffusely laminated and burrowed. Scattered mudstone turbidites with silstone laminae up to 5 mm thick (type E) in upper 1.3 m. Thin-bedded sandstone/mudstone complex (type D) with basal 2 cm sandstone defines base	2.0
Hemipelagite, above	1.9
Upper leaf of <i>percivalensis</i> Band: hemipelagite, dark grey, laminated with abundant irregular pyrite nodules up to 1 cm in diameter. Abundant percarinatus Biozone graptolites are preserved in pyrite (Table 9)	0.08
Hemipelagite, pale and medium grey, colour-banded, burrowed and diffusely laminated. Lower leaf of <i>percivalensis</i> Band: hemipelagite as in upper leaf	0.27
Hemipelagite, pale and medium grey, colour-banded, diffusely burrowed and laminated. Scattered pyrite nodules in upper part. Abundant interbedded泥质页岩与硅质页岩互层	0.17

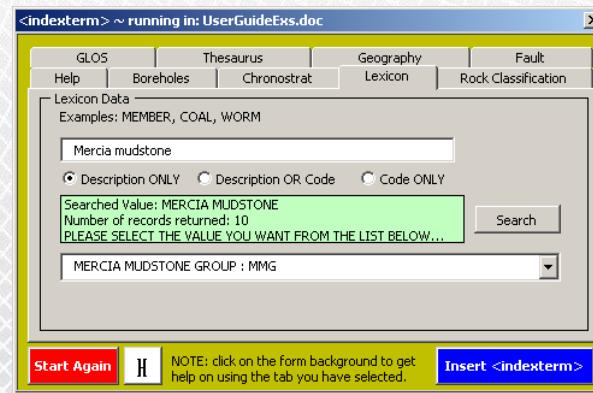
Textbase 2000

In the early 2000s, BGS ran a small project converting 30 regional reports including memoirs into XML using Word



Textbase 2000

This earlier project added not just syntactic information about the layout of the document, but also some semantic information, labelling words and phrases and linking them to BGS controlled vocabularies



Textbase 2000

- From these manually-created annotations we extracted 2500 sentences
- Every term found in the BGS Lexicon of Named Rock units is labelled LEXICON
- Every reference to a geochronological or chronostratigraphic term is labelled CHRONOSTRAT
- We used this data to train a custom Named Entity Recognition model

Stanford CoreNLP NER

Stanford CoreNLP 3.9.1 (updated 2018/04/05)

- Basic

— Text to annotate —

consistent with its palaeo-connection into the thicker and more extensive Carboniferous succession to the north-west, in the Midland Valley.

— Annotations —

named entities dependency parse openie

— Language —

English

Named Entity Recognition:

- 1 Apart from several named and extensively worked coal seams , the succession consists mainly of sandstone , siltstone , mudstone and seatearth , with ironstone ribs in places .
NATIONALITY TITLE
- 2 As with the underlying Scottish Lower Coal Measures , there is a general thinning of the succession towards the eastern part of the Sanquhar Basin , consistent with its palaeo-connection into the thicker and more extensive Carboniferous succession to the north-west , in the Midland Valley .
CITY LOCATION MISC CITY

Named Entity Recognition (NER) with default CoreNLP

CoreNLP training data - sample

```
The      O
highly   O
irregular O
middle   O
zone     O
is       O
underlain O
by       O
Permian  CHRONOSTRAT
and      O
Triassic CHRONOSTRAT
rocks    O
;
the     O
Permian  CHRONOSTRAT
strata   O
include  O
Upper    CHRONOSTRAT
Permian  CHRONOSTRAT
Zechstein LEXICON
sedimentary O
rocks    O
that     O
locally  O
crop    O
out     O
in      O
the     O
study   O
area    O
.       O
```

Train NER model with CRFClassifier

```
java -cp stanford-ner.jar  
edu.stanford.nlp.ie.crf.CRFClassifier -prop  
bgs.3class.geo.prop
```

<https://nlp.stanford.edu/software/crf-faq.shtml>

Custom NER for CoreNLP

Stanford CoreNLP

— Text to annotate —

consistent with its palaeo-connection into the thicker and more extensive Carboniferous succession to the north-west, in the Midland Valley.

— Annotations —

named entities openie

— Language —

English

Submit

Named Entity Recognition:

1 Apart from several named and extensively worked coal seams , the succession consists mainly of sandstone , siltstone , mudstone and seatearth , with ironstone ribs in places .

LEXICON

2 As with the underlying Scottish Lower Coal Measures , there is a general thinning of the succession towards the

LOCATION

eastern part of the Sanquhar Basin , consistent with its palaeo-connection into the thicker and more extensive

CHRONOSTRAT

LOCATION

Carboniferous succession to the north-west , in the Midland Valley .

Named entities recognised with custom model added

Adding additional vocabularies

- Given enough labelled examples, train a Named Entity Recognition model with a different vocabulary
- We had a few BIOZONEs in our training data but too few examples for the classifier to learn from

Rule-based approach using parts of speech

Part-of-Speech:

1 Environments representing a series of staging-posts between fully marine and limnetic settings .
2 Macrofossils and ostracods are assigned to marine , marginal marine , brackish and freshwater environments based on their faunal assemblage patterns .
3 Key brackish to freshwater ostracods are Geisina arcuata , Paraparachites circularis n. sp. , Shemonaella ornata n. sp .
4 and Silenites sp .
5 A , associated with the bivalves Anthraconalia , Carbonicola , Cardiopteridium , Curvirostrum , Naiadites , the microconchid ` Spirorbis ' , Spinicaudata and fish

CoreNLP Tools:

TokensRegex Semgrep Tregex

Enter a **TokensRegex** expression to run against the above sentence:

`(({pos>NNP|NNPS|FW})+|({pos>NN|NNS|/})?)` **TokenRegex**

1 Environments representing a series of staging-posts between fully marine and limnetic settings .
2 Macrofossils and ostracods are assigned to marine , marginal marine , brackish and freshwater environments based on their faunal assemblage patterns .
3 Key brackish to freshwater ostracods are Geisina arcuata , Paraparachites circularis n. sp. , Shemonaella ornata n. sp .
4 and Silenites sp .
5 A , associated with the bivalves Anthraconalia , Carbonicola , Cardiopteridium , Curvirostrum , Naiadites , the microconchid ` Spirorbis ' , Spinicaudata and fish

Or use manually extracted entities to recreate labelled sentences

 **Geobiodiversity Database**
Promoting collaborations based on global and regional database

Jo Walsh [Sign out](#)
Your session will be terminated after
120 minutes of inactivity.

- Home
- Add Data
- Search
- Data Subset
- Reference Management
- Specimen Identification Data
- Data Accessibility
- Analyze
- Download
- Members
- About

Links

- Choose a website ---
- Questions & Answers ---

 International Commission
on Stratigraphy

Contact us

Dr. Fan Junxuan
Phone: +86-25-83282247
E-mail: fanjunxuan@gmail.com

Free data access for anonymous users!
Register to make contributions and create!

Occurrence search result ? Consult the tip sheet for help

Occurrence search result								
Occurrence ID		Fossil name	Revised name	Collection No.	Locality	Province	Country	Section described in
<input type="checkbox"/>	436558	Palaeoneilo brevisstrom		1	Kilconquhar core	Scotland	United Kingdom of Great Britain and Northern Ireland	BENNETT et al. (2012) (ID = 82553)
<input type="checkbox"/>	436559	Lingula mytiloides		1	Kilconquhar core	Scotland	United Kingdom of Great Britain and Northern Ireland	BENNETT et al. (2012) (ID = 82553)
<input type="checkbox"/>	436560	Productus		1	Kilconquhar core	Scotland	United Kingdom of Great Britain and Northern Ireland	BENNETT et al. (2012) (ID = 82553)
<input type="checkbox"/>	436561	Productus		2	Kilconquhar core	Scotland	United Kingdom of Great Britain and Northern Ireland	BENNETT et al. (2012) (ID = 82553)
<input type="checkbox"/>	436562	Euphemites		2	Kilconquhar core	Scotland	United Kingdom of Great Britain and Northern Ireland	BENNETT et al. (2012) (ID = 82553)

Alternatives to CoreNLP

- **spacy.io** is a python-based toolkit for Natural Language Processing
- Nicer for integration with other python tools for machine learning and scientific computing
- We tried re-training spacy's NER with BGS data but couldn't get it to the same result quality
- USGS have used spacy to recognise dams and rivers in their dam removal project

Alternatives to CoreNLP

- Edinburgh Geoparser from Language Technology Group, School of Informatics, University of Edinburgh
- Uses a rule-based approach rather than a statistical model
- We used its “georesolution” phase to find coordinates for locations

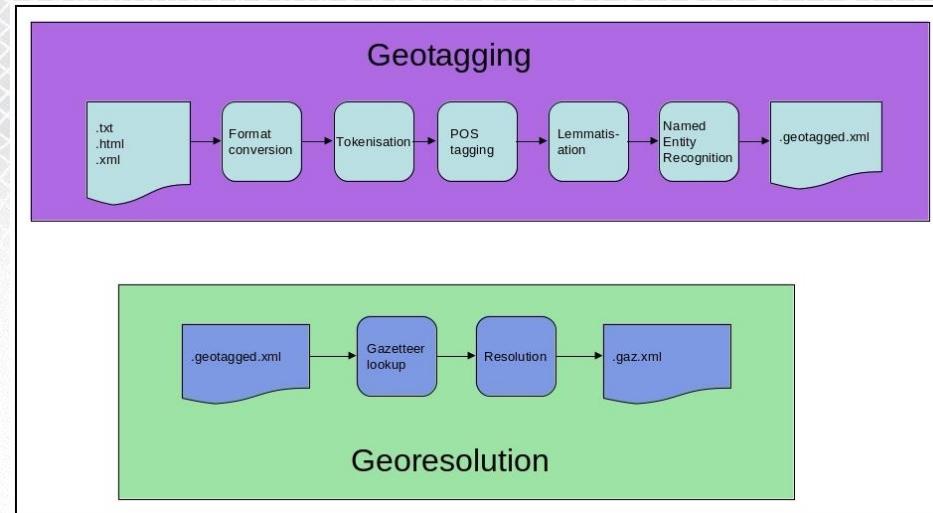


Image captions derived from sheet memoirs



Uploaded on:
14/02/2009 05:13:41

Type:
Digital Asset

File Size:
428.79 KB

Dimensions:
1001 x 791 pixels

381 views

5 downloads

P number: P002882

Old photograph number: D05094

Caption: East Kirkton Quarry, near Bathgate, West Lothian.

Description: East Kirkton Quarry, near Bathgate, West Lothian. The East Kirkton Limestone in the upper part of the West Lothian Oil Shale Formation is a laminated freshwater limestone. A unique and varied terrestrial assemblage of fossils has been collected from this site, including *Westlothiana lizziae*, a vertebrate fossil intermediate between amphibians and reptiles.

Entities found with CoreNLP

Stanford CoreNLP

— Text to annotate —
East Kirkton Quarry [NS 990 092], near Bathgate. The East Kirkton Limestone in the upper part of the West Lothian Oil Shale Formation is a laminated freshwater limestone. A unique and varied terrestrial assemblage of fossils have been collected from this site.

— Annotations —
[named entities] [openie] — Language — English

Named Entity Recognition:

1 East Kirkton Quarry [NS 990 092] , near Bathgate .

2 The East Kirkton Limestone in the upper part of the West Lothian Oil Shale Formation is a laminated freshwater limestone .

3 A unique and varied terrestrial assemblage of fossils have been collected from this site .

Open IE:

1 East Kirkton Quarry [NS 990 092] , near Bathgate .

2 The East Kirkton Limestone in the upper part of the West Lothian Oil Shale Formation is a laminated freshwater limestone .

```
graph LR; E1[The East Kirkton Limestone] -- "in" --> E2[in the upper part of the West Lothian Oil Shale Formation]; E2 -- "is a" --> E3[laminated freshwater limestone];
```

Links to BGS vocabularies

```
{  
  "_id": "B01613-001-p001-Geochronology:Division-0",  
  "term": {  
    "term": "Ordovician",  
    "vocabulary": "Geochronology:Division",  
    "url": "http://data.bgs.ac.uk/id/Geochronology/Division/0",  
    "code": "O"  
  },  
  "text": {  
    "match": "Ordovician",  
    "offset": [  
      38,  
      48  
    ],  
    "snippet": "b'The district is mainly underlain by Ordovician and Silurian mudstones and sandstones. "  
  },  
  "doc": {  
    "id": "B01613",  
    "section": [  
      "B01613",  
      "001",  
      "p001"  
    ]  
  }  
}
```

Rectangular Snip

JSON documents stored in and queried from a MongoDB document store

Lexicon of Named Rock Units

The BGS Lexicon of Named Rock Units – Result Details

Bathgate Hills Volcanic Formation

Computer Code:	BHV	Preferred Map Code:	BHV
Status Code:	Full		
Age range:	Asbian Substage (CR) — Arnsbergian Substage (CG)		
Lithological Description:	Pyroclastic rocks and related basaltic lavas, mostly alkali olivine-basalts (basanite to hawaiite) and volcaniclastic rocks.		
Definition of Lower Boundary:	First occurrence of volcanic rocks above the Two Foot Coal in the West Lothian Oil-Shale Formation.		
Definition of Upper Boundary:	Highest occurrence of volcanic rocks, below the Roman Cement Limestone, Passage Formation.		
Thickness:	To 1000m		
Geographical Limits:	Bathgate Hills, west of a line from NT0066 to NT0980, extending westwards underground to a point west of the Rashiehill Borehole (NS839730).		
Parent Unit:	Bathgate Group (BATH)		
Previous Name(s):	Bathgate Lavas (-3215) Bathgate Volcanic Formation (-58)		
Alternative Name(s):	none recorded or not applicable		
Stratotypes:			
Type Area	Natural sections in the Bathgate Hills.		
Reference Section	Rashiehill Borehole (NS87SW/22), (Anderson, 1963), 2600ft (792.48m) to 3645ft (1110.99m).		
Reference Section	Couston Borehole (NS97SW/99) 110 fathoms (201.17m) to 298 fathoms (544.98m).		

Web-based

Human-readable

Linked Data Lexicon

The screenshot shows a web browser displaying the British Geological Survey's (BGS) Linked Data Lexicon. The page title is "BGS Lexicon of Named Rock Units". The header features the BGS logo and the Natural Environment Research Council (NERC) crest. Below the title, there are links for "Home" and "Lexicon". The main content area is titled "http://data.bgs.ac.uk/ref/Lexicon" and contains a table of metadata:

Access rights	http://data.bgs.ac.uk/ref/void
Creator	http://data.bgs.ac.uk/ref/BritishGeologicalSurvey
Publisher	http://data.bgs.ac.uk/ref/BritishGeologicalSurvey
Title	"BGS Lexicon of Named Rock Units" <i>language=en</i>
See also	http://www.bgs.ac.uk/Lexicon/home.cfm http://www.bgs.ac.uk/data/services/vocabulary/1.0/home.html
Type	http://www.w3.org/2004/02/skos/core#ConceptScheme

Below the table are links for "HTML", "RDF", "Turtle", "N-Triples", and "JSON". To the left, a section titled "namespaces used" lists:

bgsr	BGS Schemas
bgsre	BGS RCS Earth Material Classes
bgsrg	BGS Geochronology
bgsrl	BGS Lexicon of Named Rock Units
bgsrm	BGS 1:625 000 Digital Geological Map
dct	Dublin Core terms
rdf	Resource Description Framework
rdfs	Resource Description Framework - Schema
skos	Simple Knowledge Organization System
xsd	XML Schema Data Type

To the right, a section titled "49 sub-concepts and/or properties" lists:

NamedRockUnit	HTML JSON N-Triples Turtle RDF
Theme	HTML JSON N-Triples Turtle RDF
Class	HTML JSON N-Triples Turtle RDF
LithologyComponent	HTML JSON N-Triples Turtle RDF
LithogeneticType	HTML JSON N-Triples Turtle RDF
DefinitionStatus	HTML JSON N-Triples Turtle RDF
RockUnitRank	HTML JSON N-Triples Turtle RDF
EquivalentName	HTML JSON N-Triples Turtle RDF

A "view all" link is located at the bottom right of this list.

Web-based

Machine-readable



“Pompidou Centre” data

A phrase coined by Peter Burnhill of EDINA, University of Edinburgh

Data infrastructure is exposed on outside of the structure



Photo credit: copyright Pascal
Poggi

“Use URIs as names for things”

- Internal code: CAL
- External code:
<http://data.bgs.ac.uk/id/Lexicon/RockUnit/CAL>



Link to other things on the semantic web

<http://data.bgs.ac.uk/id/Geochronology/Division/C>

owl:sameAs

<http://resource.geosciml.org/classifier/ics/ischart/C>

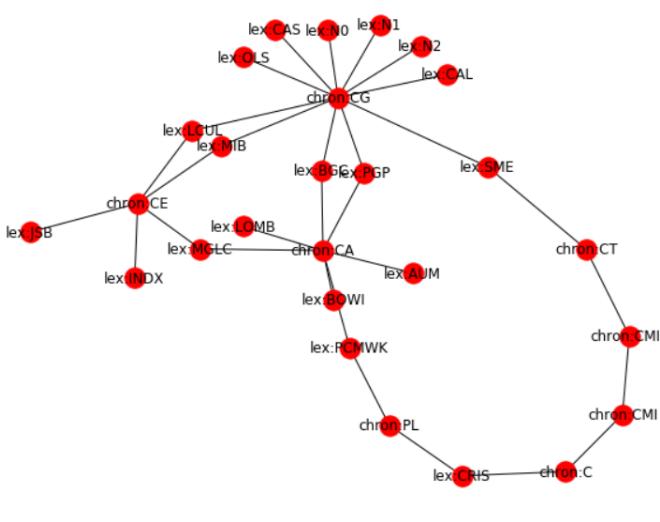
owl:sameAs

<http://dbpedia.org/page/Carboniferous>

Triples all the way down

Knowledge Graph

Fashionable term for linked data extracted from text

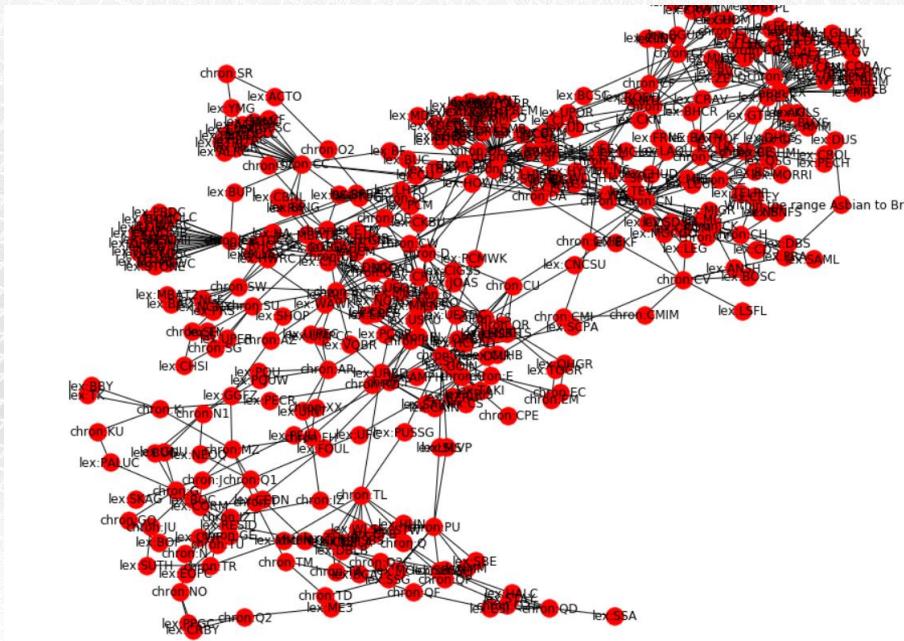


Named entities extracted from a BGS Palaeontology report

Sum of all shortest paths through a graph of connections between min/max ages of rock types and ICS hierarchy of geochronological eras

rdflib, networkx, matplotlib

Is this knowledge?



Same process;

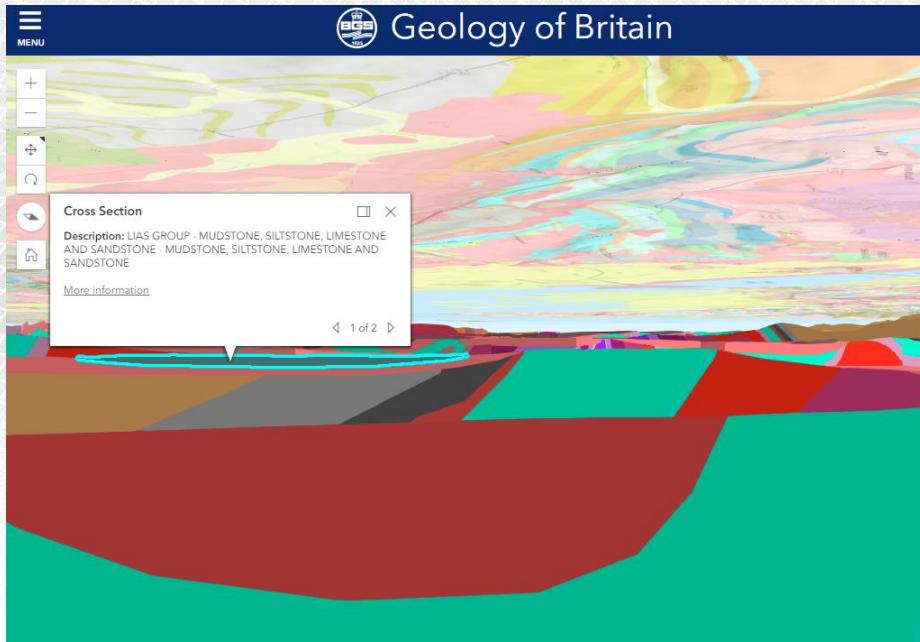
Graph inferred from
several hundred page
map sheet memoir

Links to other resources

- Other kinds of BGS data resources are already linked to the Lexicon and Geochronology
- Thin section imagery, borehole scans
- Deposited datasets have some manually-added tags
- We can automatically add more to new incoming data
- Also a lot of un-mined text content in the Lexicon itself

Other linked resources

BGS 3D geological models are linked to the rock units



BGS memoirs – entities visualised

- We ran the Named Entity Recognition over 60 BGS memoirs in cleaned-up XML format
- We extracted references to Chronostratigraphic terms, to Lexicon of Named Rock Units, and locations
- We also ran a selection of PDF scans of palaeontology reports
- Even with lower-quality OCR the results are promising

Geochronology extracted from text

32 Geochronological Divisions found in B01881

Divisions matched sort by Age (young-old) ▾

Code	Name	Min Age (MYBP)	Max Age (MYBP)	Matches	Division frequency
Q	Quaternary Period	0 ±0	2.588 ±0	22	Q 22
FH	Phanerozoic Eon	0 ±0	541 ±1	1	FH 1
QDL	Loch Lomond Stadial (Younger Dryas)	0.0118 ±0	0.01265 ±0	9	QDL 9
QD	Devensian Stage	0.0118 ±0	0.116 ±0	5	QD 5
QDD	Dimlington Stadial	0.0147 ±0	0.026 approx. ±0	2	QDD 2
N	Neogene Period	2.588 ±0	23.03 ±0	1	N 1
P	Permian Period	252.2 ±0.5	298.9 ±0.2	3	P 3
RZ	Palaeozoic Era	252.2 ±0.5	541 ±1	2	RZ 2
C	Carboniferous Period	298.9 ±0.2	358.9 ±0.4	70	C 70
CS	Stephanian Stage	302 approx. ±0	308 approx. ±0	1	CS 1
CW	Westphalian Stage	308 approx. ±0	319 approx. ±0	23	CW 23
CB	Duckmantian Substage	315.2 ±0.2	318 approx. ±0	2	CB 2
CA	Langsettian Substage	318 approx. ±0	319 approx. ±0	3	CA 3
CY	Yeadonian Substage	319 approx. ±0	320 approx. ±0	1	CY 1
CN	Namurian Stage	319 approx. ±0	329 approx. ±0	33	CN 33
CZ	Marsdenian Substage	320 approx. ±0	321.5 approx. ±0	1	CZ 1
CK	Kinderscoutian Substage	321.5 approx. ±0	322 approx. ±0	1	CK 1
CO	Alportian Substage	322 approx. ±0	323 approx. ±0	2	CO 2
CH	Chokierian Substage	323 approx. ±0	324 approx. ±0	2	CH 2
CG	Arnsbergian Substage	324 approx. ±0	328 approx. ±0	6	CG 6
CE	Pendleian Substage	328 approx. ±0	329 approx. ±0	5	CE 5
CX	Brigantian Substage	329 approx. ±0	330.9 approx. ±0	9	CX 9
CR	Asbian Substage	330.9 approx. ±0	337 approx. ±0	3	CR 3
CI	Chadian Substage	344.5 approx. ±0	346.7 ±0.4	1	

Mentions of rock unit types

BGS TextMining :: Document Report					Home	Document Report	Vocabulary Report	Location Report
Code	Name	Min Age (MYBP)	Max Age (MYBP)	Matches	Rock Unit frequency			
PSG	Passage Formation (Wenlock Series)	427.4	433.4	54	PSG			54
BHV	Bathgate Hills Volcanic Formation	324	337	37	BHV			37
GLRB	Lower Limestone	145	152.1	33	GLRB			33
LSTC	Limestone Coal	329	330.9	33	LSTC			33
ULE	Upper Limestone	145	152.1	32	ULE			32
HUR	Hurlet Limestone	329	330.9	31	HUR			31
NO	No. 0 Marine Band	324	328	29	NO			29
CAS	Castle Cary Limestone	324	328	29	CAS			29
MCMI	Middle Coal Measures (Dungannon, Northern Ireland)	318	319	24	MCMI			24
WLO	West Lothian Oil-Shale Formation	329	337	23	WLO			23
INDX	Index Limestone [Duplicate Code: Use ILS]	328	329	19	INDX			19
RSH	Raeburn Shale	329	346.7	17	RSH			17
PTLS	Petershill Limestone	329	330.9	13	PTLS			13
CAL	Calmey Limestone	324	328	13	CAL			13
EKL	East Kirkton Limestone	329	346.7	12	EKL			12
MILSL	Mill Coal (South Lancashire)	315.2	318	12	MILSL			12
BLLS	Blackhall Limestone	329	330.9	11	BLLS			11
LOMB	Lowstone Marine Band	318	319	10	LOMB			10
URLS	Under Limestone	329	330.9	10	URLS			10
KILM	Kiltongue Musselband Coal	318	319	10	KILM			10
TOHO	Top Hosie Limestone	328	330.9	9	TOHO			9

Rock unit types across documents

BGS TextMining :: Vocabulary Report

[Home](#) [Document Report](#) [Vocabulary Report](#) [Location Report](#)

[«back](#) Found 104 matches for *Hurlet Limestone*

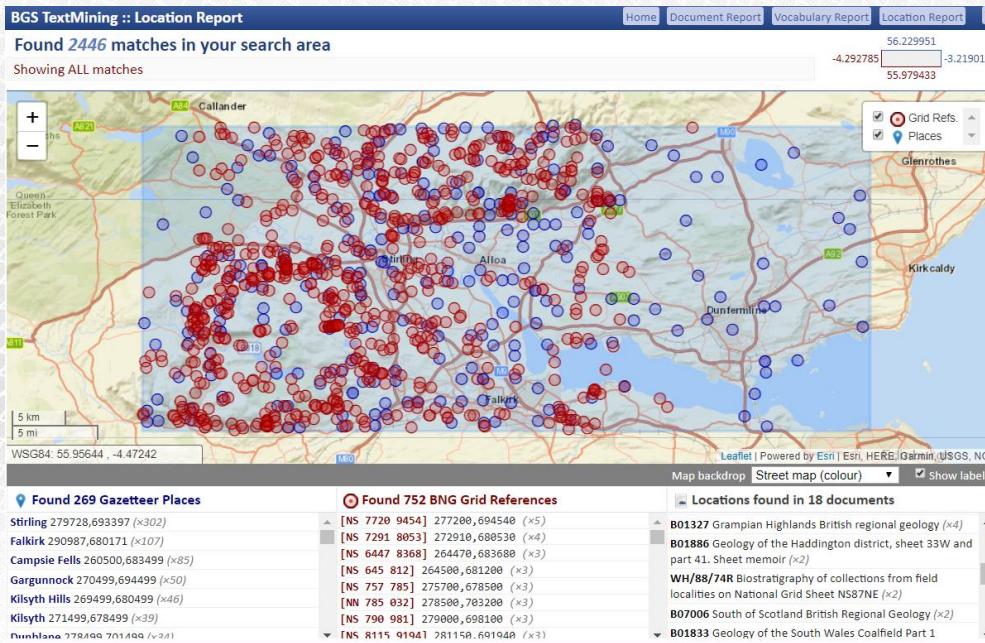
Lexicon Named Rock Unit Details

Code	HUR
Name	Hurlet Limestone
Description	Limestone with subsidiary mudstone. Stratified bedrock. Occurs onshore. Deposited during the Brigantian Substage (Carboniferous Period) (330.9-329 Ma BP).
Age Range	329 – 330.9

Matches found in 13 documents

	Filter text	Sort by
		Matches (N-0) ▾
 B01881	Geology of the Falkirk district, sheet 31E. Sheet memoir (Scotland)	x30
 B01876	Geology of the Hamilton district. Sheet 23W. Sheet memoir(Scotland)	x16
 B01879	Geology of the Glasgow district. Sheet 30E Sheet memoir (Scotland)	x14
 B01883	Geology of the Airdrie district. Sheet 31W Sheet memoir(Scotland)	x13
 B01880	Geology of the Greenock district, sheet 30W and part 29E. Sheet memoir (Scotland)	x12
 B01892	The geology of the Stirling district, sheet 39. Sheet memoir (Scotland)	x6
 B07006	South of Scotland British Regional Geology	x3
 WH/97/166R	Thin section petrology of nine samples from the archaeological site at Great Holts Farm, Boreham, Essex	x2
 WH/PD/85/4	Dinoflagellate cyst analysis of borehole 77/2, Fladen sheet, North Sea	x2
 B07005	Midland Valley of Scotland British regional geology	x2
 WH/96/268R	Faunas from and biostratigraphy of a section in upper Craig Burn, Lanark	x2
 B01886	Geology of the Haddington district, sheet 33W and part 41. Sheet memoir	x1
 B06059	Palaeogene volcanic districts of Scotland: British Regional Geology	x1

Memoir overview



Explore references to places

Coordinates found
in text:



Place-names
recognised in
text:



Correlations between entities

Lexicon Rock Unit :: Gazetteer Place		
Rock Unit	Distances	Place
Top Hosie Limestone <i>Lexicon NamedRockUnit TOHO</i>	100	Scotland 270946,768055
	100	
	100	
	100	
	100	
	100	
	100	
	100	
Johnstone Shell Bed <i>Lexicon NamedRockUnit JSB</i>	76	Scotland 270946,768055
	76	
	76	
	76	
	76	
	76	
	76	
	76	
Top Hosie Limestone <i>Lexicon NamedRockUnit TOHO</i>	78	Central Coalfield 437222,317148

Geochronological Division :: Gazetteer Place		
Division	Distances	Place
Carboniferous <i>Geochronology Division C</i>	4	South Africa 367891,445613
Namurian <i>Geochronology Division CN</i>	94	South Africa 367891,445613
Namurian <i>Geochronology Division CN</i>	18	Central Scotland 309846,725208

Next annotation task

- We can see co-occurrences of terms but want to know when they are most meaningful
- We want a model classify positive and negative correlations between named entities in sentences
- Natural History Museum in Oslo have done a similar project with excellent results
- Resulting in the ability to infer fossil ages from correlations of taxa with chronostratigraphic names

pybossa task to classify relations

pybossa.org - like an open source GalaxyZoo

The screenshot shows a web-based annotation task. At the top, there's a dark header bar with the 'pybossa' logo and navigation links for 'Community', 'Projects', 'Create', and 'About'. Below the header, the main content area has a light background. A large blue header says 'Entity Correlation Classification: Contribute'. Underneath, a question asks, 'What correlation can you identify in this phrase?'. To the right of the question, the text 'Task 212' is visible. A gray box contains the sentence: 'Early in the Late Permian, a marine incursion from the north-west into the northern part of the Cheshire Basin led to the deposition of the dolomitic and gypsiferous mudstones of the Manchester Marls Formation.' Below the sentence are three buttons labeled 'Strong', 'Weak', and 'Negative'. A large dark blue button at the bottom is partially visible.

We will need several annotators to label
several thousand sentences for a decent model

Training a TextCategorizer

- NHM in Oslo built with own classification model
- Feeding skipgram word embeddings from fasttext into an LSTM created with scikit-learn
- We are planning to start with the TextCategorizer package in spacy.io and work from there

Extracting tables and figures

- GeoDeepDive is a project at University of Madison, Wisconsin, US
- GeoDeepDive searches mainly journal publications but holds some geological survey archives
- With Stanford CoreNLP with Tesseract OCR they scan and analyse millions of short publications
- Blackstack is a new component of the GeoDeepDive project which takes a machine learning approach to identifying figures, maps and tables
- It's a work in progress!
- <https://geodeepdive.org>
- <https://github.com/UW-Deepdive-Infrastructure/blackstack>

Conclusions

“You have to have a human in the loop; what matters is where you put the human in the loop”

- *Claire Grover, Language Technology Group, University of Edinburgh*

- NER tools are accessible without NLP expertise
- Semantic web is reviving with a simpler approach
- Potential to generate scientific discoveries...
- ... but are we often discovering the obvious?

Geosemantics at BGS Informatics

Project leader: Rachel Heaven: reh@bgs.ac.uk

PhD student: Ike Nkisi-Orji: i.o.nkisi-orji@rgu.ac.uk