

Automatic Extraction of Archaeological Events from Text

Kate Byrne and Ewan Klein

Alberto González Palomo

Seminar “Unlocking the Secrets of the Past: Text Mining for Historical Documents (WS 2009/10)”

2010-02-24

Site456

SOUTH WALLS, MISBISTER, THE LOFTS
ND38NW29 centred 3325 8885

Sites recorded during an archaeological survey undertaken on the lands of the Loft, Longhope, as part of the pilot scheme for the Historic Scotland Farm Ancient Monument Survey Grant Scheme. 8885 Cairn. ND 3339 8886 Clearance cairn. ND 3342 8884 Sub-rectangular cairn. ND 3339 8883 Well Sponsors: Historic Scotland, M J Jones. N Card 1998



Clearance cairn in Germany

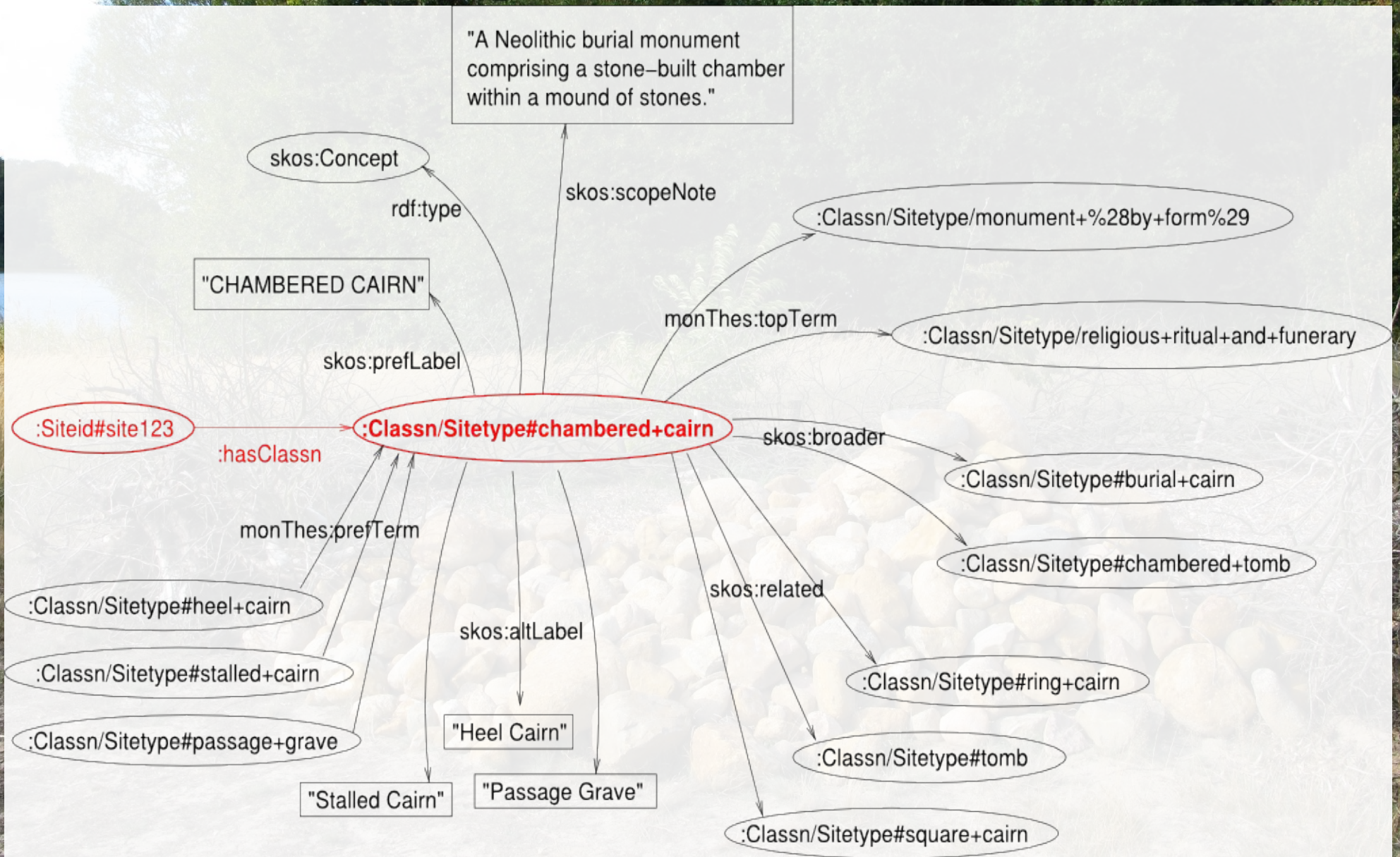


Memorial cairn in South Africa

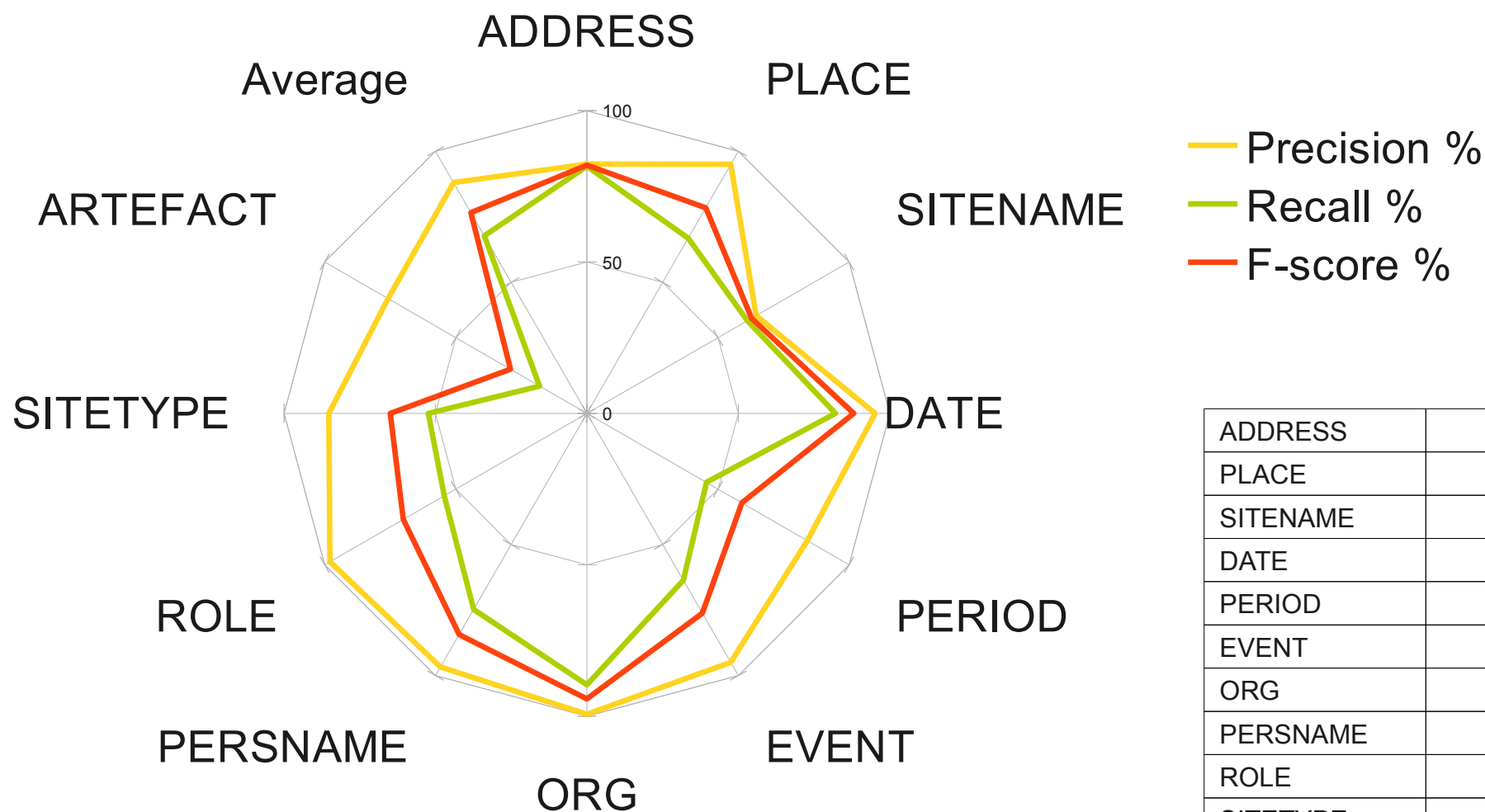
Site456

SOUTH WALLS, MISBISTER, THE LOFTS
ND38NW29 centred 3325 8885

Sites recorded during an archaeological survey undertaken on the lands of the Loft, Longhope, as part of the pilot scheme for the Historic Scotland Farm Ancient Monument Survey Grant Scheme. 8885 Cairn. ND 3339 8886 Clearance cairn. ND 3342 8884 Sub-rectangular cairn. ND 3339 8883 Well Sponsors: Historic Scotland, M J Jones. N Card 1998

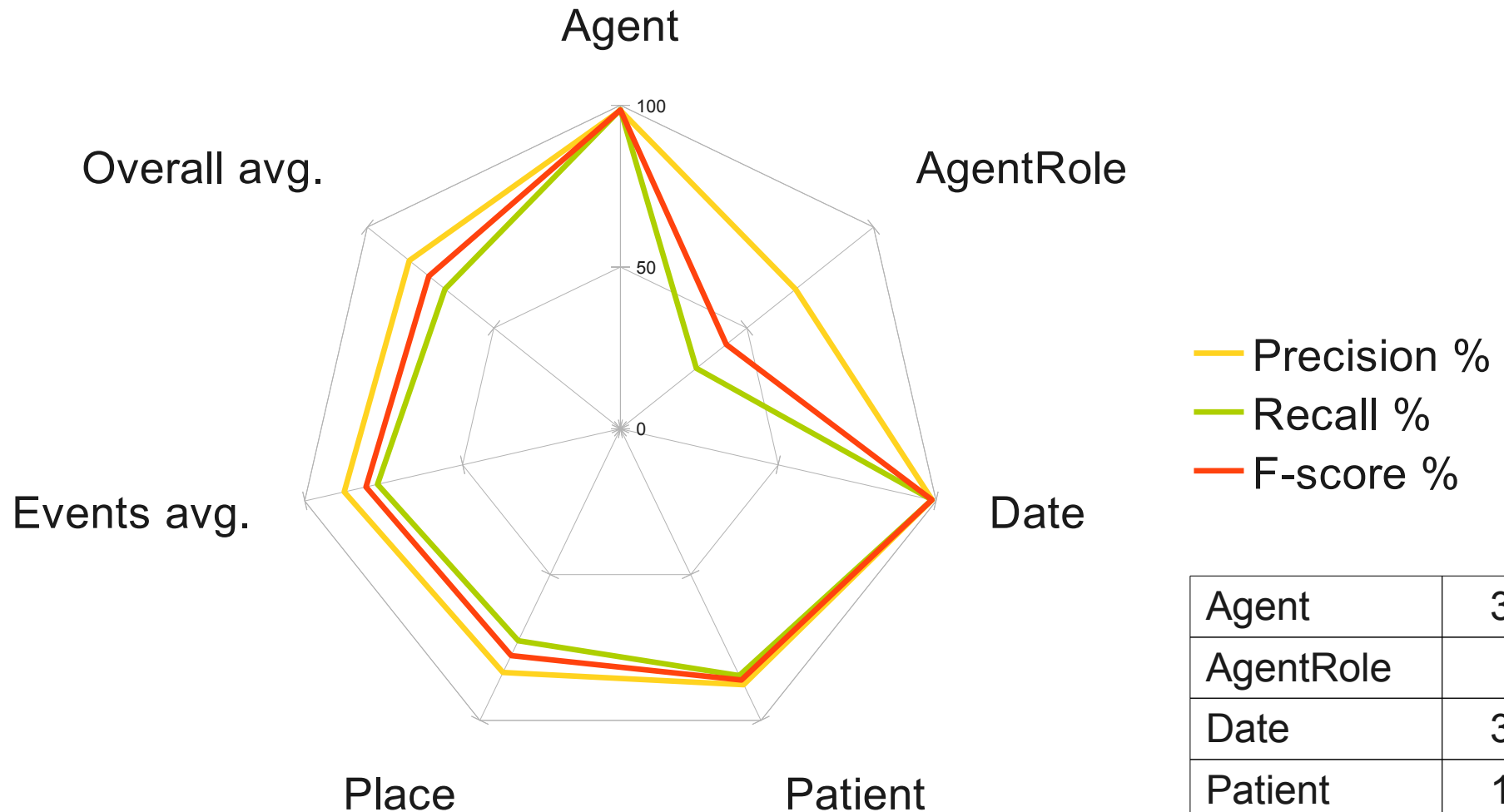


Results: Named Entity Recognition



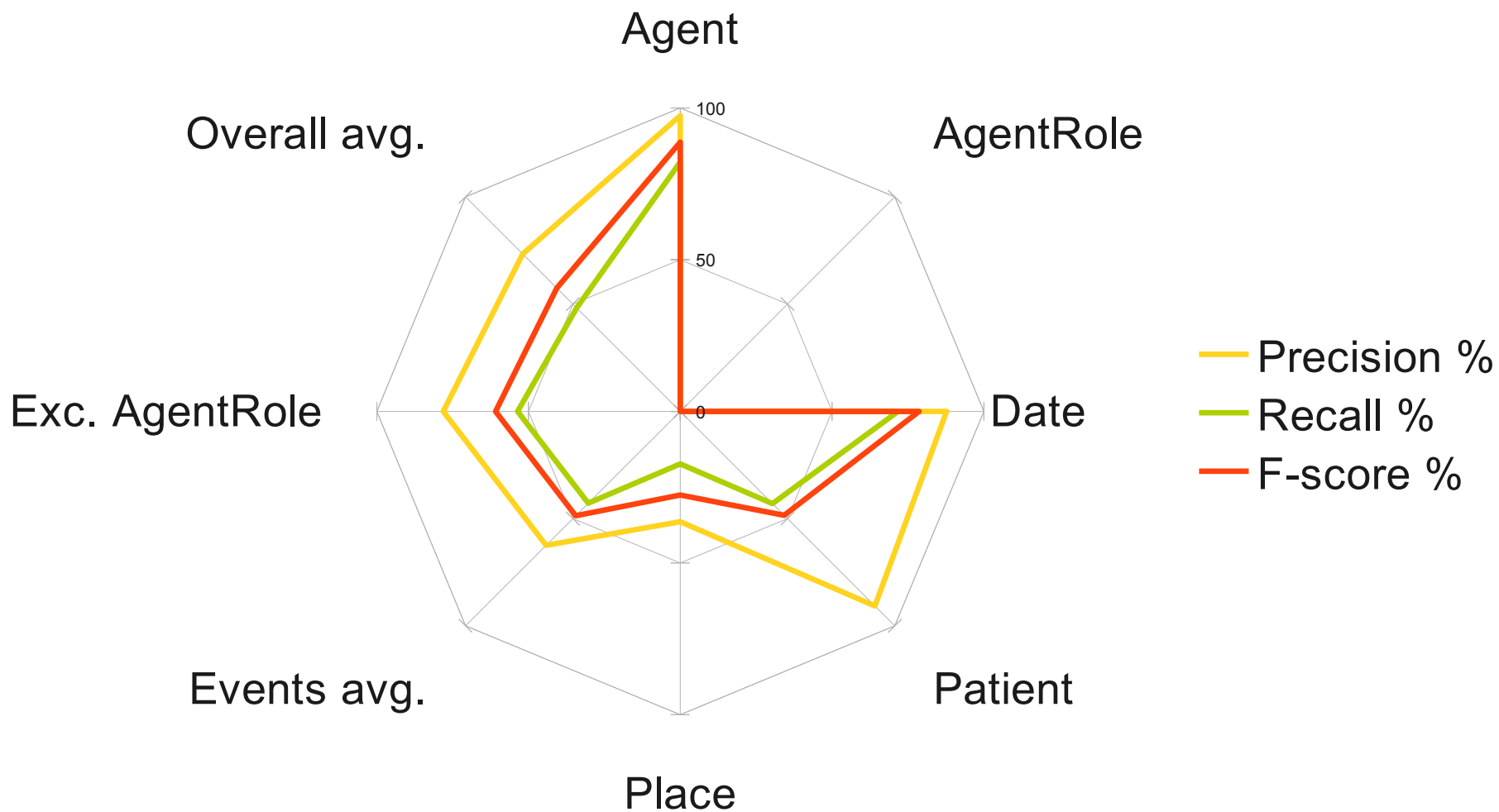
ADDRESS	3458
PLACE	2503
SITENAME	1712
DATE	3519
PERIOD	400
EVENT	3176
ORG	2730
PERSNAME	2318
ROLE	90
SITETYPE	5668
ARTEFACT	879
Total	27453

Results: Relation Extraction

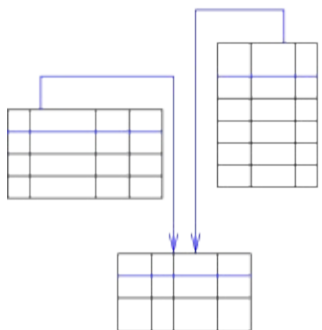


Agent	3794
AgentRole	13
Date	3189
Patient	1553
Place	341
Events	8890
Overall	21932

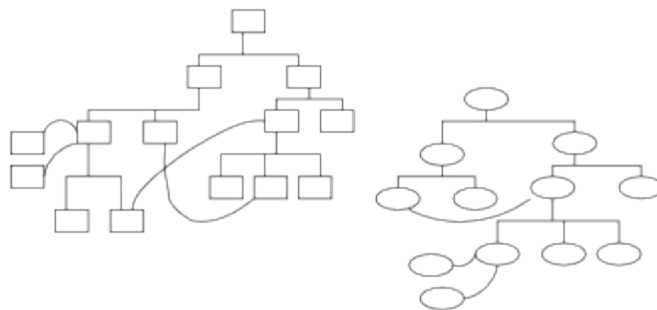
Results: whole pipeline NER + RE



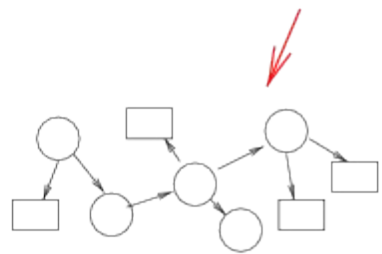
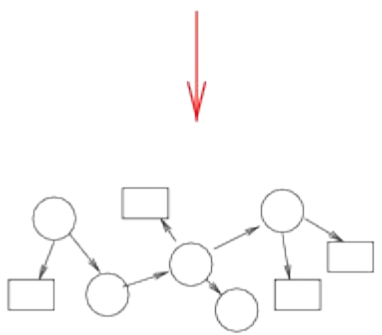
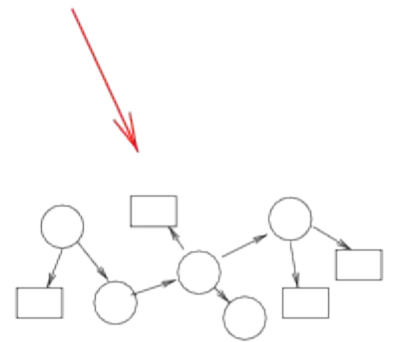
Relational database



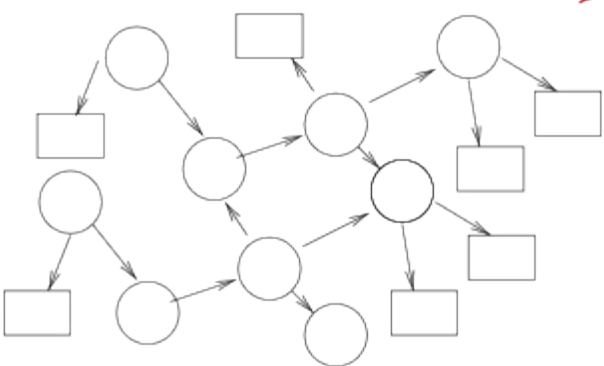
Published domain thesauri



Text documents



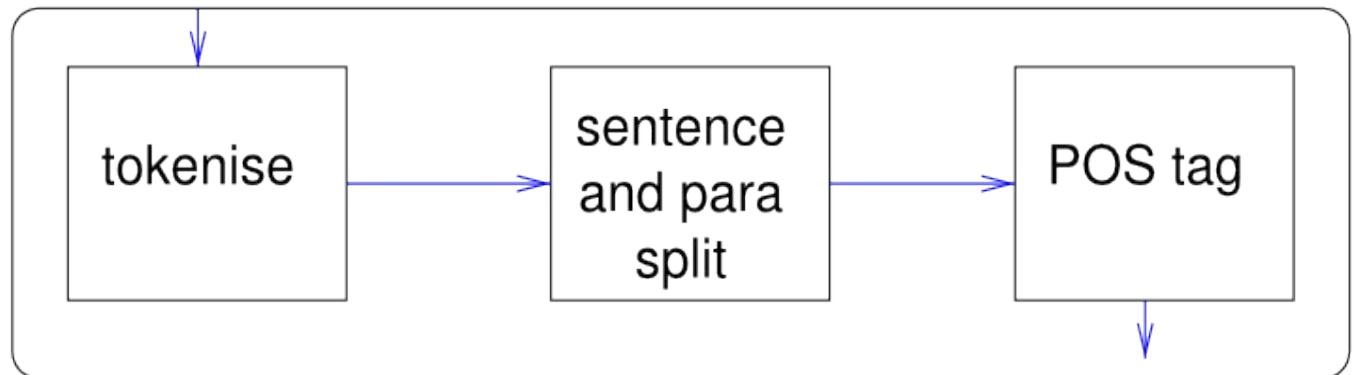
Graph of triples



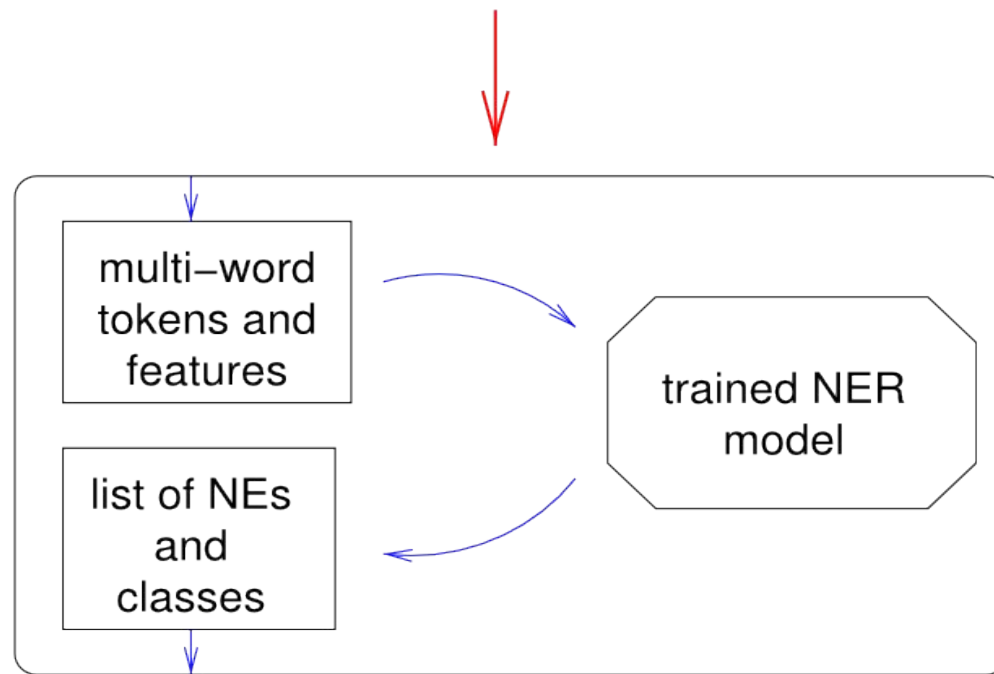
Text documents



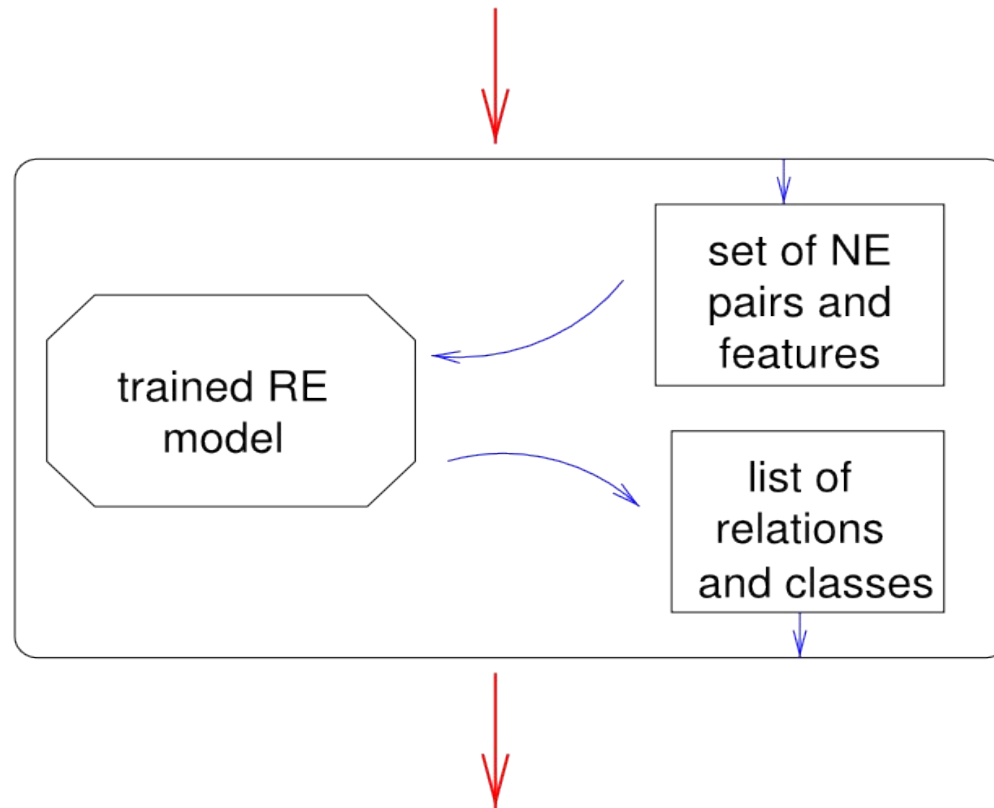
Pre-processing



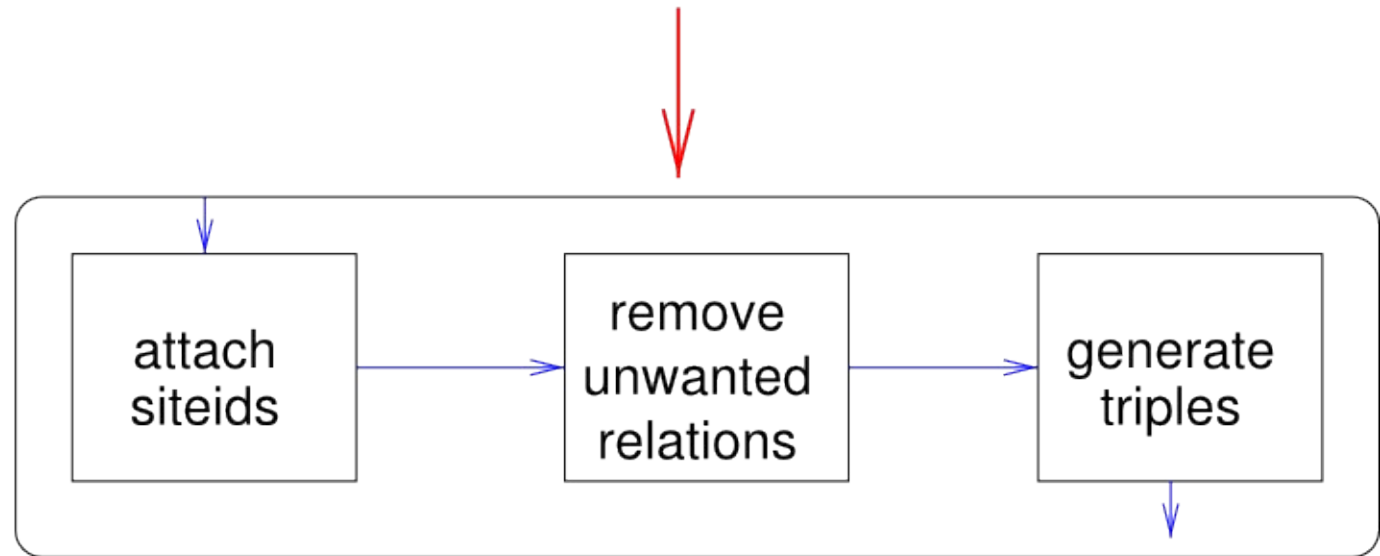
NER



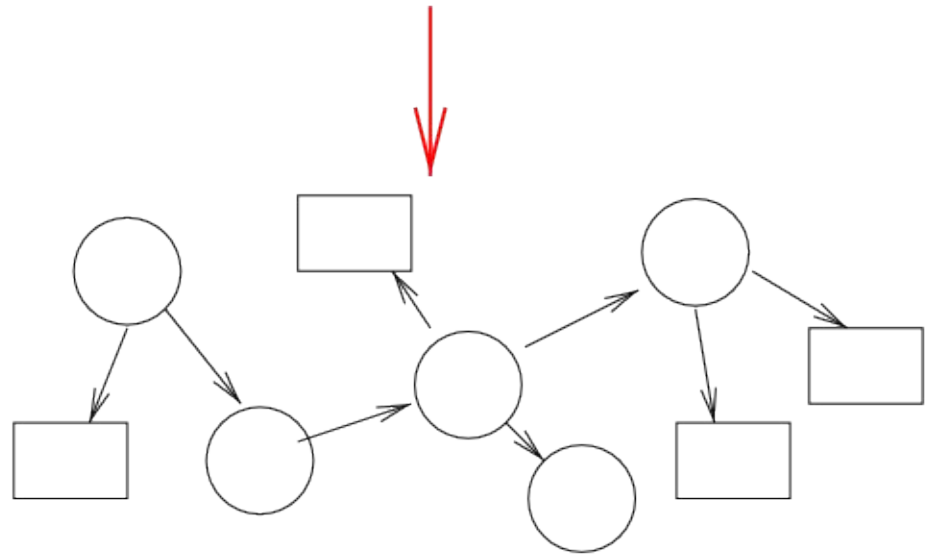
RE



***RDF
translation***



***Graph
of triples***



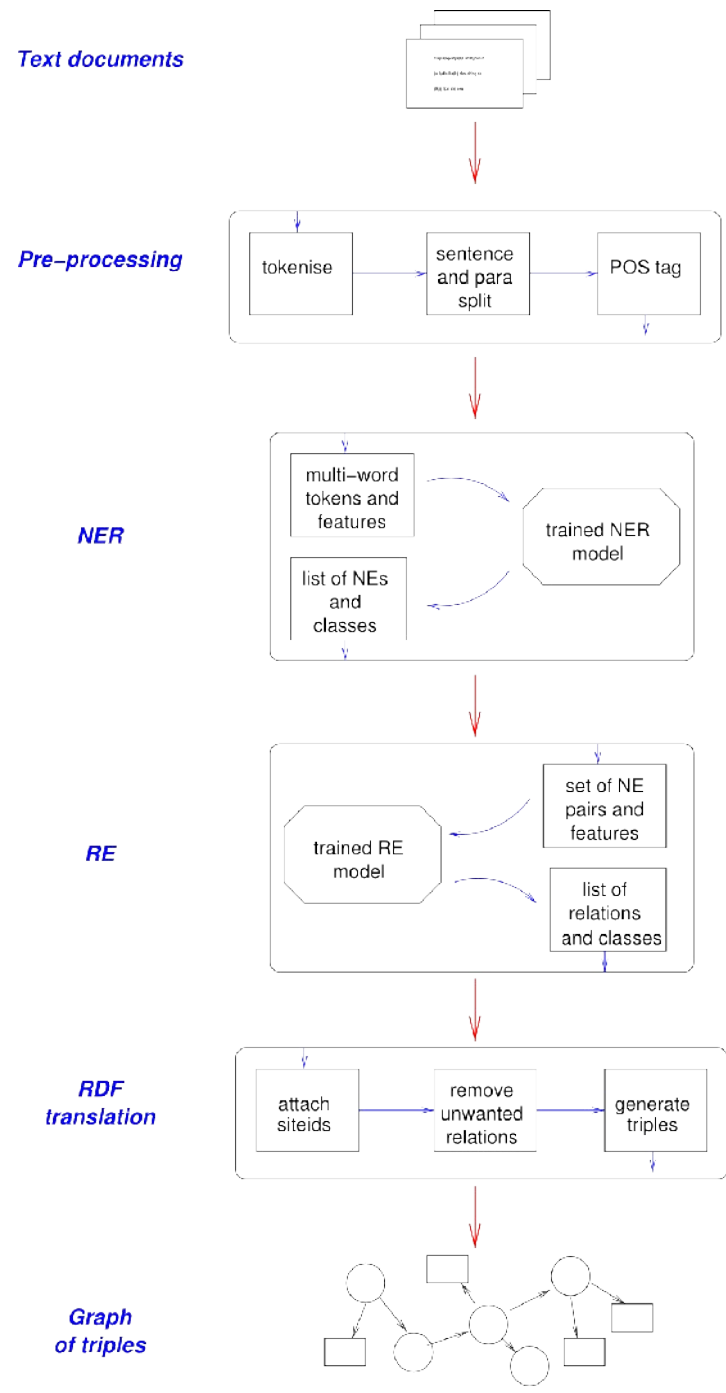


Figure 3: The text to RDF pipeline

Site456

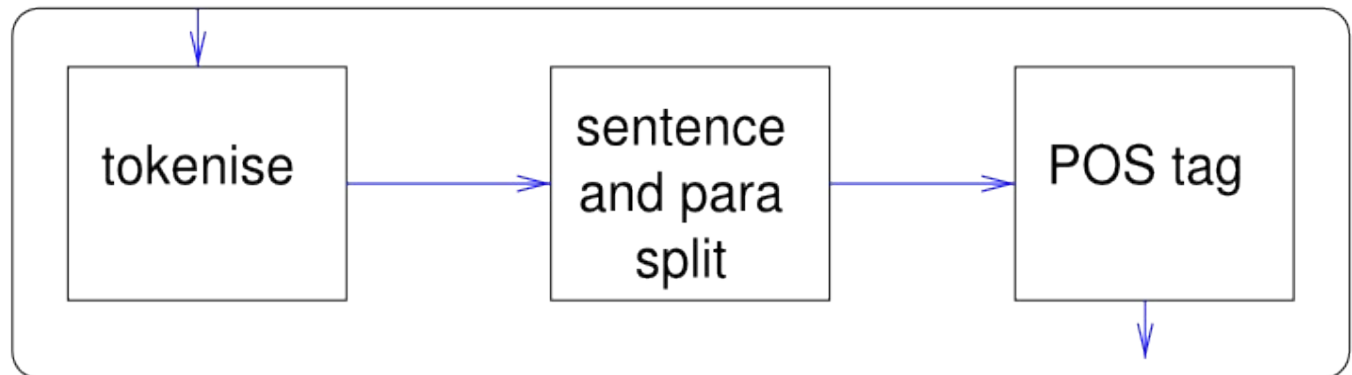
SOUTH WALLS, MISBISTER, THE LOFTS
ND38NW29 centred 3325 8885

Sites recorded during an archaeological survey undertaken on the lands of the Loft, Longhope, as part of the pilot scheme for the Historic Scotland Farm Ancient Monument Survey Grant Scheme. 8885 Cairn. ND 3339 8886 Clearance cairn. ND 3342 8884 Sub-rectangular cairn. ND 3339 8883 Well Sponsors: Historic Scotland, M J Jones. N Card 1998

Text documents



Pre-processing

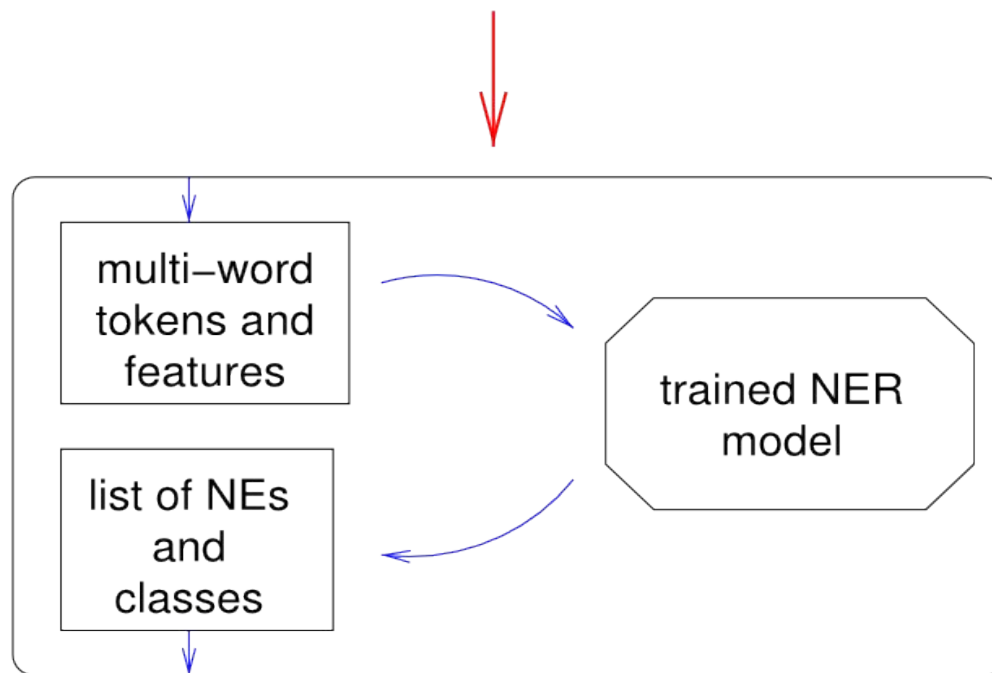


{[SOUTH WALLS]}, {[MISBISTER]}, {[[[THE {LOFTS}]]]}

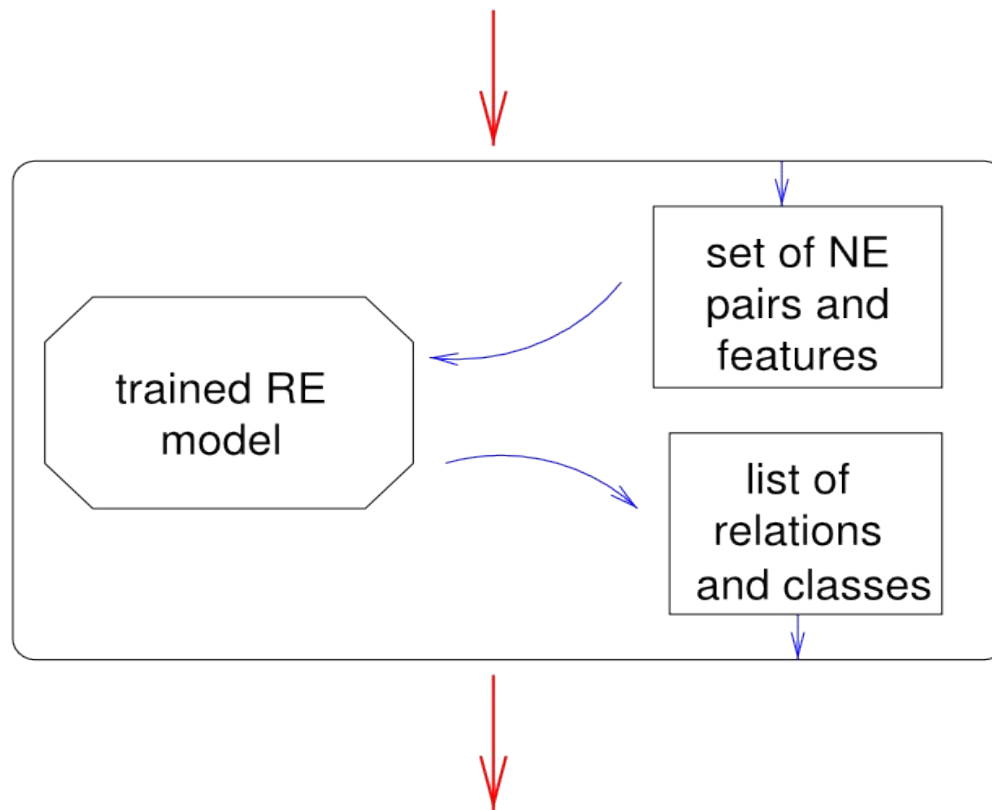
[ND38NW 29 centred 3325 8885]

Sites {[[[recorded]]]} during an ~~archaeological survey~~ undertaken on the lands of {the {Loft}}, [Longhope], as part of the pilot scheme for the {[[[Historic {Scotland}]]]} {Farm} {Ancient} {Monument} Survey Grant Scheme}. [ND 3311 8890] Two [small {cairns}]. [ND 3336 8889] {[Cairn]}. [ND 3339 8885] {[Cairn]}. [ND 3339 8886] {[Clearance cairn]}. [ND 3342 8884] [Sub-rectangular cairn]. [ND 3339 8883] {[Well]} Sponsors : {[Historic {Scotland}]} , {[M J Jones]} . {[[[N Card]]} {[1998]}

NER



RE



Named Entity Recognition

FIND EVENT *PLACE* *PERSNAME* *ARTEFACT*

The following were found in Unst by Mr A T Cluness : a steatite dish , ...

<i>eventLocation</i>	were_found	unst	<i>cls1=event</i>	<i>cls2=place</i>	<i>wdsep=+2...</i>
<i>eventAgent</i>	were_found	a_t_cluness	<i>cls1=event</i>	<i>cls2=persname</i>	<i>wdsep=+5...</i>
<i>eventPatient</i>	were_found	steatite_dish	<i>cls1=event</i>	<i>cls2=artefact</i>	<i>wdsep=+9...</i>
<i>O</i>	unst	a_t_cluness	<i>cls1=place</i>	<i>cls2=persname</i>	<i>wdsep=+9...</i>
<i>O</i>	unst	steatite_dish	<i>cls1=place</i>	<i>cls2=artefact</i>	<i>wdsep=+9...</i>
<i>O</i>	a_t_cluness	steatite_dish	<i>cls1=persname</i>	<i>cls2=artefact</i>	<i>wdsep=+9...</i>

Relation Extraction

	Form	Description
1	ne1=...	first NE string (concatenated using “-”)
2	ne2=...	second NE string
3	cls1=...	first NE type
4	cls2=...	second NE type
5	wdsep= $\circ n$	distance between NEs (+ve or -ve)
6	insent= y or n	both NEs in same sentence?
7	inpara= y or n	both NEs in same paragraph?
8	lastNEwdsame= y or n	normalised last token matches?
9	prevpos1=...	POS tag of token preceeding f irst NE
10	prevpos2=...	POS tag of token preceeding second NE
11	1begsent= y or n	f irst NE is at beginning of a sentence
12	2begsent= y or n	second NE is at beginning of a sentence
13	1endsent= y or n	f irst NE is at end of a sentence
14	2endsent= y or n	second NE is at end of a sentence
15	nest= n , 1in2 or 2in1	one NE is nested within the other
16	neBetw= n	number of NEs between this pair
17	verb=...	if insent= y , (f irst) verb between NEs; else “none”

Table 1: Textual features used for building RE model.

site456

[SOUTH WALLS] , [MISBISTER] , [THE {LOFTS}]

[ND38NW 29 centred 3325 8885] event

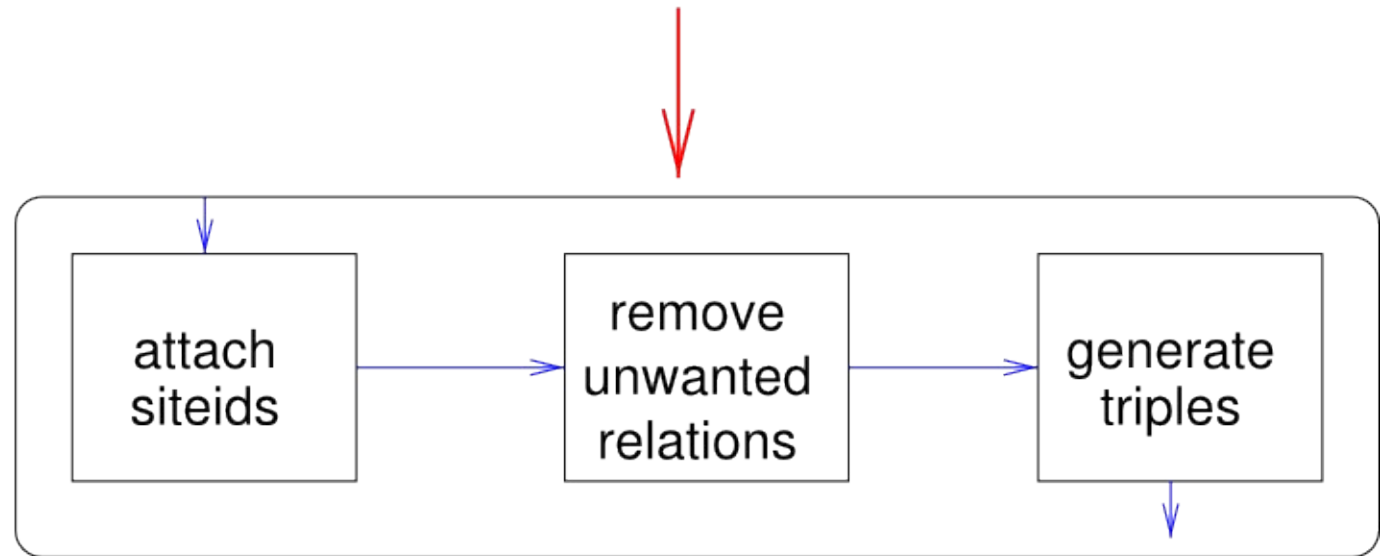
Sites [recorded] during an [archaeological survey] undertaken on the lands of [the {Loft}] , [Longhope] , as part of the pilot scheme for the [Historic {Scotland}] {Farm} {Ancient} {Monument} Survey Grant Scheme] . [ND 3311 8890] Two [small {cairns}] . [ND 3336 8889] [{Cairn}] . [ND 3339 8885] [{Cairn}] . [ND 3339 8886] [{Clearance cairn}] . [ND 3342 8884] [Sub-rectangular cairn] . [ND 3339 8883] [{Well}] Sponsors : [Historic {Scotland}] , [M J Jones] . [N Card] [1998]

eventPatient

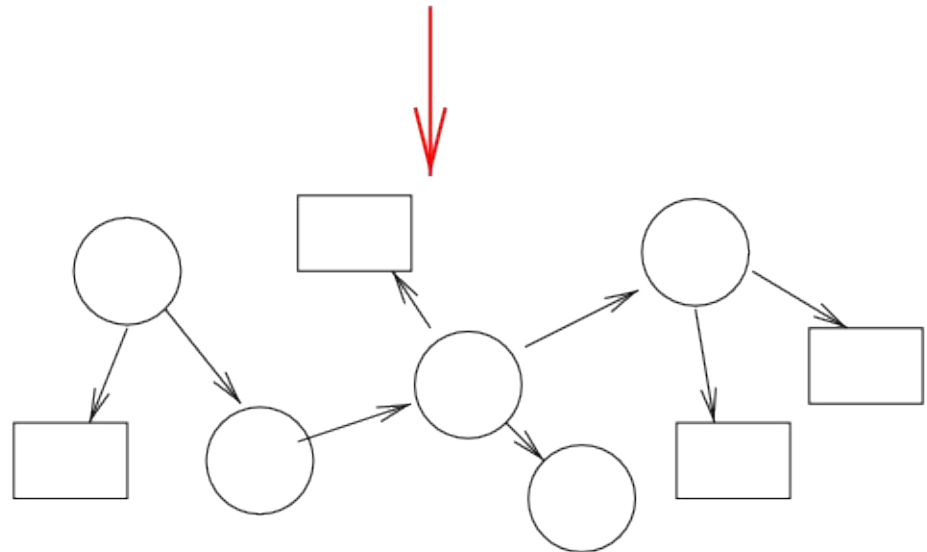
eventPlace

ADDRESS
PLACE
SITENAME
DATE
PERIOD
EVENT
ORG
PERSNAME
ROLE
SITETYPE
ARTEFACT

***RDF
translation***



***Graph
of triples***



site456

[SOUTH WALLS] , [MISBISTER] , [[THE {LOFTS}]]

[ND38NW 29 centred 3325 8885] event

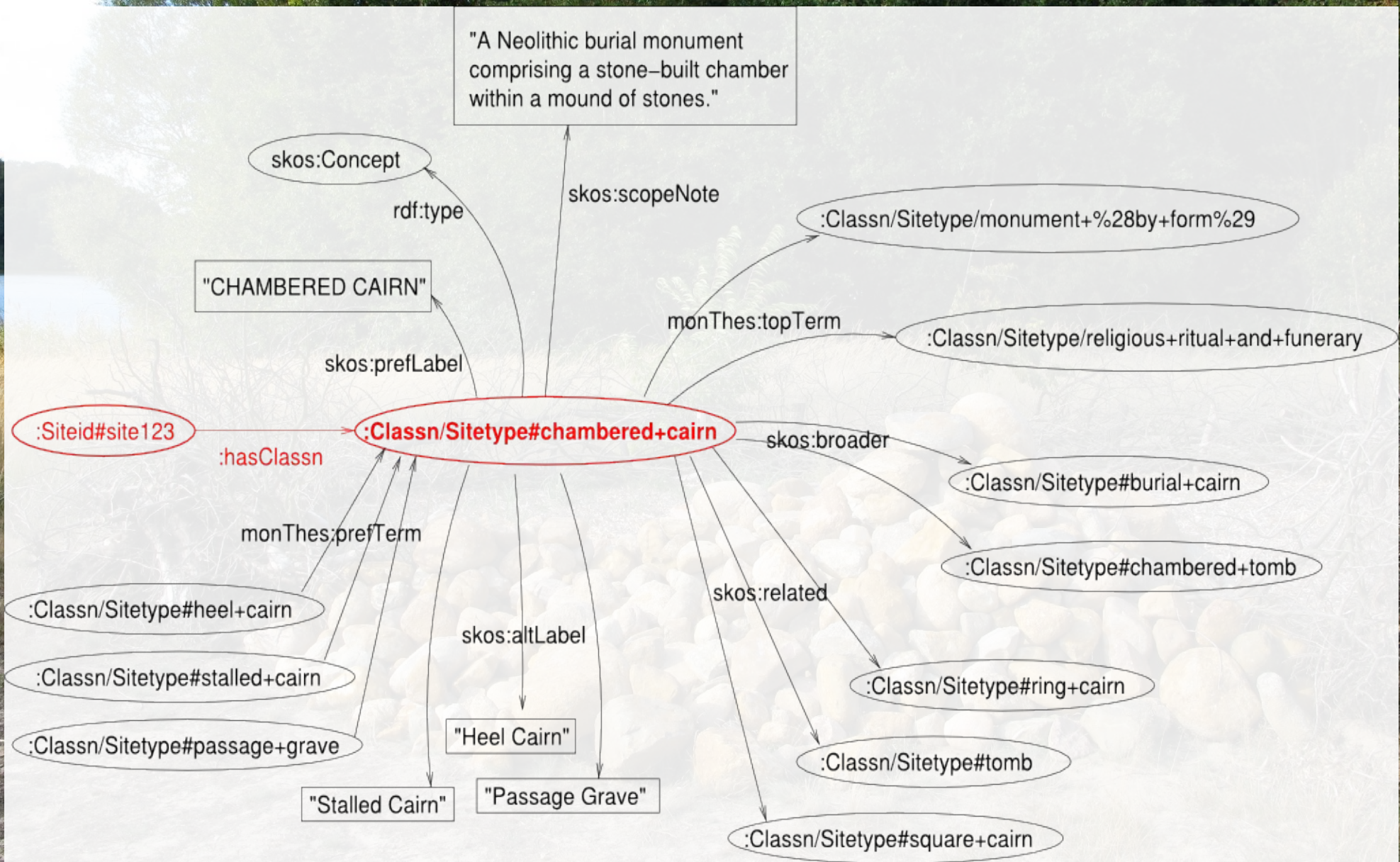
Sites [recorded] during an [archaeological survey] undertaken on the lands of [the {Loft}] , [Longhope] , as part of the pilot scheme for the [Historic {Scotland}] {Farm} {Ancient} {Monument} Survey Grant Scheme] . [ND 3311 8890] Two [small {cairns}] . [ND 3336 8889] [{Cairn}] . [ND 3339 8885] [{Cairn}] . [ND 3339 8886] [{Clearance cairn}] . [ND 3342 8884] [Sub-rectangular cairn] . [ND 3339 8883] [{Well}] Sponsors : [Historic {Scotland}] , [M J Jones] . [[N Card]] [{1998}]

eventPatient

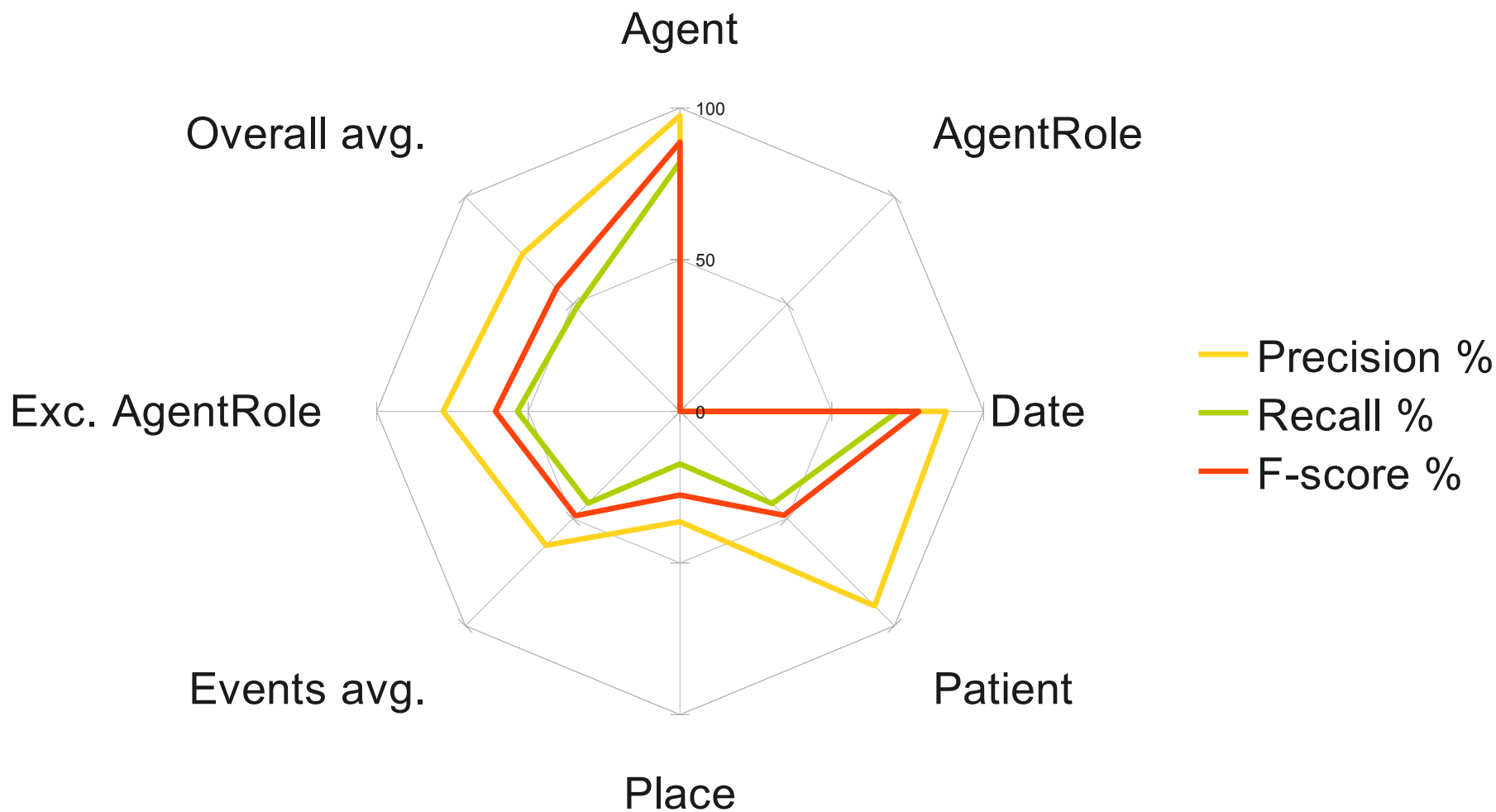
eventPlace

site456 – hasEvent – recordingX
recordingX – hasLocation – "ND 3342 8884"
recordingX – hasPatient – "Sub-rectangular cairn"

Figure 5: Mapping relations to RDF.



Results: whole pipeline NER + RE



Questions?



Image sources

- The standing stones of Callanish on the Isle of Lewis, Scotland.



- <http://commons.wikimedia.org/wiki/File:Callanish-circle.jpg>

- Clearance cairn at landscape conservation area Königswald (Potsdam).



- http://commons.wikimedia.org/wiki/File:Lesesteinhaufen,_K%C3%B6nigswald,_Potsdam.jpg

- Memorial cairn at the Isandlwana hill in Kwazulu Natal, South Africa.



- <http://en.wikipedia.org/wiki/File:Isandlwanamassgrave.JPG>