

Development of a Universal Dependencies treebank for Welsh

Johannes Heinecke¹, Francis M. Tyers²

(1) Orange Labs, TGI/Data & IA/DESKIÑ

(2) Department of Linguistics, Indiana University

Celtic Languages Technology Workshop 2019



Table of Contents

- 1 Welsh language
- 2 Welsh resources
- 3 Treebank annotation
- 4 Annotation examples
 - Periphrastic verbal constructions
 - Nonverbal predicates
 - Impersonals
 - Inflected prepositions
 - Compound numbers
- 5 Statistics
- 6 Evaluation
 - Results
 - Errors
- 7 Conclusion

Welsh language

Welsh resources

Treebank annotation

Annotation examples

Statistics

Evaluation

Conclusion

Like Breton or Irish

- initial consonant mutation
- inflected prepositions *ar* “on” *arnaf fi* “on me”
- genitive construction with single determiner (cf. Arabic):
 - *tŷ'r brenin* “the house of the king”
 - vs. *tŷ brenin* “a king’s house”
- VSO basic word order
- composed numbers
 - *tri phlentyn ar ddeugain* “three child on two twenty” → “43 children”
 - *unfed ganrif ar hugain* “first century on twenty” → “21st century”
- impersonal forms (including intransitive verbs) *gwelwyd*, *cysgir* “one saw, one will sleep”
- possession expressed with a prepositional phrase: *Mae arian gen i* “there is money with me”
- singulatives (add suffix (plus Umlaut) to have singular from basic plural):
 - *plant* “children”, *adar* “birds”, *llygod* “mice”, *caws* “cheese”
 - *plentyn* “child”, *aderyn* “bird”, *llygoden* “mouse”, *cosyn* “piece of cheese”

Welsh language

Welsh resources

Treebank annotation

Annotation examples

Statistics

Evaluation

Conclusion

Unlike Breton

- no auxiliary “to have”, thus no composed tenses with auxiliary verbs
- no participles (but verbal adjectives)
- periphrastic constructions with *bod* “to be” and tense-aspect-markers for most tenses and aspects
- pronouns not subject/object (NOM/ACC) but independent (subject position)/dependent (possessives and object position)
- no infinitives but verbnouns (direct object marked differently on verbnouns than on verbs)
 - patient marked in the same way as possessives using dependent pronouns (cf. *dy dŷ* [ti])
 - different to inflected verb forms, were nominal direct objects undergo soft mutation, and independent pronouns are used for pronominal direct objects.
 - cf. German:

inflected:	<i>ich sehe den Hund</i> _{ACC}	“I see the dog”
infinitive:	<i>den Hund sehen</i> _{ACC}	“to see the dog”
verbnoun:	<i>das Sehen des Hundes</i> _{GEN}	“the sight/seeing of the dog”
 - Welsh

inflected:	<i>Mi welodd o ti</i> _{indep}	“He saw you” (lit. “(aff) saw he you”)
verbnoun:	<i>Roedd o’n dy</i> _{dep} <i>weld</i>	“He was seeing you” (lit. “Was he (impf) your seeing”)

[Welsh language](#)[Welsh resources](#)[Treebank annotation](#)[Annotation examples](#)[Statistics](#)[Evaluation](#)[Conclusion](#)

Research on syntax

- important work on syntax using frameworks like LFG, GB, HPSG
- orthographic correction
- automatic speech recognition (ASR)
- speech synthesis (TTS)
- natural language understanding (NLU)

Welsh language

Welsh resources

Trebank annotation

Annotation examples

Statistics

Evaluation

Conclusion

(Publicly available) data

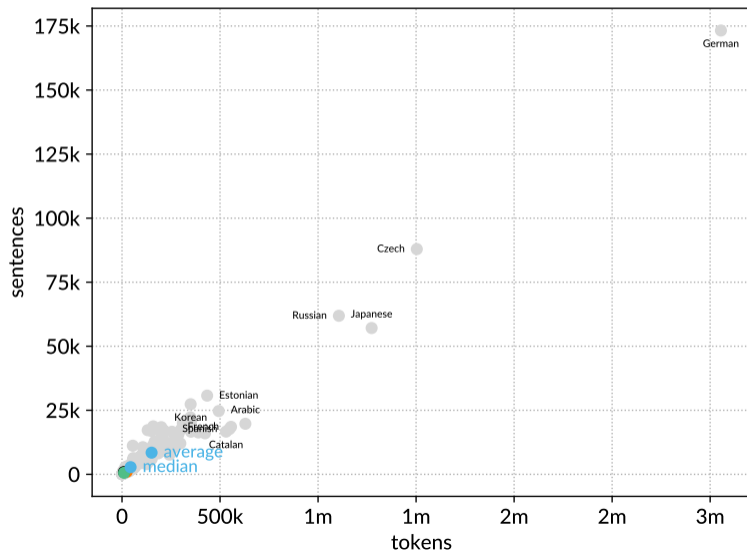
- *Cronfa Electroneg o Gymraeg*: 1 000 000 tokens annotated with lemmas and POS
- National Corpus of Contemporary Welsh
- Wikipedia (+100 000 pages), word embeddings (fastText, BERT)
- Wictionary (and Unimorph) very little data
- *Eurfa* full form dictionary (210 000 forms, 10 000 lemmas and English glosses)
- parallel Welsh/English corpus from Welsh Assembly

The Welsh language in UD

UD Welsh

Johannes Heinecke,
Francis M. Tyers

CLTW 2019



Welsh language

Welsh resources

Treebank annotation

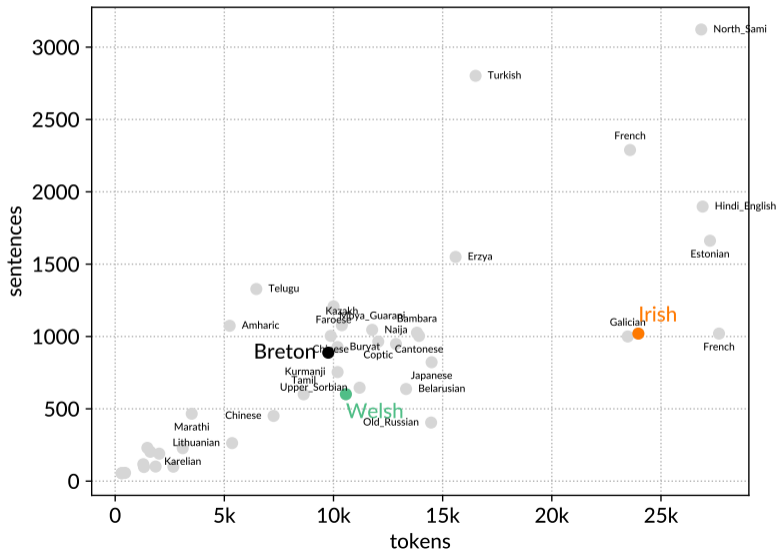
Annotation examples

Statistics

Evaluation

Conclusion

The Welsh language in UD



Welsh language

Welsh resources

Treebank annotation

Annotation examples

Statistics

Evaluation

Conclusion

Corpus of the Treebank

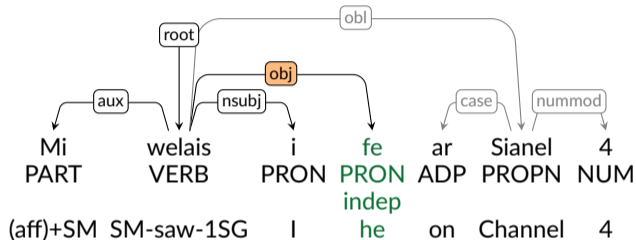
10 756 tokens, 601 sentences

- shortest sentence: 4 words, longest: 59, average length: 17.9, median length: 16
- Wikipedia (pages on items of Wales)
- Welsh assembly corpus
- Media (BBC Cymru, Y Golwg)
- Web sites of Welsh universities and organisations (*Cymdeithas yr Iaith*, *Urdd Gobaith Cymru*, county councils)
- Welsh language blogs
- Welsh Grammars and Textbooks

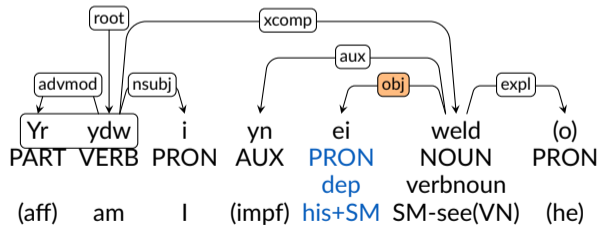
(pre)processing

- UD-isation of CEG and training UDpipe with CEG and external dictionary (*Eurfa*)
- POS annotation/lemmatisation of treebank sentences (using model trained on CEG)
- manual validation of lemmas, XPOS, UPOS and annotation of dependency relations
- validation scripts (checking (some) features, adding mutation type etc.)

Periphrastic verbal constructions

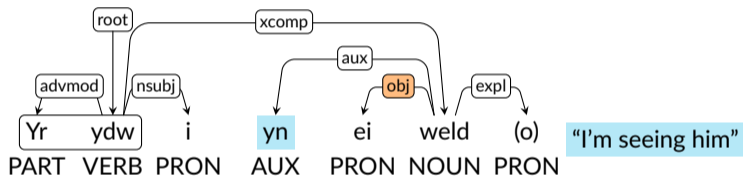


"I saw him on Channel 4"

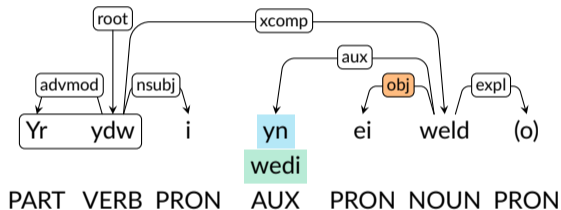


"I am seeing him"

Periphrastic verbal constructions: Tense-aspect-markers



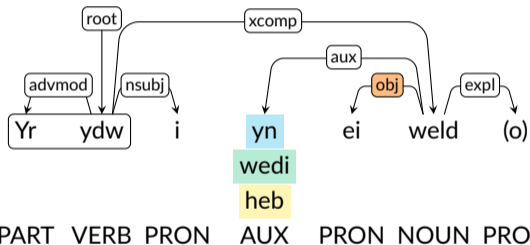
Periphrastic verbal constructions: Tense-aspect-markers



“I’m seeing him”

“I have seen him” (“I’m after seeing him”)

Periphrastic verbal constructions: Tense-aspect-markers

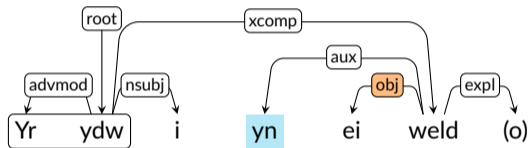


“I’m seeing him”

“I have seen him” (“I’m after seeing him”)

“I have not seen him” (“I’m without seeing him”)

Periphrastic verbal constructions: Tense-aspect-markers



PART VERB PRON AUX PRON NOUN PRON

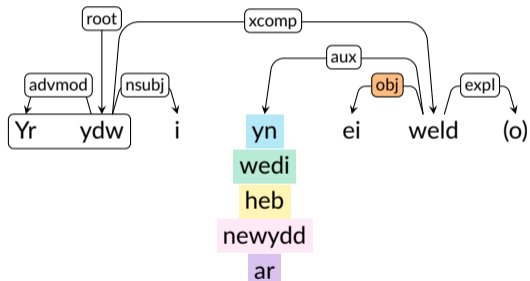
“I’m seeing him”

“I have seen him” (“I’m after seeing him”)

“I have not seen him” (“I’m without seeing him”)

“I have just seen him” (“I’m new seeing him”)

Periphrastic verbal constructions: Tense-aspect-markers



PART VERB PRON AUX PRON NOUN PRON

“I’m seeing him”

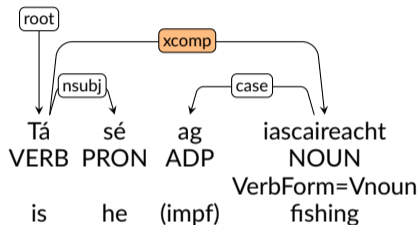
“I have seen him” (“I’m after seeing him”)

“I have not seen him” (“I’m without seeing him”)

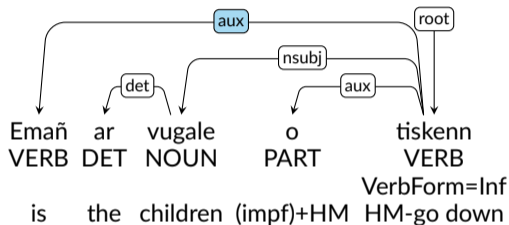
“I have just seen him” (“I’m new seeing him”)

“I’m about to see him” (“I’m on seeing him”)

Periphrastic verbal constructions in Irish and Breton



“He is fishing”

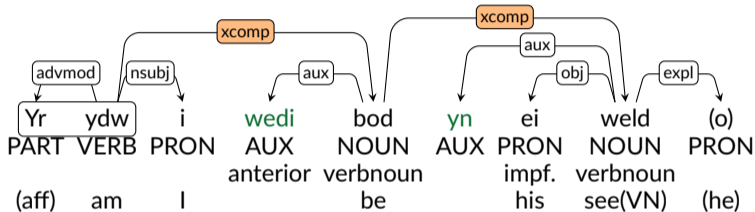


“The children are going down.”

- The Welsh treebank follows the Irish example in making “to be” head and attach the verb as *xcomp*
- Difference: Irish attaches TAM (ADP) as *case*, Welsh attaches TAM (AUX) as *aux*:
Harmonisation needed

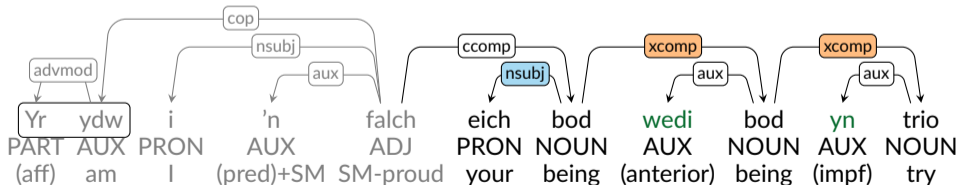
Nested periphrastic verbal constructions

■ Anterior Present, Imperfective



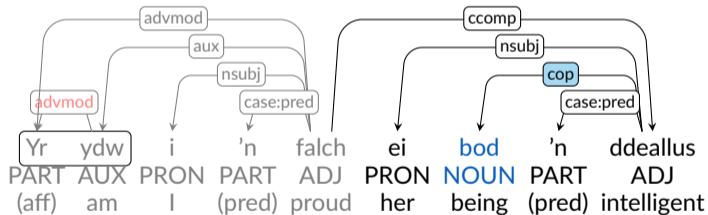
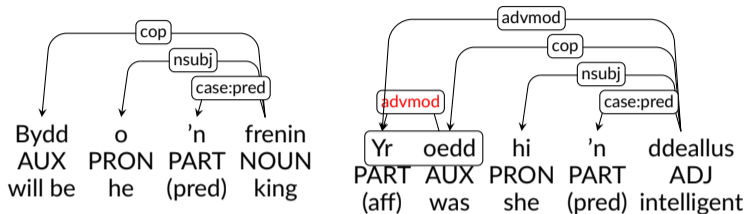
“I have been seeing him”

■ Subordinate (no finite Tense)



“I’m proud that you have been trying”

Nonverbal predicates



"I'm proud that she is intelligent"

■ N.B. *yn*: 3 syntactically distinct homographs:

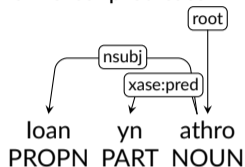
- imperfective TAM (before verbnouns, no mutation triggered)
- predicative marker (before nouns and adjectives, triggers soft mutation)
- preposition "in" (triggers nasal mutation)

Comment:

Copula can be dropped:
Liz yn frenhines "Liz
(is) queen"

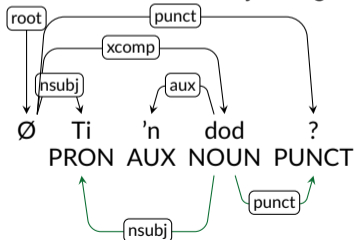
Constructions without *bod*

■ Nonverbal predicate



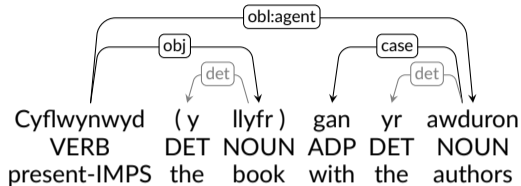
“loan (is) teacher”

■ Verbnoun: Head of subject is gone

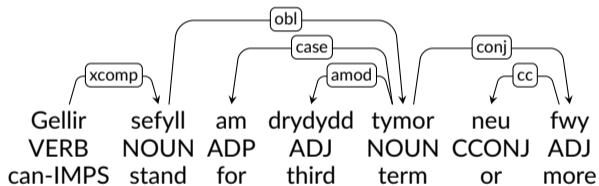


“Do you come?”

Impersonal Forms, *cael*

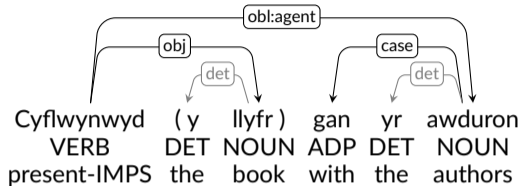


“The book was presented by the authors”
(lit. “One presented the book by the authors”)

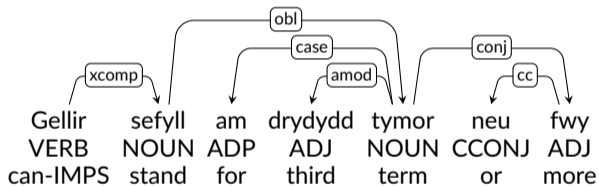


“One can stand for a third term or more”

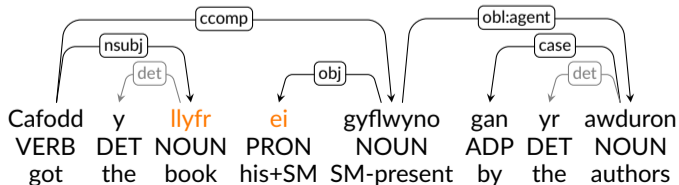
Impersonal Forms, *cael*



“The book was presented by the authors”
(lit. “One presented the book by the authors”)



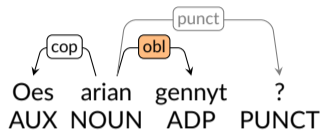
“One can stand for a third term or more”



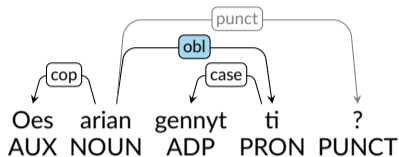
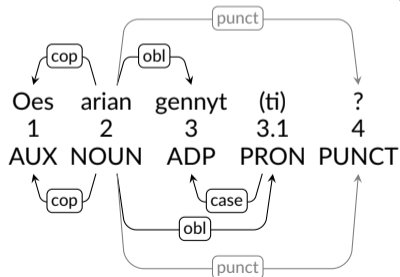
“The book was presented by the authors”
(lit. “got the book his presenting by the authors”)

[Welsh language](#)[Welsh resources](#)[Treebank annotation](#)[Annotation examples](#)[Periphrastic verbal constructions](#)[Nonverbal predicates](#)[Impersonals](#)[Inflected prepositions](#)[Compound numbers](#)[Statistics](#)[Evaluation](#)[Conclusion](#)

Inflected prepositions



“Do you have money?”
(lit. “Is money with-2SG?”)

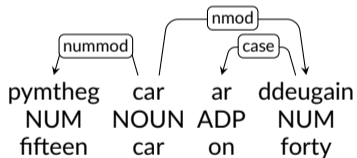


“Do you have money?”
(lit. “Is money with-2SG you-SG?”)

Compound numbers

■ vigesimal system

- 20 = *ugain*
- 30 = *deg ar hugain* “ten on twenty”
- 40 = *deugain* “two twenties”
- 60 = *trigain* “three twenties”



“55 cars” (lit. “15 cars on 2*20”)

Statistics (POS)

- 10756 tokens, 601 sentences

UPOS	%	XPOS	%	XPOS	%	XPOS	%
NOUN	30.1	noun	21.5	aff	1.3	ord	0.1
ADP	12.9	prep	12.5	person	1.3	refl	0.1
PUNCT	9.7	punct	9.7	num	1.2	intr	0.1
ADJ	6.9	verbnoun	9.1	ante	1.1	sym	0.1
DET	6.5	art	6.5	neg	0.9	int	0.0
PRON	6.3	verb	6.3	dem	0.6	work	< 0.0
VERB	6.3	pos	6.0	cprep	0.6	perf	< 0.0
AUX	5.9	cconj	2.9	sup	0.5	contr	< 0.0
PART	4.4	dep	2.7	sconj	0.5	card	< 0.0
PROPN	3.7	impf	2.5	rel	0.4	propn	< 0.0
CCONJ	2.9	indep	2.4	cmp	0.3		
ADV	1.9	place	2.3	pron	0.2		
NUM	1.3	pred	2.1	org	0.1		
SCONJ	0.5	adv	2.0	post	0.1		
SYM	0.1	aux	1.9	eq	0.1		

[Welsh language](#)[Welsh resources](#)[Treebank annotation](#)[Annotation examples](#)[Statistics](#)[Evaluation](#)[Conclusion](#)

Statistics (dependency relations)

- 10756 tokens, 601 sentences

■	<i>deprel</i>	%	<i>deprel</i>	%	<i>deprel</i>	%
	case	10.5	mark	2.1	csubj	0.1
	punct	9.7	cop	2.0	nmod:agent	0.1
	nmod	8.4	ccomp	1.7	compound	< 0.0
	det	6.9	nmod:poss	1.7	iobj	< 0.0
	obl	6.0	acl	1.6		
	nsubj	5.7	advcl	1.4		
	root	5.6	acl:relcl	1.2		
	obj	5.2	flat:name	0.8		
	advmod	5.2	flat	0.6		
	amod	4.8	nummod	0.6		
	xcomp	4.3	fixed	0.5		
	aux	3.7	appos	0.5		
	conj	3.2	expl	0.2		
	cc	2.9	obl:agent	0.2		
	case:pred	2.1	parataxis	0.2		

Welsh language

Welsh resources

Treebank annotation

Annotation examples

Statistics

Evaluation

Conclusion

Evaluation

- train: 80%, dev: 10%, test: 10%
- 10-fold cross-evaluation
- tagging and lemmatisation

	UPOS	XPOS	Lemma
baseline	89.2	87.3	86.7
+ <i>Eurfa</i>	87.9	87.5	93.5

- dependency parsing

		UAS	LAS	CLAS
tag + parse	baseline	74.3	63.9	54.8
	+ <i>Eurfa</i>	75.5	64.3	55.4
parse on gold tags	baseline	82.2	76.2	69.6
	+ <i>Eurfa</i>	81.9	75.9	69.3

- No increase in UAS, LAS, CLAS with word embeddings

Errors: UPOS

UPOS	errors	all	%	wrong UPOS
PUNCT	0	1039	0.0	
DET	7	698	1.0	PART:6 NOUN:1
ADP	91	1389	6.6	NOUN:25 PART:24 CCONJ:12 AUX:11 ADV:5 PRON:5 ADJ:4 SCONJ:2 VERB:2 PROPEN:1
CCONJ	25	312	8.0	ADP:15 PRON:4 NOUN:2 SCONJ:2 PROPEN:1 ADJ:1
NOUN	290	3242	8.9	ADJ:81 PROPEN:73 ADP:34 ADV:30 VERB:19 AUX:17 PRON:16 NUM:9 CCONJ:7 PART:2 SCONJ:1 DET:1
NUM	13	134	9.7	NOUN:8 PRON:3 ADJ:1 PROPEN:1
PRON	68	676	10.1	ADP:23 CCONJ:23 NOUN:10 PART:5 ADJ:4 AUX:2 VERB:1
PART	68	473	14.4	AUX:30 ADP:16 DET:16 PRON:4 NOUN:1 CCONJ:1
SCONJ	7	48	14.6	CCONJ:3 ADP:3 NOUN:1
VERB	111	674	16.5	AUX:50 NOUN:40 ADP:10 ADJ:7 PRON:2 PROPEN:1 CCONJ:1
ADJ	145	738	19.6	NOUN:98 PROPEN:11 ADP:9 VERB:9 ADV:8 PRON:4 AUX:3 SCONJ:1 NUM:1 CCONJ:1
AUX	126	631	20.0	VERB:80 ADP:17 NOUN:15 PART:12 ADJ:2
ADV	52	206	25.2	NOUN:27 ADJ:8 VERB:6 ADP:3 PROPEN:2 PART:2 PRON:2 CCONJ:2
SYM	2	6	33.3	PUNCT:1 NOUN:1
PROPEN	148	396	37.4	NOUN:109 ADJ:21 ADP:5 ADV:4 PART:2 PRON:2 NUM:2 AUX:2 CCONJ:1
total:	1153	10756	10.7	

[Welsh language](#)[Welsh resources](#)[Treebank annotation](#)[Annotation examples](#)[Statistics](#)[Evaluation](#)[Results](#)[Errors](#)[Conclusion](#)

Errors: dependency relations I

■ Parsing on gold UPOS

deprel	errors	all	%	wrong deprel
case:pred	0	227	0.0	
punct	0	1039	0.0	
aux	5	396	1.3	cop:2 case:1 advmod:1 case:pred:1
det	11	736	1.5	nmod:8 obl:1 advmod:1 nsubj:1
cc	6	311	1.9	mark:6
case	40	1121	3.6	mark:31 advmod:4 fixed:2 obl:1 obj:1 nsubj:1
amod	23	511	4.5	advmod:10 obl:3 det:2 nmod:2 root:2 obj:1 nsubj:1 conj:1 flat:name:1
nmod:poss	11	179	6.1	obj:4 nmod:4 obl:1 expl:1 nsubj:1
cop	19	208	9.1	root:8 acl:relcl:3 advmod:2 advcl:2 aux:2 acl:1 punct:1
advmod	52	554	9.4	amod:26 case:10 root:4 ccomp:3 obj:2 nsubj:2 obl:1 cc:1 appos:1 advcl:1 acl:1
obj	70	556	12.6	obl:23 nsubj:15 nmod:14 ccomp:4 xcomp:3 flat:name:3 root:2 case:1 cc:1 nummod:1 acl:relcl:1 conj:1 nmod:poss:1
nmod	122	899	13.6	obl:56 acl:15 flat:13 nsubj:7 conj:6 obj:3 flat:name:3 appos:3 nmod:poss:3 mark:2 acl:relcl:2 nummod:2 xcomp:1 root:1 advmod:1 det:1 case:1 amod:1 obl:agent:1
xcomp	68	461	14.8	acl:29 ccomp:12 advcl:9 nmod:7 obl:4 obj:3 root:2 amod:1 nsubj:1
mark	39	226	17.3	case:30 cc:8 advmod:1
root	109	601	18.1	acl:19 nsubj:15 nmod:12 conj:12 acl:relcl:9 cop:8 ccomp:8 advcl:7 amod:4 advmod:4 obj:3 appos:3 xcomp:2 obl:1 cc:1 flat:name:1

Errors: dependency relations II

nsubj	150	609	24.6	obj:60 nmod:23 root:20 xcomp:9 obl:8 mark:6 conj:5 amod:5 nmod:poss:3 appos:2 advmod:2 ccomp:2 case:1 det:1 flat:name:1 expl:1 nummod:1
conj	85	345	24.6	nmod:35 appos:8 acl:8 advmod:5 ccomp:5 acl:relcl:5 amod:4 obj:3 nsubj:3 xcomp:2 obl:2 cop:1 root:1 advcl:1 nmod:poss:1 nummod:1
nummod	16	64	25.0	nmod:8 obl:3 obj:2 case:1 mark:1 conj:1
acl:relcl	35	132	26.5	conj:8 acl:8 root:5 nmod:4 xcomp:3 obl:2 advcl:2 cop:1 parataxis:1 amod:1
flat:name	23	86	26.7	nmod:15 flat:2 appos:2 nummod:2 obj:1 punct:1
ccomp	77	181	42.5	xcomp:17 advcl:15 nmod:8 obl:7 obj:7 advmod:5 acl:4 amod:3 conj:3 root:3 cop:2 acl:relcl:2 nsubj:1
obl	277	639	43.3	nmod:149 obj:39 root:31 case:10 conj:10 xcomp:8 nummod:7 nsubj:5 acl:5 advcl:4 appos:2 mark:2 ccomp:2 flat:name:2 advmod:1
acl	78	167	46.7	nmod:26 xcomp:14 advcl:10 ccomp:8 acl:relcl:6 conj:4 obl:3 root:3 amod:2 obj:1 nsubj:1
obl:agent	15	22	68.2	obl:7 nmod:5 acl:1 root:1 nsubj:1
appos	39	57	68.4	nmod:12 conj:12 obl:4 nsubj:3 acl:2 flat:name:2 nummod:1 advmod:1 advcl:1 mark:1
fixed	41	56	73.2	case:24 obl:6 advmod:3 nmod:3 cc:2 conj:2 acl:1
flat	52	66	78.8	nmod:50 flat:name:2
advcl	120	150	80.0	acl:28 root:22 obl:15 xcomp:11 nmod:11 advmod:10 ccomp:9 conj:8 amod:4 obj:1 acl:relcl:1
expl	20	24	83.3	nsubj:9 nmod:6 obj:4 obl:1
parataxis	16	17	94.1	nmod:5 acl:3 ccomp:2 xcomp:1 obl:1 nsubj:1 acl:relcl:1 advcl:1 conj:1
iobj	1	1	100.0	nmod:1
compound	5	5	100.0	nmod:3 advmod:1 amod:1
nmod:agent	6	6	100.0	obl:2 nmod:2 obl:agent:1 appos:1
csubj	10	10	100.0	root:4 acl:3 xcomp:1 conj:1 nmod:1
total:	1641	10756	15.3	

Conclusion

So far ...

- First Welsh treebank, third Celtic language treebank in UD
- 10 000 tokens, allows a relatively robust POS tagging (89.2%), and parsing of simpler sentences

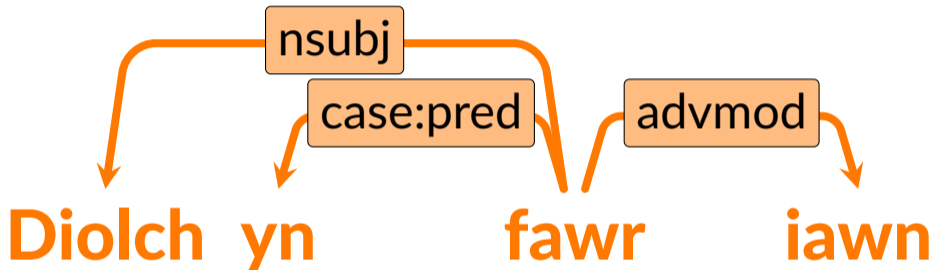
Conclusion

So far ...

- First Welsh treebank, third Celtic language treebank in UD
- 10 000 tokens, allows a relatively robust POS tagging (89.2%), and parsing of simpler sentences

Next steps

- needs to be expanded
- more annotators needed
- adding enhanced dependencies
- adding translations of sentences and glosses to words
- harmonisation with similar constructions in other languages?



[https://github.com/UniversalDependencies/UD_Welsh-CCG/
johannes.heinecke@orange.com](https://github.com/UniversalDependencies/UD_Welsh-CCG/johannes.heinecke@orange.com)