

TwittIrish: A Universal Dependencies Treebank of Tweets in Modern Irish

Anonymous ACL submission

Abstract

001 Modern Irish is a minority language lacking
002 sufficient linguistic resources for the task of
003 accurate automatic syntactic parsing of user-
004 generated content. As with other languages, the
005 linguistic style observed in Irish tweets differs,
006 in terms of orthography, lexicon and syntax,
007 to that of standard texts more commonly used
008 in Natural Language Processing (NLP) for the
009 development of language models and parsers.
010 This paper reports on the development of Twit-
011 tIrish, the first Irish Universal Dependencies
012 Twitter Treebank. We describe our bootstrap-
013 ping method, and report on preliminary parsing
014 experiments.

1 Introduction

015 User-generated content (UGC) is a valuable re-
016 source for training syntactic parsers which can
017 accurately process social media text. UGC is a
018 domain with features different from those of both
019 spoken language and standardised written language
020 more traditionally used in NLP corpora. Given that
021 the accuracy of syntactic parsing tools has been
022 shown to decline when evaluated on noisy UGC
023 data (Foster et al., 2011; Seddah et al., 2012) and
024 that domain adaptation has been shown to improve
025 parser performance in the case of dependency an-
026 notation of English tweets (Kong et al., 2014), the
027 need for domain-specific resources is clear in order
028 to reliably process this variety of data.

029 Universal Dependencies (UD) (Nivre et al.,
030 2020) provides a cross-lingually consistent frame-
031 work for dependency-based syntactic parsing.
032 UGC, especially social media text, has recently be-
033 come a popular focus within UD and NLP research
034 more broadly (Silveira et al., 2014; Luotolahti et al.,
035 2015; Albogamy and Ramsay, 2017; Wang et al.,
036 2017; Zeldes, 2017; Bhat et al., 2018; Blodgett
037 et al., 2018; Van Der Goot and van Noord, 2018;
038 Cignarella et al., 2019; Seddah et al., 2020) and has
039 encouraged active conversation around how best to
040

represent it within this framework among the UD
community.

041 This paper reports on the creation of the first
042 Irish Twitter treebank, TwittIrish. Irish is a mi-
043 nority language mostly spoken in small commu-
044 nities in Ireland called ‘Gaeltachtaí’ (CSO, 2016)
045 but social media sites, such as Twitter, provide a
046 platform for Irish speakers to communicate elec-
047 tronically from any location. Users may reach a
048 wide audience quickly, unconstrained by editors
049 and publishers. These and other socio-linguistic
050 factors contribute to the noncanonical language ob-
051 served in this domain, which we analyse through
052 the lens of orthographic, lexical and syntactic vari-
053 ation. In order to maintain optimum consistency
054 with other UD treebanks, the annotation methodol-
055 ogy employed in this research closely follows the
056 general UD guidelines and the language-specific
057 guidelines for Irish while aiming to incorporate
058 the most up-to-date recommendations (Sanguinetti
059 et al., 2020) for UGC in this evolving area of NLP.
060

061 We carry out preliminary parsing experiments
062 with TwittIrish, investigating the following two
063 questions: How effective is a parser trained on
064 the Irish Universal Dependencies Treebank (Lynn
065 and Foster, 2016), which contains only edited text
066 and no UGC, when applied to tweets? And what
067 difference do pre-trained contextualised word em-
068 beddings make? We observe a difference of approx-
069 imately 23 LAS points between TwittIrish and the
070 IUDT test set and find that the use of monolingual
071 BERT embeddings (Barry et al., 2021) improves
072 performance by over 10 LAS points.
073

074 The paper is structured as follows: Section 2
075 details the existing Irish NLP resources we use for
076 our research, Section 3 outlines the development
077 of the treebank, Section 4 describes the characteris-
078 tics of UGC evident in Irish tweets, and Section 5
079 presents parsing experiments and error analysis.

2 Irish NLP Resources

We use the following resources:

Indigenous Tweets (IT)¹ This project compiles statistics on social media data of 185 minority and indigenous languages including Irish. All tweets in the TwittIrish treebank were sourced via IT.

Lynn Twitter Corpus (LTC)² (Lynn et al., 2015) A corpus of 1,493 lemmatised and POS-tagged Irish language tweets randomly sampled from 950k tweets by 8k users posted between 2006 and 2014, identified by IT. The LTC data also contains code-switching information (Lynn and Scannell, 2019).

Irish Universal Dependencies Treebank (IUDT)³ (Lynn and Foster, 2016) A UD treebank consisting of 4910 sentences sampled from a balanced mixed-domain corpus for Irish.

gaBERT (Barry et al., 2021) A monolingual Irish BERT model, trained on approximately 7.9 million sentences, which outperforms Multilingual BERT (mBERT) (Devlin et al., 2019) and WikiBERT (Pyysalo et al., 2021) at the task of dependency parsing for Irish.

3 TwittIrish Development

We combined 700 POS-tagged tweets from the LTC with 166 tweets more recently crawled by IT in order to leverage previous linguistic annotations while also including newer tweets. This involved converting the LTC annotation scheme to that of the UD framework and then POS-tagging the new raw tweets. We provide further detail in Appendix A.

3.1 LTC Conversion

With regard to tokenisation, multiword expressions were automatically split into separate tokens following UD conventions. Only minor manual adjustments were required for lemmatisation to ensure alignment with the IUDT (to enable bootstrapping – see Section 3.3). Finally, the POS tagset used in the LTC was automatically converted to the UD tagset. Appendix A.2 describes this process.

¹<http://indigenoustweets.com/>

²<https://github.com/tlynn747/IrishTwitterPOS>

³https://github.com/UniversalDependencies/UD_Irish-IDT

3.2 Preprocessing of Newly-crawled Tweets

Due to the lack of a tokeniser designed to deal specifically with UGC in Irish, we compared two tools for this task: UDPipe (Straka et al., 2016),⁴ a language-agnostic trainable pipeline for tokenisation, tagging, lemmatisation and dependency parsing, and Tweettokenizer⁵ from NLTK (Bird, 2006), a rule-based tokeniser designed for noisy UGC. The latter proved to be more effective for tokenising UGC phenomena such as emoticons, URLs and meta-language tags. Manual corrections were then applied in order to adhere to the Irish-specific tokenisation scheme within current UD guidelines. In order to establish the best system to use for automatic lemmatising and POS-tagging, two tools, Morfette (Chrupala et al., 2008) and UDPipe (Straka et al., 2016), were analysed with Morfette achieving higher scores on both tasks.

3.3 Syntactic Annotation

As a method shown to reduce manual annotation efforts in syntactic annotation (Judge et al., 2006; Seraji et al., 2012), we carry out a bootstrapping approach to dependency parsing as recommended by UD⁶.

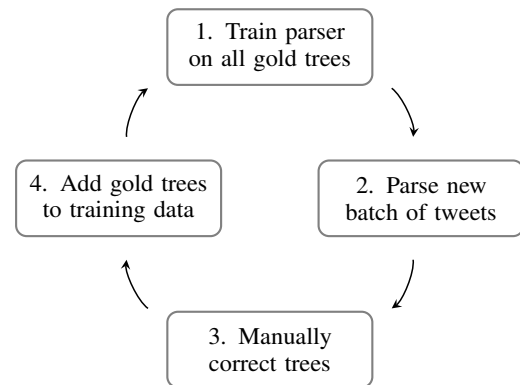


Figure 1: Bootstrapping approach to semi-automated syntax annotation.

The bootstrapping process is illustrated in Figure 1. After converting the LTC and new tweets to the CoNLL-U format, we manually annotated a small set of 166 tweets and began the bootstrapping cycle.⁷ (Step 1) A parsing model⁸ was then

⁴Trained on IUDT v2.8 with no pre-trained embeddings.

⁵<https://www.nltk.org/api/nltk.tokenize.html>

⁶https://universaldependencies.org/how_to_start.html

⁷All manual annotation and correction was performed by one linguist annotator.

⁸Biaffine Parser (Dozat and Manning, 2017) with mBERT (Devlin et al., 2019) embeddings.

trained on the IUDT in combination with the newly annotated tweets. (Step 2) The parsing model was used to automatically annotate the next batch of 100 tweets. (Step 3) These tweets were then manually corrected. (Step 4) The corrected tweets were then added to the training data. Steps 1 to 4 were repeated until all 866 tweets were fully parsed. This dataset represents the TwittIrish test set in the UD version 2.8 release.

4 Annotating Irish UGC

This section describes the linguistic features that can create challenges when parsing Irish social media text. We provide Irish language examples and discussion around the factors that influence these phenomena.

4.1 Orthographic Variation

Orthographic variation refers to deviation from the conventional spelling system of the language and is observed at the token level. Therefore, it can affect the lemmatisation of a token in an NLP pipeline, potentially affecting other downstream areas of annotation. In the TwittIrish dataset, 2.5% of tokens contained some orthographic variation. Table 1 exemplifies some frequently-occurring phenomena in Irish tweets which deviate from standard orthography.

Diacritic variation Diacritic marks are often omitted or incorrectly added to Twitter text. The acute accent or *síneadh fada* is used in Irish to indicate a long vowel and is necessary to disambiguate between certain word pairs. Example 1 shows the most probable intended word *léacht* ‘lecture’ rendered as *leacht* ‘liquid’.

- (1) *Leacht faoi stair Príosún Dún Dealgain*
‘Lecture about the history of Dundalk Prison’

Abbreviation Predictable shorthand forms can occur in standard Irish texts e.g. *lch* as an abbreviated form of *leathanach* ‘page’. While more unconventional, and thus less predictable, abbreviations are observed in Irish tweets, as per Example 2 in which the word *seachtain* ‘week’ is shortened to *seacht* ‘seven’. Abbreviations are common in tweets than standard text as the character limit and real-time, up-to-date nature of the platform encourages the user to be efficient with time and space.

- (2) *Bím de ghnáth ach sa bhaile an tseacht seo*
‘I usually am but home this week’

Lengthening This refers to the elongation of a token by repeating one or more characters. This can be thought of as an encoding of sociophonetic information (Tatman, 2015) and has been shown to be strongly linked to sentiment. Despite incentives to save time and space while tweeting, users often elongate certain words for expressive purposes (Brody and Diakopoulos, 2011). Example 3 demonstrates the lengthening of the word *buí* ‘yellow’.

- (3) *tá siad go léir buuuuuúí*
‘They are all yelloooooow’

Case Variation Nonstandard use of upper- and lowercase text is another method of encoding sociophonetic information by focusing attention or emotion on a particular word or phrase. Heath (2021) discusses the association between the use of all-caps and perceived shouting as in Example 4.

- (4) *Níl todhchaí na Gaeilge sa Ghaeltacht, ach in aon áit AR DOMHAIN*
‘The future of Irish is not in the Gaeltacht but anywhere ON EARTH’

Punctuation Variation Punctuation is used creatively in UGC to format or emphasise strings of text. However, due to the lack of standardisation, occurrences of unconventional punctuation can make text difficult to parse for both human and machine, as in Example 5 which shows a phrase from an Irish tweet appended by two punctuation characters ‘-’). It is unclear whether this should be interpreted as some form of punctuation, creative formatting or a smiley e.g. ‘;-)’.

- (5) *sin a dhóthain-)*
‘That’s enough-)’

Other Spelling Variation These are mostly slight variations very close to the intended word and may occur due to typographical error. Typos are very common in UGC due to lack of editing or proof-reading and may occur via insertion, deletion, substitution or transposition of characters.

- (6) *tus staith 6 de Imeall*
‘start of season 6 of Imeall’

Example 6 shows *sraith* (season) rendered as **staith*. Due to their phonetic dissimilarity and the fact that ‘t’ and ‘r’ are adjacent on the QWERTY keyboard layout, it is reasonable to infer that the substitution was unintentional. Less commonly, disguise or censorship of words or phrases may occur to encrypt profanity or taboo language of some variety.

Phenomenon	Example	Standard form	Gloss
Diacritic Variation	<i>níor fhoghlaím tu</i>	<i>níor fhoghlaím tú</i>	‘you did not learn’
Abbreviation	<i>fhoir rugbaí na hÉir</i>	<i>fhoireann rugbaí na hÉireann</i>	‘Irish rugby team’
Lengthening	<i>obairrrr</i>	<i>obair</i>	‘work’
Case Variation	<i>ceolchoirm DEN SCOTH</i>	<i>ceolchoirm den scoth</i>	‘excellent concert’
Punctuation Variation	<i>**folúntas**</i>	<i>folúntas</i>	‘vacancy’
Other Spelling Variation	<i>O ’Bama</i>	<i>Obama</i>	‘Obama’

Table 1: Examples of orthographic variation in Irish tweets.

4.2 Lexical Variation

Just 38.32% of the set of unique lemmata that make up the vocabulary of the TwittIrish treebank occur in the IUDT training data. Table 2 shows examples of lexical variation in Irish tweets.

Dialectal Vocabulary Irish has three major dialects; Connaught, Munster and Ulster. Distinctive features of these dialects in the form of lexical variation are evident in spoken language and informal text such as tweets. Example 7 shows the use of *domh*, the Ulster variant of *dom* meaning ‘to me’.

- (7) *Ba chóir domh rá!*
‘I should say!’

Initialism Multiword phrases are frequently represented by the initial letter of each of their constituent tokens. Example 8 shows *GRMA* ‘Thank you’ used to represent its expanded form *Go raibh maith agat*.

- (8) *Scaip an scéal! GRMA!*
‘Spread the word! Thank you!’

Pictogram Emojis, emoticons etc. can be added to text to emulate gesture (Gawne and McCulloch, 2019) or they may play a syntactic role in a phrase, replacing a word as in Example 9, in which the symbol, ♥, acts as the object of a verb. Pictograms tend not to have one-to-one correspondence with natural language words.

- (9) *Conas a deireann tú ♥?*
‘How do you say ♥?’

Truncation Due to the current limit of 280 characters per tweet, the end of the text may be truncated, sometimes mid-sentence or mid-word. Example 10 features an utterance in which the end has been unnaturally attenuated.

- (10) *Súil agam go bheas sé mar sin don...*
‘I hope it will be like that for the...’

Transliteration The practice of transliteration, in which a word in one language is written using the writing system of another, is common within the language pair of Irish and English. In the TwittIrish treebank, the English language phrase ‘fair play’ occurs twice while variations ‘fair plé’, as shown in Example 11 and ‘féar plé’ occur once each.

- (11) *Fair plé daoibh ♥,*
‘Fair play to you ♥’

Insertional Code-switching As with transliteration, code-switching occurs as a result of the high levels of contact in the Irish-English language pair. 66.74% of tokens in the TwittIrish treebank are in Irish, 4.85% of tokens are in English and the remainder are classified as neither, or indeed both in the case of intraword code-switching in which the morphologies of two languages are combined in a single word. Lynn and Scannell (2019) note the propensity of Irish speakers to conjugate an English language verb with the Irish gerund suffix *áil*.

- (12) *Eachtra i ndiaidh Happenáil i nGaoth Dobhair*
‘An event after happening in Gweedore’

Example 12 of intraword code-switching in which an Irish utterance uses the English verb root ‘happen’ instead of the Irish equivalent *tarlaigh*.

Single word code-switches and the use of loan words⁹ are a distinctive feature of informal Irish. 74.71% of the tweets in the TwittIrish treebank were considered to be entirely in Irish, the remaining 25.29% of tweets being considered bi- or multilingual. Example 13 shows a section of an Irish tweet utilizing the English word ‘Dubs’ a nickname for ‘Dubliners’.

- (13) *Roimh na Dubs*
‘Before the Dubs’

⁹Hickey (2009) and Stenson (2011) comment on the fuzzy boundary between single word switches and loan words in Irish. The TwittIrish corpus does not provide sufficient evidence to distinguish meaningfully between these terms.

Phenomenon	Example	Standard form	Gloss
Dialectal Vocabulary	anso	anseo	'here'
Initialism	BÁC	Baile Átha Cliath	'Dublin'
Pictogram	<3 mór	Grá mór	'Lots of love'
Truncation	thart fa' 53 nó...	thart fa' 53 nóiméad	'over 53 mi... (minutes)'
Transliteration	raight	ceart	'right'
Insertional Code-switching	sa town amárach	sa bhaile amárach	'in town tomorrow'
Other Nonstandard Lexical Forms	TochaltÓr	Tochaltóir óir	'Gold-digger'

Table 2: Examples of lexical variation in Irish tweets.

Other Nonstandard Lexical forms Other unfamiliar terms may occur in the form of hypercorrection and neologisms. Hypercorrection occurs when an autocorrection system is either not activated or available in a user's language of choice. As a result, their attempts to type a word are corrected to a word with a similar spelling in another language.

- (14) *Mhúscaíl mé i mo leaba féin ar maidin i ndiaidh concise mór*
'I woke up in my own bed after a big fortnight'

Example 14 shows the Irish word *coicíse* rendered as 'concise' probably due to automatic English spelling correction software. It is often difficult to distinguish between hypercorrection, neologisms, typos or other spelling variations.

- (15) *tá an teanga ag fáil bháis agus níl ach uaireanta*
'the language is dying and there are only hours'

Example 15 shows *agus* (and) rendered as *agua* which may have occurred due to hypercorrection as 'agua' (water) is a frequent token in other languages such as Portuguese and Spanish. However it could also be a simple typo.

4.3 Syntactic Variation

Grammatical phenomena observed in Irish tweets are described in this section. As these idiosyncrasies occur at the phrasal rather than token level, their effect is observed on the structure of the parse trees. Table 3 exemplifies syntactic variation in Irish tweets.

Contraction Much like abbreviation at the token level, contraction is defined here as the fusion of several tokens for the purpose of brevity sometimes mimicking spoken pronunciation. Figure 2 shows the phrase *tá a fhios agam* (lit. its knowledge is at me) reduced to *tá's agam*.

Over-splitting The inclusion of extra white space within tokens is often observed in Irish tweets e.g. *Nílím ró chinnte*. The prefix *ró-* ('too') is conventionally fused with the adjective it precedes in

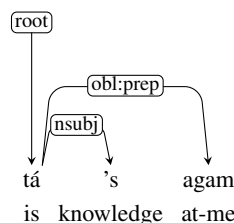


Figure 2: Contraction 'I know'.

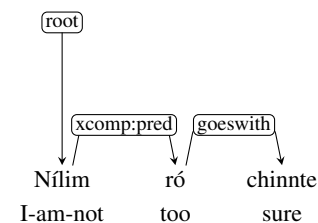


Figure 3: Over-splitting 'I am not too sure'.

standardised text and so such tokens are annotated with the *goeswith* label as shown in Figure 3.

Alternational Code-switching This refers to code-switching which alters the structure of the syntax tree, due to the differing word orders of the languages involved, thus complicating the task of dependency parsing. For example, in Irish the adjectival modifier usually *follows* the noun it modifies whereas the inverse is true for English. Figure 4 shows how the English adjective 'hippy-dippy' is incorporated in an Irish language utterance.

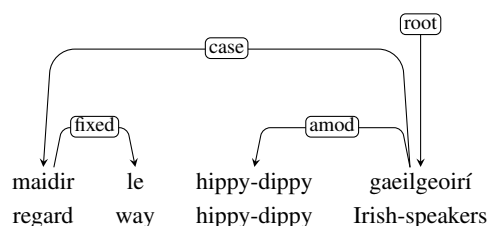


Figure 4: Alternational code-switching 'as for hippy-dippy Irish speakers'.

Dialectal Grammar Figures 5 and 6 show semantically equivalent statements rendered using the synthetic and analytic verb forms respectively. The synthetic form is more common in the Munster dialect of Irish.

Ellipsis It is also common for tweets to consist of incomplete sentences. Example 16 shows such a

Phenomenon	Example	Standard form	Gloss
Contraction	<i>go dtí'n</i>	<i>go dtí an</i>	'until the'
Over-splitting	<i>ana shuimiúil</i>	<i>an-suimiúil</i>	'very interesting'
Alternational Code-switching	<i>Tá an tweet machine ró-tapa</i>	<i>Tá inneall na tvuíte ró-tapa</i>	'The tweet machine is too fast'
Dialectal Grammar	<i>Ní fhacthas</i>	<i>ní fhaca mé</i>	'I did not see'
Ellipsis	<i>jab iontach déanta aige</i>	<i>tá jab iontach déanta aige</i>	'he has done a wonderful job'
Meta-language Tags	<i>#sonas</i>	<i>sonas</i>	'happiness'
Non-sentential Segmentation	<i>haha:) tá súil agam go raibh sé ann</i>	<i>ha ha! Tá súil agam go raibh sé ann.</i>	'haha:) I hope he was there.'
Other Grammatical Variation	<i>ce ata an athair?</i>	<i>cé hé an t-athair?</i>	'who is the father?'

Table 3: Examples of syntactic variation in Irish tweets.

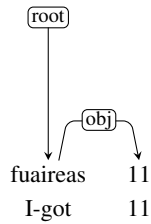


Figure 5: Synthetic verb form 'I got 11'.

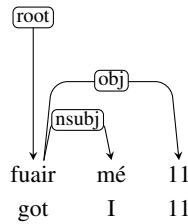


Figure 6: Analytic verb form 'I got 11'.

sentence fragment lacking a main verb. The probable inferred full phrase is *tá báisteach anseo* 'rain is here'.

- (16) *báisteach anseo*
'rain here'

Meta-language Tags Hashtags are used in tweets to render the topic searchable and at-mentions are used to address or refer to another user. Both can play syntactic roles in the sentence as shown in Figure 7.

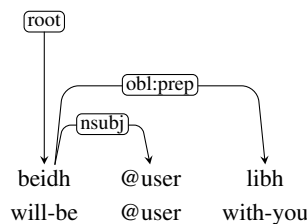


Figure 7: Syntactic meta-language tag '@user will be with you'.

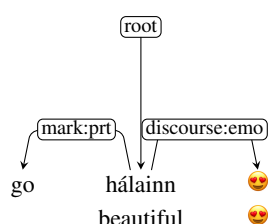


Figure 8: Non-sentential tweet using emoji as punctuation 'beautiful 😊'.

Non-sentential Structure In tweets, the sentence is not an appropriate unit of segmentation as frequently non-standard punctuation, or none at all, is used. Figure 8 exemplifies a tweet utilising an emoji instead of punctuation.

Other Grammatical Variation Grammatical variation can also occur via unintentional deviation from conventional spelling or grammar by an

L2 Irish speaker. Example 17 shows a grammatically incorrect phrase roughly translating to 'I have to *going'. In such cases, though the annotator may be able to infer the intended phrase *Caithfidh mé ag dul* 'I have to go', no corrections are made by the annotator to the surface form, however this information can be represented in the annotation via the label `CorrectForm` as described in (Sanguinetti et al., 2020).

- (17) *Caithfidh mé ag dul*
'I have to *going'

Additionally, Irish tweets contain extremely unconventional constructions. This can occur in the form of unnatural phrases that have been machine translated or generated by bots.

- (18) *Conas a Faigh tonna de Morgáiste*
'*How to get a tonne of mortgage'

Example 18 shows an ungrammatical construction that appears to have been translated automatically word by word. A more natural construction might be *conas tonna morgáiste a fháil* 'How to get a tonne of mortgage'. Some examples of this variety are easy to identify from surrounding context such as links to websites with similar content however, tweets may consist of text alone making it difficult to infer whether the author is human or machine.

5 Parsing Experiments

We compare the performance of two widely used neural dependency parsers on the TwittIrish test set, and examine the effect of using pre-trained contextualised word embeddings from a monolingual Irish BERT model (gaBERT). We report parsing performance by UPOS and dependency label and carry out a manual error analysis. Further information is detailed in Appendix B.

5.1 Parser Comparison

We experiment with two neural dependency parsing architectures: UDPipe (Straka et al., 2016), an NLP

System	LAS	
	IUDT	TwittIrish
UDPipe v1	70.58	47.33
AllenNLP	71.56	48.73
AllenNLP + gaBERT	84.25	59.34

Table 4: Comparison of parsing systems UDPipe v1, Biaffine Dependency Parser (Dozat and Manning, 2017): Biaffine dependency parser with BiLSTM encoder, and AllenNLP + gaBERT: Biaffine dependency parser where BiLSTM is replaced with pretrained BERT model of Barry et al. (2021). All were trained on IUDT version 2.8 and tested on the IUDT and TwittIrish test sets.

pipeline which includes a transition-based non-projective parser, and AllenNLP (Gardner et al., 2018), a biaffine dependency parser with a BiLSTM encoder (Dozat and Manning, 2017). Both systems are trained on IUDT version 2.8¹⁰ and tested on the IUDT and TwittIrish test sets for comparison. Gold standard tokenization is provided to the models which then predict UPOS tags and dependency relations and labels. As the TwittIrish test set is the only gold annotated treebank of Irish UGC, no UGC is used as training or development data in these experiments. We opt to preserve it as a test set so that our results and results of future research in this area will be comparable.

In order to leverage the substantial advances in accuracy achieved in dependency parsing by the use of pre-trained contextualized word representations (Che et al., 2018; Kondratyuk and Straka, 2019; Kulmizev et al., 2019), we use AllenNLP with token representations obtained from the last hidden layer of the gaBERT model (Barry et al., 2021) and passed to the biaffine parser.

Table 4 shows that, when tested on the IUDT version 2.8 test set, UDPipe achieves 70.58% labelled attachment score (LAS). In comparison, UDPipe achieves a much lower LAS score of 47.33 on the TwittIrish test set. Similarly to UDPipe, AllenNLP achieves 71.56 LAS on the IUDT test set with a similar decrease of 22.83 points on the TwittIrish test set. The highest accuracy of 84.25 LAS is achieved by gaBERT with a difference of 24.91 points when tested on the TwittIrish test set. The lower accuracy obtained by parsers on the TwittIrish test is unsurprising given the linguistic differences between the training and test sets. The

¹⁰Models were trained with and without XPOS and feature annotation. The results shown here are without XPOS and features. The addition of XPOS and features constituted a difference of approximately +/-1 LAS.

10+ LAS improvement provided by the gaBERT embeddings is seen in both test sets.

5.2 Analysis

Analysis was carried out on the AllenNLP parser with gaBERT embeddings using Dependable (Choi et al., 2015).

LAS by Number of Tokens per Sentence/Tweet

The mean sentence length of the IUDT is 23.5 tokens, whereas the mean tweet length in TwittIrish is 17.8. Figure 9 shows that, when tested

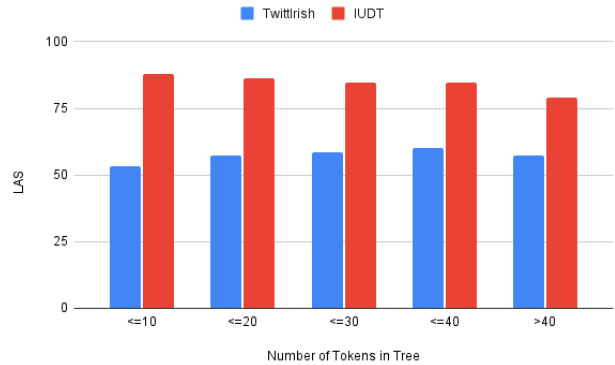


Figure 9: LAS by number of tokens per tweet achieved by AllenNLP Parser with gaBERT embeddings on the IUDT and TwittIrish test sets.

on the IUDT, parsing accuracy decreases as the length of the sentence increases. The highest accuracy of 87.92 LAS is associated with sentences of 10 tokens or fewer and the lowest accuracy is observed in sentences of 40 tokens or more. This is an unsurprising trend as a higher number of tokens increases the probability of longer dependency distances and more complex constructions within a sentence. While the range of scores is smaller and trend less pronounced, the opposite effect is observed when the same parser is tested on TwittIrish, whereby LAS tends to increase as the length of the tweet increases. The highest LAS is associated with tweets of 31 to 40 tokens in length and the lowest accuracy is associated with tweets of 10 tokens or less. This trend is also observed when gaBERT representations are not used, suggesting that, in this case, deep contextualised word embeddings do not cause this effect as observed in (Kulmizev et al., 2019). From manual inspection of the data, we observe that the genre-specific phenomena which challenge the parser such as ellipsis, meta-language tags and URLs, occur in higher proportions in shorter tweets therefore causing this

trend.

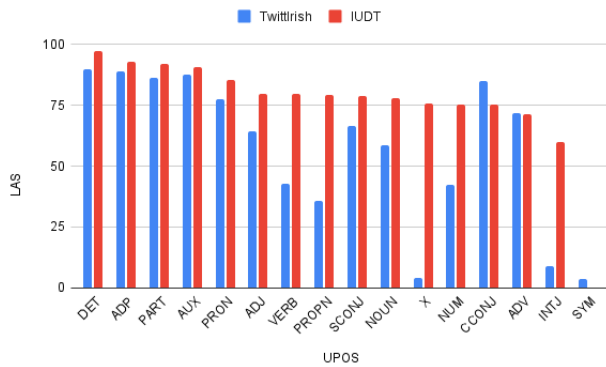


Figure 10: LAS by UPOS tag achieved by AllenNLP Parser with gaBERT embeddings on the IUDT and TwittIrish test sets.

LAS by UPOS Figure 10 shows LAS associated with each UPOS tag when tested on the IUDT and TwittIrish. LAS is higher when tested on the IUDT for all UPOS tags except CCONJ, ADV and SYM and in these cases the difference is small (<10 LAS). The most notable differences are X (71.6 LAS), INTJ (51.3 LAS), PROPN (43.5 LAS). These differences are due to 1) the divergent genres of the treebanks e.g. in the TwittIrish treebank the UPOS tag X is used for all non-syntactic hashtags, and PROPN is used for all at-mentions, neither of which occur in the IUDT and 2) differing annotation conventions e.g. in the IUDT, the tag X is used mostly for foreign-language tokens. Whereas, in TwittIrish, due to the high proportion of English language tokens, non-Irish words are annotated with their true UPOS tag where the language is known to the annotator.

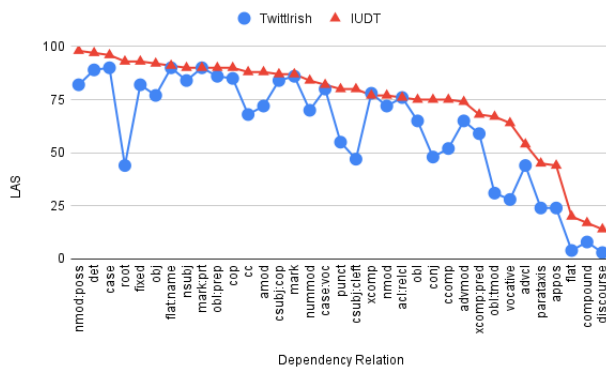


Figure 11: LAS achieved by AllenNLP Parser with gaBERT embeddings on the IUDT and TwittIrish test sets by dependency relation.

LAS by Dependency Relation Figure 11 shows parsing accuracy broken down by dependency relation. The parser obtains higher scores on the IUDT for all dependency relations except xcomp for which it is just one point higher when tested on TwittIrish. The largest differences between the accuracy of the two test sets are associated with the labels root, vocative, obl:tmod, csbj:clnt, conj, and punct.

Error Analysis In order to assess the effect of the UGC phenomena present in Irish tweets, we analyse the most and least accurate parses. 7 tweets (76 tokens) were parsed with LAS between 0 and 5. There were 15 occurrences of emojis which were most commonly incorrectly labelled punct. The 10 occurrences of code-switching were most commonly attached via flat:foreign. The 9 (2 syntactic) occurrences of usernames were most commonly labelled as root. There were 5 occurrences of ellipsis in the form of verb omission obfuscating the task of root selection. The 3 hashtags were most commonly mislabelled as nmod as were the 3 URLs. 1 occurrence of spelling variation was observed in the form of diacritic omission wherein the word *ár* ‘our’ was rendered as *ar* ‘on’ causing the parser to misinterpret the dependency label of nmod:poss as case. 7 tweets (89 tokens) were parsed with an accuracy between 95 and 100 LAS. All of these were grammatical, well-formed sentences. There were 3 usernames and 1 hashtag all of which were syntactically integrated and so they were parsed correctly. There was one occurrence of insertional single-word code-switching which was accurately parsed. There were 2 occurrences of spelling variation, both in the form of diacritic omission but, as these do not resemble any other words, they were parsed correctly.

6 Conclusion

Presented in this paper is the novel resource, TwittIrish, the first Universal Dependencies treebank for Irish UGC. Analysis of this linguistic genre and anonymised examples of Irish tweets are presented. This research facilitates the development of NLP tools such as dependency parsers for Irish by providing a test set on which future Irish language technology can be tested. Future work will involve both further annotation and exploration of semi-supervised techniques.

676	John Judge, Aoife Cahill, and Josef Van Genabith. 2006.	<i>Evaluation Conference</i> , pages 4034–4043, Marseille, France. European Language Resources Association.	733
677	Questionbank: Creating a corpus of parse-annotated		734
678	questions. In <i>Proceedings of the 21st International</i>		
679	<i>Conference on Computational Linguistics and 44th</i>	Slav Petrov, Dipanjan Das, and Ryan McDonald. 2012.	735
680	<i>Annual Meeting of the Association for Computational</i>	<i>A universal part-of-speech tagset</i> . In <i>Proceedings</i>	736
681	<i>Linguistics</i> , pages 497–504.	<i>of the Eighth International Conference on Language</i>	737
		<i>Resources and Evaluation (LREC'12)</i> , pages 2089–	738
682	Dan Kondratyuk and Milan Straka. 2019. <i>75 languages,</i>	2096, Istanbul, Turkey. European Language Re-	739
683	<i>1 model: Parsing Universal Dependencies univer-</i>	sources Association (ELRA).	740
684	<i>sally</i> . In <i>Proceedings of the 2019 Conference on</i>		
685	<i>Empirical Methods in Natural Language Processing</i>	Sampo Pyysalo, Jenna Kanerva, Antti Virtanen, and	741
686	<i>and the 9th International Joint Conference on Natu-</i>	Filip Ginter. 2021. <i>WikiBERT models: Deep trans-</i>	742
687	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	<i>fer learning for many languages</i> . In <i>Proceedings</i>	743
688	2779–2795, Hong Kong, China. Association for Com-	<i>of the 23rd Nordic Conference on Computational</i>	744
689	putational Linguistics.	<i>Linguistics (NoDaLiDa)</i> , pages 1–10, Reykjavik, Ice-	745
		land (Online). Linköping University Electronic Press,	746
690	Lingpeng Kong, Nathan Schneider, Swabha	Sweden.	747
691	Swayamdipta, Archana Bhatia, Chris Dyer, and		
692	Noah A. Smith. 2014. <i>A Dependency Parser for</i>	Manuela Sanguinetti, Cristina Bosco, Lauren Cas-	748
693	<i>Tweets</i> . In <i>Proceedings of the 2014 Conference</i>	sidy, Özlem Çetinoğlu, Alessandra Teresa Cignarella,	749
694	<i>on Empirical Methods in Natural Language Pro-</i>	Teresa Lynn, Ines Rehbein, Josef Ruppenhofer,	750
695	<i>cessing (EMNLP)</i> , pages 1001–1012, Doha, Qatar.	Djamé Seddah, and Amir Zeldes. 2020. <i>Treebanking</i>	751
696	Association for Computational Linguistics.	<i>user-generated content: A proposal for a unified rep-</i>	752
		<i>resentation in Universal Dependencies</i> . In <i>Proceed-</i>	753
697	Artur Kulmizev, Miryam de Lhoneux, Johannes	<i>ings of the 12th Language Resources and Evaluation</i>	754
698	Gontrum, Elena Fano, and Joakim Nivre. 2019. <i>Deep</i>	<i>Conference</i> , pages 5240–5250, Marseille, France. Eu-	755
699	<i>contextualized word embeddings in transition-based</i>	ropean Language Resources Association.	756
700	<i>and graph-based dependency parsing - a tale of two</i>		
701	<i>parsers revisited</i> . In <i>Proceedings of the 2019 Confer-</i>	Djamé Seddah, Farah Essaidi, Amal Fethi, Matthieu	757
702	<i>ence on Empirical Methods in Natural Language Pro-</i>	Futeral, Benjamin Muller, Pedro Javier Ortiz Suárez,	758
703	<i>cessing and the 9th International Joint Conference</i>	Benoît Sagot, and Abhishek Srivastava. 2020. <i>Build-</i>	759
704	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	<i>ing a user-generated content North-African Arabizi</i>	760
705	pages 2755–2768, Hong Kong, China. Association	<i>treebank: Tackling hell</i> . In <i>Proceedings of the 58th</i>	761
706	for Computational Linguistics.	<i>Annual Meeting of the Association for Computational</i>	762
		<i>Linguistics</i> , pages 1139–1150, Online. Association	763
707	Juhani Luotolahti, Jenna Kanerva, Veronika Laippala,	for Computational Linguistics.	764
708	Sampo Pyysalo, and Filip Ginter. 2015. Towards Uni-		
709	versal Web Parsebanks. In <i>Proceedings of the Third</i>	Djamé Seddah, Benoit Sagot, Marie Candito, Virginie	765
710	<i>International Conference on Dependency Linguistics</i>	Mouilleron, and Vanessa Combet. 2012. The french	766
711	<i>(Depling 2015)</i> , pages 211–220. Uppsala University,	social media bank: a treebank of noisy user gener-	767
712	Uppsala, Sweden.	ated content. In <i>COLING 2012-24th International</i>	768
		<i>Conference on Computational Linguistics</i> .	769
713	Teresa Lynn and Jennifer Foster. 2016. Universal		
714	Dependencies for Irish. In <i>In Proceedings of the</i>	Mojgan Seraji, Beáta Megyesi, and Joakim Nivre. 2012.	770
715	<i>2nd Celtic Language Technology Workshop</i> , Paris,	Bootstrapping a persian dependency treebank. <i>Lin-</i>	771
716	France.	<i>guistic Issues in Language Technology</i> , 7(18).	772
717	Teresa Lynn and Kevin Scannell. 2019. Code-switching	Natalia Silveira, Timothy Dozat, Marie Catherine De	773
718	in Irish tweets: A preliminary analysis. In <i>In Pro-</i>	Marneffe, Samuel R. Bowman, Miriam Connor, John	774
719	<i>ceedings of the 3rd Celtic Language Technology</i>	Bauer, and Christopher D. Manning. 2014. A Gold	775
720	<i>Workshop</i> .	Standard Dependency Corpus for English. In <i>Pro-</i>	776
		<i>ceedings of the 9th International Conference on Lan-</i>	777
721	Teresa Lynn, Kevin Scannell, and Eimear Maguire.	<i>guage Resources and Evaluation, LREC 2014</i> , pages	778
722	2015. <i>Minority Language Twitter: Part-of-Speech</i>	2897–2904. ELRA.	779
723	<i>Tagging and Analysis of Irish Tweets</i> . In <i>Proceedings</i>		
724	<i>of the Workshop on Noisy User-generated Text</i> , pages	Nancy Stenson. 2011. <i>21. Code-switching vs. borrow-</i>	780
725	1–8, Beijing, China. Association for Computational	<i>ing in Modern Irish</i> , pages 559–580. Max Niemeyer	781
726	Linguistics.	Verlag.	782
727	Joakim Nivre, Marie-Catherine de Marneffe, Filip Gin-	Nancy Stenson. 2019. <i>Modern Irish: A Comprehensive</i>	783
728	ter, Jan Hajič, Christopher D. Manning, Sampo	<i>Grammar</i> . Routledge Comprehensive Grammars.	784
729	Pyysalo, Sebastian Schuster, Francis Tyers, and	Taylor & Francis.	785
730	Daniel Zeman. 2020. <i>Universal Dependencies v2:</i>		
731	<i>An evergrowing multilingual treebank collection</i> . In	Milan Straka, Jan Hajic, and Jana Straková. 2016. UD-	786
732	<i>Proceedings of the 12th Language Resources and</i>	Pipe: trainable pipeline for processing CoNLL-U	787

files performing tokenization, morphological analysis, pos tagging and parsing. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4290–4297.

Rachael Tatman. 2015. #go awn: Sociophonetic Variation in Variant Spellings on Twitter. *Working Papers of the Linguistics Circle*, 25(2):97–108. Number: 2.

Rob Van Der Goot and Gertjan van Noord. 2018. Modeling Input Uncertainty in Neural Network Dependency Parsing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4984–4991.

Hongmin Wang, Yue Zhang, Guang Yong Leonard Chan, Jie Yang, and Hai Leong Chieu. 2017. Universal Dependencies Parsing for Colloquial Singaporean English. In *ACL 2017 - 55th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers)*, volume 1, pages 1732–1744.

Amir Zeldes. 2017. The GUM Corpus: Creating Multilayer Resources in the Classroom. *Language Resources and Evaluation*, 51(3):581–612.

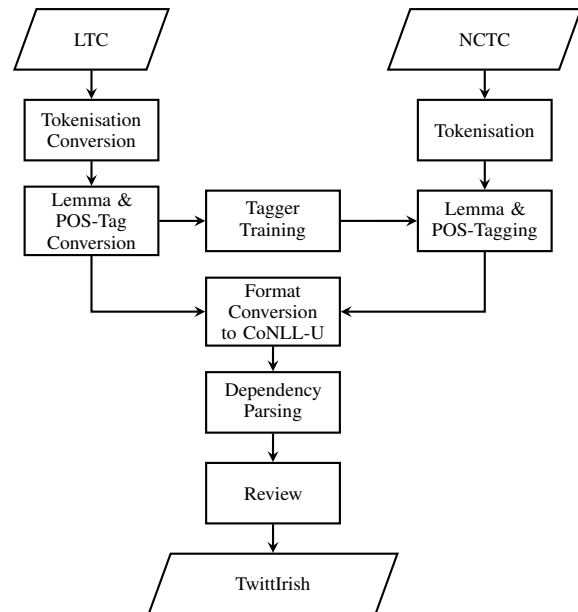


Figure 12: Diagram of the TwittIrish Creation Process. Where NCTC refers to Newly-crawled tweets.

A TwittIrish Development

A.1 LTC Tokenisation Conversion

The most notable difference in the tokenisation approach of LTC as compared to that of UD, was in the treatment of multi-word expressions (MWEs). In LTC, the individual tokens of MWEs are fused with an underscore whereas words with spaces are not allowed in UD¹¹. Several minor differences were also observed between the two tokenisation schemes such as whether or not certain symbols, abbreviations or punctuation marks should be attached to the token they follow or considered as a separate token. e.g. 5%, ama..., 1-0, 10pm. UD tends to favour the approach of separating such combinations¹² therefore we resolved to manually separate such occurrences in the TwittIrish tokenisation scheme.

A.2 LTC POS-tag Conversion

Table 5 shows the mapping of LPOS to UPOS. LPOS tags were automatically converted to the corresponding UPOS tag where a one-to-one or many-to-one mapping existed. In the case of one-to-many relationships, automatic identification and manual correction was performed.¹³

¹¹<https://universaldependencies.org/v2/mwe.html>

¹²however not all treebanks apply this consistently.

¹³Both the Gimpel et al. (2011) and UD tagsets derived from the Google Universal POS tagset (Petrov et al., 2012)

LTC POS	UPOS
N, VN	NOUN *
^, @	PROPN *
O	PRON
V	VERB, AUX †
A	ADJ
R	ADV
D	DET
P	ADP
T	PART
,	PUNCT
&	CCONJ, SCONJ †
\$	NUM
!	INTJ
U, ~, E	SYM *
#, #MWE	X *
EN	any †
G	any †

Table 5: POS tag Mapping

* Many-to-one relation

† One-to-many relation

Surface	LPOS	UPOS
@user	@	PROPN
#cutie	#	X
ca	R	ADV
bhfuil	V	VERB
an	D	DET
ghra	N	NOUN
you	EN	PRON
ask	EN	VERB

@user #cutie ca bhfuil an ghra you ask ¹⁴
 ‘@user #cutie where is the love you ask’

Table 6: Example Irish tweet with LTC and corresponding universal POS tags.

Table 6 demonstrates the mapping of a sample tweet from one scheme to the other. As all English language tokens were annotated with a single tag ‘EN’ in the LPOS scheme, these tags were converted to the appropriate UPOS tag in the TwittIrish treebank.

Table 7 shows that, using the LPOS tagset, all verbs are tagged V. According to UD, the Irish copula (e.g. *is, ní*) is tagged as AUX distinguishing it from other verbs (e.g. *tá, níl*) which are tagged VERB.

A.3 Preprocessing of Newly-crawled Tweets

Table 8 shows that hashtags and emoticons were not correctly handled by the UDPipe tokenizer trained on the IUDT. Despite being trained on Irish data, due to differences in domain, Twitter-specific features such as meta-language tags are not present in its training data.

Surface	LPOS	UPOS
Ní	V	AUX
duine	N	NOUN
cáilúil	A	ADJ
é	O	PRON
ach	&	CCONJ
táim	V	VERB
bródúil	A	ADJ
#Grá	#	X

Ní duine cáilúil é ach táim bródúil #Grá
 ‘He is not a celebrity but I’m proud #Love’

Table 7: Example Irish tweet with LTC and corresponding universal POS tags.

UDPipe (IUDT)	NLTK Tweettokenizer
Dé	Dé
Céadaoin	Céadaoin
#	#Midweek
Midweek	
#	#Beagnachann
Beagnachann	
:	:
)	:
:	:
)	:

Dé Céadaoin #Midweek #Beagnachann :) :)
 ‘Wednesday #Midweek #Almostthere :) :)’

Table 8: Example Irish tweet with UDPipe and NLTK tokenization

A.4 Lemma and POS-Tagging of Newly Crawled Tweets

A.5 Conversion to CoNLL-U Format

Table 9 shows that the Morfette format is a subset of the CoNLL-U format used by UDPipe. The LTC and CTC were thus converted automatically from the 3-column Morfette format, consisting of the token, lemma and POS-tag to the 10-column CoNLL-U format. CoNLL-U enables additional token-level annotation i.e. a token id, language-specific part-of-speech tags (XPOS), morphological features, the head of the current word, the dependency relation, an enhanced dependency graph in the form of a list of head-deprel pairs and any other miscellaneous annotation.¹⁵ CoNLL-U also requires a sentence ID and the original raw text to be included preceding the annotation. Further, in the miscellaneous column, the label ‘SpaceAfter=No’ encodes infor-

¹⁵In order to make optimum use of the time spent by the annotator, language-specific part-of-speech tags, morphological features and enhanced dependency annotation were not included in this version of the TwittIrish dataset. These elements can be automatically added in later versions of the treebank.

CoNLL-U		Morfette			CoNLL-U					
ID	FORM	LEMMA	UPOS	XPOS	FEATS	HEAD	DEPREL	DEPS	MISC	
1	Cuirfidh	cuir	VERB	–	–	0	root	–	–	
2	mé	mé	PRON	–	–	1	nsubj	–	–	
3	DM	DM	NOUN	–	–	1	obj	–	–	
4	chuici	chuig	ADP	–	–	1	obl:prep	–	–	

‘Cuirfidh mé DM chuici’
‘I will send her a DM’

Table 9: Example conversion of Irish tweet from Morfette to CoNLL-U format

mation about which tokens have a space after them in the original text for detokenisation purposes enabling automatic conversion from raw text to tree and vice versa.

A.6 Review

In order to assess the accuracy of the dependency annotation, a subset of the annotated data, consisting of 46 trees (773 tokens), was reviewed for errors by another Irish speaker trained in linguistic annotation. The task of the reviewer was to flag possible errors in the form of a token with an incorrect head and/or label. 46 possible errors were identified by the reviewer. The possible errors were then discussed by a team of two expert annotators to confirm whether the possible errors were true errors. 32 possible errors were confirmed as true errors. The overall accuracy of the treebank annotation can be estimated as 95.86% by dividing the number of correctly annotated tokens by the total number of tokens in the review.

16 tokens (2.07% of all tokens in the review) had an incorrect label and correct head. Figure 13 exemplifies one such correction. *Go* is a common particle in Irish, which can precede an adjective to create an adverb. When used for this function it is roughly equivalent to the suffix ‘-ly’ in English. e.g. *Ainnis* (‘miserable’), *go hainnis* (‘miserably’). For this reason, a parser is likely to annotate this construction as *advmod*. However, these constructions also appear as the predicate of the substantive Irish verb *bí* (‘to be’) and in this case they should be considered as *xcomp:pred*.¹⁶

12 tokens (1.55% of all tokens in the review) had an incorrect head and correct label. The most common error (5 instances) was incorrect punctuation attachment. Only 4 tokens (0.52%) were identified as having both incorrect head and label. Figure 14 shows the phrase *maith sibh* (‘good on you’)

¹⁶<https://universaldependencies.org/ga/dep/xcomp-pred.html>

incorrectly annotated with *sibh* as the *root* and *maith* as its adjectival modifier. It was identified in the review that *maith* should be considered the adjective predicate of an elided copula (Stenson, 2019).

The full phrase is thought to be *is maith sibh* and the corrected annotation is shown in Figure 15.

B Parsing Experiments

B.1 LAS by UPOS

LAS	TwittIrish High	TwittIrish Low
IUDT High	DET, ADP, PART, AUX, PRON, SCONJ	VERB, PROPN, PUNCT, X, INTJ
IUDT Low	ADJ, CCONJ, ADV	NOUN, NUM, SYM

Table 10: Confusion matrix of LAS by UPOS tag achieved by AllenNLP Parser with gaBERT embeddings on the IUDT and TwittIrish test sets

Table 10 shows which UPOS tags are associated with higher or lower than average LAS in both test sets. High accuracy is correlated with tokens which occur frequently and have low variation.¹⁷

UPOS tags DET, ADP, PART, AUX, PRON and SCONJ are associated with higher than average LAS in both the TwittIrish and IUDT test sets. In the IUDT, a high proportion, 8.87%, of tokens have the UPOS tag DET. As is common with function words, DET comprises of a closed set of lemmata and thus has the low variation of 0.21%.

The tags ADJ, CCONJ, and ADV are associated with higher than average LAS in the TwittIrish test set but lower than average LAS in the IUDT. This might be because these tags are more likely to be involved in more complex, ambiguous or long-distance attachments.

The tags VERB, PROPN, PUNCT, X, and INTJ are

¹⁷Variation is calculated by dividing the number of occurrences by then number of unique lemmata

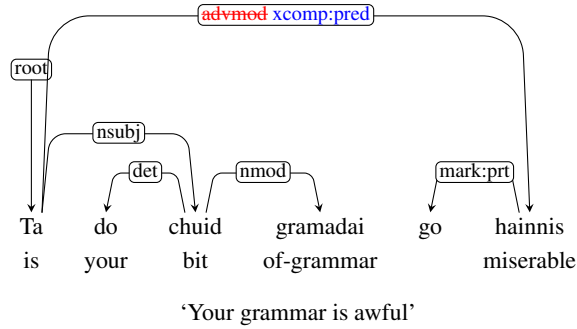


Figure 13: Reviewed tweet with corrected label.

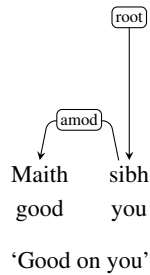


Figure 14: Example of tweet with incorrect head and label.

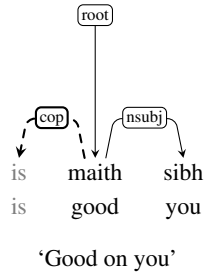


Figure 15: Reviewed tweet with corrected head and label.

936 associated with higher than average LAS in the
 937 IUDT tet set but lower than average LAS in Twit-
 938 tIrish. In the case of VERB and PUNCT, this can
 939 be attributed to the non-sentential nature of tweets.
 940 UPOS tags NOUN, NUM and SYM are associated
 941 with lower than average LAS in both the TwittIrish
 942 and IUDT test sets. In the IUDT, a low proportion,
 943 0.02%, of tokens have the UPOS tag SYM. The
 944 variation is high 83.33%.

945 B.2 LAS by Dependency Relation

946 Table 11 shows that high accuracy is associ-
 947 ated with dependency relations nmod:poss,
 948 det, case, fixed, obj, flat:name, nsubj,
 949 mark:prt, obl:prep, cop, cc, amod,
 950 csubj:cop, mark, nummod, case:voc

LAS	TwittIrish High	TwittIrish Low
IUDT High	nmod:poss, det, case, fixed, obj, flat:name, nsubj, mark:prt, obl:prep, cop, cc, amod, csubj:cop, mark, nummod, case:voc	root, csubj:cleft, punct
IUDT Low	xcomp:pred, advmod, obl, acl:relcl, nmod, xcomp	discourse, compound, flat, appos, parataxis, advcl, vocative, obl:tmod, ccomp, conj

Table 11: Confusion matrix of LAS by dependency label achieved by AllenNLP Parser with gaBERT embeddings on the IUDT and TwittIrish test sets

951 in both the IUDT and TwittIrish. root,
 952 csubj:cleft and, punct are associated with
 953 higher than average LAS in both the TwittIrish and
 954 IUDT test sets. xcomp:pred, advmod, obl,
 955 acl:relcl, nmod, and xcomp are associated
 956 with higher than average LAS in the TwittIrish tet
 957 set but lower than average LAS in the IUDT.
 958 are associated with higher than average LAS in the
 959 IUDT tet set but lower than average LAS in Twit-
 960 tIrish. discourse, compound, flat, appos,
 961 parataxis, advcl, vocative, obl:tmod,
 962 ccomp, and conj are associated with lower than
 963 average LAS in both the TwittIrish and IUDT test
 964 sets.