

# Towards better evaluation for Universal Dependencies: A typological perspective

Jimmy Callin

Uppsala University

jimmy.callin@gmail.com

## Abstract

Todo:

- Write this abstract.
- Add human judgment results (ok)
  - expand upon this in introduction and conclusion
- Eat something non-microwaved.
- Maybe add more references for UAS and LAS?

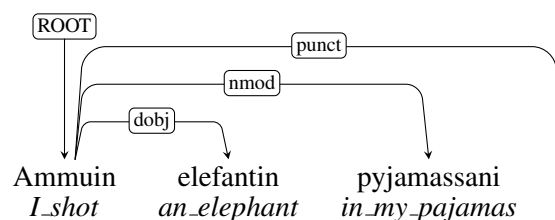


Figure 1: Finnish dependency tree for *I shot an elephant in my pajamas*. Note the edge count of Finnish being 4, while the English edge count is 8.

## 1 Introduction

The role of language processing is becoming increasingly multi-lingual, which is reflected in recent efforts into providing dependency parsing frameworks that can reliably be applied on a multitude of languages. One of the most ambitious projects in this area is the Universal Dependencies (UD) framework (Nivre, 2015), where the goal is to create a parsing framework with a cross-linguistically consistent grammatical annotation. The purpose of this work is to remove the requirement of language specific components which has up to this point been a necessity due to inconsistent annotation standards.

Data-driven evaluation metrics have been used as long as treebanks have been available (see Collins (1999) chap. 4 for a survey on early results). To evaluate statistical parsing models, two of the most ubiquitous evaluation metrics are Labeled Attachment Score (LAS) and Unlabeled Attachment Score (UAS). These have in their simplicity and intuitiveness served the research area well, but their design relies on a number of assumptions that we argue do not hold in the context of UD.

Firstly, the design of attachment scores assumes that we know next to nothing about the taxonomy and design choices of the parsing framework. This has historically been necessary since there have not been a consistently adapted framework for dependency parsing across many languages. We argue that recent progress in UD has made this assumption invalid. UD has a carefully specified framework that all UD treebanks has to adapt, and we can exploit this constraint to learn more about the performance of a model.

Secondly, parsing results are increasingly juxtaposed in a cross-linguistic manner, and this will likely continue with the establishment of UD. It is not uncommon to compare the output of e.g. English and Finnish under the assumption that equal evaluation scores is equivalent to equal parsing performance. The reason why this is problematic becomes apparent when studying grammatical morphemes in languages where these may be unbounded (i.e. function words) with languages where they typically are bounded on content words (i.e. affixes).

In figure 1, we have a Finnish sentence with four edges, while the equivalent sentence in English requires a total of eight edges. The problem appears as soon as we introduce a parsing error into the trees: one faulty edge in the Finnish sentence would result in a performance reduction of 25%, while the same error in the English sentence only would reduce the accuracy by 12.5%. Function words are regular in their appearance, and should therefore be relatively easy for a model to parse correctly. Our hypothesis is that languages with a comparatively large number of function words, like English, receive a more or less free performance boost when using classic evaluation metrics.

Based on these intuitions, we hope to contribute to these two questions:

- How much does languages with a large number of unbounded grammatical morphemes benefit from current evaluation schemes?
- Would focusing on correct classification of content word dependencies be a better evaluation scheme for cross-linguistic parsing performance?

## 2 Related work

Not much work has been done in cross-linguistic evaluation, and papers presenting evaluation scores on several languages simply use previously available metrics without analyzing their shortcomings in such a context. Before UD was publicly available, there have been several attempts at automatic normalization of dependency treebanks into a common format for a more robust evaluation (Zeman et al., 2012). In light of this work on cross-linguistically consistent annotation frameworks, Tsarfaty et al. (2011) take a separate approach with cross-framework evaluation, where they suggest an evaluation technique that is robust towards differing annotation criteria.

Since UAS arguably has been more popular of the two attachment scores, finding performance results on subsets of dependency relations has been difficult. In cases of extensive error analysis it is possible to find sections devoted to this (Plank, 2011). There have also been work looking at specific constructions, e.g. unbounded dependency evaluation (Nivre

Treebank	Token size
Arabic	282K
Basque	121K
Bulgarian	156K
Croatian	87K
Czech	1503K
Danish	100K
Dutch	200K
English	254K
Finnish	181K
Gothic	56K
Greek	59K
Hebrew	115K
Hindi	351K
Italian	252K
Norwegian	311K
Old Church Slavonic	57K
Persian	151K
Polish	83K
Portuguese	212K
Slovenian	140K
Spanish	423K
Swedish	96K

Table 1: Selected treebanks from the UD 1.2 treebank collection, with their token size.

et al., 2010), where they argue that some dependency relations are more critical for a parser to get right than others.

Plank et al. (2015) look closer at whether or not manual parsing evaluation correlate with standard dependency metrics, coming to the conclusion that none of today’s established metrics are especially well correlated with human quality judgment. One of their main findings is that humans tend to consider content dependencies to be of more importance than function dependencies, which fits well with the assumptions made in this paper.

## 3 Data

We will be using a subset of the Universal Dependencies treebank 1.2 (Nivre, 2015). To keep them as internally consistent as possible, all treebanks must adhere to the following criteria:

- They have morphological features.

Function relations	Content relations
aux auxpass case cc	acl advcl advmod
cop det expl mark	amod appos ccomp
neg mwe	compound conj
	csubj csubjpass
	dislocated dobj iobj
	list name nmod
Other	nsubj nsubjpass
list dep foreign	nummod parataxis
reparandum punct	remnant root
goeswith discourse	vocative xcomp

Table 2: Classification of content and function relations.

- They have at least 30K tokens.
- They have a small ratio of non-projective trees.
- In the case of more than one valid treebank for a language, choose the treebank with manual corrections or largest token count.

A total of 15 treebanks were removed. 5 of these had too few tokens, 4 lacked features, 2 had too many non-projective trees, while 4 treebanks were language duplicates. This leaves us with the 22 languages listed in table 1. Most notably we lost the French and German treebanks.

For measuring the correlation of metrics to manual evaluation, we will be using parts of the human judgment data as provided by Plank et al. (2015). Not all languages in the dataset are from the UD treebanks, thus only English, German, and Spanish are used.

## 4 Experimental setup

For testing our alternative evaluation metrics, we train MaltParser 1.7 using Nivre Arc-Eager with default settings on each treebank’s training data and parse the included test data (Nivre et al., 2006). We will also be using the human judgment data from Plank et al. (2015) to see how any of our proposed metrics compares against manual evaluation. Before continuing, we must split up the UD dependency relations into categories of function and content relations.

We motivate our categorization based on the spec-

ification of universal dependency relations<sup>1</sup>, linguistic intuition, and a newly developed statistical method. First, let us define what we consider to be function and content relations:

- A *function relation* is a relation that links a word with a function word.
- A *content relation* is a relation that is either the root, or links a content word with another content word.

### 4.1 Categorization by specification

Given the previous definition as well as the specification of universal dependency relations, we can categorize a relation based on how it should occur in UD treebanks. Going through each dependency relation in this manner we ended up with a classification as presented in table 2.

We chose to remove some relations where we cannot make assumptions of its content, labeled *other*. The *foreign* relation has no restrictions on what type of word it should choose as a dependent as long as it is a foreign word. *List* are used in cases where the content cannot be easily analyzed. *Reparandum* and *dep* are neither of semantic nor syntactic nature and we cannot make any assumptions of their content. *Punct* and *discourse* are removed for similar reasons, but also due to particles often being ignored in other evaluation schemes.

### 4.2 Placing dependency relations on a function–content spectrum

Next question to answer is if we can motivate our classification not only on linguistic intuition and the specification of dependency relations, but also from an empirical perspective. We do this by adhering to our previous definition on function dependencies, and what we know of the nature of function words. Since they are part of closed word classes, meaning new words rarely get introduced into their categories, we can expect the number of distinct word types to be quite small, especially when compared to the word classes shared by typical content words such as *nsubj*.

Assuming this holds, we expect that the probability of a word given a function relation to be zero for

<sup>1</sup><http://universaldependencies.github.io/docs/u/dep/index.html>

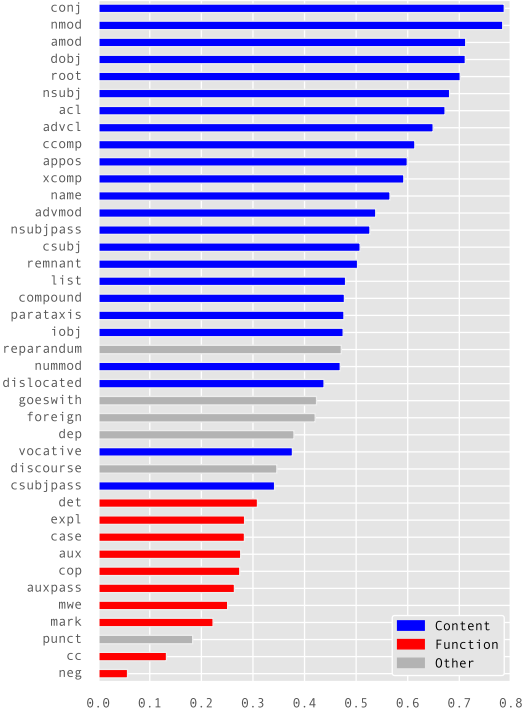


Figure 2: Averaged word dependency entropy for all universal dependency relations, with the manually created categories.

all but a few cases, while the probability of a word given a content relation should be much more evenly spread. This can be quantified by measuring a relation’s entropy given its word probabilities. We call this measure *word dependency entropy* (WDE). Calculated for all treebanks we get the averaged WDE.

Here we formally define word dependency entropy. A probability distribution  $p$  takes as input a word  $w$  conditioned on a treebank  $t \in T$ , and a dependency relation  $r$ .  $H$  is the entropy function that takes  $p(w|r, t)$  as input and calculates its entropy. The entropy function is normalized by its upper bound  $\log n_w$ , where  $n_w$  is the size of the vocabulary. This keeps the range of the function to  $[0, 1]$ . To calculate the WDE for a set of treebanks, we average WDE for all treebanks  $t \in T$ . In the case a dependency relation is not present in a treebank, we set its WDE to 0.5 to imply that it is neither a content nor a function relation.

This gives us the following mathematical functions:

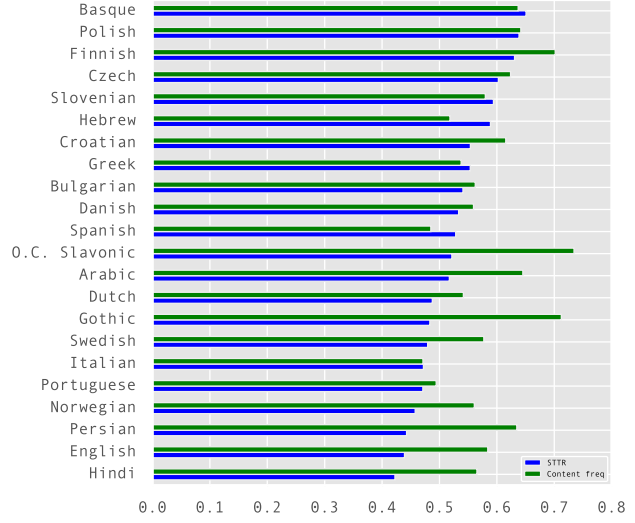


Figure 3: Standardized type/token ratio for chunks of 1000 tokens (blue) and content dependency (green).  $R = 0.27$

$$\text{WDE}(r, t) = \frac{H(p(w|r, t))}{\log n_w}$$

$$\text{Averaged WDE}(r, T) = \frac{1}{n_T} \sum_{t \in T} \text{WDE}(r, t)$$

The averaged WDE for UD 1.2 is presented in figure 2, along with our manual categorization of dependency relations. Ignoring the *other* relations, we find that the WDE gives an intuitive ordering of the dependency relations, while empirically supporting our choice of content and function relations.

### 4.3 Finding correlation with external measurements

We would expect the ratio of function words in a given language’s treebank to correlate with its degree of synthesis. Measuring degree of synthesis is not a trivial task, and there have been several proposed algorithms for this. Despite having obvious drawbacks, an often used indirect measurement of degree of synthesis is the type/token ratio (Kettunen, 2014). This assumes that synthetic languages, with their morphologically rich systems, will have fewer tokens per word than analytic languages such as English or Hindi. This is not particularly robust when comparing across corpora of

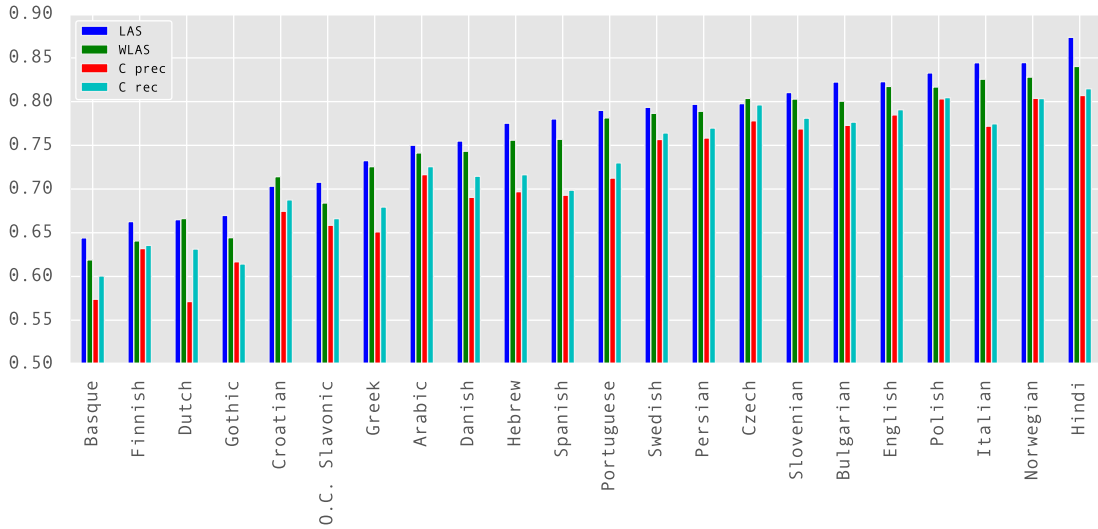


Figure 4: Overall LAS, WLAS, precision and recall for content dependencies. Sorted by LAS.

different sizes. As such, we will be using the *standardized* type/token ratio (STTR), which calculates average TTR in chunks of 1000 tokens<sup>2</sup>. Figure 3 shows that there is a weak correlation between languages’ frequency ratio of function relations and their STTR. Some languages are clear outliers such as Old Church Slavonic, Arabic, Gothic, Persian, English, and Hindi. Whether this is a result of a bad degree of synthesis measurement, a bad classification of function dependency relations, or a combination of both is up for discussion.

#### 4.4 Alternative evaluation metrics

Based on the previous findings, we propose two alternative metrics to the LAS. The first metric is based on our manual classification of content and function dependencies, while the latter is exploiting the weights outputted by the WDE.

**Precision and recall of content relations** In this metric, we look at the precision and recall for all content dependency relations, ignoring any relation that is not a part of this class. Since not all dependency relations are involved, the precision and recall can differ and thus become interesting to analyze separately. We call these *content precision* and

*content recall*.

**Weighting relations by their WDE** We weight each dependency relation by its averaged WDE as presented in figure 2. This will increase the importance of content relations, while the function relations provide less to the overall score. We call this the *Weighted Labeled Attachment Score* (WLAS). Using WLAS has the additional interesting property of also being easily calculated and deployable to non-UD frameworks.

## 5 Evaluation

Figure 4 lists the parsing results for all languages, with LAS, WLAS, and content precision and recall. We can tell that LAS is consistently providing higher scores for each output compared to WLAS, while the content precision and recall scores are substantially lower. There are small differences for precision and recall for all except the worst performing languages where, given the higher recall, the parsing model seems to have a bias towards content dependencies.

Table 3 lists the Pearson correlation coefficient between treebanks for WLAS, LAS, content relations precision and recall, function relations precision and recall, and the frequency ratio of content and function relations in the treebanks. The correlation between the various suggested measurements,

<sup>2</sup>Introduced by Mike Scott in [http://lexically.net/downloads/version6/HTML/index.html?type\\_token\\_ratio\\_proc.htm](http://lexically.net/downloads/version6/HTML/index.html?type_token_ratio_proc.htm)

	LAS	C prec	C rec	F prec	F rec	C freq	F freq
WLAS	0.99	0.96	0.98	0.81	0.85	-0.50	0.42
LAS		0.96	0.96	0.83	0.89	-0.48	0.43
C prec			0.99	0.66	0.75	-0.25	0.19
C rec				0.71	0.76	-0.33	0.23
F prec					0.97	-0.73	0.67
F rec						-0.67	0.62
C freq							-0.86

Table 3: Pearson correlation matrix between treebanks for content and function frequency ratio, content precision and recall, function precision and recall, LAS, and WLAS. Boldfaced figures are mentioned in the discussion.

	C precision	C recall	LAS	WLAS
English	0.51	<b>0.58</b>	0.55	0.55
German	0.37	0.44	<b>0.46</b>	<b>0.46</b>
Spanish	0.45	0.46	<b>0.49</b>	0.48

Table 4: Correlations of LAS, WLAS, content precision and recall against human judgment data.

as well as with LAS, are quite strong. The content relations frequency has a strong negative correlation with function relations frequency, and negative across all the other measurements as well. Function relations frequency shows more or less strong correlations across.

Figure 5 shows what happens with the variance when cumulatively adding languages in a top-scoring fashion for each measurement. For content precision and recall as well as WLAS, the variance is lower for high performing languages, but takes off for content precision and recall when you get past the first eleven treebanks. WLAS keeps a lower score compared with LAS until the end where it joins LAS in a variance of 0.0045.

Table 4 lists Spearman correlations between manual evaluation of two parser models, where the evaluators given parsed sentences in each case chose which of two parser models they consider to provide the best output. Content precision and recall are in all cases except one inferior to LAS and WLAS, where the two latter are indistinguishable.

## 6 Discussion

Looking at function frequency and its precision in table 3, they have a correlation of 0.67. This suggests that the larger rate of function words in a lan-

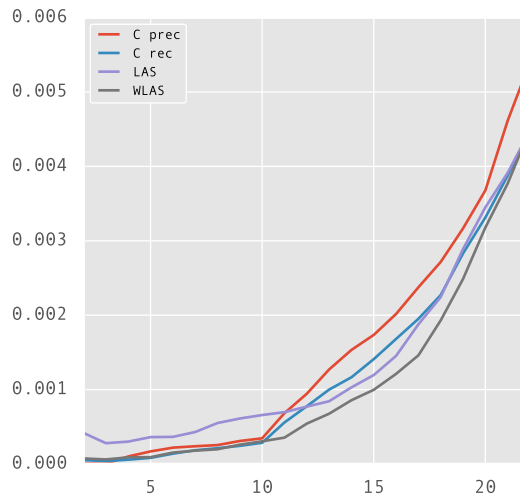


Figure 5: Cumulative variance when adding languages in a top-scoring order.

guage, the easier it is to parse its function words. What is interesting is that this does not hold when looking at content frequency and its precision where one might expect that there is a strong positive correlation which would indicate that a high degree of content relations makes it easier to parse these classes. Instead, we find a weak negative correlation of  $-0.25$ . Furthermore, given languages' LAS scores, the more function words there is in a language, the better it performs. Even when looking at the relation between the amount of function words in a language with how well it does on content words, the correlation is weakly positive. We take this to mean that our choice of parser, while not explicitly tuned for any particular language, still benefits from a context with a high rate of grammatical function words. These findings support our hypothesis, that languages with a high degree of function dependency relations has an unfair advantage when comparing attachment scores across languages.

Regarding the evaluation scores for different evaluation metrics, as presented in 4, it is difficult to tell if any metric is better than the other. One might expect that the scoring difference of LAS and WLAS, or LAS and content performance, would correlate with its STTR score, since languages like Hebrew and English have more to lose on decreasing the im-

portance of function words, but unfortunately this correlation is rather weak. We believe this has to do with what we described above, that languages with a high rate of function words provide a better context for content words for parsers.

We ran the measurements on the human judgment data with the results given in table 4. Unfortunately, none of the metrics seems to improve upon the LAS score, which was the top scoring metric reported in the original paper. Only content recall has some improvements over LAS for English, but other than that the results are either worse or equal to those of LAS. WLAS has overall very small changes compared to LAS, which is somewhat surprising given that the original paper commented on content relations being considered more important than function relations by the manual evaluators. This could possibly be explained by function relations overall performing quite well, and whenever there are erroneous function relations they are rather a result of already faulty content dependencies. This possible explanation is as of yet untested.

Another approach is to look at the variance of the metrics. If our initial hypothesis holds, this should mean that some of the differences found between languages before evens out, and thereby lowering the variance when compared to LAS. As figure 5 shows this does not hold when looking at all treebanks, but by cumulatively adding treebanks it is possible to study the effect as the performance decreases. This shows indeed that the variance is much lower among high-performing languages for WLAS and content performance. This could potentially be explained by that differences among top-scoring can much less be explained by randomness from a poor parsing model, and it is first among these that the effects of choice of metric really matters.

That brings us to the choice of parser model and its effect on the results. In the name of consistent treebanks, we chose to remove treebanks with a large number of non-projective trees. Another motivation was that our parser can not handle these types of trees especially well, and the results are thus unreliable when comparing across languages. It is quite possible, and even probable, that there are many similar factors playing a role in the discussed results. In future work, we should reproduce these results with alternative models and look for any similari-

ties or differences that might strengthen our claims or explain some of the peculiarities we have seen in this work.

## 7 Conclusion

In this paper we have presented experiments that suggest that languages with many function dependency relations are easier to parse than languages with richer morphology, given current parser models. We have motivated the necessity for new evaluation metrics that takes this into account from a typological perspective, while also referring to research that motivates this from human judgment standards. We suggested two new evaluation metrics that raise the importance of content dependency relations with some initial experiments indicating their usefulness.

## 8 Acknowledgements

We thank Jörg Tiedemann for providing his parsing output on UD 1.0 for initial experiment development.

## References

- Michael Collins. 1999. *Head-Driven Statistical Models for Natural Language Parsing*. PhD thesis, University of Pennsylvania.
- Kimmo Kettunen. 2014. Can Type-Token Ratio be Used to Show Morphological Complexity of Languages? *Journal of Quantitative Linguistics*, 21(3):223–245, July.
- Joakim Nivre. 2015. Towards a Universal Grammar for Natural Language Processing. In *Computational Linguistics and Intelligent Text Processing*, pages 3–16. Springer, editions.
- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006. Maltparser: A data-driven parser-generator for dependency parsing. In *Proceedings of LREC*, volume 6, pages 2216–2219.
- Joakim Nivre, Laura Rimell, Ryan McDonald, and Carlos Gomez-Rodriguez. 2010. Evaluation of dependency parsers on unbounded dependencies. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 833–841. Association for Computational Linguistics.
- Barbara Plank. 2011. *Domain adaptation for parsing*. PhD thesis, University of Groningen.

Barbara Plank, Héctor Martínez Alonso, Željko Agić, Danijela Merkle, and Anders Søgaard. 2015. Do dependency parsing metrics correlate with human judgments? In *Eighteenth Conference on Computational Natural Language Learning (CoNLL 2015)*.

Reut Tsarfaty, Joakim Nivre, and Evelina Andersson. 2011. Evaluating dependency parsing: robust and heuristics-free cross-annotation evaluation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 385–396. Association for Computational Linguistics.

Daniel Zeman, David Mareček, Martin Popel, Loganathan Ramasamy, Jan Štěpánek, Zdeněk Zabokrtský, and Jan Hajič. 2012. HamleDT: To Parse or Not to Parse? In *LREC*, pages 2735–2741.