

Norwegian Universal Dependencies

Lilja Øvrelid, Petter Hohle

Department of Informatics

University of Oslo

`liljao@ifi.uio.no`, `pettehoh@ifi.uio.no`

Abstract

This article describes the conversion of the Norwegian Dependency Treebank to the Universal Dependencies scheme. This paper details the mapping of PoS tags, morphological features and dependency relations and provides a description of the structural changes made to NDT analyses in order to make it compliant with the UD guidelines. We further present the PoS-tagging and dependency parsing experiments which contrast the performance of several state-of-the-art taggers and parsers on the original and converted treebank. The full converted treebank was made available with the 1.2 release of the UD treebanks.

Keywords: Dependency treebanks, Universal Dependencies, dependency parsing

1. Introduction

With the increasing popularity of dependency-based representations of syntactic structure in recent years, a wealth of different dependency annotation schemes have surfaced. It has been shown that the choice of dependency scheme influences parsing results (Schwartz et al., 2012) as well as down-stream applications (Elming et al., 2013) and even though attempts have been made to contrast different schemes theoretically (Ivanova et al., 2012), it is clear that the diversity of representation makes comparisons difficult. Cross-linguistically even more so, and it can often be difficult to tease apart aspects of annotation scheme from typological differences in cross-lingual learning (Søgaard, 2011; Skjærholt and Øvrelid, 2012).

Universal Dependencies (UD) (de Marneffe et al., 2014; Nivre, 2015) is a recent community-driven effort to create cross-linguistically consistent syntactic annotation. UD is based on the Stanford dependency scheme (de Marneffe et al., 2006) which has become a widely used dependency scheme for English in recent years. A number of existing dependency treebanks have been converted to UD (Pyysalo et al., 2015; Nivre, 2014) and new data has also been annotated from scratch in order to enable multilingual parser development, cross-lingual learning and typological studies of syntactic structure. Treebanks involved in this effort represent a diverse range of languages such as English, German, Swedish, Spanish, Italian, Persian, Japanese, and the UD release 1.1 contains treebanks for as many as 34 different languages of varying sizes.

This paper describes a fully automatic conversion procedure for the Norwegian Dependency Treebank (NDT) to UD. Due to differences both in the tag set, as well as structural analyses, the conversion requires non-trivial transformations of the dependency trees, in addition to mappings of tags and labels that make reference to a combination of various kinds of linguistic information. This paper details the mapping of PoS tags, morphological features and dependency relations and provides a description of the structural changes made to NDT analyses in order to make it compliant with the UD guidelines. We further present the PoS-tagging and dependency parsing experiments which contrast the performance of several state-of-the-art taggers

Head	Dependent
Finite verb	Complementizer
Finite auxiliary	Lexical verb
Infinitival marker	Lexical verb
Preposition	Prepositional complement
Noun	Determiner
First conjunct	Subsequent conjuncts

Table 1: Annotation choices in the NDT

and parsers on the original and converted treebank. The full converted treebank was made available with the 1.2 release of the UD treebanks (Nivre et al., 2016).

2. NDT and UD

The Norwegian Dependency Treebank (NDT) (Solberg et al., 2014) contains morphological and syntactic dependency annotation for both varieties of written Norwegian (Bokmål and Nynorsk).¹ The morphological annotation follows the Oslo-Bergen Tagger scheme (Hagen et al., 2000). The syntactic annotation scheme is, to a large extent, based on the Norwegian Reference Grammar (Faarlund et al., 1997) and the dependency representations are inspired by choices made in comparable treebanks, in particular the Swedish treebank Talbanken (Nivre et al., 2006). Skjærholt (2014) quantified inter-annotator agreement using a chance-corrected metric derived from Krippendorff’s α and showed that agreement on the NDT data is high: scoring an α of about 98%, among the highest of all the data sets studied. The annotation guidelines (Kinn et al., 2013) describe the annotation scheme in some detail and Table 1 summarizes the main annotation choices in NDT (Solberg et al., 2014).

Universal Dependencies (UD) builds on several previous initiatives for universally common morphological (Zeman, 2008; Petrov et al., 2012) and syntactic dependency (McDonald et al., 2013; Rosa et al., 2014) annotation. Among

¹The current Norwegian UD treebank contains the data from the Bokmål section of the treebank which consists of 311,000 tokens. There are plans to include the Nynorsk data in the next UD release.

NDT	UD
adj	ADJ
adv	ADV
clb	PUNCT, SYM
det	DET, NUM
konj	CONJ
interj	INTJ
inf-merke	PART
prep	ADP, ADV
pron	PRON
<komma>	PUNCT
sbu	SCONJ
<strek>	PUNCT
subst	NOUN, PROP
<anf>	PUNCT
<parentes-slutt>	PUNCT
<parentes-beg>	PUNCT
symb	SYM
ukjent	X
verb	AUX, VERB

Table 2: Mapping between NDT and UD parts-of-speech

its main tenets are the primacy of content-words, i.e. content words, as opposed to function words, are syntactic heads wherever possible. It is intended to be a universal annotation scheme, i.e. applicable to any language, however also offers some possibilities for language-specific information. With reference to the NDT annotation choices in Table 1, the UD scheme adopts the reverse attachment for auxiliaries, infinitival markers and prepositions.

3. Parts-of-speech

The part-of-speech tag set used in the UD scheme is based on the Universal PoS tag set of Petrov et al. (2012) and contains 17 tags. The NDT tag set contains 19 tags. The conversion of the part-of-speech information in NDT to the UD pos tag set is fairly straightforward and largely relies on a direct mapping presented in Table 2. A few parts-of-speech require conversion rules which make reference to additional information in the treebank, represented by disjunction in the mapping. Below we will discuss a few of these cases.

The universal scheme makes a distinction between proper and common nouns at the part-of-speech level. This information can be found among the morphological features in NDT (*prop*), hence the mapping is straightforward.

For verbs UD distinguishes auxiliaries (AUX) from main verbs (VERB). This distinction is not explicitly made in NDT, hence our conversion procedure must make use of the syntactic structure of the verbs in order to implement this distinction. Verbs that have a direct, non-finite dependent (a dependent with the NDT dependency relation *INF*) are marked as auxiliaries and all other verbs as regular verbs.

The relative pronoun *som* 'who/that' is analyzed as a subjunction in NDT, whereas the universal scheme, and thus our conversion, uses the pronominal tag *PRON* for these words.

NDT	UD
mask, fem, nøytt	Gender=Masc, Fem, Neut
ent, fl	Number=Sing, Plur
be, ub	Definite=Def, Ind
pres, pret	Mood=Ind, Tense=Pres, Past, VerbForm=Fin
perf-part	VerbForm=Part
imp	Mood=Imp, VerbForm=Fin
pass	Voice=Pass
inf	VerbForm=Inf
1, 2, 3	Person=1, 2, 3
nom, akk, gen	Case=Nom, Acc, Gen
pos, komp, sup	Degree=Pos, Cmp, Sup
hum	Animacy=Anim
pers	PronType=Prs
dem	PronType=Dem
sp	PronType=Int
res	PronType=Rcp
poss	Poss=Yes
refl	Refl=Yes

Table 3: Mapping between NDT and UD morphological features

NDT explicitly marks headings using the | symbol, tagged as a clause boundary (*clb*) along with information in the morphological features *<overskrift>*. These are converted to the UD symbol tag *SYM*.

Numerical expressions, e.g. *to* 'two' or *45* are tagged as determiners in NDT, also when they do not explicitly modify a nominal, e.g. *de to gjorde en imponerende innsats* 'the two did an impressive job'. These are however marked as 'kvant' (quantificational) in the morphological features, hence the conversion mapping makes reference to this property in order to distinguish *NUM* from *DET*.

NDT implements a somewhat broader definition of prepositions than the UD scheme and includes several cases which are counted as adverbs in the universal schemes. In particular, this is the case for demonstrative adverbs such as *her* 'here' and *der* 'there'. These are followingly converted to adverbs (*ADV*) in the PoS conversion, whereas all other prepositions are converted to *ADP*.

4. Morphological information

In addition to part-of-speech information, NDT contains a rich inventory of morphological features, e.g. information about properties like gender, definiteness, tense, voice, etc. The UD guidelines specify a universal set of morphological features and the conversion between the two does not require reference to information in addition to the feature information. The feature mapping is described in Table 3. Note that since the number of UD features is larger than the NDT features, some of the NDT features correspond to a set of UD features, e.g. the NDT features for verbs (*pres*, *pret*) which instantiate both the *Mood*, *Tense* and *VerbForm* features.

5. Structural conversion

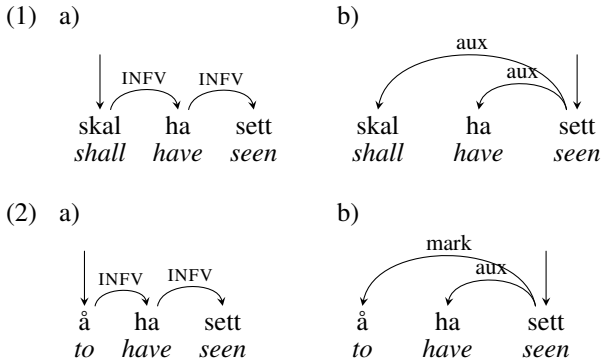
The NDT annotation scheme differs structurally from the UD scheme in a number of important ways. The conversion is therefore non-trivial and requires a set of structural rules which operate on the dependency graphs in addition to a mapping procedure over the dependency labels. The conversion is implemented as a cascade of structural rules followed by a relation conversion procedure over the modified graph structures. The structural rules employ a small set of graph operations that reverse, reattach, delete and add arcs.

5.1. Root

In NDT, there is no designated root label. Rather, the root of the dependency graph may have different labels, e.g. FINV (finite verb) or FRAG (fragment), depending on the structure of the sentence. In the conversion, every node with the dummy node (0) as head receives the `root` relation.

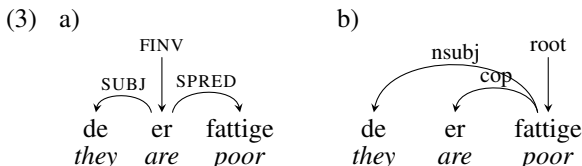
5.2. Verbal groups

NDT consistently marks the finite verb as head of a clause, with other non-finite verbs as dependents (INFV), see example (1a). In a parallel manner, infinitival markers are also annotated as heads with the infinitival verb as its dependent, see (2a). UD on the other hand annotates the lexical, main verb as head of the verbal group and various finite and non-finite auxiliaries receive an auxiliary relation (`aux`, `auxpass`), see (1b) and (2b) below. The conversion rule locates the main verb within the chain of nonfinite dependents of the finite verb and makes this node the head of the other verbs in the chain.



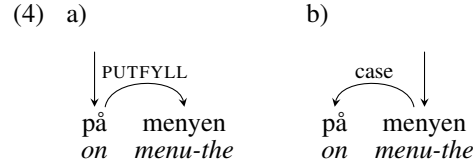
5.3. Copula constructions

The treatment of copula constructions within the UD scheme differs markedly from that of the NDT by appointing the predicative element as head of the entire construction and attaching the copula verb with a special relation `cop`, see (3b). Our conversion thus reverses the arc between the copula and its complement and reattaches all its dependents to the predicative element.



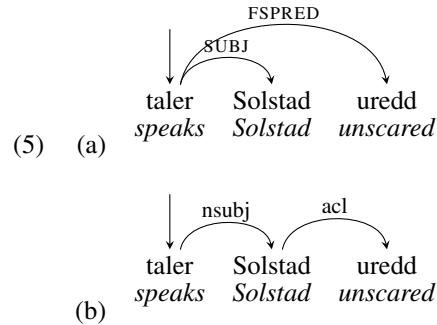
5.4. Prepositions and their complements

In NDT, prepositions are heads of their prepositional complements which receive the `PUTFYLL` label, see (4a). Seeing that languages differ greatly in their use of pre/postpositions, the UD scheme annotates the prepositional complement as head and attaches the preposition using the `case` relation, see (4b). In conversion this is once again a matter of reversing the arc and reattaching dependents to the new head.



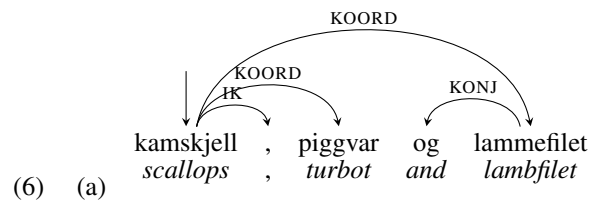
5.5. Predicatives

NDT distinguishes several types of predicatives, both predicatives that are arguments of verbs (subject predicative `SPRED` and object predicative `OPRED`) and “free predicatives” which are not arguments of the verb, but nonetheless characterize either a subject or an object in the preceding context (free subject predicative `FSPRED` and free object predicative `FOPRED`). Both of these are attached to the finite verb in NDT, as we can see in (5a). In a similar manner, UD distinguishes between obligatory and optional predicatives, where the former are analyzed using the `xcomp` relation and attached to the main predicate, whereas the optional predicatives are attached as adverbial clauses (`acl`) modifying the argument they characterize, see (5b). Our conversion thus attaches the `FSPRED` argument to its sibling subject argument, `FOPRED` to an object sibling.



5.6. Coordination

Coordination is a phenomenon which exhibits considerable variation in terms of dependency representation across various annotation schemes (Popel et al., 2013). As we can see from the example in (6), the analyses chosen in the NDT and UD schemes are fairly similar in their choice of the first conjunct as head of the coordinate structure. They differ mainly in the attachment of the conjunction and the relation names.



NDT	UD
ADV	advcl, advmod, compound:prt, neg, nmod
ATR	acl:relcl, amod, nmod
DET	nmod, nummod, det
FINV	aux, auxpass, root
FLAT	foreign, name
OBJ, POBJ	dobj, ccomp, xcomp
SBU	nsubj, nsubjpass, dobj, iobj, mark
SUBJ, PSUBJ	nsubj, nsubjpass, csbj, csubjpass

Table 5: Non-direct mapping between NDT and UD dependency relations; requires additional constraints with reference to PoS, morphological features or dependency context.

Data set	Tokens	Sentences
no-ud-train	244766	15696
no-ud-dev	36467	2410
no-ud-test	30034	1939
Total	20045	311277

Table 6: Overview of the Norwegian UD train, dev and test data sets.

7. The converted treebank

Following the conversion of the Norwegian Dependency Treebank to Universal Dependencies scheme, 51.5% of the tokens in the original treebank were reattached. The resulting treebank contains 17 PoS tags and 35 different morphological features for 311,277 tokens of Norwegian Bokmål. All UD treebanks consist of three data sets: a training, development and test set. In creating these data splits for Norwegian, care has been taken to preserve contiguous texts in the different splits and also to keep a balance of genres in each of the splits. Table 6 shows an overview of the Norwegian UD data sets in terms of tokens and sentences.

8. Tagging and parsing

Without a gold standard for Norwegian UD dependencies it is difficult to evaluate our conversion directly. We may however, evaluate PoS-tagging and dependency parsing performance for the converted treebank. In the following we report on a set of experiments which investigate the performance of a set of state-of-the-art PoS-taggers and parsers trained and evaluated on the converted Norwegian UD treebank.

As noted in previous work (de Marneffe et al., 2014), several of the design choices of the UD scheme, such as the attachment of auxiliaries and prepositions, are known to cause a drop in parsing accuracy (Schwartz et al., 2012).

Data set	Tag set	Accuracy	
		NDT	UD
Dev	Coarse/original	97.90%	96.96%
Dev	Fine/full	93.74%	94.59%
Test	Coarse/original	97.82%	96.82%
Test	Fine/full	93.19%	94.15%

Table 7: Overview of the results from tagging NDT and UD with the various data sets and tag sets

Danish parsing experiments: (Johannsen et al., 2015) orig 84.38 -> ud 81.56 Bulgarian (Osenova and Simov, 2015): orig 89.14 -> ud 83.5

The tagging was performed by SVMTool³ employing strategy 1, which proved to be optimal for parsing NDT in previous work ((Hohle, 2016; forthcoming)). The original tag set of NDT comprises 19 tags, 12 of which are morphosyntactic tags, while its fine-grained counterpart (concatenation of tag and set of morphological features) totals 368 tags. The corresponding fine-grained tag set of UD comprises 169 tags. We see that tagging accuracy is higher on NDT with the original tag set, while UD surpasses NDT when using the fine-grained tag set. This is as expected, as the fine-grained tag set of NDT contains almost 200 more tags than that of UD, which markedly complicates the tagging.

9. Bibliographical References

- de Marneffe, M.-C., MacCartney, B., and Manning, C. D. (2006). Generating typed dependency parses from phrase structure parses. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, pages 449–454.
- de Marneffe, M.-C., Dozat, T., Silveira, N., Haverinen, K., Ginter, F., Nivre, J., and Manning, C. D. (2014). Universal Stanford dependencies. A cross-linguistic typology. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, pages 4585–4592.
- Elming, J., Johannsen, A., Klerke, S., Lapponi, E., Martinez, H., and Søgaaard, A. (2013). Down-stream effects of tree-to-dependency conversions. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 617–626, Atlanta, Georgia, USA.
- Faarlund, J. T., Lie, S., and Vannebo, K. I. (1997). *Norsk referansegrammatikk*. Universitetsforlaget, Oslo, Norway.
- Hagen, K., Johannessen, J. B., and Nøklestad, A. (2000). A constraint-based tagger for norwegian. In *17th Scandinavian Conference in Linguistics*, pages 31–48.
- Ivanova, A., Oepen, S., Øvrelid, L., and Flickinger, D. (2012). Who Did What to Whom? A contrastive study of syntacto-semantic dependencies. In *Proceedings of the Sixth Linguistic Annotation Workshop*, pages 2–11, Jeju, Republic of Korea.

lund et al., 1997)

³cs.upc.edu/~nlp/SVMTool

- Johannsen, A., Alonso, H. M., and Plank, B. (2015). Universal dependencies for danish. In *Proceedings of Treebanks and Linguistic Theories (TLT14)*.
- Kaplan, R. and Bresnan, J. (1982). Lexical Functional Grammar. A formal system for grammatical representation. In Joan Bresnan, editor, *The Mental Representation of Grammatical Relations*, pages 173–281. MIT Press, Cambridge, MA, USA.
- Kinn, K., Eriksen, P. K., and Solberg, P. E. (2013). Retningslinjer for morfologisk og syntaktisk annotasjon i språkbankens gullkorpus. Technical report, National Library of Norway.
- McDonald, R., Nivre, J., Quirmbach-Brundage, Y., Goldberg, Y., Das, D., Ganchev, K., Hall, K., Petrov, S., Zhang, H., and Täckström, O. (2013). Universal dependency annotation for multilingual parsing. In *Proceedings of Association for Computational Linguistics (ACL)*, pages 92–97.
- Nivre, J., Nilsson, J., and Hall, J. (2006). Talbanken05: A swedish treebank with phrase structure and dependency annotation. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)*, Genoa, Italy.
- Nivre, J., de Marneffe, M.-C., Ginter, F., Goldberg, Y., Hajič, J., Manning, C. D., McDonald, R., Petrov, S., Pyysalo, S., Silveira, N., Tsarfaty, R., and Zeman, D. (2016). Universal dependencies v1: A multilingual treebank collection. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*.
- Nivre, J. (2014). Universal Dependencies for Swedish. In *Proceedings of the Swedish Language Technology Conference (SLTC)*.
- Nivre, J. (2015). Towards a universal grammar of natural language processing. In *Proceedings of the International Conference on Intelligent Text Processing and Computational Linguistics (CICLING)*.
- Osenova, P. and Simov, K. (2015). Universalizing BulTreeBank: a linguistic tale about glocalization. In *Proceedings of the 5th Workshop on Balto-Slavic Natural Language Processing*.
- Petrov, S., Das, D., and McDonald, R. (2012). A universal part-of-speech tagset. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, pages 2089–2096.
- Popel, M., Mareček, D., Štěpánek, J., Zeman, D., and Žabokrtský, Z. (2013). Coordination Structures in Dependency Treebanks. In *Proceedings of Association for Computational Linguistics (ACL)*, pages 517–527.
- Pyysalo, S., Kanerva, J., Missilä, A., Laippala, V., and Ginter, F. (2015). Universal dependencies for finnish. In *Proceedings of the Nordic Conference of Computational Linguistics (NODALIDA)*.
- Rosa, R., Mašek, J., Mareček, D., Popel, M., Zeman, D., and Žabokrtský, Z. (2014). Hamledt 2.0: Thirty dependency treebanks stanfordized. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC 2014)*, pages 2334–2341.
- Schwartz, R., Abend, O., and Rappoport, A. (2012). Learnability-based syntactic annotation design. In *Proceedings of the International Conference on Computational Linguistics COLING*, Mumbai, India.
- Skjærholt, A. and Øvrelid, L. (2012). Impact of treebank characteristics on cross-lingual parser adaptation. In *Proceedings of Treebanks and Linguistic Theories (TLT11)*.
- Skjærholt, A. (2014). A chance-corrected measure of inter-annotator agreement for syntax. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Søgaard, A. (2011). Data point selection for cross-language adaptation of dependency parsers. In *Proceedings of the meeting of the Association for Computational Linguistics*.
- Solberg, P. E., Skjærholt, A., Øvrelid, L., Hagen, K., and Johannessen, J. B. (2014). The Norwegian Dependency Treebank. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*.
- Zeman, D. (2008). Reusable tagset conversion using tagset drivers. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*.