# UPOS-DEPREL Mismatches:
# Detecting Annotation Errors and Improving UD Guidelines based on Linguistic Knowledge

**Tsy Yih[1]**
School of Foreign Languages
Tongji University
`yezi_leafy@hotmail.com`

**Zheyuan Dai**✉
College of Foreign Languages
Zhejiang University of Technology
`zydai329@163.com`

## Abstract

The current study explores the power of the UPOS-DEPREL combination in UD. By examining its distribution in 20 PUD treebanks, we found that most parts of speech have a matching dependency relation, which is supported by linguistic rationales. Based on the matches and minor combinations denoting the phenomenon of transcategorization, we set up a Possible Combination Zone. Those falling out of the zone are then considered errors or inconsistencies in annotation, or underdescribed special constructions in UD. Our findings thus provide a new way to detect annotation errors based on linguistic knowledge. In addition, we discuss two case studies on improving the current UD guideline enlightened by the combination rules.

## 1 Introduction

Universal Dependencies [2] (henceforth, UD) is either understood as a treebank annotation project or a specific annotation guideline presented in a sheet format with ten columns called CoNLL-U (Nivre et al. 2016, 2020, de Marneffe et al. 2021). Apart from the two columns with numerical values (ID and HEAD) that constitute the skeleton of dependency structures, most other columns are annotated with tags carrying distinctive information, among which the universal part of speech (UPOS) and the universal dependency relations (DEPREL) are two most basic yet

important layers. The tags are distinctive in the sense that each tag is used to represent a certain phenomenon while different phenomena are well distinguished and annotated different tags.

It was inevitable that the proposer(s) could not carefully look into every linguistic phenomenon and provide a perfect guideline when it first came out. Therefore, follow-up researchers have contributed to the improvement of annotations by clarifying the meaning of certain tags and how concrete linguistic phenomena should be related to them (e.g. Ahrenberg 2019), and refining the tag system with respect to certain field (e.g. Schneider & Zeldes 2021).

However, previous studies often focus on one layer at a time, but few pay attention to the combination of tags from different layers. In fact, the combination of UPOS and dependency relations DEPREL [3] has always been a powerful tool for researchers. For instance, to investigate nouns with certain grammatical roles in a sentence, one only needs to search the combinations of UPOS = `NOUN` & DEPREL [4], such as `nsubj`, `obj`, etc. to locate wanted results accurately, a task that cannot be done in a raw or POS tagged corpus. In addition, for any researcher who has some slightest knowledge about UD, it is easy to notice that certain UPOS-DEPREL combinations, such as `ADJ` & `amod`, `DET` & `det`, are quite frequent in UD treebanks. Hence, there must be some regularities if we do a cross-tabulation.

In this paper, we aim to investigate to what extent are UPOS and DEPREL correlated cross-

---

[1] Tsy Yih is the transliteration of the name of the first author in his mother tongue, Shanghai Wu Chinese. He is also known as ZI YE in Mandarin pinyin.
[2] https://universaldependencies.org/
[3] In the rest of this paper, **combination**, without further notice, simply refers to the UPOS-DEPREL combination.

[4] UPOS tags are presented in capitalized letters in Courier New font, such as `NOUN`, `ADJ`, `ADV`, etc. DEPREL tags are in lowercase letters of Courier New font, such as `amod`, `advcl`, `nsubj`, etc. The combination of UPOS-DEPREL thus goes as capitals & lowercase letters, like `ADJ` & `amod`.

linguistically, explore how they can be further employed as an annotation error detector, and see what insights it can provide for improving the current UD guideline.

## 2   Methods

The Parallel Universal Dependencies (PUD) treebank, in UD version 2.10, was chosen as the data source for this study. The PUD treebanks comprise of a set of parallel UD-annotated treebanks with 20 languages. Each treebank is composed of 1,000 sentences in each language, always in the same order. The genres include news and Wikipedia. The sizes of PUD treebanks range from 15,813 (Finnish) to 28,788 (Japanese), mostly around 20,000 tokens.

Currently there are 17 UPOS tags (Petrov et al. 2012) and 37 DEPREL tags with more sub-types (de Marneffe et al. 2014) in the UD guideline version 2. All UPOS tags are roughly attested in all languages except for a few missing (such as CCONJ in Korean). For DEPREL, the number ranges from 25 (Japanese) to 58 (Polish). The combinations of UPOS-DEPREL were extracted and then arranged into a $17 \times n$ table for further data processing and illustration.

We proposed a measure called **utilization rate (UR)**, which is defined as follows. With this measure are we able to find the proportion of combinations in actual use.

$$UR = \frac{attested\ combinations\ of\ UPOS-DEPREL}{all\ theoretical\ combinations\ of\ UPOS-DEPREL} \quad (1)$$

We then differentiated between two-way and one-way matches. A **two-way match** means for the UPOS tag X, its most frequent corresponding DEPREL is Y, whereas for DEPREL tag Y, the most frequent corresponding UPOS is X. Since the number of DEPREL is larger than that of UPOS, it is impossible for each DEPREL to have a two-way match. Therefore **one-way match** of DEPRELs is also needed. That is, UPOS tags co-occur with these DEPREL most frequently but not vice versa. **Mismatch** then in this paper is understood as all those combinations falling out of the scope of our restrictions, that is, neither a two-way match, nor a one-way match.

After identifying the frequent matches between UPOS and DEPREL, we manually searched through the mismatches and then divided them into different types.

| Language | Attested | Theoretical | UR |
|---|---|---|---|
| Arabic | 155 | 40×17=680 | 22.8% |
| Chinese | 169 | 44×17=748 | 22.6% |
| Czech | 177 | 42×17=714 | 24.8% |
| English | 209 | 47×17=799 | 26.2% |
| Finnish | 161 | 43×17=731 | 22.0% |
| French | 156 | 44×17=748 | 20.9% |
| German | 160 | 42×17=714 | 22.4% |
| Hindi | 181 | 38×17=646 | 28.0% |
| Icelandic | 188 | 35×17=595 | 31.6% |
| Indonesian | 140 | 46×17=782 | 17.9% |
| Italian | 170 | 39×17=663 | 25.6% |
| Japanese | 86 | 25×17=425 | 20.2% |
| Korean | 118 | 34×17=578 | 20.4% |
| Polish | 183 | 58×17=986 | 18.6% |
| Portuguese | 149 | 40×17=680 | 21.9% |
| Russian | 186 | 39×17=663 | 28.1% |
| Spanish | 161 | 39×17=663 | 24.3% |
| Swedish | 178 | 41×17=697 | 25.5% |
| Thai | 166 | 42×17=714 | 23.2% |
| Turkish | 164 | 39×17=663 | 24.7% |

Table 1:  The utilization of UPOS-DEPREL combinations in PUD.



Figure 1: The rank-frequency distributions of UPOS-DEPREL combinations in 20 PUD treebanks**.**

2

## 3 Correlation between UPOS and DEPREL

Extracted data show that the UPOS-DEPREL combination matrices in most languages are sparse (See Appendix A for the case in the English PUD). Table 1 lists the utilization rate of the combinations in each language. It can be seen that the utilization rates are generally low, ranging roughly from 20% to 30%. A low UR means most of the actual combinations must be found in cells.

We then investigated the distribution of attested combinations. As shown in Figure 1, the rank-frequency distributions of UPOS-DEPREL combinations in all languages manifest a long-tail, Zipfian, or power-law pattern. Put another way, the majority of the attested combinations is low frequency terms. It is thus indicated that the actual utilization rate may be even lower.

For those cells with zero or a few examples, it might be due to several factors including the size of the treebank, the genre and content of the texts, or the wrong annotation. However, the overall non-uniform pattern reveals innate properties. Certain combinations are indeed more frequent than others, and this finding needs a further look and probably a linguistic explanation.

Hence in what follows, we analyze the matches between UPOS and DEPREL. Table 2 shows the most frequent DEPREL co-occurring with each UPOS with its proportion in 20 languages in the last column. The denominator lower than 20 means that DEPREL is not existent in every treebank. From Table 2 can we read a number of things. First, three nominal UPOS (NOUN, PRON, PROPN) do not have two-way matches with single DEPREL, but if we merge several argument relations, such as nsubj, obj, iobj, and obl, then the regularity manifest itself. For NOUN, the modifier relations and argument relations together take 10% of proportion. In addition, PRON and PROPN are overwhelmingly matched for dependencies expl (expletive or dummy subjects) and flat (multiword proper names) from a one-way perspective. What might be surprising is that in half of languages nouns

| UPOS | DEPREL | Proportions |
|---|---|---|
| Two-way matches | | |
| NOUN | argument (nsubj, obj, iobj, obl) | 9/20 |
| | modifier (nmod, compound) | 11/20 |
| VERB | root | 18/20 |
| DET | det | 17/20 |
| ADJ | amod | 19/20 |
| ADV | advmod | 19/20 |
| NUM | nummod | 20/20 |
| AUX | cop | 7/20 |
| | aux, aux:pass | 12/20 |
| CCONJ | cc | 20/20 |
| SCONJ | mark | 14/20 |
| ADP | case | 18/20 |
| INTJ | discourse | 5/20 |
| PUNCT | punct | 20/20 |
| One-way matches | | |
| VERB | acl | 16/19 |
| VERB | acl:relcl | 18/19 |
| VERB | advcl | 19/20 |
| VERB | xcomp | 19/19 |
| VERB | ccomp | 18/20 |
| PROPN | flat, flat:name | 19/19 |
| PRON | expl | 11/11 |
| NOUN, PROPN | appos | 18/20 |
| NOUN, PRON | nmod:poss | 10/10 |

Table 2: Most frequent matching DEPREL for each UPOS in 20 languages.

appear most in the modifier position[5], which might be counterintuitive at first sight but indeed an interesting phenomena noted in the literature (Croft, 1991: 91-92; 2001: 104). Second, similar to NOUN, AUX basically matches two dependency types, cop and aux, the total amount of which add up to 95% languages. This probably indicates that AUX could have been divided into two

---

[5] The compound relation is defined as one of the three types of multiword expressions in UD (the other two being flat and fixed). However, the working definition or delimitation between nmod and compound, at least in

English PUD, lies in whether this noun modifiers is preposed barely or postposed with the help of adpositions. According to different languages, treebanks, and annotators, this definitely causes diversions. We simply report the current usage here.

categories. Thirdly, the UPOS `PART` (particle) does not have a universal matching dependency across languages. This might be due to the fact that its definition is a hotchpotch. According to the UD guideline[6], PART contains four cases, namely, negation particle, possessive marker, sentence-final particle as in Asian languages and sentential mood marker. Hence, this definition does not manifest a unique semantic prototype, but plays the role of a "garbage" among parts of speech. Moreover, in practice, annotators of different languages and treebanks have varied understandings, and they might not follow the guideline perfectly. The situation apparently explains why it does not have a consistent match. Thus, the UPOS tag `PART` is indeed problematic, evidenced by our data and in need of retrospection. Fourth and finally, while only 25% of the languages have a two-way match between `INTJ` & `discourse`, yet that is due to the low frequency in treebanks. However, according to the description of UD framework, these two seem to be designed to represent the same phenomenon, which do form a two-way match.

Overall, Table 2 reveals that most parts of speech match only one dependency relation, indicating that the DEPREL tag covers the UPOS tag to a large extent.

What is the linguistic rationale underlying this correlation?

For the major parts of speech, we find Croft's model of parts of speech appropriate to explain this distribution (1991, 2001, 2022). His model is spanned by the two dimension of semantic classes and propositional acts[7] as shown in Figure 2, and

| ARG | PRED | MOD |
|---|---|---|
| nsubj<br>obj<br>iobj<br>obl | main clause:<br>`root`<br>`parataxis`<br>subordinate clause:<br>finite:<br>`ccomp`<br>`csubj`<br>`acl:relcl`<br>`advcl`<br>non-finite:<br>`xcomp`<br>`acl`<br>`advcl` | nominal:<br>`amod`<br>`nummod`<br>`nmod`<br>verbal or sentential:<br>`advmod` |

Table 3: The DEPREL tags grouped by their functions

prototypical parts of speech falls at the diagonal of this space.

In addition to the primary ones, in fact Croft also analyzed several minor parts of speech and proposed several minor propositional acts (categorizing, situating, and selecting), although these are less noted by other scholars.

Relating this to UD, the universal parts of speech can be seen as a manifestation of semantic classes, whereas the universal dependency relations stands for propositional acts here except for the different terms or labels. If we combine related dependency relations into group as shown in Table 3, the situation becomes even clearer. According to Croft's theory, the terms on the diagonal of the model are predicted to be most frequent. This has been corroborated by our findings, that is, nouns mostly co-occur with argument DEPREL, verbs appear mostly with predicative DEPREL and adjectives appear mostly as modifiers, apart from one exception (nouns as modifiers being also numerous) which has also been noticed by himself as mentioned above.

However, do the findings above mean that combinations which are not a match are logically impossible? If in an extreme case where UPOS and DEPREL are perfectly correlated, then DEPREL becomes but the sub-classes of UPOS, and UPOS seems redundant and can be completely discarded. But that is not the case. In



|  | Reference | Modification | Predication |
|---|---|---|---|
| **Objects** | UNMARKED NOUNS | genitive, adjectivalizations, PPs on nouns | predicate nominals, copulas |
| **Properties** | deadjectival nouns | UNMARKED ADJECTIVES | predicate adjectives, copulas |
| **Actions** | action nominals, complements, infinitives, gerunds | participles, relative clauses | UNMARKED VERBS |

Figure 2: Croft's model of major parts of speech (excerpted from Croft 2001: 88).

the next section, we explore the second tier of combinations.

## 4 Causes of mismatches

The findings above indicate strong correlation between certain UPOS and DEPREL. In this section, we pay special attention to the English PUD due to the researchers' intelligibility. We manually searched through all the mismatches and divided them into the following types.

**Denoting transcategorization.** The first type of mismatches is also theoretically possible but simply less frequent than typical matches. However, their frequencies are still not so small that approach zero. They can be viewed as the second tier of combinations, describing a phenomenon called transcategorization (Malchukov, 2004; Ježek and Ramat, 2009). The morphological conversion with zero marker is often considered a common strategy, such as in *water*$_N$ > *water*$_V$. The following are examples of some other transcategorial processes realized with syntactic strategies:

| | |
|---|---|
| N > V: *be a teacher* | **NOUN** & `pred` |
| N > A: *a linguistic textbook* | **NOUN** & `mod` |
| V > N: *his leaving* | **VERB** & `arg` |
| V > A: *a crying boy* | **VERB** & `mod` |
| A > N: *the old* | **ADJ** & `arg` |
| A > V: *be smart* | **ADJ** & `pred` |

Recall that in Croft's model, except for the cells on the diagonals, there are also six cells left which represent these non-prototypical phenomena, thereby being less frequent.

This set of combinations plus typical matches found in the last section are together called **Possible Combination Zone**. With this can we set up combining rules for UPOS and DEPREL. In other words, only those combinations falling within the zone are legitimate, while all the rest falling outside should be considered wrong annotations. Therefore, the UPOS-DEPREL combination rule can serve as an error detector.

**Wrong or inconsistent annotations.** A second type regarding low frequency combinations is that they reflect wrong or inconsistent annotations. Empirically the numbers are small compared with the previous type. Here we present a set of examples:

(1a) During this time, Marcelle was often left alone in the room while Piaf and Mômone were out **on** the streets or at the club singing. (En_PUD sent_id = w01138045)

wrong combination: ADP & `obj`

right combination: ADP & `case`

(1b) Mr Osborne signed **up** with a US speakers agency after being sacked in July. (En_PUD sent_id = n01013005)

wrong combination: PART & `compound:prt`

right combination: ADV & `compound:prt`

(1c) However, they were intercepted and had to do battle in Freeman, **close** to the Hudson River. (En_PUD sent_id = w05005087)

wrong combination: ADV & `amod`

right combination: ADJ & `acl`

(1a) shows an apparent annotation error that the function word *on* improperly carries the typical dependency relation of nominals. In (1b), the UPOS of *up* is annotated as PART. Although it is not our intent here to discuss the validity of this analysis, it is at least inconsistent compared with similar examples in English PUD, and it is the UPOS-DEPREL combination that helps us to identity this inconsistency. Finally in (1c), the UPOS and DEPREL tags are both wrong, where the case is a reduced relative clause. This might be due to the dual status of *close* which could either be an adjective or adverb, but here the determination is unfortunately incorrect.

Note that previous attempts to error detecting are often computationally or technically motivated by way of alignment, sentence regeneration, etc. (van Halteren, 2000; Wisniewski, 2018; Lapalme, 2021). Our findings, on the contrary, suggest a linguistically motivated method to do so. This idea can also further be used to produce automatic error detectors. Moreover, it provides a systematic method for researchers who are in need of highly accurate treebanks to facilitate manual checking simply based on the information contained in the CoNLL-U data itself without resorting to other software or application. Otherwise, they might have to go through all the tokens, which is time consuming.

Note that this method is not almighty in that it simply detects those falling out of the matching range but is unable to find out the wrong annotations within the matches. Yet sometimes inconsistencies discovered by this method would help to further find out those hiding within the safe zone. Taking *due* in *due to* in English PUD as an instance, 7 in 8 take the ADP & `case` combination, while only one takes the uncommon ADJ & `case`.

As a matter of fact, since this is a `fixed` expression, the former which takes the majority is contrary to the UD guideline, according to which the first token is supposed to be annotated with its original part of speech rather than that of the whole expression. Yet the word *due* per se is obviously not an adposition, thereby being wrongly annotated. It is due to the mismatch that are we able to locate this error.

**To-be-determined annotations.** The last type of mismatches concerns special constructions, whose treatment are still in dispute and needs further discussion in UD. The UPOS tags `SYM` and `NUM` are mostly involved. A few examples in English include ellipsis constructions, headless constructions, complex quantifiers (e.g. *one of, some of, all of*), money structures with symbols (e.g. $500), scales and units (5 meters), spatio-temporal ranges (20–30 degrees), dates (e.g. March 20), years (treated as proper noun or numeral), etc. Mismatches in the treebanks help us to identify them.

These are beyond the scope of this paper, which await to be settled in future studies. We have seen a number of attempts to extend the UD guideline to well analyze rare phenomena, such as Hassert et al. (2021) on special expressions in technical documents, Höhn (2021) on adnominal pronoun constructions, Rueter et al. (2021) on numerals, Schneider and Zeldes (2021) on "mischievous" nominal constructions, Zeman (2021) on temporal expressions, Tyers and Mishchenkova (2020) on noun incorporation, Przepiórkowski and Patejuk (2019) on nested coordination, Bouma et al. (2018) on expletives, Droganova et al. (2018) and Droganova and Zeman (2017) on elliptical constructions, Schuster et al. (2017) on gapping constructions, to name a few if not all.

To summarize, the mismatches in the treebanks are primarily due to three types of reasons. In the next section, we discuss how UD could be improved in one respect in light of the combination constraint that we put forward.

# 5 Improving the UD annotation—insights from combinations

In this section, we discuss two case studies for improving the current UD guideline, the insights of which come from previous findings.

## 5.1 Multiword expressions

The first and foremost concern is given to multiword expressions (MWE), which intersect with a number of research fields such as Construction Grammar, formulaic language, lexical bundles, etc., and have drawn extensive attention in both linguistics and computational linguistics. In UD, MWE is taken as the cover term and comprises three types of dependency relations. The tag `fixed` primarily represents complex adpositions (e.g. *as well as*), conjunctions (e.g. *in order to*), and adverbs (e.g. *in addition*), while `flat` stands for exocentric proper nouns. The last one `compound` is endocentric, which generally denotes complex content words with clear internal structures. The former two take a head-initial approach, which makes the first token head and all others dependent on it, and at the same time leaves all non-initial tokens carrying a `fixed or flat` DEPREL but keeping their original parts of speech.

The controversial annotations that we look into in this section are generally correct in terms of the current UD guideline, while the results disobey the matches in our sense. According to our counts, there are 211 tokens related to MWE, taking up 33 combinations cells and a large part among mismatches. If the UD guideline could be changed regarding MWE, more combinations would be consistent.

There are two kinds of problems relating to MWE. The first type is that they could have been annotated as MWE but not. For `flat`, it contains the words within a long proper names such as the title of a song, a movie, a book and the like. They are not marked as `flat` according to the current UD guideline and thus likely to have an unnatural combination.

(2) Rafferty recorded a new version of his Humblebums song "Her Father Didn't **Like** Me Anyway" on the album Over My Head (1994). (En_PUD sent_id = w01130102)

    current combination: `VERB` & `appos`
    ideal combination: `VERB` & `flat`

In this example, the head of the song name *like* lies at the nexus between the internal structure of and the outside world, carrying an `appos` tag. Since generally the `appos` is realized by a nominal, here it provides evidence for such cases to be analyzed as a whole with the flat analysis.

To solve this issue, we suggest enlarging the scope of MWEs and include more in the basic,

standard UD guideline. In effect, on Construction Grammarians' view, which basically also holds for MWE, there might be much more constructions in real language than we might have thought (Hilpert, 2013). Therefore, long proper names such as the titles of songs, movies, and books should be analyzed in light of the `flat` construction since these are in fact unchangeable expressions. In addition, expressions such as *in theory*, *in practice* should also be considered complex adverbs as a whole because the nouns inside are not referential, being different from a prototypical noun.

The second type of problem is that the current MWE analysis per se is problematic. For instance, the current analysis of `fixed` gives rise to an unnatural combination for the initial token in the structure, such as *in* (ADP & advmod) in *in addition*. The UPOS tag `ADP` reflects the original or regular part of speech of *in*, while the DEPREL tag `advmod` is played by the whole complex adverb. In addition, in current UD, only the DEPREL of the non-initial tokens in the `flat/fixed` construction is crossed off (replaced by a `flat/fixed` label), but not their parts of speech. The same problem exist that the original part of speech of each word makes little sense since here it does not have full grammatical capacity or potential compared with the prototypical members in those categories.

One of the designing principle of UD is to achieve maximum cross-linguistic correspondence. In linguistic typology, it has long been hold that only functional terms are primarily cross-linguistic valid (Haspelmath, 2010; Croft et al., 2017). Since the tokens within MWEs are meaningless or have different behavior from their normal counterpart, thus the construction should function as a whole with their internal structure masked. Otherwise, the above-mentioned principle could by no means be achieved.

Therefore, for this second case, our solution is to keep the head-initial approach to the `fixed` type but mask the original parts of speech of all internal, non-initial tokens. Meanwhile, the UPOS of the first word should be assigned the category of that MWE as a whole. For the case of *in addition*, that means the UPOS of *in* should be `ADV` and that of *addition* (including any other word in case of a

MWE longer than two words) should be "_" (null). The final result looks as follows:

| FORM | UPOS | … | DEPREL |
|---|---|---|---|
| in | ADV | … | advmod |
| addition | _ | … | fixed |

A further improvement is to concatenate all the tokens into the LEMMA column of the first word. In doing so, when we try to extract certain data, it would avoid the awkward combination of *addition* annotated as `advmod` at first sight and make it easier for language researchers to extract the exact information they want. It is an idea found and realized in Yih (2022)'s Construction-based Universal Dependencies (CUD). The present analysis also holds for the case of `flat`.

| FORM | LEMMA | UPOS | … | DEPREL |
|---|---|---|---|---|
| in | in addition | ADV | … | advmod |
| addition | _ | _ | … | fixed |

By this adjustment, the UPOS-DEPREL combination becomes consistent, and it is easier for researchers to distinguish between different lexical items, identify them exactly, and conduct quantitative surveys more accurately. If the relation of *in* is to be kept with a normal, prototypical *in* as an adposition, we suggest that be shown in the language-specific part of speech, i.e., the XPOS column.

## 5.2 Modifiers

In the preceding sections, we have shown that the utilization rates in all treebanks are low so far. One possible solution to increase the utilization is to eliminate fine-grained difference within similar dependencies whose information are already contained in UPOS. This is an approach taken in the Surface-syntactic Universal Dependencies framework [8] (SUD, Gerdes et al. 2018), which merged `amod`, `nmod`, `nummod`, `advmod` into one tag `mod`. In doing so, the contingency table of UPOS and DEPREL becomes more dense. If there is no further distinction between these modifiers except for their categorical information, then we suggest UD follow SUD in this regard.

Alternatively, if this fine distinction is to be kept, we suggest more to be absorbed from the existing literature on the internal layered structure of noun phrases (e.g. Halliday, 1985; Davidse & Breban,

---

[8]
https://surfacesyntacticud.github.io/

2019), which has roughly reached a consensus on dividing the nominal modifiers into at least three layers or subtypes with distinct properties. For instance, `det`, `nummod`, `amod`, and `nmod` could be adjusted to represent deictics, numeratives, epithets and classifiers. Certain intermediate cases lying between two categories, such as the secondary determiner (e.g. *other*) can be annotated with the combination of `ADJ` & `det`. In addition, there is a well-known difference between descriptive adjectives, such as *big*, and associative adjectives, such as *political*. One way to distinguish between these while maximally keeping the current setting of UD is to resort to mismatches. The associative adjective could be annotated as `ADJ` & `nmod`, which sometimes behave in the same way as a noun. Likewise, fuzzy quantifiers, such as *many* and *much*, which on the one hand behaves like adjectives, and on the other like numerals, could be marked as `ADJ` & `nummod`. In doing so, more combinations regarding modifiers become legitimate in theory and the utilization rate increases.

## 6   Concluding remarks

The present contribution explores the distribution of UPOS-DEPREL combinations in 20 parallel UD treebanks and take a closer look at the mismatches in English PUD. The findings and suggestions include:

• The distributions of combinations are all highly non-uniform cross-linguistically, indicating that a small group of combinations take a rather large part.

• There are two-way matches for most UPOS, which indicates that the information contained in these two tags are redundant to some extent. We also provide a linguistic rationale for this result.

• UPOS-DEPREL mismatches fall into three types: those denoting transcategorization, wrong or inconsistent annotations, and representations of special constructions. The first type together with previously found matches form a Possible Combination Zone, the second calls for a systematic retroflection and adjustment, and the last group is left for further discussion in the future.

Our findings provide implications for producing automatic error detectors or manual checking guidelines based on linguistic knowledge. An integrated parser, which contains both a statistical annotator and a post-hoc rule-based error checker is likely to achieve higher accuracy. As Feng (2017) pointed out, the NLP would better be improved with the involvement of linguistic knowledge to bear more fruitful results. By way of two case studies, this research also sheds light on how UD annotation guidelines can be improved to better represent linguistic phenomena and to provide convenience for researchers to easily locate their wanted constructions.

## References

Lars Ahrenberg. 2019. Towards an adequate account of parataxis in Universal Dependencies. In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)*, pages 94–100, Paris, France. Association for Computational Linguistics.

Gosse Bouma, Jan Hajic, Dag Haug, Joakim Nivre, Per Erik Solberg, and Lilja Øvrelid. 2018. Expletives in Universal Dependency Treebanks. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 18–26, Brussels, Belgium. Association for Computational Linguistics. https://doi.org/1010.18653/v1/W18-6003.

William Croft. 1991. *Syntactic Categories and Grammatical Relations*. The University of Chicago Press, Chicago, US.

William Croft. 2001. *Radical Construction Grammar: Syntactic Theory in Typological Perspective*. Oxford University Press, Oxford, UK.

William Croft. 2022. *Morphosyntax: Constructions of the World's Languages*. Cambridge University Press, Cambridge, UK.

William Croft , Dawn Nordquist, Katherine Looney, and Michael Regan. 2017. Linguistic Typology meets Universal Dependencies. In *Proceedings of the 15th International Workshop on Treebanks and Linguistic Theories (TLT15)*, pages 63–75, Bloomington, IN, USA, January 20-21.

Kristin Davidse and Tine breban 2019. A cognitive-functional approach to the order of adjectives in the English noun phrase. *Linguistics*, 57(2): 327–371. https://doi.org/10.1515/ling-2019-0003.

Marie-Catherine de Marneffe, Timothy Dozat, Natalia Silveira, Katri Haverinen, Filip Ginter, Joakim Nivre, and Christopher D. Manning. 2014. Universal Stanford dependencies: A cross-linguistic

typology. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4585–4592, Reykjavik, Iceland. European Language Resources Association (ELRA).

Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. Universal Dependencies. *Computational Linguistics*, 47(2):255–308, 07.

Kira Droganova and Daniel Zeman. 2017. Elliptic Constructions: Spotting Patterns in UD Treebanks. In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 48–57, Gothenburg, Sweden. Association for Computational Linguistics.

Kira Droganova, Filip Ginter, Jenna Kanerva, and Daniel Zeman. 2018. Mind the Gap: Data Enrichment in Dependency Parsing of Elliptical Constructions. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 47–54, Brussels, Belgium. Association for Computational Linguistics. https://doi.org/10.18653/v1/W18-6006.

Zhiwei Feng. 2017. *Theory and Method for Formal Analysis of Natural Language by Computer*. Press of University of Science and Technology of China, Anhui, China.

Kim Gerdes, Bruno Guillaume, Sylvain Kahane, and Guy Perrier. 2018. SUD or Surface-Syntactic Universal Dependencies: An annotation scheme near-isomorphic to UD. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 66–74, Brussels, Belgium. Association for Computational Linguistics.

M.A.K. Halliday. 1985. *An Introduction to Functional Grammar*. Edward Arnold, London, UK.

Naïma Hassert, Pierre André Ménard, and Edith Galy. 2021. UD on Software Requirements: Application and Challenges. In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021)*, pages 62–74, Sofia, Bulgaria. Association for Computational Linguistics.

Martin Hilpert. 2013. *Construction Grammar and its Application to English*. Edinburgh University Press, Edinburgh, UK.

Georg F.K. Höhn. 2021. Towards a consistent annotation of nominal person in Universal Dependencies. In *Proceedings of the Fifth Workshop on Universal Dependencies* (UDW, SyntaxFest 2021), pages 75–83, Sofia, Bulgaria. Association for Computational Linguistics.

Elisabetta Ježek and Paolo Ramat. 2009. On parts-of-speech transcategorization. *Folia Linguistica*, 43(2): 391–416. https://doi.org/10.1515/FLIN.2009.011.

Martin Haspelmath. 2010. Comparative concepts and descriptive categories in crosslinguistic studies. *Language*, 86(3): 663–687. https://doi.org/10.1353/lan.2010.0021.

Guy Lapalme. 2021. Validation of Universal Dependencies by regeneration. In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021)*, pages 109–120, Sofia, Bulgaria. Association for Computational Linguistics.

Andrej L. Malchukov. 2004. *Nominalization/verbalization: Constraining a typology of transcategorial operations*. LINCOM Europa, Muenchen, Germany.

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. Universal Dependencies v1: A multilingual treebank collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1659–1666. European Language Resources Association (ELRA), May.

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France, May. European Language Resources Association.

Slav Petrov, Dipanjan Das, and Ryan McDonald. 2012. A Universal Part-of-Speech Tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2089–2096, Istanbul, Turkey. European Language Resources Association (ELRA).

Adam Przepiórkowski and Agnieszka Patejuk. 2019. Nested Coordination in Universal Dependencies. In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)*, pages 58–69, Paris, France. Association for Computational Linguistics. https://doi.org/10.18653/v1/W19-8007.

Jack Rueter, Niko Partanen, and Flammie A. Pirinen. 2021. Numerals and what counts. In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021)*, pages 151–159, Sofia, Bulgaria. Association for Computational Linguistics.

Nathan Schneider and Amir Zeldes. 2021. Mischievous nominal constructions in Universal Dependencies. In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW,*

*SyntaxFest 2021)*, pages 160–172, Sofia, Bulgaria. Association for Computational Linguistics.

Sebastian Schuster, Matthew Lamm, and Christopher D. Manning. 2017. Gapping Constructions in Universal Dependencies v2. In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 123–132, Gothenburg, Sweden. Association for Computational Linguistics.

Francis Tyers and Karina Mishchenkova. 2020. Dependency annotation of noun incorporation in polysynthetic languages. In *Proceedings of the Fourth Workshop on Universal Dependencies (UDW 2020)*, pages 195–204, Barcelona, Spain (Online). Association for Computational Linguistics.

Hans van Halteren. 2000. The Detection of Inconsistency in Manually Tagged Text. In *Proceedings of the COLING-2000 Workshop on Linguistically Interpreted Corpora*, pages 48–55, Centre Universitaire, Luxembourg. International Committee on Computational Linguistics.

Guillaume Wisniewski. 2018. Errator: a Tool to Help Detect Annotation Errors in the Universal Dependencies Project. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Tsy Yih. 2022. The combination of dependency relations and constructions: the annotation scheme of Construction-based Universal Dependencies (CUD). In Wei Huang, editor, *Quantitative Studies on Vocabulary and Syntax*, pages 285–296. Zhejiang University Press, Hangzhou, China.

Daniel Zeman. 2021. Date and Time in Universal Dependencies. In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021)*, pages 173–193, Sofia, Bulgaria. Association for Computational Linguistics.

**A  Cross-tabulation between UPOS and DEPREL in English PUD (in alphabetical order)**

| | ADJ | ADP | ADV | AUX | CCONJ | DET | INTJ | NOUN | NUM | PART | PRON | PROPN | PUNCT | SCONJ | SYM | VERB | X |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| acl | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 184 | 0 |
| acl:relcl | 9 | 0 | 0 | 3 | 0 | 0 | 0 | 8 | 0 | 0 | 1 | 3 | 0 | 0 | 0 | 187 | 0 |
| advcl | 11 | 0 | 3 | 1 | 0 | 1 | 0 | 13 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 262 | 0 |
| advmod | 11 | 13 | 770 | 0 | 0 | 1 | 0 | 0 | 0 | 57 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| amod | 1221 | 0 | 1 | 0 | 0 | 0 | 0 | 15 | 1 | 0 | 0 | 10 | 0 | 0 | 0 | 88 | 0 |
| appos | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 46 | 6 | 0 | 0 | 86 | 0 | 0 | 1 | 2 | 0 |
| aux | 0 | 0 | 0 | 410 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| aux:pass | 0 | 0 | 0 | 274 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| case | 18 | 2339 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 104 | 0 | 0 | 0 | 4 | 4 | 22 | 1 |
| cc | 0 | 0 | 9 | 0 | 565 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| cc:preconj | 0 | 0 | 1 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ccomp | 13 | 0 | 0 | 3 | 0 | 0 | 0 | 14 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 103 | 0 |
| compound | 13 | 0 | 1 | 0 | 0 | 1 | 0 | 433 | 26 | 0 | 0 | 330 | 0 | 0 | 2 | 2 | 2 |
| compound:prt | 0 | 62 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| conj | 44 | 0 | 8 | 1 | 0 | 2 | 0 | 265 | 12 | 0 | 3 | 96 | 0 | 0 | 3 | 199 | 1 |
| cop | 0 | 0 | 0 | 316 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| csubj | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 27 | 0 |
| csubj:pass | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| dep | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| det | 0 | 0 | 0 | 0 | 0 | 2047 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| det:predet | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| discourse | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| dislocated | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| expl | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 59 | 0 | 0 | 0 | 0 | 0 | 0 |
| fixed | 4 | 60 | 16 | 0 | 1 | 2 | 0 | 8 | 1 | 3 | 0 | 2 | 0 | 5 | 0 | 1 | 0 |
| flat | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 7 | 0 | 0 | 212 | 0 | 0 | 0 | 0 | 7 |

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| goeswith | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| iobj | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 3 | 3 | 0 | 0 | 0 | 0 | 0 |
| mark | 0 | 9 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 261 | 0 | 0 | 0 | 280 | 0 | 0 | 0 |
| nmod | 8 | 1 | 2 | 0 | 0 | 1 | 0 | 723 | 31 | 0 | 18 | 284 | 0 | 0 | 6 | 1 | 1 |
| nmod:npmod | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 16 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| nmod:poss | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 36 | 0 | 0 | 260 | 68 | 0 | 0 | 0 | 0 | 1 |
| nmod:tmod | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 26 | 10 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 |
| nsubj | 15 | 0 | 1 | 0 | 0 | 4 | 0 | 599 | 10 | 0 | 477 | 280 | 0 | 1 | 4 | 2 | 0 |
| nsubj:pass | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 134 | 2 | 0 | 61 | 39 | 0 | 0 | 0 | 1 | 0 |
| nummod | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 254 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| obj | 13 | 1 | 1 | 0 | 0 | 10 | 0 | 689 | 7 | 0 | 72 | 73 | 0 | 0 | 9 | 1 | 0 |
| obl | 14 | 5 | 6 | 0 | 0 | 1 | 0 | 829 | 81 | 0 | 57 | 235 | 0 | 0 | 6 | 1 | 2 |
| obl:npmod | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 16 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| obl:tmod | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| orphan | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 0 |
| parataxis | 5 | 0 | 1 | 0 | 0 | 0 | 0 | 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 80 | 0 |
| punct | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2451 | 0 | 0 | 0 | 0 |
| reparandum | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| root | 73 | 1 | 4 | 6 | 0 | 0 | 0 | 93 | 3 | 0 | 6 | 12 | 0 | 0 | 4 | 797 | 1 |
| vocative | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| xcomp | 47 | 2 | 2 | 0 | 0 | 0 | 0 | 31 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 186 | 0 |

Table 4: The cross-tabulation between UPOS and DEPREL in English PUD (in alphabetical order). The numbers in violet represent two-way matches, those in red and in blue one-way matches.

**B  Cross-tabulation between UPOS and DEPREL in English PUD rearranged according to the function of dependencies**

| | | PROPN | NOUN | PRON | DET | ADJ | NUM | ADV | VERB | AUX | ADP | SCONJ | CCONJ | INTJ | PUNCT | PART | X | SYM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ARG | nsubj | 319 | 733 | 538 | 6 | 15 | 12 | 1 | 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 4 |
| | obj | 73 | 689 | 72 | 10 | 13 | 7 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 9 |
| | iobj | 3 | 4 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | obl | 235 | 863 | 57 | 2 | 15 | 81 | 7 | 1 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 2 | 7 |
| MOD | nmod | 354 | 801 | 281 | 1 | 8 | 41 | 3 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 6 |
| | det | 0 | 0 | 0 | 2056 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | amod | 10 | 15 | 0 | 0 | 1221 | 1 | 1 | 88 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | nummod | 0 | 0 | 0 | 0 | 0 | 254 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | advmod | 0 | 0 | 0 | 1 | 11 | 0 | 770 | 0 | 0 | 13 | 0 | 0 | 0 | 0 | 57 | 0 | 0 |
| PRED | root | 12 | 93 | 6 | 0 | 73 | 3 | 4 | 797 | 6 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 4 |
| | parataxis | 0 | 11 | 0 | 0 | 5 | 0 | 1 | 80 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | acl | 0 | 1 | 0 | 0 | 7 | 0 | 0 | 184 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| | acl:relcl | 3 | 8 | 1 | 0 | 9 | 0 | 0 | 187 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | advcl | 0 | 13 | 0 | 1 | 11 | 1 | 3 | 262 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| | csubj | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 29 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | ccomp | 2 | 14 | 0 | 0 | 13 | 0 | 0 | 103 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | xcomp | 2 | 31 | 0 | 0 | 47 | 1 | 2 | 186 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| OTHER | aux | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 684 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | cop | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 316 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | case | 0 | 0 | 0 | 0 | 18 | 0 | 7 | 22 | 0 | 2339 | 4 | 0 | 0 | 0 | 104 | 1 | 4 |
| | mark | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 9 | 280 | 0 | 0 | 0 | 261 | 0 | 0 |
| | cc | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 575 | 0 | 0 | 0 | 0 | 0 |
| | discourse | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| | punct | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2451 | 0 | 0 | 0 |
| MWE | compound | 330 | 433 | 0 | 1 | 13 | 26 | 8 | 2 | 0 | 62 | 0 | 0 | 0 | 0 | 1 | 2 | 2 |
| | fixed | | 8 | 0 | 2 | 4 | 1 | 16 | 1 | 0 | 60 | 5 | 1 | 0 | 0 | 3 | 0 | 0 |
| | flat | | 2 | 0 | 1 | 1 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 0 |

Table 5: Rearranged to better show the matches (yellow), the possible zone (orange), controversial examples (green), and presumably problematic cells (white).