

A Description of Turkish Discourse Bank 1.2 and an Examination of Common Dependencies in Turkish Discourse

Deniz Zeyrek¹, Mustafa Erolcan Er¹

¹Middle East Technical University, Graduate School of Informatics, Cognitive Science Department, Dumlupınar Boulevard, No:1, 06800, Ankara, Turkey

Abstract

We describe Turkish Discourse Bank 1.2, the latest version of a discourse corpus annotated for explicitly or implicitly conveyed discourse relations, their constitutive units, and senses in the Penn Discourse Treebank style. We present an evaluation of the recently added tokens and examine three commonly occurring dependency patterns that hold among the constitutive units of a pair of adjacent discourse relations, namely, shared arguments, full embedding and partial containment of a discourse relation. We present three major findings: (a) implicitly conveyed relations occur more often than explicitly conveyed relations in the data; (b) it is much more common for two adjacent implicit discourse relations to share an argument than for two adjacent explicit relations to do so; (c) both full embedding and partial containment of discourse relations are pervasive in the corpus, which can be partly due to subordinator connectives whose preposed subordinate clause tends to be selected together with the matrix clause rather than being selected alone. Finally, we briefly discuss the implications of our findings for Turkish discourse parsing.

Keywords

Turkish, discourse connectives, converbial suffixal connectives, postpositions, dependencies in discourse

1. Introduction

Turkish is a language of more than 80M speakers and belongs to the Turkic sub-family of the Altaic language family. It is a free word-order, agglutinating language with a complex morphology, where suffixation is a major tool of both derivation and inflection.

The existing Natural Language Processing (NLP) methods for Turkish have been developed primarily targeting its morphology and syntax, lately extending to semantics [1], [2]. But there is also need for discourse processing research, i.e. NLP beyond the boundaries of the sentence, which would inform systems such as information retrieval, dialogue systems, summarization. The first annotated discourse corpus of Turkish, Turkish Discourse Bank, or TDB [3] has been developed to fill the gap in the discourse processing of Turkish and is expected to support language technology applications that need information at the discourse level. It is a manually

The International Conference and Workshop on Agglutinative Language Technologies as a challenge of Natural Language Processing, ALTNLP'22, June 7-8, Koper, Slovenia


✉ dezeyrek@metu.edu.tr (D. Zeyrek); erolcan.er@metu.edu.tr (M. E. Er)

🌐 <http://users.metu.edu.tr/dezeyrek/> (D. Zeyrek); <https://github.com/erolcan-er> (M. E. Er)

🆔 0000-0001-9248-0141 (D. Zeyrek); 0000-0002-3009-4517 (M. E. Er)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

annotated corpus of modern Turkish that follows the rules and principles of the Penn Discourse Bank (PDTB) [4] annotating discourse relations over texts from various genres (fiction, biography, newspaper editorials, popular magazines, etc.).¹

While the PDTB still remains the largest resource, the creation of PDTB-style discourse corpora in languages such as Turkish, Hindi, Arabic and Chinese (see [5] and the references therein) has been significant for empirical purposes and for discourse processing studies on those languages. The empirical value of new resources is high because they underscore both the variability and similarity of discourse-related phenomena across languages and enable researchers to reach a better understanding of discourse structure.

The goal of the current paper was twofold: (a) To describe the latest version of TDB, namely TDB 1.2, a 40.000-word corpus, and evaluate the recently added tokens, (b) to highlight three commonly occurring discourse dependencies found in TDB 1.2; i.e., *shared arguments*, *full embedding* and *proper containment of a discourse relation*, and discuss the issues revolving around these dependencies from the viewpoint of a morphologically rich language.

The layout of the paper is as follows. We start with an overview of the notions that underlie TDB, describe major annotation categories, and offer an evaluation of the new discourse relation tokens (§2). In §3 we present the most common dependencies in the corpus and discuss the linguistic issues surrounding them, and in §4, we conclude the paper.

2. Turkish Discourse Bank 1.2

2.1. What is discourse and what are discourse relations?

Discourse is the level of language above the sentence and can be found even within a sentence. The assumption in discourse research is that a stretch of text is not an arbitrary sequence of sentences but a structured, coherent unit that has a meaning more than the sum of its parts. Discourse structure can be discovered by examining the patterns in multi-sentence or multi-clausal texts and by finding the constitutive units of these patterns. This is essential for correctly interpreting the text [6] and for the first step of discourse processing, i.e. discourse segmentation, known as discourse parsing. One of the key aspects of discourse structure is *discourse relations (DRs)*, which denote the semantic relatedness of two text pieces at the local level, such as contrast, additive, condition.²

Following the PDTB’s lexicalized approach to discourse relations, it is assumed that there is lexico-syntactic evidence for the existence of discourse relations. Thus, connectives are seen as a primary source of evidence for the occurrence of a discourse relation. These are expressions such as conjunctions and adverbs (*or*, *although*, *moreover*) linking clauses that have an *abstract object* interpretation (propositions or eventualities) [7]. They are referred to as (*explicit*) *discourse connectives (DCs)* signalling the presence of discourse relations (see example (7) in Appendix A).

¹The earliest version, TDB 1.0, is a \sim 400.000-word corpus available at <https://github.com/disrpt/sharedtask2019/tree/master/data/tur.pdtb.tdb>. TDB 1.1, a 40.000-word-version with fewer annotations, is available at: <https://github.com/disrpt/sharedtask2021/tree/main/data/tur.pdtb.tdb>.

²Although discourse relations can also express the pragmatic relatedness of discourse units (e.g. claim-evidence), they are not annotated in TDB.

But readers do not necessarily need discourse connectives, because they can easily infer the relation from the adjacency of textual units, lexical relations, anaphoric links, etc. These have been known as implicit relations. Furthermore, readers can add a discourse connective to an implicitly conveyed relation to make it salient – called “*implicit discourse connectives*” [4] – and can specify the textual parts of an implicit relation (see example (8) in Appendix A). Finally, implicit relations may be realized by other means, namely, through Alternative Lexicalization (AltLex), or as Hypophora, Entity Relation, as well as No Relation (more explanation and examples are provided for each relation type in Appendix A).

2.2. What is annotated in TDB?

Based on the notions described in §2, three major aspects of discourse are annotated in TDB 1.2: (a) Discourse relations conveyed explicitly or implicitly as well as by other means, (b) constitutive units of discourse relations, which are known as *arguments*, (c) the sense of explicitly and implicitly conveyed relations and AltLexes. There are always two textual units that constitute a relation. The textual unit syntactically hosting the discourse connective is called Argument 2, the other argument is named as Argument 1.³

Although all languages have elements that function as discourse connectives, the syntactic class to which they belong may differ. For example, Turkish not only has lexical connectives (*and, but, so*) as most languages do but also converbial and postpositional connectives, grouped as subordinators. These connectives relate a non-finite subordinate adverbial clause to the matrix clause. In converbial structures, the marker of the relation is merely a suffix, called suffixal connectives here, which generally correspond to subordinating conjunctions in English. In postpositional structures, the marker of the relation has two parts, a postposition and a nominalization suffix on the subordinate verb. Converbial suffixes and postpositions are annotated as explicit discourse connectives in TDB.

Importantly, the neutral order of arguments to subordinators is Argument2-Argument1 (i.e. the argument that hosts the connective, which is the second argument, is normally **preposed**). Both subordinator types are typically translated to English with a **postposed** subordinate clause). Example (1) presents a suffixal connective, *-ince* ‘when’ (2), while (2) illustrates the use of a postposition *-diği gibi* ‘as’ used as a discourse connective. Both connectives relate a preposed non-finite subordinate clause to the matrix clause. In the examples throughout the paper, the discourse connective is underlined, the inferred implicit discourse connective is both underlined and put between parentheses. Argument 1 is shown in italic fonts, Argument 2 in bold fonts. Each Turkish example is translated into English and shown between single quotation marks.

- (1) **Öğrenciler gel-ince** aşağı *indi*.
‘He came down when **the students arrived**’.

³At least two annotators, who were graduate students at Middle East Technical University, Cognitive Science Department, were involved in each annotation cycle. The annotations were regularly checked and adjudicated by the research team.

- (2) **Ali'nin göster-diği gibi resim yaptım.**
'I drew as Ali showed'.

In the rest of the paper, the patterns that involve subordinator connectives will be in focus as their syntactic behaviour is peculiar to Turkish and their analysis could highlight the differences between Turkish and other languages annotated in the same style.

2.3. Evaluation and the finalization of TDB 1.2

TDB 1.2 currently has a total of 3870 relations, surpassing TDB 1.1 by 2014 relations (see Appendix B for the tokens recently added to the corpus). Since earlier versions of TDB 1.1 have already been evaluated, it appeared meaningful to evaluate the recently added tokens. A group of three expert annotators worked on a randomly chosen $\sim 42\%$ of the new relations (849 tokens in total) annotated since [8]. They were told to accept the annotations, revise them where needed, or reject them, suggesting a new relation token where possible. All decisions were made unanimously by them independently of the annotators who created and adjudicated the recent tokens. In calculating inter-annotator agreement (IAA) statistics, we considered the already adjudicated tokens as created by Annotator1, and the unanimously revised tokens as created by Annotator2. Thus IAA was measured between two annotators. We measured various types of IAA as described below and obtained a high degree of agreement in each case.

- *Agreement on the DRs' type of realization*: This is defined as the number of common discourse relations (pairs of clauses specified as a discourse relation by both annotators) over the number of unique relations, where all relations have the same type of realization [8, 9]. We used the exact match criterion [10] and present the results of this analysis in Table 1 in Appendix C.
- *Agreement on senses*: The PDTB introduces a hierarchically organized semantic categorization used to tag the sense(s) of Explicit and Implicit relations and AltLexes. The sense hierarchy has four Level-1 senses (Expansion, Contingency, Comparison, Temporal), which are further refined by Level-2 senses. A third level specifies the semantic contribution of each argument [4]. Thus, a temporal relation anchored by *then* would be annotated as Temporal:Asynchronous:Precedence, while a temporal relation expressed by *after* would be annotated as Temporal:Asynchronous:Succession. Following [9], we calculated sense agreement on all three sense levels of the PDTB 3.0 sense hierarchy among common discourse relations using the exact match criterion. The results are listed in Table 2 in Appendix C.
- *Agreement on argument spans*: TDB 1.2 asks the annotators to observe the PDTB's minimality principle, which states that the extent of the arguments to a discourse connective should be as minimal as possible as needed by the sense of the relation. The annotators are not encouraged to select distant arguments to a discourse connective but they should leave out certain expressions specified in the annotation manual (e.g. attribution phrases such as *he said* should be excluded).

To evaluate the stability of the argument span annotations, we measured IAA using Cohen's Kappa [11].

The first step involves determining the boundaries of arguments, both Argument1 and Argument2. This is known as unitization of the data ([12, 13]). In earlier work on TDB 1.0, the data was unitized with respect to words [3]. In the current work, we unitized the data with respect to *characters* by encoding each of them as 1 or 0 (selected/excluded); that is, we recorded the number of judgements a character receives for each category and calculated agreement over the data unitized in this manner. This encoding method has been considered more advantageous than the word-based encoding as it suits the agglutinating nature of Turkish better, enabling for example, the calculation of the agreement on argument spans to suffixal connectives. The agreement on each argument was measured separately. The results are given in Table 3 in Appendix C.

All disagreements were resolved by the research team and the remaining discourse relation tokens checked and updated where needed. The results were recorded in the data and TDB 1.2 was created. Recently added tokens yielded a corpus with the annotation categories distributed as shown in Table 4 in Appendix C. The table reveals that the majority of the relations are implicit amounting to 62.09% of the total number of annotated tokens as opposed to explicit relations that constitute 37.91% of the data.

3. Common dependencies in TDB 1.2

TDB annotation style reflects the incremental interpretation of texts by humans. The annotators are asked to read the text sentence by sentence and annotate different realizations of discourse relations as they appear in the text, also tagging the constituents of discourse relations along with the relations' senses. Although they are not required to annotate any dependencies among discourse units, by examining the annotation files produced by this annotation style, certain dependencies can be detected, which in turn would inform us about discourse structure, ultimately supporting discourse parsers and other language technology applications. Discourse-level dependencies have been examined in PDTB 2.0 for English [14], over TDB 1.0 for Turkish [15, 16], and recently for Czech [17]. In this paper, we continue this line of research started by Lee et al. [14]. Examining TDB 1.2 with a Python script, we investigate the dependencies among three discourse units belonging to two consecutive discourse relations related by explicit or implicit discourse connectives (other discourse relations are out of scope of our analysis).

The object of our investigation can be represented as: $DU_1 - DC1 - DU_2 - DC2 - DU_3$. That is, we deal with the dependencies among three linearly ordered discourse units (DUs), where DU means any text span selected as an argument by one or both of the discourse connectives. The major dependency types that we find are listed in Table 5 in Appendix C together with the number of times each type occurs in the data.

3.1. Shared arguments

Shared arguments refer to multiple parenthood, a kind of dependency where DU_2 is shared by the right side and the left side discourse connectives without any part of the argument span being excluded (in the examples, the shared argument is shown in a double-lined frame box to distinguish it from other DUs, which are placed in a frame box). Table 5 shows that 632

tokens (72.48% of the total number of shared arguments) are an argument to an *implicit* DC_1 shared by an *implicit* DC_2 (the Implicit-Implicit pattern in Table 5). Given the high number of implicit discourse relations in the corpus, the common occurrence of shared arguments in the Implicit-Implicit pattern is not unexpected. Also, recall that TDB is a multi-genre corpus including works of fiction, where few discourse connectives tend to occur. So, the inclusion of fiction in our corpus could be one of the reasons why implicit relations occur more frequently than explicit ones, eventually leading to the frequent occurrence of arguments shared by implicit relations.

Example (3) illustrates an Implicit-Implicit dependency structure, where DU_2 is shared by two implicit relations and the shared argument is syntactically a finite clause just like other DUs in the example. Each DU of this example is a main clause expressing an independent eventuality that can take the discourse forward. This appears to be a valid reason to make them available for reselection.

- (3)

Bu ben değildim

, (çünkü)

ben yere bakmazdım

, (bilakis)

gözüne gözüne bakardım insanların

.

This was not me

 (because)

I would not look down

, (rather)

I would look into people's eyes

.

Given the saliency of main clauses in discourse [18], their reselection is no surprise, but are subordinate clauses shared? As already mentioned, in Turkish, postpositional and suffixal connectives anchor non-finite (preposed) subordinate clauses. Are such clauses shared or not? We found that such subordinate clauses can be shared, though very rarely. For example, we found only 6 instances where the subordinate clause of a postposition is shared. Sentence (4) presents a causal postpositional connective *için* (DC_2), and its subordinate clause (**görüştürmeyi kabul ettiđi** ‘accepting to meet us’) reselected without its matrix clause. Although a detailed analysis is needed to reveal the conditions under which a preposed subordinate clause (DU_2) is shared, it appears that in (4), annotators have interpreted the eventuality described in DU_2 as semantically independent possibly co-occurring with the event described in the matrix clause (DU_3). This could have triggered the subordinate clause to be reselected.

- (4)

Bizi aray-

-arak

görüştürmeyi kabul ettiđi

 için

çok teşekkür ediyoruz

.
‘(By)

Calling us

he accepted to meet with us

, it’ for this reason that

we are thankful to him

’.

To summarize, our analysis shows that while it is common for two adjacent implicit discourse relations to share an argument, it is much less common for two adjacent explicit relations to share an argument, and subordinate clauses of subordinators are shared on rare occasions.

3.2. Full embedding

Full embedding refers to cases where a discourse relation is *totally* realized as the argument to the connective. It is similar to embedding in syntax and expected to occur in TDB 1.2, too.

Indeed, it is common in the corpus, as Table 5 (Appendix C) reveals.

Most of the fully embedded discourse relations appear in patterns where DC_2 is an explicit

discourse connective, either lexical or suffixal. The Implicit-Explicit pattern, for example, occurs in 59.77% of all fully embedded instances in Table 5. This is where the second argument to an *implicit* DC_1 is a fully embedded relation anchored by an *explicit* DC_2 .

Example (5) is chosen from the Explicit-Explicit pattern. It presents a suffixal discourse connective -ip ‘after’ and its binary arguments being fully embedded as an argument to a suffixal connective on the left side, -arak ‘once’. In other words, the subordinate clause of -ip (**anneannesinin yanına gel-** ‘move to her grandmother’s’) is selected together with the matrix clause, as the translation also shows.

- (5)

Hukuk Fakültesini yarım bırak

 -arak

anneannesinin yanına gel

 -ip

Ankara’ya yerleşmesinin

nedeni ...
‘the reason why

after	moving to her grandmother’s	she settled in Ankara
-------	------------------------------------	-----------------------

once

she quitted the Law School

’ ...

Different from example (4), this subordinate clause is not selected alone and a shared argument structure does not arise. The selection of the entire discourse relation seems due to a semantic reason: rather than being an independent eventuality, the event in the subordinate clause is in a sense dependent on the event described in the matrix clause: it brings about the event in the matrix clause. The preposed position of the subordinate clause and possibly its non-finiteness coupled with its semantics appears to block its selection alone as an argument. Although our annotation guidelines do not have rules regarding such subtle issues, the annotators opted to select most of the preposed non-finite subordinate clauses together with their matrix clauses (i.e. the entire discourse relation) as an argument, leading to fully embedded clauses or properly contained discourse relations, which is the next topic below.

3.3. Properly contained discourse relations

Properly contained discourse relations are a subtype of fully embedded ones except that some material is left out (shown with three dots in Ex. (6)) (the examination of the excluded part is left for further research). Similar to fully embedded relations, properly contained relations tend to occur in the patterns where DC_2 is an explicit discourse connective. For example, the Implicit-Explicit pattern comprises 55.25% of all properly contained relations.

- (6)

çarşafarla gecedem giderek terasa saklandı
--

 (sonra) ...

çarşafı giy

 -erek

terastan indi

 .

he hid at the terrace with the hijab

 (then) ... after

wearing the hijab

he came down

 .

In Ex. (6), chosen from the Implicit-Explicit pattern, the preposed subordinate clause (DU_2) and its matrix clause (DU_3) are selected entirely as the second argument to DC_1 rather than the subordinate clause being selected alone, which would have resulted in a shared argument structure. Once again, this seems to be due to the position of the subordinate clause as well as its semantics: the event described by the preposed subordinate clause **çarşafı giy-** ‘wear the hijab’ engenders the main event *terastan indi* ‘he came down’; the man wears the hijab and only

then, he comes down from where he is hiding (otherwise, he would be noticed by the women, as the narrative describes). These events are not inferred as independently (co-)occurring, which seems a good reason why we find a properly contained dependency structure.

In short, preposed (non-finite) subordinate clauses in Turkish seem to trigger full embedding or proper containment structures, which could be explained not only by the position and non-finiteness of the subordinate clauses but also by their semantics in relation to the matrix clauses.

4. Summary and conclusion

We introduced TDB 1.2, a corpus that annotates different realizations of discourse relations, their arguments and senses in the PDTB style, and found that the corpus contains more implicit relations than explicit ones. Then, we zoomed in three types of dependency, which revealed an asymmetry between the occurrence patterns of shared arguments on one hand and fully embedded and properly contained discourse relations on the other. Our analyses showed that arguments are shared frequently by two adjacent implicit discourse relations, but much less so by two adjacent explicit discourse relations. Instead, discourse relations conveyed by explicit connectives such as suffixal ones or postpositions tend to be selected totally as an argument to another discourse relation, mostly an implicit one.

Our findings have implications both for discourse parsing and the theoretical understanding of Turkish paving the way for comparisons with other languages towards a better understanding of discourse. While there is room for more research on both sides, the findings minimally show that the implicit discourse relation recognition task can be improved by considering shared arguments, which demonstrate, among others, that three adjacent implicit discourse relations is a highly likely sequence in Turkish discourse. Also, automatic argument span detection can be improved by considering the availability of an entire discourse relation anchored by postpositions or suffixal connectives as an argument, as fully embedded and properly contained dependency patterns reveal.

What we have not examined in this paper is whether there are other factors involved in the formation of the dependency structures described, e.g. the sense of DC_1 and/or DC_2 . The investigation of such factors is left for further research.

Acknowledgments

We acknowledge the partial support of Middle East Technical University (BAP-07-04-2017-001) and thank Salih Fırat Canpolat, Deniz Dilek Bilgiç, Ozan Deniz, Ali Can Serhan Yılmaz, Zeynep Başer, Özgür Şen Bartan, Aytaç Çeltek and Murathan Kurfalı for their assistance at various stages of the development of TDB 1.2. Any remaining errors are our own.

References

- [1] G. Eryiğit, J. Nivre, K. Oflazer, Dependency Parsing of Turkish, *Computational Linguistics* 34 (2008) 357–389. doi:10.1162/coli.2008.34.4.627.

- [2] R. Çakıcı, M. Steedman, C. Bozsahin, Wide-coverage parsing, semantics, and morphology, in: *Turkish Natural Language Processing*, Springer, 2018, pp. 153–174. doi:10.1007/978-3-319-90165-7_8.
- [3] D. Zeyrek, I. Demirşahin, A. B. S. Çallı, Turkish Discourse Bank: Porting a discourse annotation style to a morphologically rich language, *Dialogue & Discourse* 4 (2013) 174–184.
- [4] R. Prasad, N. Dinesh, A. Lee, E. Miltsakaki, L. Robaldo, A. K. Joshi, B. L. Webber, The Penn Discourse TreeBank 2.0., in: *LREC*, 2008. URL: <https://www.aclweb.org/anthology/L08-1093/>.
- [5] R. Prasad, B. Webber, A. Joshi, Reflections on the Penn Discourse TreeBank, comparable corpora, and complementary annotation, *Computational Linguistics* (2014). doi:10.1162/COLI_a_00204.
- [6] B. Webber, M. Egg, V. Kordoni, Discourse structure and language technology, *Natural Language Engineering* 18 (2011) 437–490. doi:10.1017/S1351324911000337.
- [7] N. Asher, *Reference to Abstract Objects in Discourse*, Kluwer, Dordrecht, 1993.
- [8] D. Zeyrek, M. Kurfalı, TDB 1.1: Extensions on Turkish Discourse Bank, in: *Proceedings of the 11th Linguistic Annotation Workshop*, 2017, pp. 76–81. doi:10.18653/v1/W17-0809.
- [9] K. Forbes-Riley, F. Zhang, D. Litman, Extracting PDTB Discourse Relations from Student Essays, in: *Proc. of the SIGDIAL*, 2016, pp. 117–127.
- [10] E. Miltsakaki, R. Prasad, A. K. Joshi, B. L. Webber, The Penn Discourse TreeBank., in: *LREC*, 2004.
- [11] J. Cohen, A Coefficient of Agreement for Nominal Scales, *Educational and psychological measurement* 20 (1960) 37–46.
- [12] R. Artstein, M. Poesio, Inter-coder agreement for computational linguistics, *Computational Linguistics* 34 (2008) 555–596.
- [13] Ş. İ. Yalçınkaya, An inter-annotator agreement measurement methodology for the Turkish Discourse Bank (TDB), Master's thesis, Middle East Technical University, 2010.
- [14] A. Lee, R. Prasad, A. Joshi, N. Dinesh, B. Webber, Complexity of dependencies in discourse: Are dependencies in discourse more complex than in syntax, in: *Proceedings of the 5th International Workshop on Treebanks and Linguistic Theories*, Citeseer, 2006, pp. 12–23.
- [15] B. Aktaş, C. Bozsahin, D. Zeyrek, Discourse Relation Configurations in Turkish and an Annotation Environment, in: *Proc. of the 4th Linguistic Annotation Workshop*, ACL, 2010, pp. 202–206. URL: <https://www.aclweb.org/anthology/W10-1832.pdf>.
- [16] I. Demirsahin, A. Ozturel, C. Bozsahin, D. Zeyrek, Applicative Structures and Immediate Discourse in the Turkish Discourse Bank, in: *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, 2013, pp. 122–130. URL: <https://www.aclweb.org/anthology/W13-2315.pdf>.
- [17] L. Poláková, J. Mírovský, Š. Zikánová, E. Hajičová, Discourse Relations and Connectives in Higher Text Structure, *Dialogue & Discourse* 12 (2021) 1–37.
- [18] W. C. Mann, S. A. Thompson, Rhetorical Structure Theory: Toward a functional theory of text organization, *Text-Interdisciplinary Journal for the Study of Discourse* 8 (1988) 243–281.

A. Appendix: Major annotation categories and examples in TDB 1.2

TDB 1.2 annotates implicitly and explicitly conveyed discourse relations that hold between adjacent verb phrases, clauses, and sentences. This section illustrates major annotation categories together with examples.

Explicit relations - An explicit discourse relation holds when the relation is encoded through an overt discourse connective.

- (7) **Ali uzun boylu** ama *kız kardeşi kısa boyludur.*
'**Ali is tall**, but *his sister is short.*'

Implicit relations - In cases where an overt discourse connective is absent, an implicit discourse relation is inferred and shown by inserting a discourse connective in the relation.

- (8) *Yol kaygandı, (Imp=o yüzden)* **Ali arabayı dikkatli kullandı.**
'*The road was slippery, (Imp=due to that)* **Ali was driving carefully.**'

Alternative Lexicalization (AltLex) - When a discourse relation is alternatively lexicalized through linguistic expressions such as *despite this, because of this, the reason is*, the relation is called and AltLex.

- (9) *Ali Latince öğrendi. Bundan sonra* **Fransızca kitap okumak çok kolay oldu.**
'*Ali learnt Latin. After that,* **reading books in French has been so easy.**'

Entity Relation (EntRel) - This is where the text spans express a relation with an entity.

- (10) *Dr. Ahmet bey yeni bir hastahane* **işe başladı. Rahmetli Dr. Ali bey'in yerini aldı.**
'*Dr. Ahmet Beg has started to work in a new hospital. He succeeds the late Dr. Ali Beg.*'

Hypophora - These are questions and meaningful answers given to the questions.

- (11) *Fıkra hoşuna gitti mi?* **Evet bayıldım.**
'*Did you like the joke?* **Yes I loved it.**'

No Relation (NoRel) - A NoRel involves cases where no relation can be inferred between adjacent text spans.

- (12) '*Okul yakında tatile girecek. Öğretmenler okula gönderilmeyen öğrencilerle uğraşamaz.*
Children will have a break soon. Teachers can't deal with students not sent to school.

Explicit and Implicit relations and AltLexes are annotated both within and across sentences, while Hypophora tokens, EntRels, and NoRels are annotated only between adjacent sentences.

B. Appendix: Tokens recently added to TDB 1.2

The most recent additions to the corpus involve implicit verb phrase conjunctions (Ex. (13)) and multiple relations (examples (14) - (16)).

- (13) *Çabuk değişen (Imp=ve) yaşlanan bir nüfusumuz var.*
 ‘We have a population that *rapidly changes (Imp=and) ages.*’

Multiple relations comprise:

- the implicit senses of explicitly conveyed verb phrase conjunctions (only the senses of relations marked by the conjunction *ve* ‘and’ were considered) (Ex. (14)).
- multiple relations between the same argument spans conveyed by co-occurring explicit connectives, such as *ve böylece* ‘and hence’ (Ex. (15)).
- multiple relations between the same argument spans conveyed by an explicit connective and an AltLex, such as *ve buna rağmen* ‘and despite this’ (Ex. (16)).⁴

- (14) *Okulu bıraktı ve (Imp=sonra) evlendi.*

‘She *left school and (Imp=then) got married.*’

- (15) *Ayşe sevdiğiyle evlendi ve böylece dünyanın en mutlu kızı oldu.*

‘Ayşe married her beloved one *and so she became the happiest women in the world.*’

- (16) *Ali okuldan nefret etti ve buna rağmen liseden mezun olmayı başardı.*

‘*Ali hated school and despite this he managed to finish high school.*’

Multiple relations were annotated separately on each token as in the PDTB, then linked with the same index value in their link fields.

C. Appendix: Summarization tables

Table 1

IAA results for agreement on DRs’ type of realization

Realization Type	Agreement
Implicit	0.97
Explicit	0.99
AltLex	0.98
EntRel	0.95
Hypophora	1.00
NoRel	0.97

Table 2

IAA results for sense agreement

Sense	Explicit (%)	Implicit (%)	AltLex (%)
Level-1	99.02	99.75	100
Level-2	98.66	99.75	100
Level-3	80.11	79.44	79.83

⁴PDTB 3.0 annotates multiple senses for explicit or implicit relations if annotators infer more than one sense as holding between a pair of spans. In TDB 1.2, multiple senses were not annotated systematically.

Table 3

IAA results for argument span selection

Arg Type	Cohen's κ
Argument1	0.90
Argument2	0.85

Table 4

Number of different realizations of discourse relations and their Level-1 sense tags in TDB 1.2

	Expansion	Temporal	Comparison	Contingency	DRs with no sense tag	Total
Implicit	1090	158	162	333	0	1743
Explicit	540	400	259	268	0	1467
AltLex	33	32	14	67	0	146
EntRel	0	0	0	0	233	233
Hypophora	0	0	0	0	78	78
NoRel	0	0	0	0	203	203
Total	1663	590	435	668	514	3870

Table 5

Number of common dependencies in TDB 1.2

	Explicit	Explicit	Implicit	Sub Total	Implicit	Total
DC1	Explicit	Implicit	Explicit		Implicit	
DC2						
Shared Arguments	41	105	96	240	632	872
Fully embedded DRs	117	85	471	673	115	788
Properly Contained DRs	145	82	521	748	195	943
Total	303	272	1088	1663	942	2605