

Large-scale evidence of dependency length minimization in 37 languages

Richard Futrell¹, Kyle Mahowald, and Edward Gibson

Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139

Edited by Barbara H. Partee, University of Massachusetts Amherst, Amherst, MA, and approved June 2, 2015 (received for review February 2, 2015)

Explaining the variation between human languages and the constraints on that variation is a core goal of linguistics. In the last 20 y, it has been claimed that many striking universals of cross-linguistic variation follow from a hypothetical principle that dependency length—the distance between syntactically related words in a sentence—is minimized. Various models of human sentence production and comprehension predict that long dependencies are difficult or inefficient to process; minimizing dependency length thus enables effective communication without incurring processing difficulty. However, despite widespread application of this idea in theoretical, empirical, and practical work, there is not yet large-scale evidence that dependency length is actually minimized in real utterances across many languages; previous work has focused either on a small number of languages or on limited kinds of data about each language. Here, using parsed corpora of 37 diverse languages, we show that overall dependency lengths for all languages are shorter than conservative random baselines. The results strongly suggest that dependency length minimization is a universal quantitative property of human languages and support explanations of linguistic variation in terms of general properties of human information processing.

language universals | language processing | quantitative linguistics

Finding explanations for the observed variation in human languages is the primary goal of linguistics and promises to shed light on the nature of human cognition. One particularly attractive set of explanations is functional in nature, holding that language universals are grounded in the known properties of human information processing. The idea is that grammars of languages have evolved so that language users can communicate using sentences that are relatively easy to produce and comprehend. Within the space of functional explanations, a promising hypothesis is dependency length minimization (DLM).

Dependency lengths are the distances between linguistic heads and dependents. In natural language syntax, roughly speaking, heads are words that license the presence of other words (dependents) modifying them (1). For example, the verb “throw” in sentence C in Fig. 1 licenses the presence of two nouns, “John”—its subject—and “trash”—its object. Subject and object relations are kinds of dependency relations where the head is a verb and the dependent is a noun. Another way to think about dependency is to note that heads and dependents are words that must be linked together to understand a sentence. For example, to correctly understand sentence C in Fig. 1, a comprehender must determine that a relationship of adjectival modification exists between the words “old” and “trash”, and not between, say, the words “old” and “kitchen”. In typical dependency analyses, objects of prepositions (“him” in “for him”) depend on their prepositions, articles depend on the nouns they modify, and so on. Most aspects of dependency analysis are generally agreed on, although the analysis of certain relations is not settled, primarily those relations involving function words such as prepositions, determiners, and conjunctions. Fig. 1 shows the dependencies involved in some example sentences according to the analysis we adopt.

The DLM hypothesis is that language users prefer word orders that minimize dependency length. The hypothesis makes two broad predictions. First, when the grammar of a language provides

multiple ways to express an idea, language users will prefer the expression with the shortest dependency length (2). Indeed, speakers of a few languages have been found to prefer word orders with short dependencies when multiple options are available (3, 4) (Fig. 1 provides English examples). Second, grammars should facilitate the production of short dependencies by not enforcing word orders with long dependencies (5, 6).

Explanations for why language users would prefer short dependencies are various, but they all involve the idea that short dependencies are easier or more efficient to produce and comprehend than long dependencies (7, 8). The difficulty of long dependencies emerges naturally in many models of human language processing. For example, in a left-corner parser or generator, dependency length corresponds to a timespan over which a head or dependent must be held in a memory store (9–11); because storing items in memory may be difficult or error prone, short dependencies would be easier and more efficient to produce and parse according to this model. In support of this idea, comprehension and production difficulty have been observed at the sites of long dependencies (8, 12).

If language users are motivated by avoiding difficulty, then they should avoid long dependencies. Furthermore, if languages have evolved to support easy communication, then they should not enforce word orders that create long dependencies. The DLM hypothesis thus provides a link between language structure and efficiency through the idea that speakers and languages find ways to express meaning while avoiding structures that are difficult to produce and comprehend.

Over the last 20 y, researchers have proposed DLM-based explanations of some of the most pervasive properties of word order in languages. We can see the word order in a sentence as a particular linearization of a dependency graph, where a linearization is an arrangement of the words of the dependency graph in a certain linear order. For instance, sentences A and B in Fig. 1 are two linearizations of the same graph. Below we give examples of applications of the DLM idea.

Significance

We provide the first large-scale, quantitative, cross-linguistic evidence for a universal syntactic property of languages: that dependency lengths are shorter than chance. Our work supports long-standing ideas that speakers prefer word orders with short dependency lengths and that languages do not enforce word orders with long dependency lengths. Dependency length minimization is well motivated because it allows for more efficient parsing and generation of natural language. Over the last 20 y, the hypothesis of a pressure to minimize dependency length has been invoked to explain many of the most striking recurring properties of languages. Our broad-coverage findings support those explanations.

Author contributions: R.F. designed research; R.F. performed research; R.F. and K.M. analyzed data; and R.F., K.M., and E.G. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. Email: futrell@mit.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1502134112/-DCSupplemental.

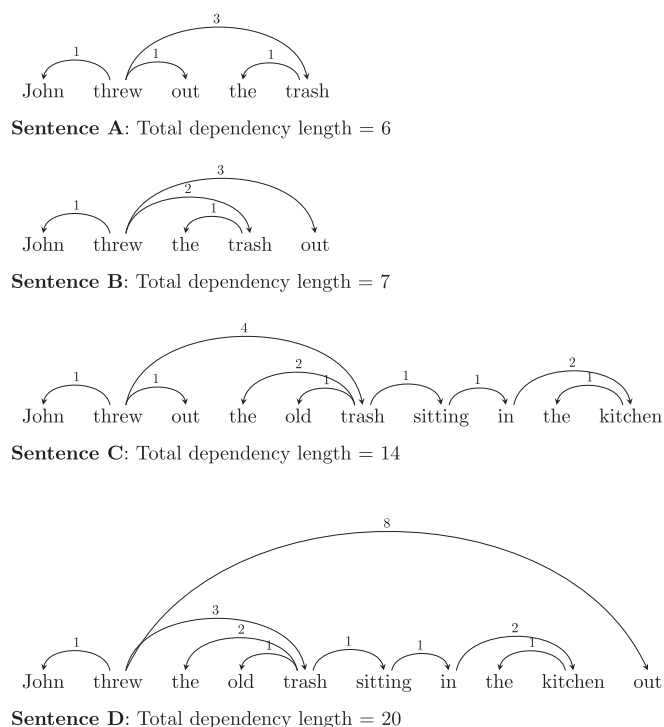


Fig. 1. Four sentences along with their dependency representations. The number over each arc represents the length of the dependency in words. The total dependency length is given below each sentence. Sentences A and B have the same semantics, and either word order is acceptable in English; English speakers typically do not find one more natural than the other. Sentences C and D also both have the same semantics, but English speakers typically find C more natural than D.

Languages constrain what linearizations are possible; for example, some languages require that a noun depending on a preposition come after the preposition, and some require that it come before. Greenberg (13) found striking correlations between different ordering constraints in languages, such that languages tend to be consistent in whether heads come before dependents or vice versa (14, 15). Both this generalization and exceptions to it have been explained as linearizations that minimize dependency length (7, 16). Hawkins (17) documents that the basic grammatical word orders for many constructions in many languages minimize dependency length over alternatives.

Another pervasive property of languages is projectivity, the property that, in linearizations of dependency graphs, the lines connecting heads and dependents do not cross (18). Ferrer i Cancho (19) has argued that this ubiquitous property of languages arises from dependency length minimization, because orders that minimize dependency length have a small number of crossing dependencies on average.

Minimal dependency length has also been widely assumed as a reliable generalization in the field of natural language processing. For example, most state-of-the-art models for natural language grammar induction incorporate a bias toward positing short dependencies, and their performance is greatly improved by this assumption (20, 21). Influential practical parsing algorithms also incorporate this assumption (22).

The studies mentioned above, for the most part, use categorical descriptions of the most common word orders in languages or examine small numbers of languages. Therefore, a crucial question remains open: is dependency length actually minimized overall in real utterances, considering the full range of possible syntactic constructions and word orders as they are used, or is the effect confined to the constructions and languages that have been studied? If indeed there is a universal preference to

minimize dependency lengths, then utterances in all natural languages should have shorter dependency lengths than would be expected by chance. On the other hand, if observed dependency lengths are consistent with those that would be produced by chance, then this would pose a major challenge to DLM as an explanatory principle for human languages.

Here, we answer that question using recently available dependency-parsed corpora of many languages (23–25). We obtained hand-parsed or hand-corrected corpora of 37 languages, comprising 10 language families. Thirty-six of the corpora follow widely recognized standards for dependency analysis (25, 26); the remaining corpus (Mandarin Chinese) uses its own system that is nonetheless similar to the standards [see Table S1 for details on each corpus]. The texts in the corpora are for the most part written prose from newspapers, novels, and blogs. Exceptions are the corpora of Latin and Ancient Greek, which include a great deal of poetry, and the corpus of Japanese, which consists of spoken dialogue. Previous comprehensive corpus-based studies of DLM cover seven languages in total, showing that overall dependency length in those languages is shorter than various baselines (16, 27–30). However, these studies find only weak evidence of DLM in German, raising the possibility that DLM is not a universal phenomenon. Noji and Miyao (31) use dependency corpora to show that memory use in a specific parsing model is minimized in 18 languages, but they do not directly address the question of dependency length minimization in general.

We compare real language word orders to counterfactual baseline orders that experience no pressure for short dependencies. These baselines serve as our null hypotheses. Our baselines represent language users who choose utterances without regard to dependency length, speaking languages whose grammars are not affected by DLM. We do not distinguish between DLM as manifested in grammars and DLM as manifested in language users' choice of utterances; the task of distinguishing grammar and use in a corpus study is a major outstanding problem in linguistics, which we do not attempt to solve here. In addition to the random baselines, we present an optimal baseline for the minimum possible dependency length in a projective linearization for each sentence. This approach allows us to evaluate the extent to which different languages minimize their dependency lengths compared with what is possible. We do not expect observed dependency lengths to be completely minimized, because there are other factors influencing grammars and language use that might come into conflict with DLM.

Results

Free Word Order Baseline. Our first baseline is fully random projective linearizations of dependency trees. Random projective linearizations are generated according to the following procedure, from Gildea and Temperley (28), a method similar to one developed by Hawkins (32). Starting at the root node of a dependency tree, collect the head word and its dependents and order them randomly. Then repeat the process for each dependent. For each sentence in our corpora, we compare real dependency lengths to dependency lengths from 100 random linearizations produced using this algorithm. Note that the 100 random linearization all have the same underlying dependency structure as the original sentence, just with a potentially different linear order. Under this procedure, the random linearizations do not obey any particular word order rules: there is no consistency in whether subjects precede or follow verbs, for example. In that sense, these baselines may most closely resemble a free word order language as opposed to a language like English, in which the order of words in sentences are relatively fixed.

Fig. 2 shows observed and random dependency lengths for sentences of length 1–50. As the figure shows, all languages have average dependency lengths shorter than the random baseline, especially for longer sentences. To test the significance of the effect, for each language, we fit regression models predicting dependency length as a function of sentence length. The models show a significant effect where the dependency length of real sentences

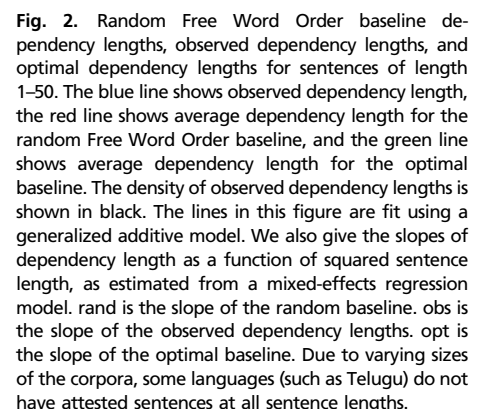
Fig. 3 shows histograms of observed and random dependency lengths for sentences of length 12, the shortest sentence length to show a significant effect in all languages ($P < 0.01$ for Latin, $P < 0.001$ for Telugu, and $P < 0.0001$ for all others, by Stouffer’s method). In languages for which we have sufficient data, there is a significant DLM effect for all longer dependency lengths.

Fig. 4 shows observed dependency lengths compared with the random fixed-order baselines. The results are similar to the comparison with the free word order baselines in that all languages have dependencies shorter than chance, especially for longer sentences. We find that this random baseline is more conservative than the free word order baseline in that the average dependency lengths of the fixed word order random baselines are shorter than those of the free word order random baselines (with significance $P < 0.0001$ by a t test in each language). For this baseline, the DLM effect as measured in the regression model is significant at $P < 0.0001$ in all languages

Discussion

Fig. 2 also reveals that, whereas observed dependency lengths are always shorter than the random baselines, they are also longer than the minimal baselines (although some languages such as Indonesian come quite close). In part, this is due to the unrealistic nature of the optimal baseline. In particular, that baseline does not have any consistency in word order [see ref. 16 for attempts to develop approximately optimal baselines which address this issue].

In general, we believe dependency length should not be fully minimized because of other factors and desiderata influencing languages that may conflict with DLM. For example, linearizations should allow the underlying dependency structure to be recovered incrementally, to allow incremental understanding of utterances. In a sequence of two words A and B , when the comprehender receives B , it would be desirable to be able to determine immediately and correctly whether A is the head of B , B is the head of A , or A and B are both dependents of some as-yet-unheard word. If the order of dependents around a head is determined only by minimizing dependency length, then there is no guarantee that word orders will facilitate correct incremental inference. More generally, it has been argued that linearizations should allow the comprehender to quickly identify the syntactic and semantic properties of each word [see Hawkins (17) for



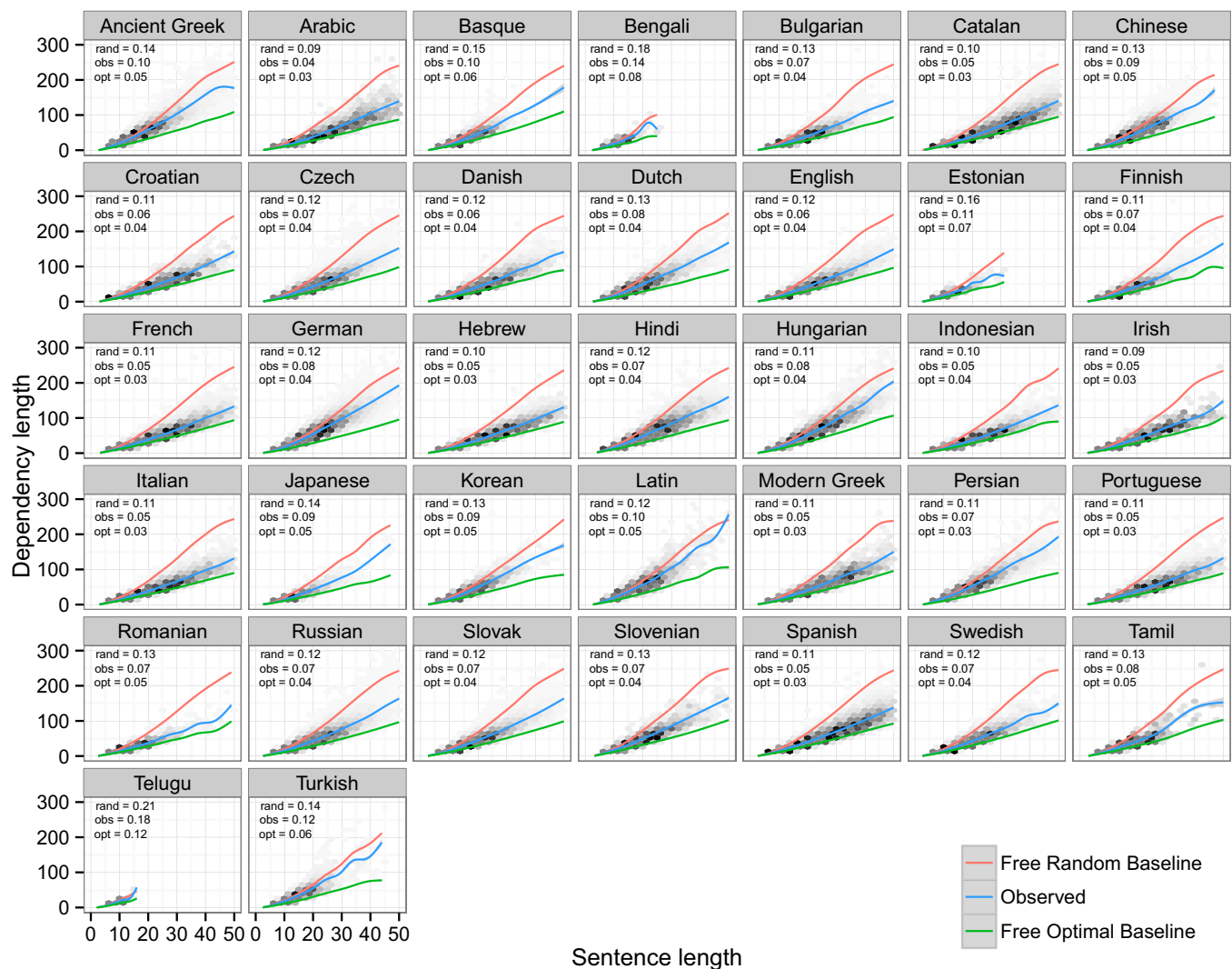


Fig. 3. Histograms of observed dependency lengths and Free Word Order random baseline dependency lengths for sentences of length 12. m_{rand} is the mean of the free word order random baseline dependency lengths; m_{obs} is the mean of observed dependency lengths. We show P values from Stouffer's Z-transform test comparing observed dependency lengths to the dependency lengths of the corresponding random linearizations.

detailed discussion of the interaction of this principle with DLM]. The interactions of DLM with these and other desiderata for languages are the subject of ongoing research.

The results presented here also show great variance in the effect size of DLM across languages. In particular, the head-final languages such as Japanese, Korean, and Turkish show much less minimization than more head initial languages such as Italian, Indonesian, and Irish, which are apparently highly optimized. This apparent relationship between head finality and dependency length is a new and unexpected discovery. Head final languages typically have highly informative word morphology such as case marking on dependents (33), and morphology might give languages more freedom in their dependency lengths because it makes long dependencies easier to identify. In line with this idea, long dependencies in German (a language with case marking) have been found to cause less processing difficulty than in English (34). In general, explaining in general why dependency lengths in some languages are shorter than in others is an interesting challenge for the DLM hypothesis.

This work has shown that the preference for short dependencies is a widespread phenomenon that not confined to the limited languages and constructions previously studied. Therefore, it lends support to DLM-based explanations for language universals. Inasmuch

as DLM can be attributed to minimizing the effort involved in language production and comprehension, this work joins previous work showing how aspects of natural language can be explained by considerations of efficiency (17, 35–39).

Materials and Methods

Data. We use the dependency trees of the HamleDT 2.0, Google Universal Treebank 2.0, and Universal Dependencies 1.0 corpora (23–25); these are projects that have aimed to harmonize details of dependency analysis between dependency corpora. In addition, we include a corpus of Mandarin, the Chinese Dependency Treebank (40). See the Table S1 for details on the source and annotation standard of each corpus. We normalize the corpora so that prepositional objects depend on their prepositions (where the original corpus has a case relation) and verbs depend on their complementizers (where the original corpus has a mark relation). For conjunctions, we use Stanford style. We also experimented with corpora in the original content-head format of HamleDT and Universal Dependencies; the pattern of results and their significance was the same.

Measuring Dependency Length. We calculate the length of a single dependency arc as the number of words between a head and a dependent, including the dependent, as in Fig. 1. For sentences, we calculate the overall dependency length by summing the lengths of all dependency arcs. We do

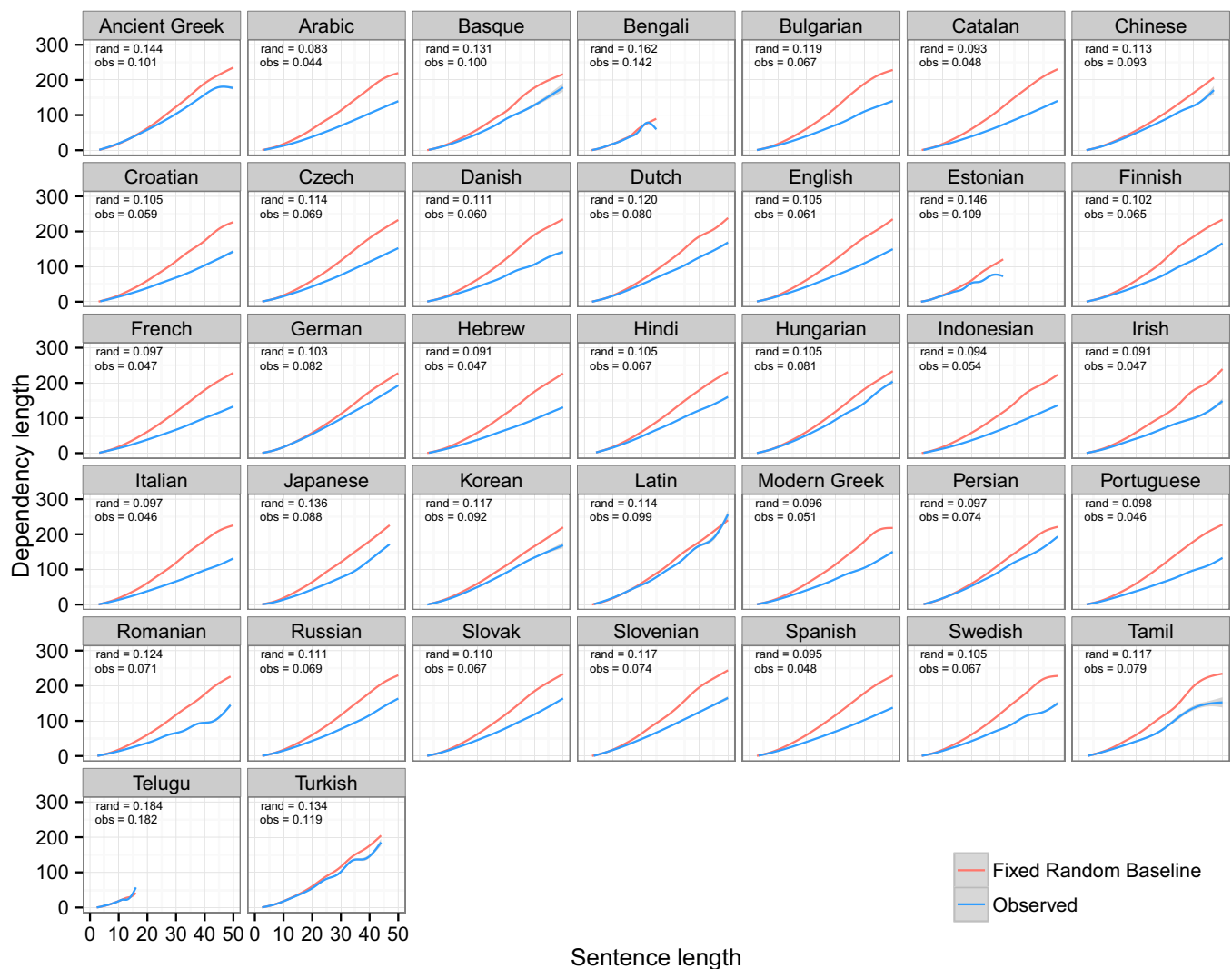


Fig. 4. Real dependency lengths as a function of sentence length (blue) compared with the Fixed Word Order Random baseline (red). GAM fits are shown. rand and obs are the slopes for random baseline and observed dependency length as a function of squared sentence length, as in Fig. 2.

not count any nodes representing punctuation or root nodes, nor arcs between them; sentences that are not singly rooted are excluded.

Fixed Word Order Random Baseline. Fixed word order random linearizations are generated according to the following procedure per sentence. Assign each relation type a random weight in $[-1, 1]$. Starting at the root node, collect the head word and its dependents and order them by their weight, with the head receiving weight 0. Then repeat the process for each dependent, keeping the same weights. This procedure creates consistency in word order with respect to relation types.

This linearization scheme can capture many aspects of fixed order in languages, but cannot capture all of them; for example, linearization order in German depends on whether a verb is in a subordinate clause or not. The fixed linearization scheme is also inaccurate in that it produces entirely deterministic orders. In contrast, many languages permit the speaker a great deal of freedom in choosing word order. However, creating a linearization model that can handle all possible syntactic phenomena is beyond the scope of this paper.

Generalized Additive Models. For the figures, we present fits from generalized additive models predicting dependency length from sentence length using cubic splines as a basis function. This model provides a line that is relatively close to the data for visualization.

Regression Models. For hypothesis testing and comparison of effect sizes, we use regression models fit to data from each language independently. For these regressions, we only consider sentences with length < 100 words. For each sentence

s in a corpus, we have $N + 1$ data points: 1 for the observed dependency length of the sentence and $N = 100$ for the dependency lengths of the random linearizations of the sentence's dependency tree. We fit a mixed-effects regression model (41) with the following equation, with coefficients β representing fixed effects and coefficients S representing random effects by sentence:

$$\hat{y}_i = \beta_0 + S_0 + \beta_1 l_s^2 + (\beta_2 + S_2)r_i + \beta_3 r_i l_s^2 + \epsilon_i, \quad [1]$$

where \hat{y}_i is the estimated total dependency length of data point i , β_0 is the intercept, l_s^2 is the squared length of sentence s in words, r_i is an indicator variable with value 1 if data point i is a random linearization and 0 if it is an observed linearization, and m_i is an indicator variable with value 1 if data point i is a minimal linearization and 0 if it is an observer linearization. We use l_s^2 rather than l_s because we found that a model using squared sentence length provides a better fit to the data for 33 of 37 languages, as measured by the Akaike information criterion and Bayesian information criterion; the pattern and significance of the results are the same for a model using plain sentence length rather than squared sentence length. The coefficient β_3 determines the extent to which dependency length of observed sentences grows more slowly with sentence length than dependency length of randomly linearized sentences. This growth rate is the variable of interest for DLM; summary measures that are not a function of length fall prey to inaccuracy due to mixing dependencies of different lengths (30). For significance testing comparing the real dependencies and random baselines, we performed a likelihood ratio test comparing models with and without β_3 . We fit the model using the lme4 package in R (42).

ACKNOWLEDGMENTS. We thank David Temperley, Gary Marcus, Ernie Davis, and the audience at the 2014 Conference on Architectures and Mechanisms in Language Processing for comments and discussion and

Dan Popel for help accessing data. K.M. was supported by the Department of Defense through the National Defense Science and Engineering Graduate Fellowship Program.

1. Corbett GG, Fraser NM, McGlashan S, eds (1993) *Heads in Grammatical Theory* (Cambridge Univ Press, Cambridge, UK).
2. Behaghel O (1932) [*Deutsche Syntax: Eine geschichtliche Darstellung* (Wortstellung)] (Carl Winter, Heidelberg), Vol IV. German.
3. Yamashita H, Chang F (2001) "Long before short" preference in the production of a head-final language. *Cognition* 81(2):B45–B55.
4. Wasow T (2002) *Postverbal Behavior* (CSLI Publications, Stanford, CA).
5. Rijkhoff J (1990) Explaining word order in the noun phrase. *Linguistics* 28(1):5–42.
6. Hawkins JA (1990) A parsing theory of word order universals. *Linguist Inq* 21(2): 223–261.
7. Hawkins JA (1994) *A Performance Theory of Order and Constituency* (Cambridge Univ Press, Cambridge, UK).
8. Gibson E (1998) Linguistic complexity: Locality of syntactic dependencies. *Cognition* 68(1):1–76.
9. Abney SP, Johnson M (1991) Memory requirements and local ambiguities of parsing strategies. *J Psycholinguist Res* 20(3):233–250.
10. Gibson E (1991) A computational theory of human linguistic processing: Memory limitations and processing breakdown. PhD thesis (Carnegie Mellon Univ, Pittsburgh).
11. Resnik P (1992) Left-corner parsing and psychological plausibility. *Proceedings of the 14th International Conference on Computational Linguistics*, ed Boileau C (Association for Computational Linguistics, Nantes, France), pp 191–197.
12. Grodner D, Gibson E (2005) Consequences of the serial nature of linguistic input for sentential complexity. *Cogn Sci* 29(2):261–290.
13. Greenberg J (1963) *Universals of Language*, ed Greenberg J (MIT Press, Cambridge, MA), pp 73–113.
14. Vennemann T (1974) Theoretical word order studies: Results and problems. *Papier Linguistik* 7:5–25.
15. Dryer MS (1992) The Greenbergian word order correlations. *Language* 68(1):81–138.
16. Gildea D, Temperley D (2010) Do grammars minimize dependency length? *Cogn Sci* 34(2):286–310.
17. Hawkins JA (2014) *Cross-Linguistic Variation and Efficiency* (Oxford Univ Press, Oxford, UK).
18. Kuhlmann M (2013) Mildly non-projective dependency grammar. *Comput Linguist* 39(2):507–514.
19. Ferrer i Cancho R (2006) Why do syntactic links not cross? *Europhys Lett* 76(6):1228.
20. Klein D, Manning CD (2004) Corpus-based induction of syntactic structure: Models of dependency and constituency. *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics* (Association for Computational Linguistics, Barcelona), pp 478–485.
21. Smith NA, Eisner J (2006) Annealing structural bias in multilingual weighted grammar induction. *Proceedings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, eds Calzolari N, Cardie C, Isabelle P (Association for Computational Linguistics, Sydney), pp 569–576.
22. Collins M (2003) Head-driven statistical models for natural language parsing. *Comput Linguist* 29(4):589–637.
23. McDonald RT, et al. (2013) Universal dependency annotation for multilingual parsing. *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, eds Fung P, Poesio M (Association for Computational Linguistics, Sofia, Bulgaria), pp 92–97.
24. Zeman D, et al. (2014) HamleDT: Harmonized multi-language dependency treebank. *Lang Resour Eval* 48(4):601–637.
25. Nivre J, et al. (2015) *Universal Dependencies 1.0* (LINDAT/CLARIN Digital Library at Institute of Formal and Applied Linguistics, Charles University in Prague, Prague).
26. de Marneffe MC, et al. (2014) Universal Stanford Dependencies: A cross-linguistic typology. *Proceedings of the Ninth International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Declerck T, Loftsson H, Maegaard B, Mariani J, Moreno A, Odijk J, Piperidis S (European Language Resources Association, Reykjavik, Iceland).
27. Ferrer i Cancho R (2004) Euclidean distance between syntactically linked words. *Phys Rev E Stat Nonlin Soft Matter Phys* 70(5 Pt 2):056135.
28. Gildea D, Temperley D (2007) Optimizing grammars for minimum dependency length. *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, eds Zaenen A, van den Bosch A (Association for Computational Linguistics, Prague), pp 184–191.
29. Park YA, Levy R (2009) Minimal-length linearizations for mildly context-sensitive dependency trees. *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, (Association for Computational Linguistics, Boulder, CO), pp 335–343.
30. Ferrer i Cancho R, Liu H (2014) The risks of mixing dependency lengths from sequences of different length. *Glottology* 5(2):143–155.
31. Noji H, Miyao Y (2014) Left-corner transitions on dependency parsing. *Proceedings of the 25th International Conference on Computational Linguistics*, eds Tsujii J, Hajič J (Association for Computational Linguistics, Dublin), pp 2140–2150.
32. Hawkins JA (1998) *Constituent Order in the Languages of Europe*, ed Siewierska A (Mouton de Gruyter, Berlin), pp 729–781.
33. Dryer MS (2002) Case distinctions, rich verb agreement, and word order type. *Theoretical Linguistics* 28(2):151–157.
34. Konieczny L (2000) Locality and parsing complexity. *J Psycholinguist Res* 29(6):627–645.
35. Zipf GK (1949) *Human Behavior and the Principle of Least Effort* (Addison-Wesley Press, Oxford).
36. Jaeger TF (2006) Redundancy and syntactic reduction in spontaneous speech. PhD thesis (Stanford Univ, Stanford, CA).
37. Piantadosi ST, Tily H, Gibson E (2011) Word lengths are optimized for efficient communication. *Proc Natl Acad Sci USA* 108(9):3526–3529.
38. Fedzechkina M, Jaeger TF, Newport EL (2012) Language learners restructure their input to facilitate efficient communication. *Proc Natl Acad Sci USA* 109(44):17897–17902.
39. Kemp C, Regier T (2012) Kinship categories across languages reflect general communicative principles. *Science* 336(6084):1049–1054.
40. Che W, Li Z, Liu T (2012) *Chinese Dependency Treebank 1.0 LDC2012T05* (Linguistic Data Consortium, Philadelphia).
41. Gelman A, Hill J (2007) *Data Analysis Using Regression and Multilevel/Hierarchical Models* (Cambridge Univ Press, Cambridge, UK).
42. Bates D, Maechler M, Bolker B, Walker S (2014) lme4: Linear mixed-effects models using Eigen and S4. R package version 1.1-7. Available at CRAN.R-project.org/package=lme4. Accessed June 18, 2015.

Supporting Information

Futrell et al. 10.1073/pnas.1502134112

Further Baselines

Consistent Head Direction Baseline. DLM has been advanced as an explanation for the consistency of head direction in languages: the fact that whether heads come before dependents or vice versa is typically consistent across dependency types within a language. The argument is that consistent head direction leads to lower dependency lengths, given that a language is not head medial, i.e., that heads do not appear between their dependents. However, there are plausible alternative explanations for consistency in head direction. Most compelling is the argument that the grammar of a language with consistent head direction is less complex than the grammar of a language with inconsistent head direction, because describing head direction in such a language requires only a single parameter. If there is a simplicity bias in grammars, then we would expect consistent head order independent of DLM. There is the possibility that our findings actually reflect independently motivated consistency in head order rather than DLM per se.

Here we test this idea by comparing languages to random and optimal baselines where head direction is fixed for all relation types. In this case, the only way that dependency length can be minimized is by choosing an optimal ordering of the dependents of a single head; this is accomplished by ordering constituents from short to long in the case of a head initial language or from long to short in the case of a head final language.

Fig. S1 shows real dependency lengths compared with the consistent head direction baselines. We find that all languages have shorter dependencies than we would expect by chance given consistent head direction. The difference between real and ran-

dom slopes is significant at $P < 0.001$ for all languages. The baseline is especially interesting in the case of the overwhelmingly head final languages in our sample, such as Japanese, Korean, Turkish, Telugu, Tamil, and Hindi. For these languages, which are similar to the baselines in the consistency of their head direction, the fact that they have dependency lengths shorter than the random baseline indicates that they accomplish dependency length minimization through long before short order.

Fixed Head Position Baseline. To what extent is DLM accomplished by choosing an optimal position of the head relative to its dependents and to what extent is it accomplished by choosing an optimal ordering of the dependents? To address this question, we compare real dependency lengths to random and optimal baselines where the position of the head and the direction of each dependent with respect to the head is fixed at the observed values. For example, given an observed head H with left dependents A , B , and C , and right dependents D , E , and F , we consider random orderings such as $[C, A, B, H, E, F, D]$, $[A, C, B, H, D, F, E]$, etc., where A , B , and C and D , E , and F are shuffled but maintain their direction with respect to the head.

Fig. S2 shows real dependency lengths compared with the random and optimal fixed head position baselines. We find that all languages have dependency lengths shorter than this baseline. The difference between real and random slopes is significant at $P < 0.001$ for all languages. The finding suggests that given a fixed head position, the ordering of dependents of the head is optimized across all languages, i.e., there is long before short order before heads and short before long order after heads.

Fig. S1. Real dependency lengths as a function of sentence length (blue) compared with the Consistent Head Direction Free Word Order Random baseline (red) and the Consistent Head Direction Free Word Order Optimal baseline (green). GAM fits are shown. rand, obs, and opt are the slopes for random, observed, and optimal dependency length as a function of squared sentence length, as in Fig. 3.

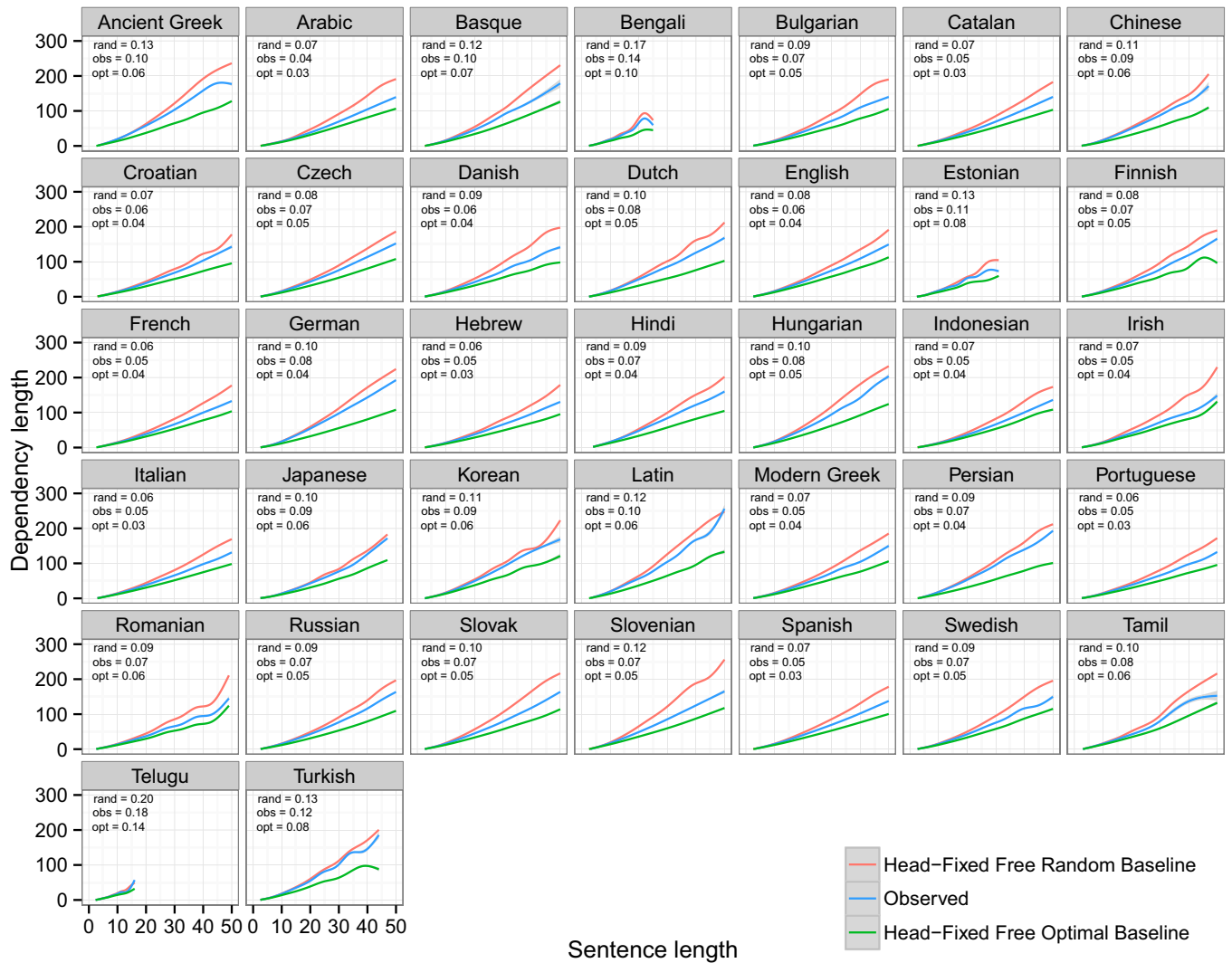


Fig. S2. Real dependency lengths as a function of sentence length (blue) compared with the Head-Fixed Free Word Order Random baseline (red) and the Head-Fixed Free Word Order Optimal baseline (green). GAM fits are shown. rand, obs, and opt are the slopes for random, observed, and optimal dependency length as a function of squared sentence length, as in Fig. 3.

Other Supporting Information Files

[Table S1 \(PDF\)](#)

Table S1. We draw data from three projects aimed at standardizing dependency annotation schemes and reformatting existing dependency corpora according to them: the Google Universal Dependency Treebank (23) (UDT), HamleDT 2.0 (24), and Universal Dependencies 1.0 (25) (UD)

Language	Family/subfamily	Standardization	Source
Arabic	Semitic	HamleDT	Prague Arabic Dependency Treebank (43)
Bulgarian	IE/Slavic	HamleDT	BulTreeBank (44)
Bengali	IE/Indo-Aryan	HamleDT	Hyderabad Dependency Treebank (45)
Catalan	IE/Romance	HamleDT	AnCorà-CA (46)
Chinese	Sino-Tibetan	—	Chinese Dependency Treebank[40]
Croatian	IE/Slavic	HamleDT+	Croatian Dependency Treebank (47)
Czech	IE/Slavic	UD	Prague Dependency Treebank (48)
Danish	IE/Germanic	HamleDT	Danish Dependency Treebank (49)
German	IE/Germanic	HamleDT	TIGER (50)
Modern Greek	IE/Greek	HamleDT	Greek Dependency Treebank (51)
English	IE/Germanic	UD	English Web Treebank (52)
Spanish	IE/Romance	HamleDT	AnCorà-ES (46)
Estonian	Uralic/Finnic	HamleDT	Eesti keele puudepank (53)
Basque	Isolate	HamleDT	Basque Dependency Treebank (54)
Persian	IE/Indo-Aryan	HamleDT	Persian Dependency Treebank (55)
Finnish	Uralic/Finnic	UD	Turku Dependency Treebank (56)
French	IE/Romance	UD	French Dependency Treebank (57)
Ancient Greek	IE/Greek	HamleDT	Perseus Project (58)
Hebrew	Semitic	HamleDT+	MILA HaAretz Treebank (59)
Hindi	IE/Indo-Aryan	HamleDT	Hyderabad Dependency Treebank (45)
Hungarian	Uralic/Ugric	HamleDT	Szeged Treebank (60)
Indonesian	Austronesian	UDT	Indonesian Dependency Treebank (61)
Irish	IE/Celtic	UD	Irish Dependency Treebank (62)
Italian	IE/Romance	UD	Italian Stanford Dependency Treebank (63)
Japanese	Isolate	HamleDT	Tüba-J/S (64)
Korean	Isolate	UDT	Universal Dependency Treebank 1.0
Latin	IE/Romance	HamleDT	Perseus Project (58)
Dutch	IE/Germanic	HamleDT	Alpino Treebank (65)
Portuguese	IE/Romance	HamleDT	Bosque/Floresta Sintá(c)tica (66)
Romanian	IE/Romance	HamleDT	Resurse pentru Gramaticile de Dependenta (67)
Russian	IE/Slavic	HamleDT	SynTagRus (68)
Slovak	IE/Slavic	HamleDT	Slovak Treebank (69)
Slovenian	IE/Slavic	HamleDT	Slovene Dependency Treebank (70)
Swedish	IE/Germanic	UD	Talbanken05 (71)
Tamil	Dravidian	HamleDT	TamilTB (72)
Telugu	Dravidian	HamleDT	Hyderabad Dependency Treebank (45)
Turkish	Turkic	HamleDT	METU-Sabancı Treebank (73)

The HamleDT project corpora are automatically converted from their original hand-parsed form to the Stanford Dependencies standard. We also used the scripts provided by the HamleDT project to convert some corpora not in the official HamleDT 2.0 release to HamleDT format (Hebrew and Croatian). These are labeled as HamleDT+ in the table.

43. Smrž O, et al. (2008) Prague Arabic dependency treebank: A word on the million words. *Proceedings of the Workshop on Arabic and Local Languages (LREC 2008)*, eds Choukri K, Diab M, Maegaard B, Rosso P, Soudi A, Farghaly A (European Language Resources Association, Marrakech, Morocco), pp 16–23.
44. Simov K, Osenova P (2005) Extending the annotation of BulTreeBank: Phase 2. *The Fourth Workshop on Treebanks and Linguistic Theories (TLT 2005)*, eds Civit M, Kuebler S, Martí MA (Universitat de Barcelona, Barcelona), pp 173–184.
45. Husain S, Mannem P, Ambati B, Gadde P (2010) The ICON-2010 tools contest on Indian language dependency parsing. *Proceedings of ICON-2010 Tools Contest on Indian Language Dependency Parsing*, eds Husain S, Mannem P, Ambati B, Gadde P (IIIT-Hyderabad, Kharagpur, India).
46. Taulé M, Martí MA, Recasens M (2008) AnCorà: Multilevel annotated corpora for Catalan and Spanish. *Proceedings of the 6th International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Maegaard B, Mariani J, Odijk J, Piperidis S, Tapias D (European Language Resources Association, Marrakech, Morocco), pp 96–101.
47. Berović D, Agić Ž, Tadić M (2012) Croatian dependency treebank: Recent development and initial experiments. *Proceedings of the 8th International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Declerck T, Uğur Doğan M, Maegaard B, Mariani J, Moreno A, Odijk J, Piperidis S (European Language Resources Association, Istanbul, Turkey), pp 1902–1906.
48. Hajič J, et al. (2006) *Prague Dependency Treebank 2.0* (Linguistic Data Consortium, Philadelphia).
49. Kromann M (2003) The Danish dependency treebank and the DTAG treebank tool. *Proceedings of the Second Workshop on Treebanks and Linguistic Theories*, eds Nivre J, Hinrichs E (Växjö University Press, Växjö, Sweden), pp 217–220.
50. Brants S, Dipper S, Hansen S, Lezius W, Smith G (2002) The TIGER treebank. *Proceedings of the Workshop on Treebanks and Linguistic Theories*, eds Hinrichs E and Simov K (Sopotol, Bulgaria), pp 24–42.
51. Prokropidis P, Desipri E, Koutsombogera M, Papageorgiou H, Piperidis S (2005) Theoretical and practical issues in the construction of a Greek dependency treebank. *Proceedings of the 4th Workshop on Treebanks and Linguistic Theories (TLT)*, eds Hajič J, Nivre J (Charles University in Prague, Department of Formal and Applied Linguistics, Prague), pp 149–160.
52. Silveira N, et al. (2014) A gold standard dependency corpus for English. *Proceedings of the 9th International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Declerck T, Loftsson H, Maegaard B, Mariani J, Moreno A, Odijk J, Piperidis S (European Language Resources Association, Reykjavik, Iceland), pp 2897–2904.
53. Bick E, Uibo H, Müürisep K (2004) Arboret: A VISL-style treebank derived from an Estonian constraint grammar corpus. *Proceedings of Treebanks and Linguistic Theories*, eds Kuebler S, Nivre J, Hinrichs E, Wunsch H (Eberhard-Karls-Universität Tübingen, Tübingen, Germany), pp 9–20.
54. Aduriz I, et al. (2003) Construction of a Basque dependency treebank. *Proceedings of the 2nd Workshop on Treebanks and Linguistic Theories*, eds Nivre J, Hinrichs E (Växjö University Press, Växjö, Sweden), pp 201–204.

55. Rasooli MS, Moloodi A, Kouhestani M, Minaei-Bidgoli B (2011) A syntactic valency lexicon for Persian verbs: The first steps towards Persian dependency treebank. *5th Language & Technology Conference (LTC): Human Language Technologies as a Challenge for Computer Science and Linguistics*, ed Vetulani Z (Adam Mickiewicz University, Poznań, Poland), pp 227–231.
56. Haverinen K, et al. (2010) Treebanking Finnish. *Proceedings of the Ninth International Workshop on Treebanks and Linguistic Theories (TLT9)*, eds Dickinson M, Müürisep K, Passarotti M (Northern European Association for Language Technology, Tartu, Estonia), pp 79–90.
57. Candito M, Crabbé B, Denis P (2010) Statistical French dependency parsing: Treebank conversion and first results. *Proceedings of the 7th International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Maegaard B, Mariani J, Odijk J, Piperidis S, Rosner M, Tapias D (European Language Resources Association, Valletta, Malta), pp 1840–1847.
58. Bamman D, Crane G (2011) The Ancient Greek and Latin dependency treebanks. *Language Technology for Cultural Heritage: Selected Papers from the LaTeCH Workshop Series, Theory and Applications of Natural Language Processing*, eds Sporleder C, von den Bosch A, Zervanou K (Springer, Berlin), pp 79–98.
59. Sima'an K, Itai A, Winter Y, Altman A, Nativ N (2001) Building a tree-bank of Modern Hebrew text. *Traitement Automatique des Langues* 42(2):347–380.
60. Csendes D, Csirik J, Gyimóthy T, Kocsor A (2005) The Szeged Treebank (TSD). *Text, Speech and Dialogue, 8th International Conference, Lecture Notes in Computer Science*, eds Matousek V, Mautner P, Pavelka T (Springer, Berlin), Vol 3658, pp 123–131.
61. Green N, Larasati SD, Žabokrtský Z (2012) Indonesian dependency treebank: Annotation and parsing. *Proceedings of the 26th Pacific Asia Conference on Language, Information and Computation*, eds Manurung R, Bond F (Faculty of Computer Science, Universitas Indonesia, Bali, Indonesia), pp 137–145.
62. Lynn T, et al. (2012) Irish treebanking and parsing: A preliminary evaluation. *Proceedings of the 8th International Conference on Linguistic Resources and Evaluation*, eds Calzolari N, Choukri K, Declerck T, Uğur Doğan M, Maegaard B, Mariani J, Moreno A, Odijk J, Piperidis S (European Language Resources Association, Istanbul, Turkey), pp 1939–1946.
63. Bosco C, Montemagni S, Simi M (2013) Converting Italian treebanks: Towards an Italian Stanford dependency treebank. *The 7th Linguistic Annotation Workshop Interoperability with Discourse*, eds Pareja-Lora A, Liakata M, Dipper S (Association for Computational Linguistics, Sofia, Bulgaria), pp 61–69.
64. Kawata Y, Bartels J (2000) *Stylebook for the Japanese Treebank in VerbMobil*. Report 240 (Eberhard-Karls-Universitaet Tuebingen, Tuebingen, Germany).
65. van der Beek L, et al. (2002) The Alpino dependency treebank. *Algorithms for Linguistic Processing NWO PIONIER Progress Report*, eds van der Beek L, Bouma G, Daciuk J, Gaustad T, Malouf R, van Noord G, Prins R, Villada V (Graduate School for Behavioral and Cognitive Neurosciences Alfa-informatica, Groningen, The Netherlands), Chap 5.
66. Afonso S, Bick E, Haber R, Santos D (2002) "Floresta sintá(c)tica": A treebank for Portuguese. *Proceedings of the 3rd International Conference on Language Resources and Evaluation*, eds Zampolli A, (European Language Resources Association, Las Palmas, Spain), pp 1968–1703.
67. Căläcean M (2008) Data-driven dependency parsing for Romanian. MS thesis (Uppsala Univ, Uppsala).
68. Boguslavsky I, Grigorieva S, Grigoriev N, Kreidlin L, Frid N (2000) Dependency treebank for Russian: Concept, tools, types of information. *Proceedings of the 18th Conference on Computational Linguistics*, eds Kay M, (Association for Computational Linguistics, Saarbruecken, Germany), Vol 2, pp 987–991.
69. Šimková M, Garabik R (2006) Sintaksičeskaja razmetka v slovackom nacional'nom korpusse. Trudy meždunarodnoj konferencii Korpusnaja lingvistika, (University of St. Petersburg, St. Petersburg, Russia), pp 389–394.
70. Džeroski S, et al. (2006) Towards a Slovene dependency treebank. *Proceedings of the Fifth International Language Resources and Evaluation Conference*, eds Calzolari N, Maegaard B, Choukri K, Marconi L (European Language Resources Association, Genoa, Italy), pp 1388–1391.
71. Nivre J, Nilsson J, Hall J (2006) Talbanken05: A Swedish treebank with phrase structure and dependency annotation. *Proceedings of the 5th International Conference on Language Resources and Evaluation*, eds Calzolari N, Maegaard B, Choukri K, Marconi L (European Language Resources Association, Genoa, Italy), pp. 1392–1395.
72. Ramasamy L, Žabokrtský Z (2012) Prague dependency style treebank for Tamil. *Proceedings of the 8th International Conference on Language Resources and Evaluation*, eds Calzolari N, Choukri K, Declerck T, Uğur Doğan M, Maegaard B, Mariani J, Moreno A, Odijk J, Piperidis S (European Language Resources Association, Istanbul, Turkey), pp 1888–1894.
73. Atalay NB, Oflazer K, Say B (2003) The annotation process in the Turkish treebank. *Proceedings of the 4th International Workshop on Linguistically Interpreted Corpora (LINC)*, eds Abeillé A, Hansen-Schirra S, Uszkoreit H (Association for Computational Linguistics, Budapest), pp 33–38.