

Intro

- Starting point: a classical philologist interested in Greek and Latin word order

Intro

- Starting point: a classical philologist interested in Greek and Latin word order
 - I need data → I need a treebank
 - I need an annotation scheme → I need linguistic theory

Intro

- Starting point: a classical philologist interested in Greek and Latin word order
 - I need data → I need a treebank
 - I need an annotation scheme → I need linguistic theory
 - What can linguistic theory do for treebanks?

Greek and Latin word order: state of the art

- The characteristic feature of languages like Greek and Latin are their free word order:
 - All permutations of S, V and O found with reasonable frequency
 - Discontinuous constituency is common
- The agreement stops there. . .

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause
 - The verb assigns accusative, genitive, or dative to an argument that is a patient or theme

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause
 - The verb assigns accusative, genitive, or dative to an argument that is a patient or theme
 - The V consists of one word (no periphrastic forms, modal embeddings or light verbs)

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause
 - The verb assigns accusative, genitive, or dative to an argument that is a patient or theme
 - The V consists of one word (no periphrastic forms, modal embeddings or light verbs)
 - S and O are determiner phrases (this includes nominalizations) or quantifier phrases, and not clausal

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause
 - The verb assigns accusative, genitive, or dative to an argument that is a patient or theme
 - The V consists of one word (no periphrastic forms, modal embeddings or light verbs)
 - S and O are determiner phrases (this includes nominalizations) or quantifier phrases, and not clausal
 - S and O are continuous strings
- Admirably explicit, but what are the chances of getting things right at first try?

Being specific about it

- The criteria of Kirk (2012)
 - The clause contains at least an S(ubject), V(erb) and O(bject)
 - The clause is continuous
 - S and O are not embedded in a participial clause
 - The verb assigns accusative, genitive, or dative to an argument that is a patient or theme
 - The V consists of one word (no periphrastic forms, modal embeddings or light verbs)
 - S and O are determiner phrases (this includes nominalizations) or quantifier phrases, and not clausal
 - S and O are continuous strings
- Admirably explicit, but what are the chances of getting things right at first try?
- Even if everything that should be included is included, things may have been excluded that should not have been

The cure

- We need treebanks and linguistic theory
 - Treebanks for replicability and iterated search
 - Linguistic theory for annotation schemes, and to know what we are looking for
- “annotations are no substitute for the understanding of a phenomenon. They are an encoding of that understanding.” Zaenen (2006)

Creating a treebank

- There are obvious upsides for a linguist in working within an established linguistic theory
 - Build on earlier work

Creating a treebank

- There are obvious upsides for a linguist in working within an established linguistic theory
 - Build on earlier work
 - Compare and contrast analyses for specific languages

Creating a treebank

- There are obvious upsides for a linguist in working within an established linguistic theory
 - Build on earlier work
 - Compare and contrast analyses for specific languages
 - Compare and contrast with other theories

Creating a treebank

- There are obvious upsides for a linguist in working within an established linguistic theory
 - Build on earlier work
 - Compare and contrast analyses for specific languages
 - Compare and contrast with other theories
- Unfortunately there are also downsides in treebanking with an established linguistic theory
 - Requires linguistically (very) sophisticated annotators
 - Hard to use for outsiders

Creating a treebank

- There are obvious upsides for a linguist in working within an established linguistic theory
 - Build on earlier work
 - Compare and contrast analyses for specific languages
 - Compare and contrast with other theories
- Unfortunately there are also downsides in treebanking with an established linguistic theory
 - Requires linguistically (very) sophisticated annotators
 - Hard to use for outsiders
 - Risks circular confirmation of biases in the theory
- So we need “theory-neutral” treebanks

Principles of annotation

Two conflicting constraints on annotation

- 1 Encode enough structure to enable reconstruction of theoretically motivated structures

Principles of annotation

Two conflicting constraints on annotation

- 1 Encode enough structure to enable reconstruction of theoretically motivated structures
 - 2 Encode no more structure than is common to all frameworks
- 1 is obvious if the goal is theoretical adequacy

Principles of annotation

Two conflicting constraints on annotation

- 1 Encode enough structure to enable reconstruction of theoretically motivated structures
 - 2 Encode no more structure than is common to all frameworks
- 1 is obvious if the goal is theoretical adequacy
 - 2 is desirable to minimize the assumptions that go into the annotation and hence cannot be tested using the corpus

A naturally occurring parallel corpus

- The New Testament in its Greek original and Latin, Gothic, Classical Armenian and OCS translations
- The NT translations are the oldest attestations of Armenian and OCS, and virtually the only attestation of Gothic

A naturally occurring parallel corpus

- The New Testament in its Greek original and Latin, Gothic, Classical Armenian and OCS translations
- The NT translations are the oldest attestations of Armenian and OCS, and virtually the only attestation of Gothic
- The gospels constitute the core of the OCS text canon

A naturally occurring parallel corpus

- The New Testament in its Greek original and Latin, Gothic, Classical Armenian and OCS translations
- The NT translations are the oldest attestations of Armenian and OCS, and virtually the only attestation of Gothic
- The gospels constitute the core of the OCS text canon
- So these are important texts, and they are parallel texts

Later extensions

- Classical Greek and Latin:
 - Herodotus
 - Gallic War, Letters to Atticus, De officiis

Later extensions

- Classical Greek and Latin:
 - Herodotus
 - Gallic War, Letters to Atticus, De officiis
- Post-classical Greek and Latin
 - Sphrantzes' Chronicles
 - Peregrinatio Aetheriae
 - Palladius' De Agricultura

Other corpora in the same syntactic annotation scheme

- Poetic Edda (Greinir skáldskapar)

Many-layered annotation

- Morphology
- Syntax

Many-layered annotation

- Morphology
- Syntax
- Semantics and other customised annotation (e.g. animacy)
- Givenness

Many-layered annotation

- Morphology
- Syntax
- Semantics and other customised annotation (e.g. animacy)
- Givenness
- Experimental discourse structure annotation

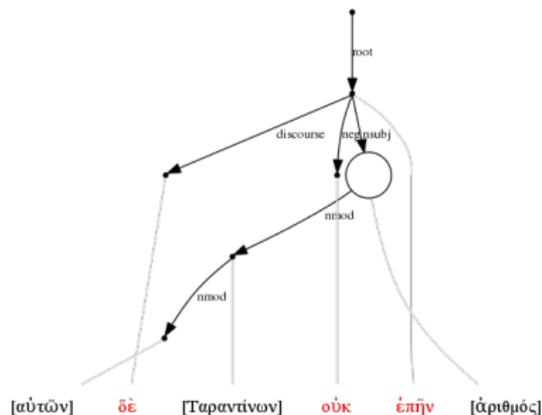
Syntactic annotation scheme

- Dependency grammar
 - practicable, easy to teach
 - no assumptions about word order/constituency

Syntactic annotation scheme

- Dependency grammar
 - practicable, easy to teach
 - no assumptions about word order/constituency
- But some aspects were felt too limiting already at the outset
 - We violate the unique head principle with secondary edges in e.g. control structures
 - We use empty nodes in e.g. ellipsis structures
- Experience now vindicates these choices
- For standardization a UD version is available

Gap degree: the concept



Gap degree (trees) without punctuation

| | 0 | 1 | 2 | 3 | | 0 | 1 | 2 | 3 |
|----------------------|-------|-------|-------|-------|---------------------|-------|-------|-------|-------|
| Ancient_Greek | 0.372 | 0.547 | 0.076 | 0.004 | Hungarian | 0.815 | 0.172 | 0.010 | 0.002 |
| Ancient_Greek-PROIEL | 0.605 | 0.351 | 0.038 | 0.004 | Indonesian | 0.994 | 0.006 | 0.000 | 0.000 |
| Arabic | 0.945 | 0.055 | 0.000 | 0.000 | Irish | 0.883 | 0.117 | 0.000 | 0.000 |
| Basque | 0.771 | 0.215 | 0.013 | 0.000 | Italian | 0.988 | 0.011 | 0.000 | 0.000 |
| Bulgarian | 0.972 | 0.028 | 0.000 | 0.000 | Japanese-KTC | 1.000 | 0.000 | 0.000 | 0.000 |
| Croatian | 0.936 | 0.063 | 0.001 | 0.000 | Latin | 0.551 | 0.409 | 0.038 | 0.002 |
| Czech | 0.876 | 0.121 | 0.003 | 0.000 | Latin-ITT | 0.631 | 0.350 | 0.018 | 0.001 |
| Danish | 0.750 | 0.250 | 0.000 | 0.000 | Latin-PROIEL | 0.699 | 0.274 | 0.024 | 0.002 |
| Dutch | 0.831 | 0.156 | 0.012 | 0.000 | Norwegian | 0.926 | 0.074 | 0.000 | 0.000 |
| English | 1.000 | 0.000 | 0.000 | 0.000 | Old_Church_Slavonic | 0.784 | 0.203 | 0.012 | 0.000 |
| Estonian | 0.992 | 0.008 | 0.000 | 0.000 | Persian | 0.955 | 0.045 | 0.000 | 0.000 |
| Finnish | 0.924 | 0.074 | 0.001 | 0.000 | Persian | 0.997 | 0.003 | 0.000 | 0.000 |
| Finnish-FTB | 0.933 | 0.065 | 0.002 | 0.000 | Portuguese | 0.838 | 0.149 | 0.013 | 0.000 |
| French | 0.960 | 0.039 | 0.001 | 0.000 | Romanian | 0.892 | 0.100 | 0.006 | 0.002 |
| German | 0.926 | 0.074 | 0.001 | 0.000 | Slovenian | 0.869 | 0.122 | 0.008 | 0.001 |
| Gothic | 0.761 | 0.224 | 0.012 | 0.002 | Spanish | 0.964 | 0.036 | 0.000 | 0.000 |
| Greek | 0.782 | 0.212 | 0.007 | 0.000 | Swedish | 0.973 | 0.027 | 0.000 | 0.000 |
| Hebrew | 1.000 | 0.000 | 0.000 | 0.000 | Tamil | 0.987 | 0.013 | 0.000 | 0.000 |
| Hindi | 0.864 | 0.132 | 0.003 | 0.000 | | | | | |

Gap degree (trees) without punctuation

| | 0 | 1 | 2 | 3 | | 0 | 1 | 2 | 3 |
|----------------------|-------|-------|--------------|--------------|---------------------|-------|-------|-------|-------|
| Ancient_Greek | 0.372 | 0.547 | 0.076 | 0.004 | Hungarian | 0.815 | 0.172 | 0.010 | 0.002 |
| Ancient_Greek-PROIEL | 0.605 | 0.351 | 0.038 | 0.004 | Indonesian | 0.994 | 0.006 | 0.000 | 0.000 |
| Arabic | 0.945 | 0.055 | 0.000 | 0.000 | Irish | 0.883 | 0.117 | 0.000 | 0.000 |
| Basque | 0.771 | 0.215 | 0.013 | 0.000 | Italian | 0.988 | 0.011 | 0.000 | 0.000 |
| Bulgarian | 0.972 | 0.028 | 0.000 | 0.000 | Japanese-KTC | 1.000 | 0.000 | 0.000 | 0.000 |
| Croatian | 0.936 | 0.063 | 0.001 | 0.000 | Latin | 0.551 | 0.409 | 0.038 | 0.002 |
| Czech | 0.876 | 0.121 | 0.003 | 0.000 | Latin-ITT | 0.631 | 0.350 | 0.018 | 0.001 |
| Danish | 0.750 | 0.250 | 0.000 | 0.000 | Latin-PROIEL | 0.699 | 0.274 | 0.024 | 0.002 |
| Dutch | 0.831 | 0.156 | 0.012 | 0.000 | Norwegian | 0.926 | 0.074 | 0.000 | 0.000 |
| English | 1.000 | 0.000 | 0.000 | 0.000 | Old_Church_Slavonic | 0.784 | 0.203 | 0.012 | 0.000 |
| Estonian | 0.992 | 0.008 | 0.000 | 0.000 | Persian | 0.955 | 0.045 | 0.000 | 0.000 |
| Finnish | 0.924 | 0.074 | 0.001 | 0.000 | Persian | 0.997 | 0.003 | 0.000 | 0.000 |
| Finnish-FTB | 0.933 | 0.065 | 0.002 | 0.000 | Portuguese | 0.838 | 0.149 | 0.013 | 0.000 |
| French | 0.960 | 0.039 | 0.001 | 0.000 | Romanian | 0.892 | 0.100 | 0.006 | 0.002 |
| German | 0.926 | 0.074 | 0.001 | 0.000 | Slovenian | 0.869 | 0.122 | 0.008 | 0.001 |
| Gothic | 0.761 | 0.224 | 0.012 | 0.002 | Spanish | 0.964 | 0.036 | 0.000 | 0.000 |
| Greek | 0.782 | 0.212 | 0.007 | 0.000 | Swedish | 0.973 | 0.027 | 0.000 | 0.000 |
| Hebrew | 1.000 | 0.000 | 0.000 | 0.000 | Tamil | 0.987 | 0.013 | 0.000 | 0.000 |
| Hindi | 0.864 | 0.132 | 0.003 | 0.000 | | | | | |

Immediate lessons

- Careful with your commas!
- Results in line with reports on Arabic (Kuhlmann, 2010), Czech (Kuhlmann & Nivre, 2006), Danish (Kuhlmann & Nivre, 2006), Slovene (Kuhlmann, 2010), Swedish (Gómez-rodríguez et al., 2009): gap degree 1 gives very little loss (<0.5%)

Immediate lessons

- Careful with your commas!
- Results in line with reports on Arabic (Kuhlmann, 2010), Czech (Kuhlmann & Nivre, 2006), Danish (Kuhlmann & Nivre, 2006), Slovene (Kuhlmann, 2010), Swedish (Gómez-rodríguez et al., 2009): gap degree 1 gives very little loss (<0.5%)
- The same holds for Croatian, English, Estonian, Finnish, French, German, Hindi, Indonesian, Irish, Italian, Japanese, Norwegian, Persian, Polish, Spanish and Tamil

Immediate lessons

- Careful with your commas!
- Results in line with reports on Arabic (Kuhlmann, 2010), Czech (Kuhlmann & Nivre, 2006), Danish (Kuhlmann & Nivre, 2006), Slovene (Kuhlmann, 2010), Swedish (Gómez-rodríguez et al., 2009): gap degree 1 gives very little loss (<0.5%)
- The same holds for Croatian, English, Estonian, Finnish, French, German, Hindi, Indonesian, Irish, Italian, Japanese, Norwegian, Persian, Polish, Spanish and Tamil
- Ancient Greek (Gothic, OCS), Basque, Dutch, Greek, Hungarian, Latin, Portuguese, Romanian and Slovenian: more than 0.5% loss

Immediate lessons

- Careful with your commas!
- Results in line with reports on Arabic (Kuhlmann, 2010), Czech (Kuhlmann & Nivre, 2006), Danish (Kuhlmann & Nivre, 2006), Slovene (Kuhlmann, 2010), Swedish (Gómez-rodríguez et al., 2009): gap degree 1 gives very little loss (<0.5%)
- The same holds for Croatian, English, Estonian, Finnish, French, German, Hindi, Indonesian, Irish, Italian, Japanese, Norwegian, Persian, Polish, Spanish and Tamil
- Ancient Greek (Gothic, OCS), Basque, Dutch, Greek, Hungarian, Latin, Portuguese, Romanian and Slovenian: more than 0.5% loss
- Ancient Greek stands out (Mambrini & Passarotti, 2013)

Theoretical implications

- ‘Standard’ (linguistically well-developed) mildly context-sensitive grammar formalisms such as TAG and CCG do not generate structures with gap degree > 1

Theoretical implications

- 'Standard' (linguistically well-developed) mildly context-sensitive grammar formalisms such as TAG and CCG do not generate structures with gap degree > 1
- There are extensions (e.g. multiple-component tree-adjoining grammars), but we lose the benefits of working with well-developed theories for cross-linguistic comparison
- So it is reasonable to "step up" to unification-based formalisms

Theoretical implications

- ‘Standard’ (linguistically well-developed) mildly context-sensitive grammar formalisms such as TAG and CCG do not generate structures with gap degree > 1
- There are extensions (e.g. multiple-component tree-adjoining grammars), but we lose the benefits of working with well-developed theories for cross-linguistic comparison
- So it is reasonable to “step up” to unification-based formalisms
- Lexical-functional grammar (LFG) is particularly well-developed when it comes to studies of discontinuous syntax

Lexical-functional grammar

- LFG analyses sentences in terms of a surface-oriented c(onstituent)-structure, a more abstract f(eature)-structure and the mapping between them

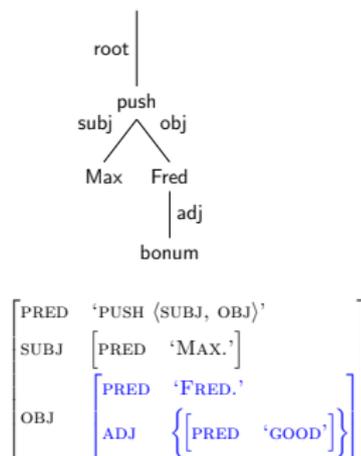
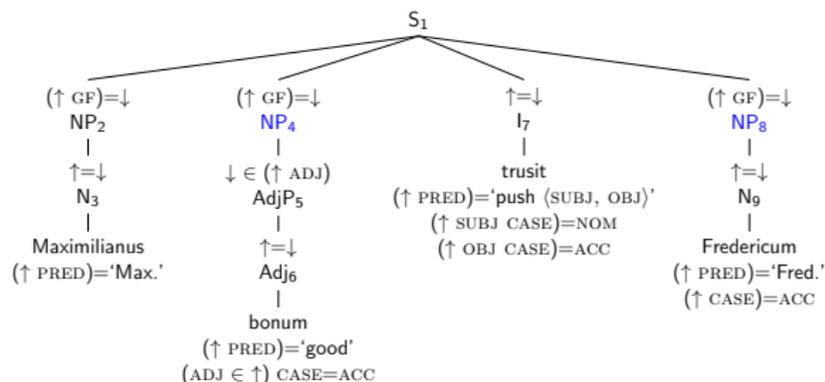
Lexical-functional grammar

- LFG analyses sentences in terms of a surface-oriented c(onstituent)-structure, a more abstract f(eature)-structure and the mapping between them
- Bröker et al. (1994): Dependency grammar is an LFG that only knows f-structure
- A DG treebank is a perfect match if we want to study c-structures given particular f-structures!

Lexical-functional grammar

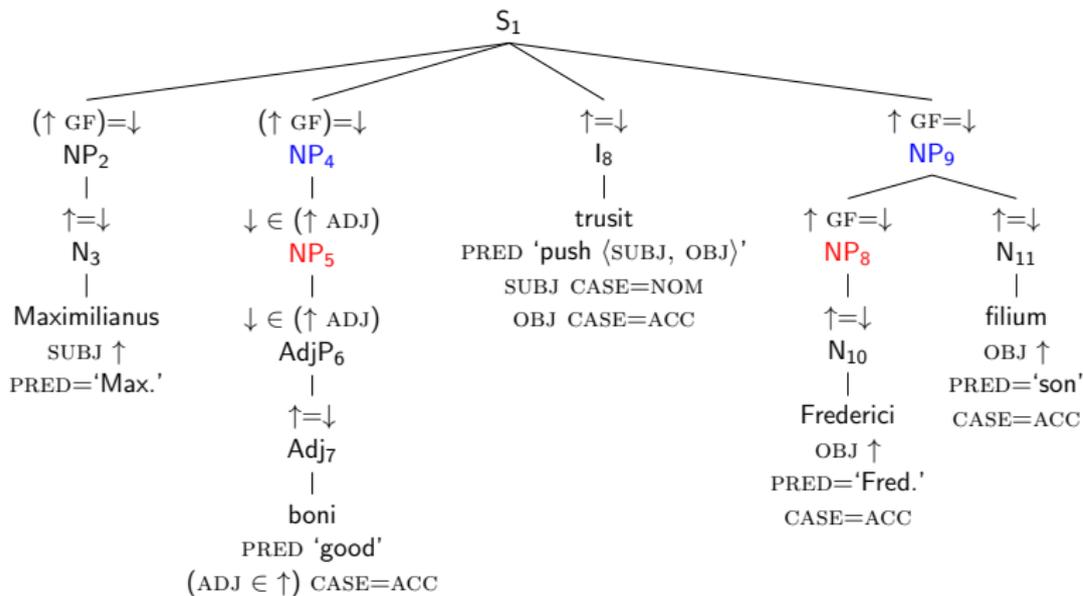
- LFG analyses sentences in terms of a surface-oriented c(onstituent)-structure, a more abstract f(eature)-structure and the mapping between them
- Bröker et al. (1994): Dependency grammar is an LFG that only knows f-structure
- A DG treebank is a perfect match if we want to study c-structures given particular f-structures!
- In practice, DG analyses are in between c- and f-structures:
 - Typically surface-oriented: the tokens are the nodes

An LFG analysis



- Dependency graph \approx f-structure
- Discontinuities correspond to reentrancies (one f-structure corresponding to multiple phrase structure nodes)

LFG analysis: gap depth



- Deeper gaps require more reentrancies

Depth vs. degree (edges)

| Universal dependencies | Degree | Depth | | | | | | |
|------------------------|---------|-------|------|-----|-----|----|----|----|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6+ |
| 0 | 1416015 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 58667 | 4833 | 592 | 128 | 30 | 43 | |
| 2 | 0 | 3475 | 294 | 49 | 8 | 2 | 0 | |
| 3 | 0 | 266 | 15 | 2 | 0 | 0 | 0 | |
| 4 | 0 | 40 | 3 | 0 | 0 | 0 | 0 | |
| 5 | 0 | 18 | 3 | 1 | 0 | 0 | 0 | |
| 6 | 0 | 4 | 1 | 0 | 0 | 0 | 0 | |
| 7 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | |

af-dauþs* (Ger.) Pl. Pr. zu *af-
 dauþ, éckauþog geschunden,
 gepagt; N.M. über M. 1300.
 dauþs Pl. Pr. zu *af-
 dauþ, éckauþog geschunden,
 wöþir eðstíðauþerod k 2.15
 E.5.2; N. R. 1217 k 2.1.16; J.
 k 2.14; G. J. 1200; D. E. vomnþoc
 dauþeins Pl. (152) vékþwic das
 Absterben A. k 4.10; éy þauvá-
 þwic in einim in Todesnöten
 k 1.23.
 dauþjan sic. V.1 vékþvön tóten
 C.5.5.
 af-dauþjan tóten (perfektiv, 291 ff.)

Depth vs. degree (edges)

| Universal dependencies | Degree | Depth | | | | | | |
|-------------------------|---------|-------|------|-----|-----|----|----|----|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6+ |
| 0 | 1416015 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 58667 | 4833 | 592 | 128 | 30 | 43 | |
| 2 | 0 | 3475 | 294 | 49 | 8 | 2 | 0 | |
| 3 | 0 | 266 | 15 | 2 | 0 | 0 | 0 | |
| 4 | 0 | 40 | 3 | 0 | 0 | 0 | 0 | |
| 5 | 0 | 18 | 3 | 1 | 0 | 0 | 0 | |
| 6 | 0 | 4 | 1 | 0 | 0 | 0 | 0 | |
| 7 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | |
| <hr/> | | | | | | | | |
| UD-Ancient Greek | | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| 0 | 65707 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 16182 | 1548 | 187 | 21 | 3 | 1 | |
| 2 | 0 | 1259 | 116 | 16 | 3 | 0 | 0 | |
| 3 | 0 | 73 | 2 | 0 | 0 | 0 | 0 | |
| 4 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | |
| <hr/> | | | | | | | | |
| UD-Ancient Greek-PROIEL | | 0 | 1 | 2 | 3 | 4 | 5 | 6+ |
| 0 | 75129 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 8900 | 596 | 65 | 13 | 2 | 8 | |
| 2 | 0 | 742 | 57 | 5 | 2 | 0 | 0 | |
| 3 | 0 | 88 | 5 | 0 | 0 | 0 | 0 | |
| 4 | 0 | 18 | 1 | 0 | 0 | 0 | 0 | |
| 5 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | |
| 6 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |
| 7 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |

Aiming too high

- For practical parsing purposes we could limit ourselves to gap degree 1 + gap degree 2 with depth 1 only
- LFG as a theory as a theory could derive any gap degree and depth, but
 - reflects the low complexity in a low number of reentrancies in the LFG analyses
 - offers a body of theoretical and cross-linguistic work to lean on

Aiming too high

- For practical parsing purposes we could limit ourselves to gap degree 1 + gap degree 2 with depth 1 only
- LFG as a theory as a theory could derive any gap degree and depth, but
 - reflects the low complexity in a low number of reentrancies in the LFG analyses
 - offers a body of theoretical and cross-linguistic work to lean on
- To connect we need to derive c-structures from the dependencies (Haug, 2012)

Order domains (Adapted from Bröker 1998)

Definition

The order domain \mathcal{D}_w of a word w is the largest subset of \mathcal{W} such that

- 1 $w \in \mathcal{D}_w$

Order domains (Adapted from Bröker 1998)

Definition

The order domain \mathcal{D}_w of a word w is the largest subset of \mathcal{W} such that

- 1 $w \in \mathcal{D}_w$
 - 2 all words in \mathcal{D}_w are dominated by w
 - 3 \mathcal{D}_w is continuous, i.e. for any two words in \mathcal{D}_w , all words in between are also contained in \mathcal{D}_w
- Intuitively, the order domain corresponds to all of the node's dependents that are not 'displaced'

Order domain structures

Definition

The order domain structure \mathcal{O} of a sentence S with the words \mathcal{W} is the set of order domains of all words $w \in \mathcal{W}$.

Order domain structures

Definition

The order domain structure \mathcal{O} of a sentence S with the words \mathcal{W} is the set of order domains of all words $w \in \mathcal{W}$.

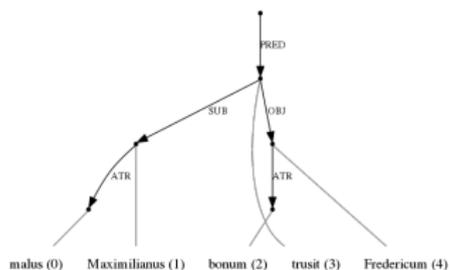
Order domain structures

Definition

The order domain structure \mathcal{O} of a sentence S with the words \mathcal{W} is the set of order domains of all words $w \in \mathcal{W}$.

- Subset inclusion corresponds to phrase structure dominance
- Each order domain is continuous, so we have a total precedence relation
- \mathcal{O} is an ordered tree

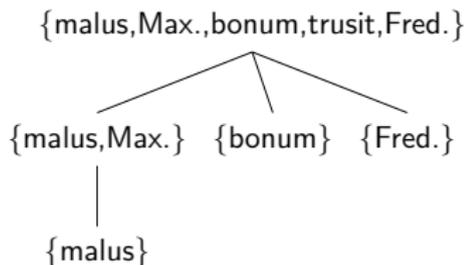
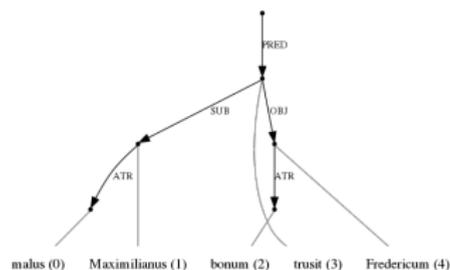
Example



{malus,Max.,bonum,trusit,Fred.},
{malus,Max.} {malus}, {bonum}, {Fred.}

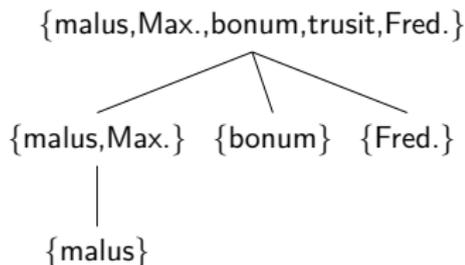
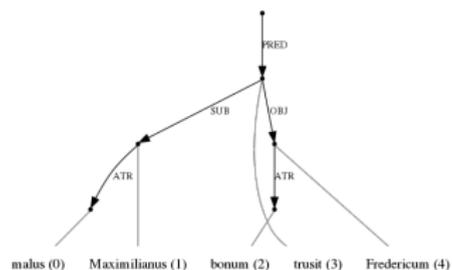
- We order the order domains by subset inclusion and precedence

Example



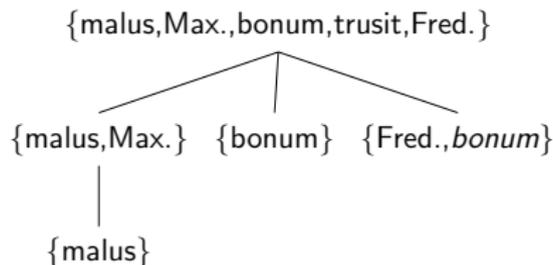
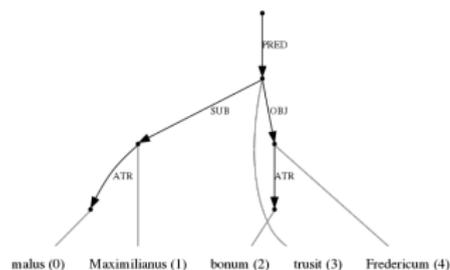
- We order the order domains by subset inclusion and precedence
- Problem: no way to retrieve the dependency of *bonum* on *Fredericum*

Example



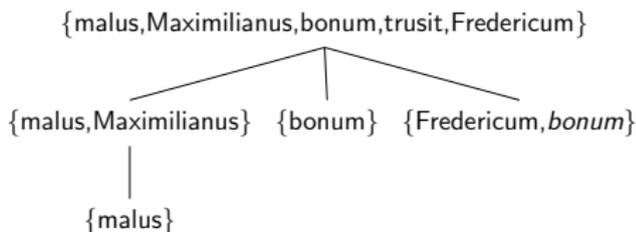
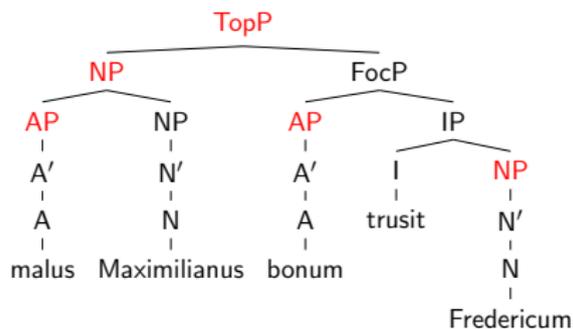
- We order the order domains by subset inclusion and precedence
- Problem: no way to retrieve the dependency of *bonum* on *Fredericum*
- Solution: add a trace keeping track of the discontinuity

Example

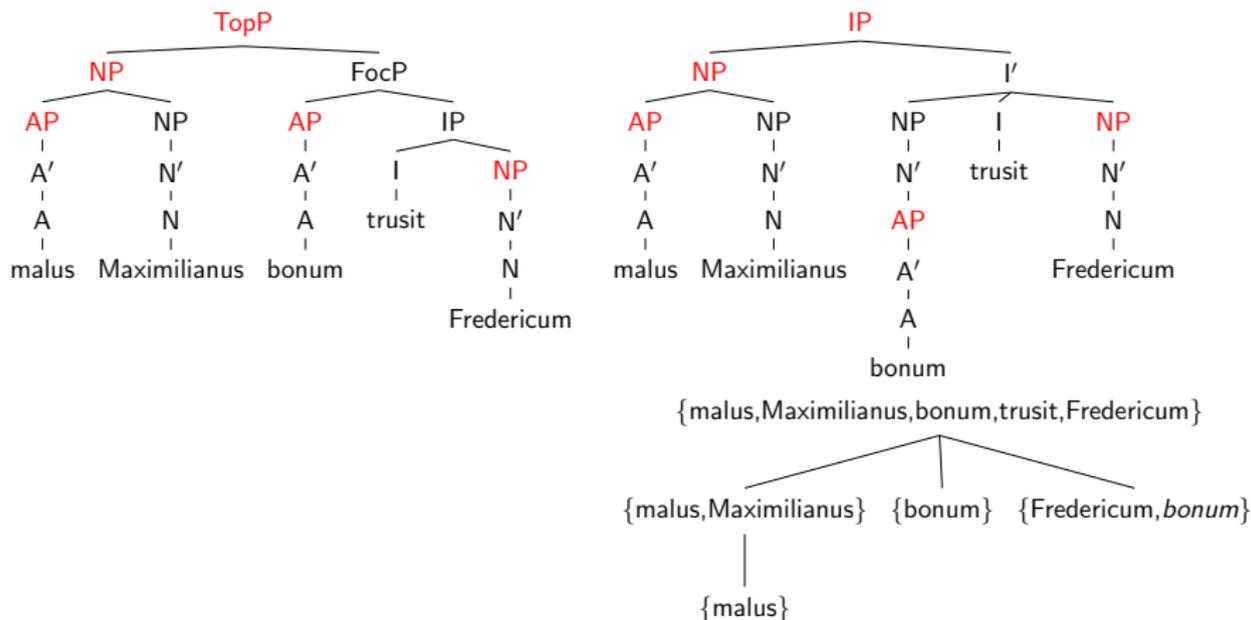


- We order the order domains by subset inclusion and precedence
- Problem: no way to retrieve the dependency of *bonum* on *Fredericum*
- Solution: add a trace keeping track of the discontinuity
- We get a structure that is implicitly present in the dependency graph. but is isomorphic to the expected phrase structure

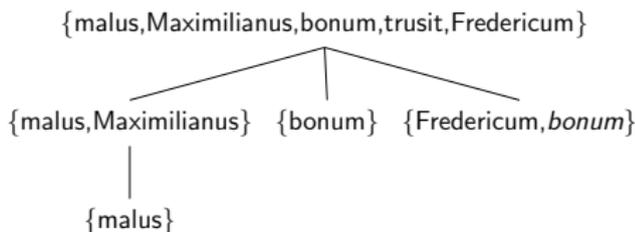
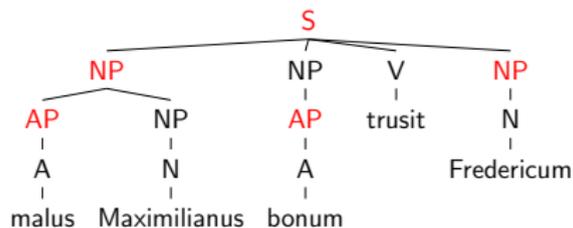
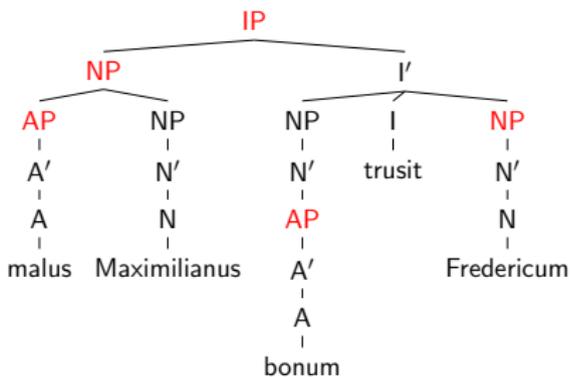
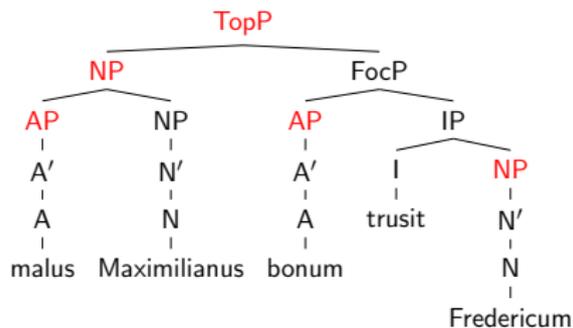
Some alternatives



Some alternatives



Some alternatives



Hypothesis testing

- The configurations of the maximal projections in these trees are all isomorphic to the order domain structure
- So we can view the creation of phrase structures from the dependency structure as an expansion of the order domain structure

Hypothesis testing

- The configurations of the maximal projections in these trees are all isomorphic to the order domain structure
- So we can view the creation of phrase structures from the dependency structure as an expansion of the order domain structure
- The task then is to determine the internal structure of each word's projection

Why does this work?

- In Greek and Latin syntax, we know more about grammatical relations than about phrase structure
- This may be true for discontinuous syntax in general

Why does this work?

- In Greek and Latin syntax, we know more about grammatical relations than about phrase structure
- This may be true for discontinuous syntax in general
- DG here achieves the ideal situation: the treebank encodes our linguistic understanding, but does not make presuppositions about uncertain things

Why does this work?

- In Greek and Latin syntax, we know more about grammatical relations than about phrase structure
- This may be true for discontinuous syntax in general
- DG here achieves the ideal situation: the treebank encodes our linguistic understanding, but does not make presuppositions about uncertain things
- Function words are typically challenging, as they are often c-structure heads, but typically taken as dependents in DG analyses

Participle clauses

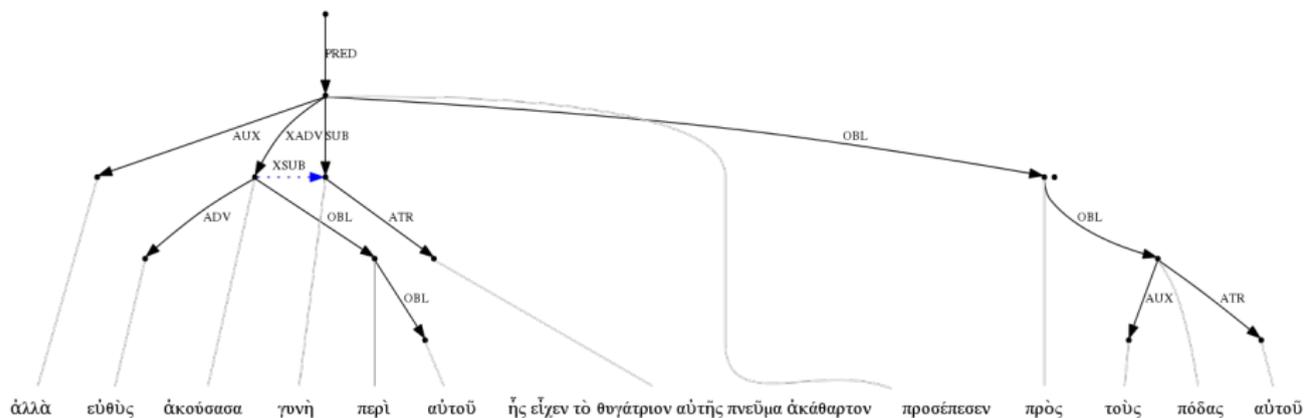
- Let us now look closer at one of Kirk's exclusion cases:
 - S and O are not embedded in a participial clause

Participle clauses

- Let us now look closer at one of Kirk's exclusion cases:
 - S and O are not embedded in a participial clause

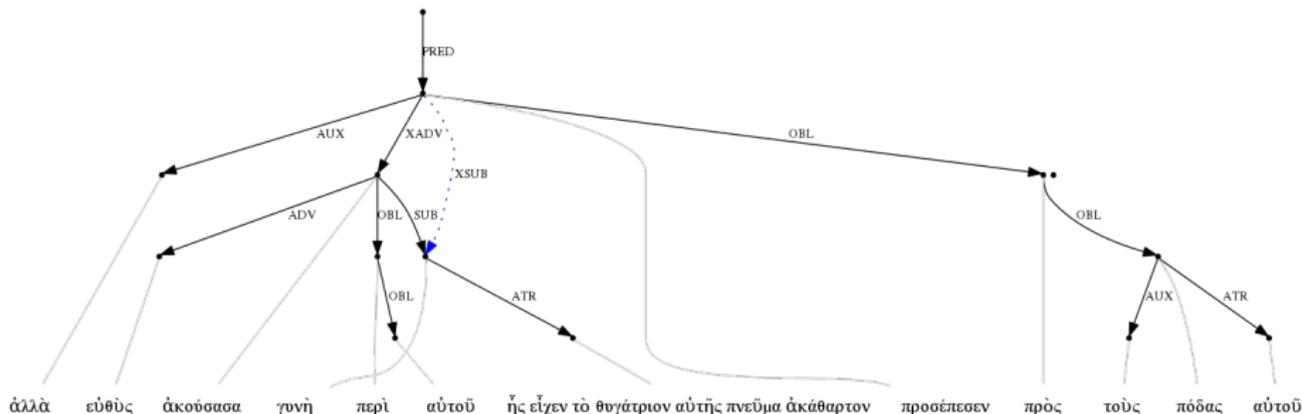
- (1) akousasa gunê peri autou
 hearing woman about him
 hês eikhen to thugatron autês pneuma akatharton
 whose daughter had an evil spirit
 prosepesen pros tous podas autou
 fell down before his feet (Mark 7.25)

The analysis



- On this analysis, it is legitimate to analyze the clause as S - V - OBL

The alternative



- We get one discontinuity (the relative clause) and no illnesting
- On this analysis it is **not** legitimate to count the matrix clause as S - V - OBL because there is no overt S

How to choose?

- The question is better framed in terms of constituents than dependencies, i.e. as questions about properties of the participle's projection

How to choose?

- The question is better framed in terms of constituents than dependencies, i.e. as questions about properties of the participle's projection
- We can connect with known typologies of discontinuities
 - unbounded dependencies
 - extraposition
 - scrambling

How free is the word order?

- The question, then, is how free the word order in participle clauses is

How free is the word order?

- The question, then, is how free the word order in participle clauses is
- We know that unbounded dependencies can break up clausal categories
- Can participle clauses be broken up by scrambling too?
- We therefore need to count both analyses and test their predictions
 - The internal subject hypothesis predicts only the subject can interrupt the participle clause

How free is the word order?

- The question, then, is how free the word order in participle clauses is
- We know that unbounded dependencies can break up clausal categories
- Can participle clauses be broken up by scrambling too?
- We therefore need to count both analyses and test their predictions
 - The internal subject hypothesis predicts only the subject can interrupt the participle clause
 - The discontinuous clause hypothesis predicts other things to intervene as well
- We generate c-structures based on both the SUB and XSUB relation

Why does this work?

- Crucially, the answer is not directly encoded in the representation

Why does this work?

- Crucially, the answer is not directly encoded in the representation
- There is a core of “facts” (incl. multiple subjecthood) encoded as theory-neutral as possible

Pullum's conjecture (Pullum, 1982, p. 8)

“no constituent of a recursive category (one that can immediately dominate itself) can scramble out of that category.”

The obvious problem: infinitives

- It is well-known that German has long-distance scrambling across (coherent) infinitive constructions

(2) dass [die Witwen]_j [der Opfer]_i [dem Pfarrer]_k der Rat
 that the widows the victims the priest the council
 gedenken zu lassen versprochen hat.
 to commemorate let promised have
 ‘... that the council has promised the priest to let the widows
 commemorate the victims’ (Becker et al., 1991)’

A challenge for linguistic theories: Illnestedness

- As we saw, the correct analysis of participle clauses in Ancient Greek (and Latin) removes one source of illnestedness in the language
- In fact, there are strikingly few illnested constructions in AG and Latin, given the overall nonprojectivity in the languages

Illnestedness in UD

| | illnested | wellnested |
|----------------------|-----------|------------|
| Ancient_Greek | 0.019 | 0.981 |
| Ancient_Greek-PROIEL | 0.023 | 0.977 |
| Arabic | 0.001 | 0.999 |
| Basque | 0.016 | 0.984 |
| Bulgarian | 0.000 | 1.000 |
| Croatian | 0.007 | 0.993 |
| Czech | 0.001 | 0.999 |
| Danish | 0.009 | 0.991 |
| Dutch | 0.001 | 0.999 |
| English | 0.002 | 0.998 |
| Estonian | 0.000 | 1.000 |
| Finnish | 0.006 | 0.994 |
| Finnish-FTB | 0.006 | 0.994 |
| French | 0.001 | 0.999 |
| German | 0.001 | 0.999 |
| Gothic | 0.019 | 0.981 |
| Greek | 0.005 | 0.995 |
| Hebrew | 0.000 | 1.000 |
| Hindi | 0.002 | 0.998 |

| | illnested | wellnested |
|---------------------|-----------|------------|
| Hungarian | 0.022 | 0.978 |
| Indonesian | 0.001 | 0.999 |
| Irish | 0.000 | 1.000 |
| Italian | 0.004 | 0.996 |
| Japanese-KTC | 0.000 | 1.000 |
| Latin | 0.042 | 0.958 |
| Latin-ITT | 0.003 | 0.997 |
| Latin-PROIEL | 0.015 | 0.985 |
| Norwegian | 0.001 | 0.999 |
| Old_Church_Slavonic | 0.018 | 0.982 |
| Persian | 0.000 | 1.000 |
| Polish | 0.000 | 1.000 |
| Portuguese | 0.001 | 0.999 |
| Romanian | 0.019 | 0.981 |
| Slovenian | 0.003 | 0.997 |
| Spanish | 0.000 | 1.000 |
| Swedish | 0.000 | 1.000 |
| Tamil | 0.000 | 1.000 |

“Real” illnesting (without remnants and punctuation)

| | illnested | wellnested |
|----------------------|--------------|------------|
| Ancient_Greek | 0.014 | 0.986 |
| Ancient_Greek-PROIEL | 0.003 | 0.997 |
| Arabic | 0.000 | 1.000 |
| Basque | 0.001 | 0.999 |
| Bulgarian | 0.000 | 1.000 |
| Croatian | 0.000 | 1.000 |
| Czech | 0.001 | 0.999 |
| Danish | 0.001 | 0.999 |
| Dutch | 0.001 | 0.999 |
| English | 0.000 | 1.000 |
| Estonian | 0.000 | 1.000 |
| Finnish | 0.000 | 1.000 |
| Finnish-FTB | 0.000 | 1.000 |
| French | 0.000 | 1.000 |
| German | 0.000 | 1.000 |
| Gothic | 0.002 | 0.998 |
| Greek | 0.000 | 1.000 |
| Hebrew | 0.000 | 1.000 |
| Hindi | 0.001 | 0.999 |

| | illnested | wellnested |
|---------------------|--------------|------------|
| Hungarian | 0.000 | 1.000 |
| Indonesian | 0.000 | 1.000 |
| Irish | 0.000 | 1.000 |
| Italian | 0.000 | 1.000 |
| Japanese-KTC | 0.000 | 1.000 |
| Latin | 0.037 | 0.963 |
| Latin-ITT | 0.002 | 0.998 |
| Latin-PROIEL | 0.002 | 0.998 |
| Norwegian | 0.001 | 0.999 |
| Old_Church_Slavonic | 0.002 | 0.998 |
| Persian | 0.000 | 1.000 |
| Polish | 0.000 | 1.000 |
| Portuguese | 0.000 | 1.000 |
| Romanian | 0.000 | 1.000 |
| Slovenian | 0.000 | 1.000 |
| Spanish | 0.000 | 1.000 |
| Swedish | 0.000 | 1.000 |
| Tamil | 0.000 | 1.000 |

Conclusions

- Designing treebanks for theoretical linguistics research requires careful attention to annotation schemes
 - The knowns go in the annotation

Conclusions

- Designing treebanks for theoretical linguistics research requires careful attention to annotation schemes
 - The knowns go in the annotation
 - The known unknowns go in post-annotation experimental enrichments
 - Careful with the unknown unknowns!
- Linguistically motivated structures can reflect complexity in a way that corresponds to frequency (e.g. gap degrees and depths)
- Motivated structures can provide evidence for non-standard analyses which avoid exponential blowup (e.g. reentrancies between recursive categories)
- But sometimes there is no connection between treebank data and theoretical structures, so there is more to do

- Haug, Dag Trygve Truslew. 2012. From dependency structures to LFG representations. In Miriam Butt & Tracy Holloway King (eds.), *Proceedings of the LFG12 conference*, 271–291. CSLI Publications.
- Kirk, Allison. 2012. *Word order and information structure in new testament greek*: Universiteit Leiden dissertation.
- Kuhlmann, Marco. 2010. Mildly non-projective dependency grammar. *Computational Linguistics* 39. 355–387.
- Kuhlmann, Marco & Joakim Nivre. 2006. Mildly non-projective dependency structures. In *Proceedings of the coling/acl on main conference poster sessions COLING-ACL '06*, 507–514. Stroudsburg, PA, USA: Association for Computational Linguistics.
<http://dl.acm.org/citation.cfm?id=1273073.1273139>.

- Seki, Hiroyuki, Ryuichi Nakanishi, Yuichi Kaji, Sachiko Ando & Tadao Kasami. 1993. Parallel multiple context-free grammars, finite-state translation systems, and polynomial-time recognizable subclasses of lexical-functional grammars. In *Proceedings of the 31st annual meeting on association for computational linguistics* ACL '93, 130–139. Stroudsburg, PA, USA: Association for Computational Linguistics. doi:10.3115/981574.981592.
<http://dx.doi.org/10.3115/981574.981592>.
- Zaenen, Annie. 2006. Mark-up barking up the wrong tree. *Comput. Linguist.* 32(4). 577–580. doi:10.1162/coli.2006.32.4.577.
<http://dx.doi.org/10.1162/coli.2006.32.4.577>.