

From dependency structures to LFG representations

Dag Haug

Seminar in computational linguistics
 April 18

The texts

- Core – parallel corpus of New Testament translations:
 - Ancient Greek (original, 1st century AD)
 - Gothic (4th century AD)
 - Latin (ca. 400 AD)
 - Classical Armenian (ca. 400 AD)
 - Old Church Slavic (9th century AD)

The texts

- Core – parallel corpus of New Testament translations:
 - Ancient Greek (original, 1st century AD)
 - Gothic (4th century AD)
 - Latin (ca. 400 AD)
 - Classical Armenian (ca. 400 AD)
 - Old Church Slavic (9th century AD)
- Extensions:
 - Herodotus' Histories (Greek 5th century BC)
 - Caesar's Gallic War (Latin, 1st century BC)
 - Cicero's Letters to Atticus (Latin, 1st century BC)
 - Peregrinatio Aetheriae (Vulgar Latin, ca. 400 AD)
 - Hagiographies (The Slavic Codex Suprasliensis, 11th century AD)

The texts

- Core – parallel corpus of New Testament translations:
 - Ancient Greek (original, 1st century AD)
 - Gothic (4th century AD)
 - Latin (ca. 400 AD)
 - Classical Armenian (ca. 400 AD)
 - Old Church Slavic (9th century AD)
- Extensions:
 - Herodotus' Histories (Greek 5th century BC)
 - Caesar's Gallic War (Latin, 1st century BC)
 - Cicero's Letters to Atticus (Latin, 1st century BC)
 - Peregrinatio Aetheriae (Vulgar Latin, ca. 400 AD)
 - Hagiographies (The Slavic Codex Suprasliensis, 11th century AD)
- Ultimate goal: a representative corpus of early IE languages

Small but beautiful

language	tokens
chu	64031
got	56315
grc	137750
lat	120253
xcl	22614
total	400963

On the languages

- Old languages → no native speakers

On the languages

- Old languages → no native speakers
- But fairly well-understood and much-studied texts

On the languages

- Old languages → no native speakers
- But fairly well-understood and much-studied texts
- Morphologically rich

On the languages

- Old languages → no native speakers
- But fairly well-understood and much-studied texts
- Morphologically rich
- Non-configurational, grammatical functions indicated by case rather than word order

On the languages

- Old languages → no native speakers
- But fairly well-understood and much-studied texts
- Morphologically rich
- Non-configurational, grammatical functions indicated by case rather than word order
- All in all quite different from English, which creates lots of problems. . .

Workflow for annotation

- International team of student annotators

Workflow for annotation

- International team of student annotators
- Manual disambiguation of morphology and lemmatization

af-dauifš* (*Germ.* Pl. Pl. zu *af-
 doni, éckvauioç geschunden,
 gepagt; N.M. iñva-ñ 9.39)
 dauþeins *F.ñ* (152) vékpwic *das*
 Alsterben *A.* k4.10; éy θavá-
 vovc in eínim in Todesnöten
 k4.10)
 dauþjan *gr.* V.1 vékpwñ tōten
 C.5.5)
 af-dauþjan *tōten* (perfektiv, 291 ff.)

Workflow for annotation

- International team of student annotators
- Manual disambiguation of morphology and lemmatization
- Syntactic annotation

Workflow for annotation

- International team of student annotators
- Manual disambiguation of morphology and lemmatization
- Syntactic annotation
- Review by project members

Workflow for annotation

- International team of student annotators
- Manual disambiguation of morphology and lemmatization
- Syntactic annotation
- Review by project members
- Advanced annotation done by project members

Morphology

- Verbs inflect for tense, mood, voice, person, number

Morphology

- Verbs inflect for tense, mood, voice, person, number
- Nominals inflect for case, number, gender + possibly grade and definiteness

Morphology

- Verbs inflect for tense, mood, voice, person, number
- Nominals inflect for case, number, gender + possibly grade and definiteness
- All in all this makes for 1817 unique MSD-tags

Morphology

- Verbs inflect for tense, mood, voice, person, number
- Nominals inflect for case, number, gender + possibly grade and definiteness
- All in all this makes for 1817 unique MSD-tags
- In addition there are 25 POS-tags (fairly traditional, with some subdivisions especially in the pronouns)

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer
- Ignores the context and offers spurious ambiguities

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer
- Ignores the context and offers spurious ambiguities
- When we have enough data within a domain, we now use TnT to pretag the text

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer
- Ignores the context and offers spurious ambiguities
- When we have enough data within a domain, we now use TnT to pretag the text
- MDSs are supplemented with lemmatization from the transducer

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer
- Ignores the context and offers spurious ambiguities
- When we have enough data within a domain, we now use TnT to pretag the text
- MDSs are supplemented with lemmatization from the transducer
- Skjærholt (2011, 2012):

Experiment	Token accuracy
Cross-validation on BG	84.3%
Vulgate → BG	62.8%

Morphological annotation

- Started out with manual disambiguation of alternatives from a transducer
- Ignores the context and offers spurious ambiguities
- When we have enough data within a domain, we now use TnT to pretag the text
- MDSs are supplemented with lemmatization from the transducer
- Skjærholt (2011, 2012):

Experiment	Token accuracy
Cross-validation on BG	84.3%
Vulgate → BG	62.8%

- Annotation accuracy goes up and time goes down

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements
- Inherent problems of: (asyndetic) coordination, structure sharing

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements
- Inherent problems of: (asyndetic) coordination, structure sharing
- Dependency grammar with LFG adjustments

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements
- Inherent problems of: (asyndetic) coordination, structure sharing
- Dependency grammar with LFG adjustments
 - Limited set of empty nodes (for asyndetic coordination and ellipsis)

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements
- Inherent problems of: (asyndetic) coordination, structure sharing
- Dependency grammar with LFG adjustments
 - Limited set of empty nodes (for asyndetic coordination and ellipsis)
 - Secondary dependencies (for structure sharing, incl. control/raising)

The syntactic annotation scheme: dependency grammar

- Information about syntactic relations and word order stored separately
- Reliance on overt elements
- Inherent problems of: (asyndetic) coordination, structure sharing
- Dependency grammar with LFG adjustments
 - Limited set of empty nodes (for asyndetic coordination and ellipsis)
 - Secondary dependencies (for structure sharing, incl. control/raising)
 - More granular syntactic relations than usual

Syntactic relations

Label	Function
PRED	Predicate
SUB	Subject
OBJ	Object
OBL	Oblique
AG	Agent
ADV	Adverbial
ATR	Attribute
APOS	Apposition
NARG	Nominal argument

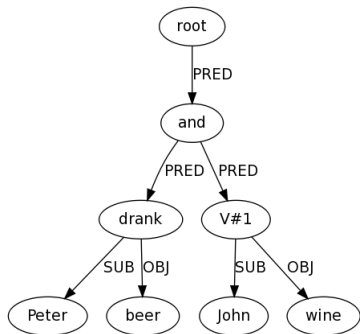
Label	Function
XADV	Free predicative
XOBJ	Open complement
Aux	Auxiliary
XOBJ	Open complement clause
COMP	Complement clause
PART	Partitive
PARPRED	Parenthetical
VOC	Vocative

▶ example

Empty nodes

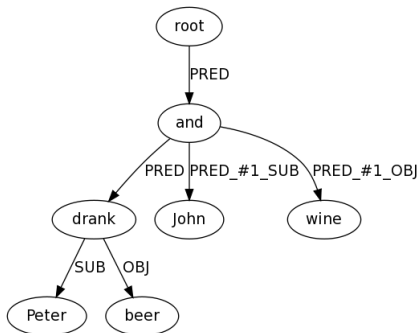
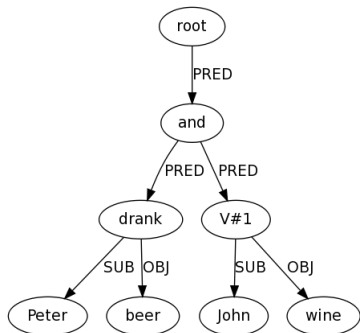
- Null conjunctions for asyndetic parataxis
- Null verbs for null copulas and elided verbs

Eliminability of empty nodes

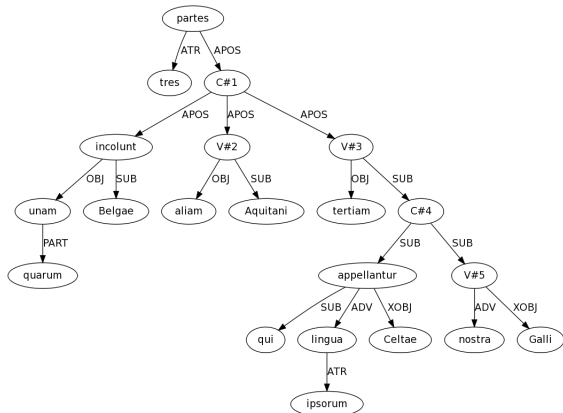


af-daujɨs* (152) Pl.Pf. zu *af-dorai, ékxatúōc geschunden, gefügt; N.M. idem M.9.36; dauheis F.6 (152) vékpwic das Absterben A. k.4.10; éy θavá-voic in einim in Todesnoten k.1.23; dauhjan sic.V.1 vékpoúv tóten C.5.5; af-dauhjan tóten (perfektiv, 291 ff.) k.1.23; D. éyavmnoç

Eliminability of empty nodes

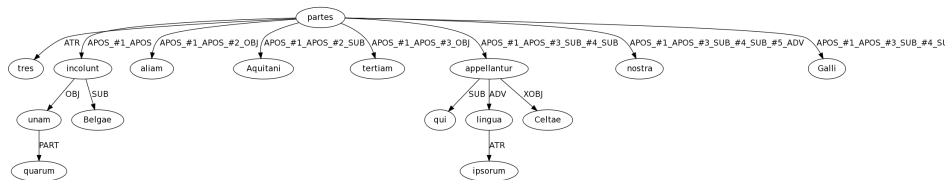


Human processing



of which the Belgians inhabit one, the Aquitani V another, C those who are called Celts in their own language – C Gauls V in our – V the third.

Human processing



of which the Belgians inhabit one, the Aquitani V another, C those who are called Celts in their own language – C Gauls V in our – V the third.

Structure sharing

- Subject control: ▶ Example
- Object control: ▶ Example
- Various other possibilities
- Could also be encoded in the label but typically not with the same precision

Projectivity

language	source	nonprojective	projective
Latin	Gallic War	1887	22717
	Letters to Atticus	2006	20416
	Vulgate	4217	92186
	Per. Aeth.	1279	14890
Greek	Herodotus	6606	56175
	NT	4377	103418
OCS	Zographensis	36	1034
	Suprasliensis	416	7780
	Marianus	1828	47731
Gothic	NT	1886	46884
Armenian	NT	1231	59556
	Koriwn	48	1556

Token alignments

- The translations of the NT have been aligned with the Greek original

Token alignments

- The translations of the NT have been aligned with the Greek original
- A 'dictionary' based on likelihood of occurring in the same bible verse

Token alignments

- The translations of the NT have been aligned with the Greek original
- A 'dictionary' based on likelihood of occurring in the same bible verse
- Information from the annotation: syntax, morphology, word order

Token alignments

- The translations of the NT have been aligned with the Greek original
- A 'dictionary' based on likelihood of occurring in the same bible verse
- Information from the annotation: syntax, morphology, word order
- Manual correction of the Slavic indicates very good results (and a very literal translation)

Precision	Recall	F-score
95.27%	92.97%	94.11%

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD
 - Situation (deixis) → ACC-SIT

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD
 - Situation (deixis) → ACC-SIT
 - Scenarios (inferences) → ACC-INF

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD
 - Situation (deixis) → ACC-SIT
 - Scenarios (inferences) → ACC-INF
 - Encyclopedic knowledge → ACC-GEN

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD
 - Situation (deixis) → ACC-SIT
 - Scenarios (inferences) → ACC-INF
 - Encyclopedic knowledge → ACC-GEN
 - No context (no extra-NP information) → NEW

Givenness

- Givenness tags based on which context the hearer uses to establish reference
 - Discourse (anaphora) → OLD
 - Situation (deixis) → ACC-SIT
 - Scenarios (inferences) → ACC-INF
 - Encyclopedic knowledge → ACC-GEN
 - No context (no extra-NP information) → NEW

▶ example

Modal subordination

Luke 5:39

Und niemand ist, der vom alten trinkt und wolle bald den neuen; denn er spricht: Der alte ist milder.

- The subject and the old and the new wine are embedded under subordination
- Should be inaccessible (Karttunen, COLING 69) but they aren't
- We ignore recursive embeddings and use a special tagset for all embedded referents

Tagset for embedded referents

- NONSPEC (but QUANT for quantification)
- NONSPEC_INF
- NONSPEC_OLD

Tagset for embedded referents

- NONSPEC (but QUANT for quantification)
- NONSPEC_INF
- NONSPEC_OLD

No counterparts to ACC-GEN or ACC-SIT as these belong in the main DRS by definition

Interannotator agreement

- Towards the end of the NT tagging projects, kappa values were around 0.8 (after long periods of weekly meetings)

Interannotator agreement

- Towards the end of the NT tagging projects, kappa values were around 0.8 (after long periods of weekly meetings)
- New project: Caesar's Gallic War

Interannotator agreement

- Towards the end of the NT tagging projects, kappa values were around 0.8 (after long periods of weekly meetings)
- New project: Caesar's Gallic War
- Supervised tagging of 8 chapters (ca. 400 taggables)

Interannotator agreement

- Towards the end of the NT tagging projects, kappa values were around 0.8 (after long periods of weekly meetings)
- New project: Caesar's Gallic War
- Supervised tagging of 8 chapters (ca. 400 taggables)
- Unsupervised tagging of 5 chapters (ca. 250 taggables)
 - $\kappa = 0.66$ counting divergences in taggables
 - $\kappa = 0.75$ on tags set by both annotators

Interannotator agreement

- Towards the end of the NT tagging projects, kappa values were around 0.8 (after long periods of weekly meetings)
- New project: Caesar's Gallic War
- Supervised tagging of 8 chapters (ca. 400 taggables)
- Unsupervised tagging of 5 chapters (ca. 250 taggables)
 - $\kappa = 0.66$ counting divergences in taggables
 - $\kappa = 0.75$ on tags set by both annotators
- Decent; but much potential for more agreement, especially in taggables

Size of IS corpus

Tag	Freq
old	34430
old_inact	1395
acc_gen	3755
acc_inf	2634
acc_sit	883
new	5768
kind	1178
non_spec	4485
non_spec_inf	408
non_spec_old	1799
quant	2021
total	58756

edge type	freq
coreference	36650
bridging	2847
total	39497

Storing linguistic analyses

- Theory-neutrality →
 - data for larger audiences

Storing linguistic analyses

- Theory-neutrality →
 - data for larger audiences
 - widening gulf between corpus linguistics and linguistic theory

Storing linguistic analyses

- Theory-neutrality →
 - data for larger audiences
 - widening gulf between corpus linguistics and linguistic theory
- DG corpora (Prague, PROIEL) → DG not really in use as a linguistic theory



Storing linguistic analyses

- Theory-neutrality →
 - data for larger audiences
 - widening gulf between corpus linguistics and linguistic theory
- DG corpora (Prague, PROIEL) → DG not really in use as a linguistic theory
- PS corpora (Penn, NEGRA) typically use flatter tree structures than anyone believes in

Storing linguistic analyses

- Theory-neutrality →
 - data for larger audiences
 - widening gulf between corpus linguistics and linguistic theory
- DG corpora (Prague, PROIEL) → DG not really in use as a linguistic theory
- PS corpora (Penn, NEGRA) typically use flatter tree structures than anyone believes in
- On the other hand, LFG and HPSG corpora can be hard to use for people who do not share the theoretical assumptions of these theories

Our take

Principles

- 1 Encode no more structure than is common to all frameworks

Our take

Principles

- 1 Encode no more structure than is common to all frameworks
- 2 Enoded structure could be seen as derived/secondary in some frameworks

Our take

Principles

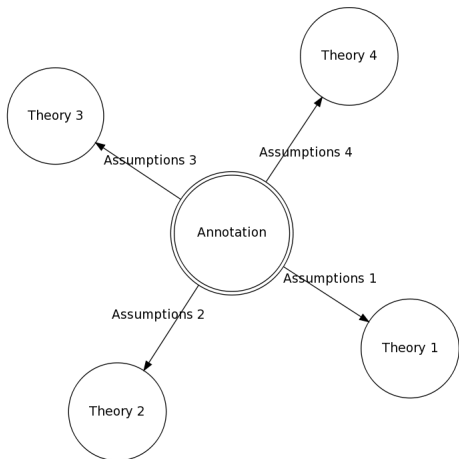
- 1 Encode no more structure than is common to all frameworks
- 2 Enoded structure could be seen as derived/secondary in some frameworks
- 3 Encode enough structure to allow reconstruction of theoretically motivated structures

Our take

Principles

- 1 Encode no more structure than is common to all frameworks
 - 2 Enoded structure could be seen as derived/secondary in some frameworks
 - 3 Encode enough structure to allow reconstruction of theoretically motivated structures
- In the ideal situation, the information in the annotation can be (monotonically) expanded to structures conforming to a particular theory by adding information from the assumptions of that theory

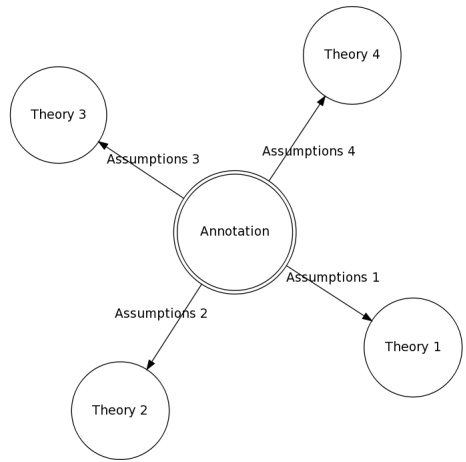
The ideal situation



- The added assumptions will typically be about phrase structure, such as various versions of X' theory

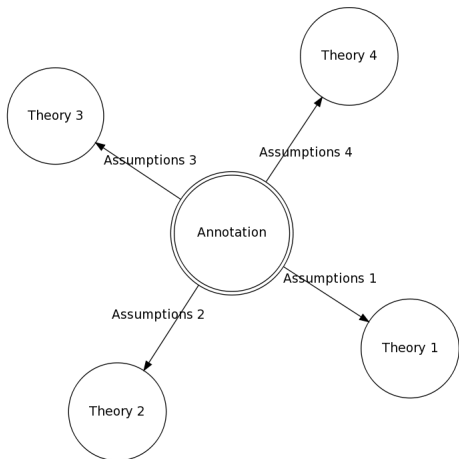
af-daupjans* (Gen.) Pl. Pl. zu *af-
dojan, ékkaútioc geschunden,
geplagt; N. Pl. ána M 9,31
dauns Pl éan (Genos) ánauns
wópi: énaús (AU) ánaús k 2,15
E 5,2; N. K 12,17 k 2,1; 6; J.
k 2,14; G. J 12,1; D. éavomroc
danbeins Pl. (15?) vékpwic das
Absterben A. k 4,10; éy θανά-
vovc in einim in Todesnöten
k 1,23
avupjan sic. V. I vékpoúv töten
G 5,5
af-daupjan töten (perfektiv, 291 ff.)

The ideal situation



- The added assumptions will typically be about phrase structure, such as various versions of X' theory
- Given information about what the subject is, it will be possible to create a structure where the subject has a specific position if the theory requires that (unless the data contradict the theory)

The ideal situation



- The added assumptions will typically be about phrase structure, such as various versions of X' theory
- Given information about what the subject is, it will be possible to create a structure where the subject has a specific position if the theory requires that (unless the data contradict the theory)
- Useful for hypothesis testing

Basic principles

- Modular: several levels of grammatical description connected by projections (functions)



Basic principles

- Modular: several levels of grammatical description connected by projections (functions)



- The c-structure is a tree structure described by a CFG

Basic principles

- Modular: several levels of grammatical description connected by projections (functions)



- The c-structure is a tree structure described by a CFG
- The f-structure is a set of ordered attribute-value pairs

Basic principles

- Modular: several levels of grammatical description connected by projections (functions)



- The c-structure is a tree structure described by a CFG
- The f-structure is a set of ordered attribute-value pairs
- the attribute is a grammatical function or feature and the value is
 - a symbol
 - a semantic form
 - an f-structure
 - a set of f-structures (for adjuncts)

Basic principles

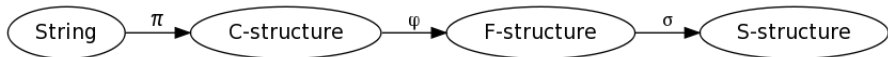
- Modular: several levels of grammatical description connected by projections (functions)



- The c-structure is a tree structure described by a CFG
- The f-structure is a set of ordered attribute-value pairs
- the attribute is a grammatical function or feature and the value is
 - a symbol
 - a semantic form
 - an f-structure
 - a set of f-structures (for adjuncts)
- Lexical items and CFG rules can contribute f-descriptions

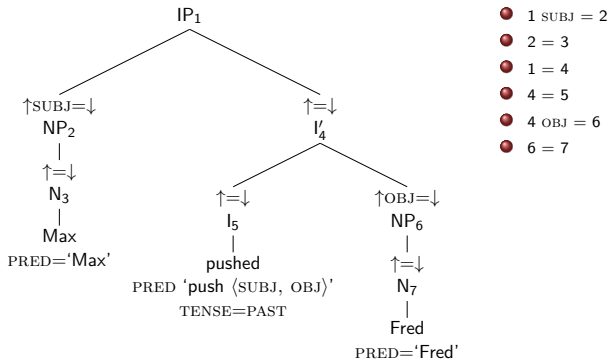
Basic principles

- Modular: several levels of grammatical description connected by projections (functions)

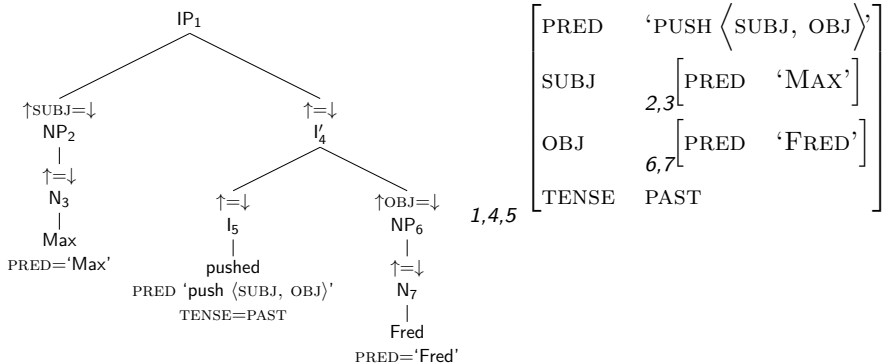


- The c-structure is a tree structure described by a CFG
- The f-structure is a set of ordered attribute-value pairs
- the attribute is a grammatical function or feature and the value is
 - a symbol
 - a semantic form
 - an f-structure
 - a set of f-structures (for adjuncts)
- Lexical items and CFG rules can contribute f-descriptions
- Lexical-functional languages \in context-sensitive languages

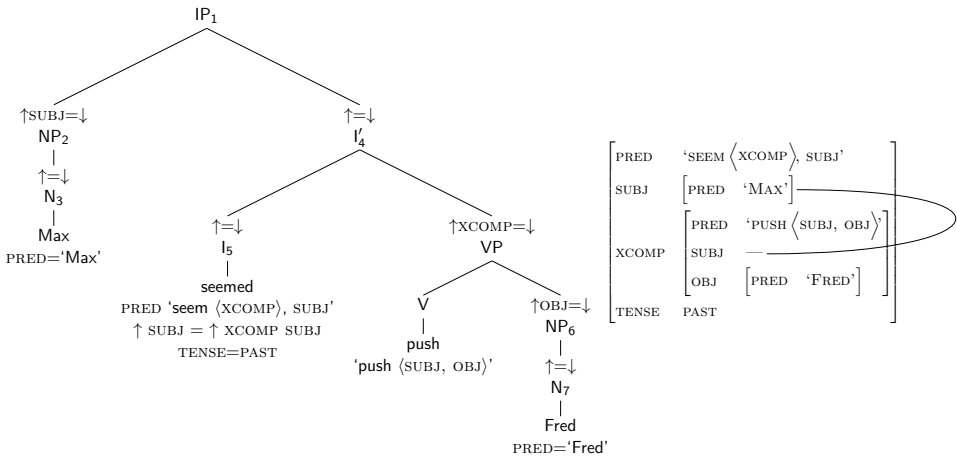
Configurational encoding



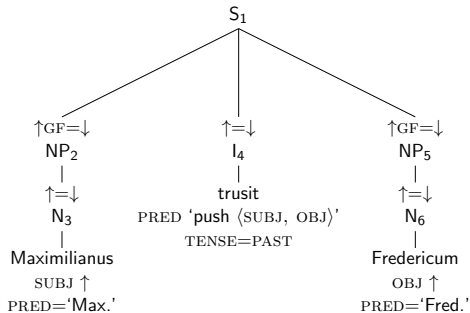
Configurational encoding



Structure sharing



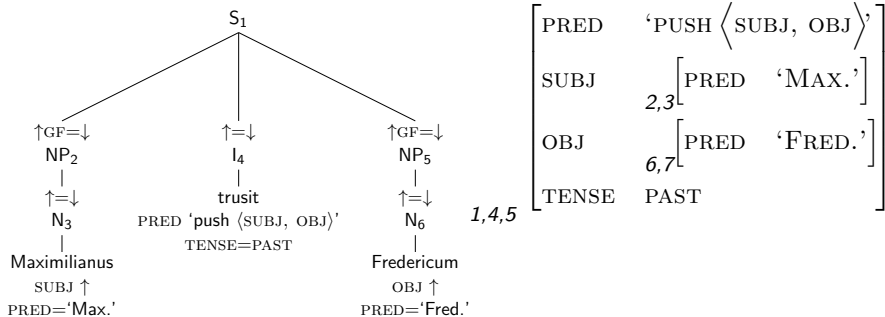
Non-configurational encoding



- 1 GF = 2
- 2 = 3
- $\exists f.f$ SUBJ = 3
- 1 = 4
- 4 GF = 5
- 5 = 6
- $\exists f.f$ OBJ = 6

af-dauþs* (291), Pl. PP. zu *af-doum, éckuauþs geschunden, gepugt; N.M. 291. M. 9.30. dauþs Pl. 6m. Germ. 291. dauþs wópi: eð. 291. dauþs k 2.15. E 5.2; N. R. 12.17 k 2.1. 16; J. k 2.14; G. 1. 13. 14. D. 2. 13. 14. 15. 16. 17. 18. 19. 20. 21. 22. 23. 24. 25. 26. 27. 28. 29. 30. 31. 32. 33. 34. 35. 36. 37. 38. 39. 40. 41. 42. 43. 44. 45. 46. 47. 48. 49. 50. 51. 52. 53. 54. 55. 56. 57. 58. 59. 60. 61. 62. 63. 64. 65. 66. 67. 68. 69. 70. 71. 72. 73. 74. 75. 76. 77. 78. 79. 80. 81. 82. 83. 84. 85. 86. 87. 88. 89. 90. 91. 92. 93. 94. 95. 96. 97. 98. 99. 100. 101. 102. 103. 104. 105. 106. 107. 108. 109. 110. 111. 112. 113. 114. 115. 116. 117. 118. 119. 120. 121. 122. 123. 124. 125. 126. 127. 128. 129. 130. 131. 132. 133. 134. 135. 136. 137. 138. 139. 140. 141. 142. 143. 144. 145. 146. 147. 148. 149. 150. 151. 152. 153. 154. 155. 156. 157. 158. 159. 160. 161. 162. 163. 164. 165. 166. 167. 168. 169. 170. 171. 172. 173. 174. 175. 176. 177. 178. 179. 180. 181. 182. 183. 184. 185. 186. 187. 188. 189. 190. 191. 192. 193. 194. 195. 196. 197. 198. 199. 200. 201. 202. 203. 204. 205. 206. 207. 208. 209. 210. 211. 212. 213. 214. 215. 216. 217. 218. 219. 220. 221. 222. 223. 224. 225. 226. 227. 228. 229. 230. 231. 232. 233. 234. 235. 236. 237. 238. 239. 240. 241. 242. 243. 244. 245. 246. 247. 248. 249. 250. 251. 252. 253. 254. 255. 256. 257. 258. 259. 260. 261. 262. 263. 264. 265. 266. 267. 268. 269. 270. 271. 272. 273. 274. 275. 276. 277. 278. 279. 280. 281. 282. 283. 284. 285. 286. 287. 288. 289. 290. 291. 292. 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 303. 304. 305. 306. 307. 308. 309. 310. 311. 312. 313. 314. 315. 316. 317. 318. 319. 320. 321. 322. 323. 324. 325. 326. 327. 328. 329. 330. 331. 332. 333. 334. 335. 336. 337. 338. 339. 340. 341. 342. 343. 344. 345. 346. 347. 348. 349. 350. 351. 352. 353. 354. 355. 356. 357. 358. 359. 360. 361. 362. 363. 364. 365. 366. 367. 368. 369. 370. 371. 372. 373. 374. 375. 376. 377. 378. 379. 380. 381. 382. 383. 384. 385. 386. 387. 388. 389. 390. 391. 392. 393. 394. 395. 396. 397. 398. 399. 400. 401. 402. 403. 404. 405. 406. 407. 408. 409. 410. 411. 412. 413. 414. 415. 416. 417. 418. 419. 420. 421. 422. 423. 424. 425. 426. 427. 428. 429. 430. 431. 432. 433. 434. 435. 436. 437. 438. 439. 440. 441. 442. 443. 444. 445. 446. 447. 448. 449. 450. 451. 452. 453. 454. 455. 456. 457. 458. 459. 460. 461. 462. 463. 464. 465. 466. 467. 468. 469. 470. 471. 472. 473. 474. 475. 476. 477. 478. 479. 480. 481. 482. 483. 484. 485. 486. 487. 488. 489. 490. 491. 492. 493. 494. 495. 496. 497. 498. 499. 500. 501. 502. 503. 504. 505. 506. 507. 508. 509. 510. 511. 512. 513. 514. 515. 516. 517. 518. 519. 520. 521. 522. 523. 524. 525. 526. 527. 528. 529. 530. 531. 532. 533. 534. 535. 536. 537. 538. 539. 540. 541. 542. 543. 544. 545. 546. 547. 548. 549. 550. 551. 552. 553. 554. 555. 556. 557. 558. 559. 560. 561. 562. 563. 564. 565. 566. 567. 568. 569. 570. 571. 572. 573. 574. 575. 576. 577. 578. 579. 580. 581. 582. 583. 584. 585. 586. 587. 588. 589. 590. 591. 592. 593. 594. 595. 596. 597. 598. 599. 600. 601. 602. 603. 604. 605. 606. 607. 608. 609. 610. 611. 612. 613. 614. 615. 616. 617. 618. 619. 620. 621. 622. 623. 624. 625. 626. 627. 628. 629. 630. 631. 632. 633. 634. 635. 636. 637. 638. 639. 640. 641. 642. 643. 644. 645. 646. 647. 648. 649. 650. 651. 652. 653. 654. 655. 656. 657. 658. 659. 660. 661. 662. 663. 664. 665. 666. 667. 668. 669. 670. 671. 672. 673. 674. 675. 676. 677. 678. 679. 680. 681. 682. 683. 684. 685. 686. 687. 688. 689. 690. 691. 692. 693. 694. 695. 696. 697. 698. 699. 700. 701. 702. 703. 704. 705. 706. 707. 708. 709. 710. 711. 712. 713. 714. 715. 716. 717. 718. 719. 720. 721. 722. 723. 724. 725. 726. 727. 728. 729. 730. 731. 732. 733. 734. 735. 736. 737. 738. 739. 740. 741. 742. 743. 744. 745. 746. 747. 748. 749. 750. 751. 752. 753. 754. 755. 756. 757. 758. 759. 760. 761. 762. 763. 764. 765. 766. 767. 768. 769. 770. 771. 772. 773. 774. 775. 776. 777. 778. 779. 780. 781. 782. 783. 784. 785. 786. 787. 788. 789. 790. 791. 792. 793. 794. 795. 796. 797. 798. 799. 800. 801. 802. 803. 804. 805. 806. 807. 808. 809. 810. 811. 812. 813. 814. 815. 816. 817. 818. 819. 820. 821. 822. 823. 824. 825. 826. 827. 828. 829. 830. 831. 832. 833. 834. 835. 836. 837. 838. 839. 840. 841. 842. 843. 844. 845. 846. 847. 848. 849. 850. 851. 852. 853. 854. 855. 856. 857. 858. 859. 860. 861. 862. 863. 864. 865. 866. 867. 868. 869. 870. 871. 872. 873. 874. 875. 876. 877. 878. 879. 880. 881. 882. 883. 884. 885. 886. 887. 888. 889. 890. 891. 892. 893. 894. 895. 896. 897. 898. 899. 900. 901. 902. 903. 904. 905. 906. 907. 908. 909. 910. 911. 912. 913. 914. 915. 916. 917. 918. 919. 920. 921. 922. 923. 924. 925. 926. 927. 928. 929. 930. 931. 932. 933. 934. 935. 936. 937. 938. 939. 940. 941. 942. 943. 944. 945. 946. 947. 948. 949. 950. 951. 952. 953. 954. 955. 956. 957. 958. 959. 960. 961. 962. 963. 964. 965. 966. 967. 968. 969. 970. 971. 972. 973. 974. 975. 976. 977. 978. 979. 980. 981. 982. 983. 984. 985. 986. 987. 988. 989. 990. 991. 992. 993. 994. 995. 996. 997. 998. 999. 1000.

Non-configurational encoding

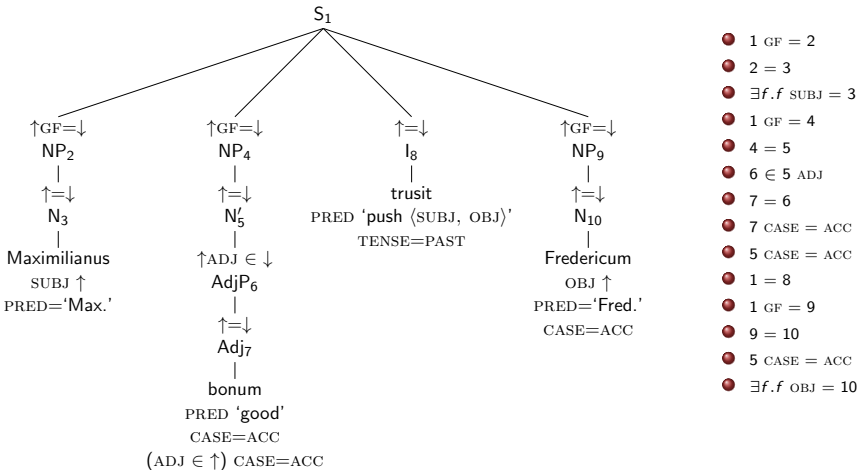


Non-projectivity

A mock Latin example

Maximilianus	bonum	trusit	Fredericum
Maximilian.NOM	good.ACC	pushed	Frederick.ACC

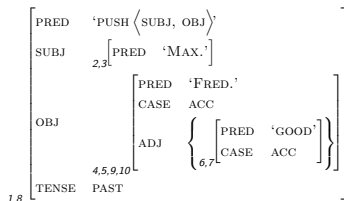
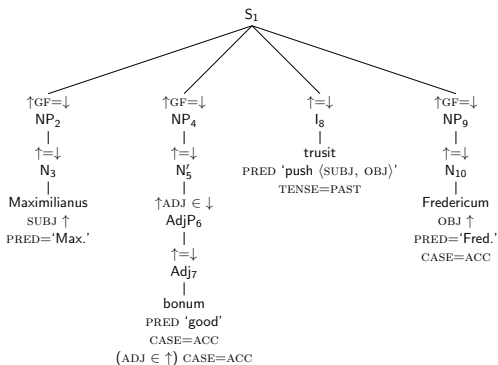
Non-projectivity



- 1 GF = 2
- 2 = 3
- ∃f.f SUBJ = 3
- 1 GF = 4
- 4 = 5
- 6 ∈ 5 ADJ
- 7 = 6
- 7 CASE = ACC
- 5 CASE = ACC
- 1 = 8
- 1 GF = 9
- 9 = 10
- 5 CASE = ACC
- ∃f.f OBJ = 10

af-dauipš* (Ger.) Pl. PP zu *af-doum, éckauuOC geschunden, gepugt; N.M. idem M 9361
 danheins F.f. (152) vékpwic das Absterben A. k4.10; év θανά-
 dauns Ff éon, Gémo, AFR dains wópi: éon, AFR, wépoúv k 2,15 E 5,2; N. K 1217 k 2,1; 16; J. k 2,14; G. J 13; D. éwonnroc
 af-dauipjan k 1, 23
 af-dauipjan sic.V.1 vékpoúv tóten C 5,5
 af-dauipjan tóten (perfektiv, 291 ff.)

Non-projectivity





Relationship to DG

- F-structures and DGs both encode labelled syntactic dependencies

Relationship to DG

- F-structures and DGs both encode labelled syntactic dependencies
- Two major differences
 - LFG's structure sharing runs against DG's unique head principle
 - In DG, every word introduces depth in the graph, whereas multiple words can contribute to the same F-structure (without nesting)

Relationship to DG

- F-structures and DGs both encode labelled syntactic dependencies
- Two major differences
 - LFG's structure sharing runs against DG's unique head principle
 - In DG, every word introduces depth in the graph, whereas multiple words can contribute to the same F-structure (without nesting)
- We have already given up the unique head principle in our DG

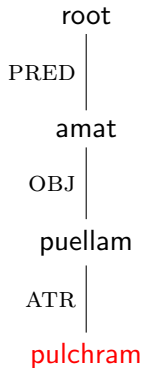
Relationship to DG

- F-structures and DGs both encode labelled syntactic dependencies
- Two major differences
 - LFG's structure sharing runs against DG's unique head principle
 - In DG, every word introduces depth in the graph, whereas multiple words can contribute to the same F-structure (without nesting)
- We have already given up the unique head principle in our DG
- The words that do not introduce separate layers of f-structures are typically function words, so they can be identified from the labels

Label mapping

Function	Label	LFG	Function	Label	LFG
Adverbial	ADV	ADJ	Oblique	OBL	OBJ _θ /OBL
Agent	AG	OBL _{AG}	Parenthetical	PARPRED	—
Apposition	APOS	ADJ	Partitive	PART	ADJ
Attribute	ATR	ADJ	Predicate	PRED	—
Auxiliary	AUX	—	Subject	SUB	SUBJ
Complement	COMP	COMP	Vocative	VOC	—
Argument of noun	NARG	≈ OBL	Free predicative	XADV	XADJ
Object	OBJ	OBJ	Open complement	XOBJ	XCOMP

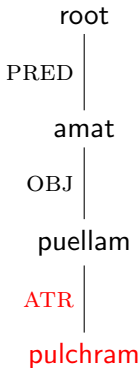
A simple example



- Each node maps to an attribute-value matrix with morphological features and a semantic form

PRED	'PULCHER'
CASE	ACC
GEND	FEM

A simple example

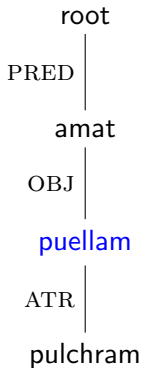


- The relations are translated to attributes with the dependents' AVM as value

$$\left[\text{ADJ} \left\{ \begin{array}{l} \left[\begin{array}{ll} \text{PRED} & \text{'PULCHER'} \end{array} \right] \\ \left[\begin{array}{ll} \text{CASE} & \text{ACC} \end{array} \right] \\ \left[\begin{array}{ll} \text{GEND} & \text{FEM} \end{array} \right] \end{array} \right\} \right]$$

af-dauīþs* (Gen. 1) Pl. PP. zu *af-
 doan, éckuauioç geschunden,
 gepagt; N.M. idm. M. 9.30.
 dauþeins F. Sg. (152) vékrwic *das*
 Absterben A. k.4.10; éy θανά-
 voc in einim in Todesnöten
 k.1.25.
 dauþjan sic.V.1 vékrwōn *töten*
 G. S. S.
 af-dauþjan *töten* (perfektiv, 291 ff.)
 k.2.14; G. S. S. D. éyovmtoç

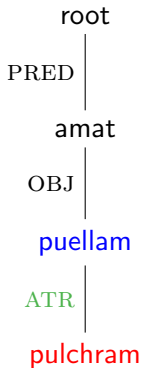
A simple example



- We do this for all nodes in the structure

PRED	'PUELLA'
CASE	ACC
GEND	FEM

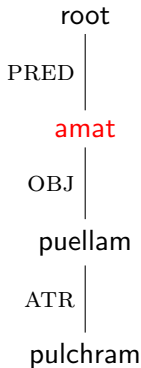
A simple example



- The AVMs of the head and the relation+dependent are unified

PRED	'PUELLA'						
CASE	ACC						
GEND	FEM						
ADJ	<table border="1"> <tr> <td>PRED</td> <td>'PULCHER'</td> </tr> <tr> <td>CASE</td> <td>ACC</td> </tr> <tr> <td>GEND</td> <td>FEM</td> </tr> </table>	PRED	'PULCHER'	CASE	ACC	GEND	FEM
PRED	'PULCHER'						
CASE	ACC						
GEND	FEM						

A simple example



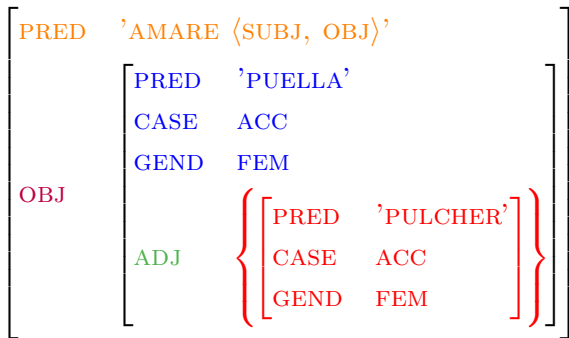
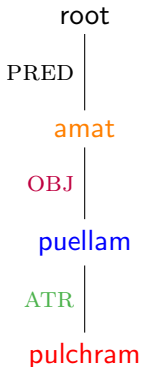
- The process terminates with the main verb
- NB PRED \neq PRED

[PRED 'AMARE < SUBJ, OBJ >']

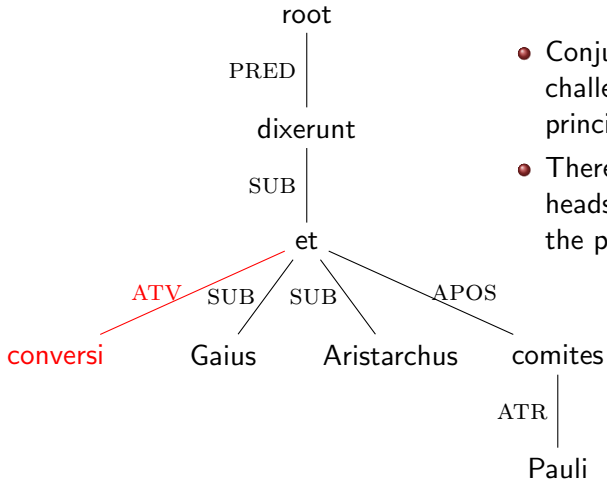
af-dauþs* (Gen.) Pl. PP. zu *af-
 doan, éckuauþs geschunden,
 geflagt; N.M. idem M 936.
 dauus Pl. Gen. Germ. ATR dauus
 wópi: ed. 15. 16. 17. k 2. 15.
 E 5. 2; N. R 12. 17 k 2. 1. 16; J.
 k 2. 14; G. J. 13. 14. D. E. vomroc
 dauþeins Pl. (15?) vékruic das
 Absterben A. k 4. 10; év θανά-
 τωic in einim in Todesnöten
 k 1. 23.
 dauþjan sic. V. 1 vékruþv tóten
 C 5. 5.
 af-dauþjan tóten (perfektiv, 291 ff.)

A simple example

- The final result

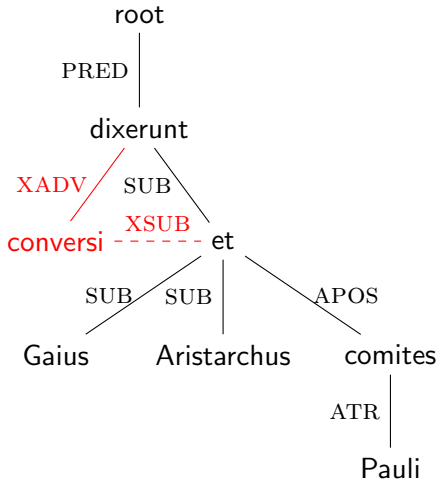


Structure sharing 1



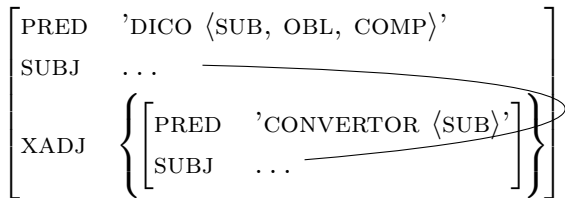
- Conjunct participles challenge the unique head principle
- There are two candidate heads: the main verb and the participle subject

Structure sharing 2

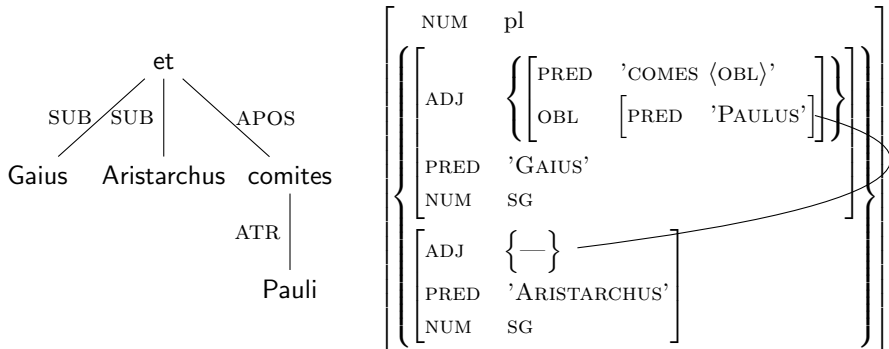


- With secondary edges we can represent both dependencies

F-structure representation



Features in coordination



- The adjunct is a distributive feature
- Non-distributive features are computed from the set members
- Number, gender and person are such features

Preliminaries

- C-structures and DGs contain very different information

Preliminaries

- C-structures and DGs contain very different information
- Instead of syntactic dependencies, c-structures contain information about
 - category
 - word order
 - word groupings (constituents)

Preliminaries

- C-structures and DGs contain very different information
- Instead of syntactic dependencies, c-structures contain information about
 - category
 - word order
 - word groupings (constituents)
- Of these, only word order is present in a DG (assuming there is a precedence order on terminals)

Preliminaries

- C-structures and DGs contain very different information
- Instead of syntactic dependencies, c-structures contain information about
 - category
 - word order
 - word groupings (constituents)
- Of these, only word order is present in a DG (assuming there is a precedence order on terminals)
- We will see how we can enrich DGs with ‘projections’ that include the other information

Preliminaries

- C-structures and DGs contain very different information
- Instead of syntactic dependencies, c-structures contain information about
 - category
 - word order
 - word groupings (constituents)
- Of these, only word order is present in a DG (assuming there is a precedence order on terminals)
- We will see how we can enrich DGs with ‘projections’ that include the other information
- The makeup of constituents is a matter of theoretical debate, so we need to introduce theoretical assumptions from LFG

Basic DG

What's in a DG?

A DG is a tuple $\langle \mathcal{W}, r, R_{\mathcal{D}} \rangle$ where

- \mathcal{W} is the set of words totally ordered by \prec
- $R_{\mathcal{D}}$ is a set of dependency relations that forms a tree over \mathcal{W} rooted in $r (\in \mathcal{W})$

DG with categories

- The basic point is to note that category constraints are in principle independent of other constraints

DG with categories

- The basic point is to note that category constraints are in principle independent of other constraints
- The classic case is the German Mittelfeld (Bröker 1998)

DG with categories

- The basic point is to note that category constraints are in principle independent of other constraints
- The classic case is the German Mittelfeld (Bröker 1998)
- We can simply extend our model with a class of categories \mathcal{C} and a function $V_{\mathcal{C}} : \mathcal{W} \mapsto \mathcal{C}$

DG with categories

- The basic point is to note that category constraints are in principle independent of other constraints
- The classic case is the German Mittelfeld (Bröker 1998)
- We can simply extend our model with a class of categories \mathcal{C} and a function $V_{\mathcal{C}} : \mathcal{W} \mapsto \mathcal{C}$
- In practice we will use the morphological annotations on the words and map them to a set of theoretically motivated categories

DG with categories

- The basic point is to note that category constraints are in principle independent of other constraints
- The classic case is the German Mittelfeld (Bröker 1998)
- We can simply extend our model with a class of categories \mathcal{C} and a function $V_{\mathcal{C}} : \mathcal{W} \mapsto \mathcal{C}$
- In practice we will use the morphological annotations on the words and map them to a set of theoretically motivated categories
- Notice that if we conceive of $V_{\mathcal{C}}$ as a projection, it is different from LFG projections since it embodies linguistic knowledge (the ϕ function is not similarly restricted)

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

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

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

Order domain structure

The set of order domains of all words $w \in \mathcal{W}$ is a semi-lattice ordered by set inclusion. The join/meet of the semi-lattice is \mathcal{W} .

- Every order domain is included in exactly one other order domain, and the order domains are ordered by precedence so the order domain structure is in effect an ordered tree

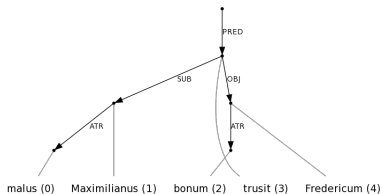
Order domain structures

Order domain structure

The set of order domains of all words $w \in \mathcal{W}$ is a semi-lattice ordered by set inclusion. The join/meet of the semi-lattice is \mathcal{W} .

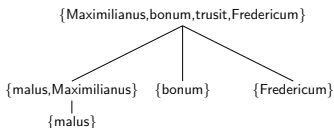
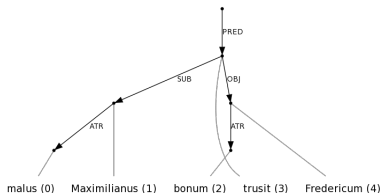
- Every order domain is included in exactly one other order domain, and the order domains are ordered by precedence so the order domain structure is in effect an ordered tree
- Similar to those generated by CFGs but without the categorial information

Example

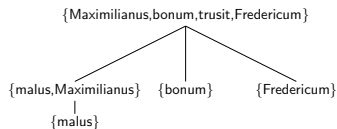
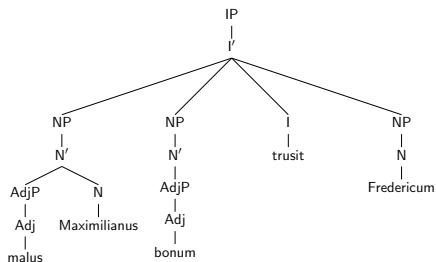


af-dauþs* (291), Pl. Pr. zu *af-doum, éckuauþs geschunden, gepugt; N.M. idem M. 936.
 dauþs Pl. Pr. zu *af-doum, éckuauþs geschunden, gepugt; N.M. idem M. 936.
 dauþs Pl. Pr. zu *af-doum, éckuauþs geschunden, gepugt; N.M. idem M. 936.
 dauþs Pl. Pr. zu *af-doum, éckuauþs geschunden, gepugt; N.M. idem M. 936.

Example

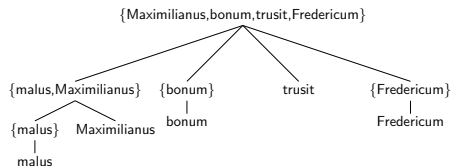
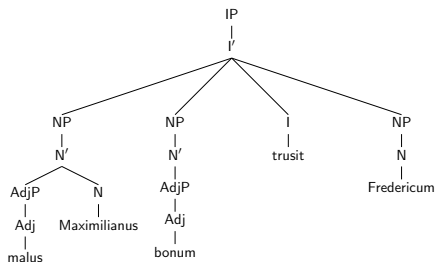


Example



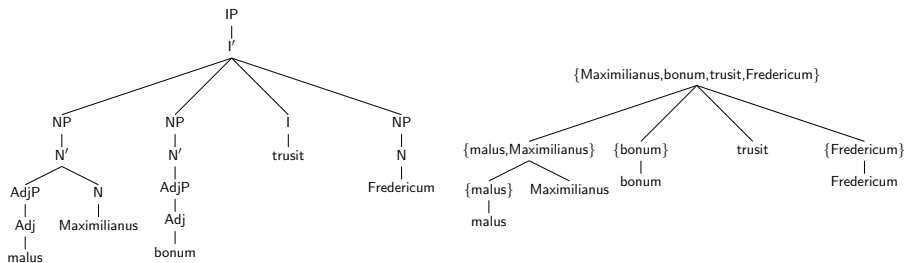
- Each Bröker node corresponds to a $X'' - X' - X$ spine

Example



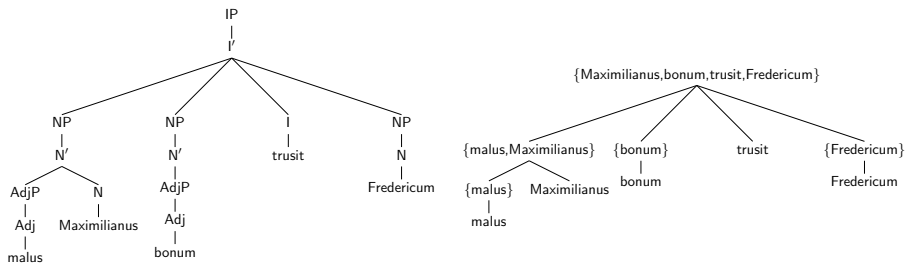
- Each Bröker node corresponds to a $X'' - X' - X$ spine
- We can add explicit heads (each w is the head of \mathcal{D}_w)

Example



- Each Bröker node corresponds to a $X'' - X' - X$ spine
- We can add explicit heads (each w is the head of \mathcal{D}_w)
- Probably as close as we can come to a pure projection from the DG

Example



- Each Bröker node corresponds to a $X'' - X' - X$ spine
- We can add explicit heads (each w is the head of \mathcal{D}_w)
- Probably as close as we can come in a pure projection from the DG
- What we are lacking is a theory of the internal structure of phrases

Internal structure of phrases

Questions (from Xia 2001)

- ① for a category X, what kind of projections can X have?
- ② if a category Y depends on a category Y in a dependency structure, how far should Y project before it attaches to Xs projection?
- ③ if a category Y depends on a category X in a dependency structure, to what position on X's projection chain should Y's projection attach?

Internal structure of phrases

Answers

- 1 all categories X project two levels X' and XP.
- 2 a dependent Y always projects to Y' then YP and the YP attaches to the head's projection
- 3 dependents are divided into three types using a set of handwritten rules: specifiers, modifiers and arguments. Specifiers are made sisters of X' and arguments are made daughters of X. Modifiers Chomsky-adjoin to either X' or XP depending on whether they are restrictive, as indicated by the dependency edge label (ATR or APOS).

An algorithm

 $\mathcal{L} = \{\}$
function CREATEPROJECTION(n)

 $\mathcal{D} = \{\}$
for all d : daughters of n **do**

 put CREATEPROJECTION(d) in \mathcal{D}
end for
for all $d \in \mathcal{D} \cup \mathcal{L}$ **do**

 if d is in n 's order domain **then**

 put/leave d' in \mathcal{D}

else

 put/leave d in \mathcal{L}

end if

end for

 make the elements in \mathcal{D} daughters of n
end function

af-dauþs* (Gothic) Pl. Pr. zu *af-
doan, êckauþog geschunden,
geþagt; N.M. idan 3. 11.
dauns Pl. Pr. Germ. 1.1.
wôpi: eð 1. 1. 1. 1. 1. 1.
E 5.2; N. K 12.17 k 2.1; 16; J.
k 2.14; G. 1. 1. 1. 1. 1. 1.

dauþeins Pl. Pr. 152 vekrowic *das*
Absterben A. k 4.10; êy þavâ-
yoic in einum in *Todesnôten*
k 1. 1. 1.
daupjan 3rd Pl. Pr. vekroþv *töten*
G 5.5
af-daupjan *töten (perfektiv, 291 ff.)*

Adding linguistic knowledge

- This algorithm gives us the Bröker trees

Adding linguistic knowledge

- This algorithm gives us the Bröker trees
- We can enrich these with linguistic knowledge

Adding linguistic knowledge

- This algorithm gives us the Bröker trees
- We can enrich these with linguistic knowledge
- We will use our X' assumptions, the category mapping and handwritten phrase structure rules

A sample rule

N:	
:phrase_adjuncts:	- NP - AdjP
:specifier:	- DP
:bar_adjuncts:	- AdjP - NP
:complements:	- NP

Adding linguistic knowledge

- This algorithm gives us the Bröker trees
- We can enrich these with linguistic knowledge
- We will use our X' assumptions, the category mapping and handwritten phrase structure rules
- We can recursively embed loose nodes under headless structures to achieve the LFG analysis of non-projectivity

Where to add linguistics

 $\mathcal{L} = \{ \}$
function CREATEPROJECTION(n)

 $\mathcal{D} = \{ \}$
for all d : daughters of n **do**

 put CREATEPROJECTION(d) in \mathcal{D}
end for
for all $d \in \mathcal{D} \cup \mathcal{L}$ **do**

 if d is in n 's order domain **then**

 put/leave d' in \mathcal{D}

else

 put/leave d in \mathcal{L}

end if

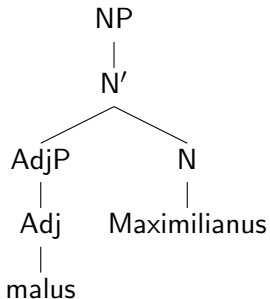
end for

 make the elements in \mathcal{D} daughters of n
end function

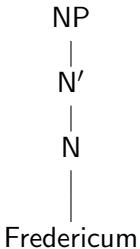
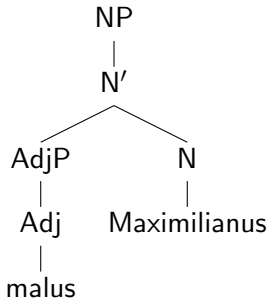
af-dauþs* (291), Pl. PP. zu *af-
 dauþu, éckuauþs geschunden,
 gepagt; N.M. 291 M 936
 dauþeins F.6 (152) vékrwic das
 Absterben A. k 4.10; év þau-
 dauþuic in einim in Todesnöten
 k 1.23
 dauþjan sic.V.1 vékröuþ töten
 G.5.5
 af-dauþjan töten (perfektiv, 291 ff.)
 k 2.14; G. 1.13; D. þauþmöt

D

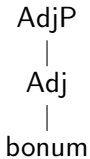
L



\mathcal{D}

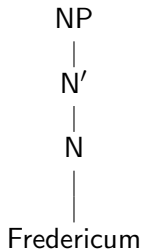
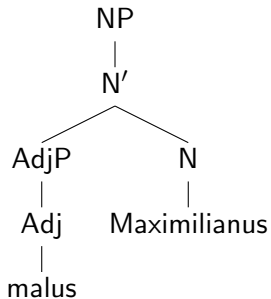


\mathcal{L}



af-dauþs* (Gen.) Pl. M. zu *af-
 dau, éckuauþs geschunden,
 gepagt; N.M. 1. Pl. M. 1. Pl. M.
 dauþs Pl. Gen. Germ. Pl. M. 1. Pl. M.
 wóþi: eð. Pl. M. 1. Pl. M. k 2,15
 E 5,2; N. K 12,17 k 2,17; J.
 k 2,14; G. J 1,1; D. Pl. M. 1. Pl. M.
 dauþs Pl. (152) vékrwic das
 Absterben A. k 4,10; éy þauá-
 þaic in einim in Todesnöten
 k 1,23
 dauþjan sic. V. 1 vékröuþ tóten
 G 5,5
 af-dauþjan tóten (perfektiv, 291 ff.)

D

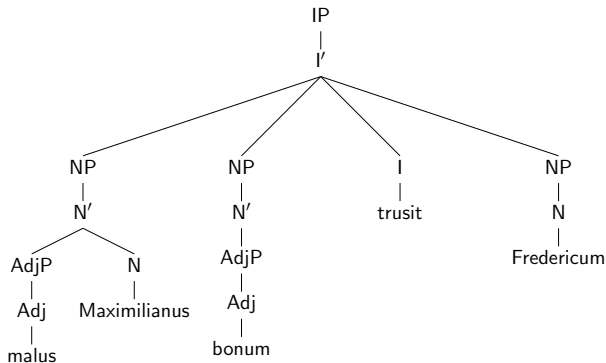


L



af-dauþs* (Ger.) Pl. PP. zu 'af-
 dauþi, éckuauþs geschunden,
 geflagt; N.M. einu M 936i
 dauþs Pl. þauþi Germ. Pl. dauþs
 wóþi: eð. N. Pl. dauþs k 2,15
 E 5,2; N. R 12,17 k 2,1; 6; A.
 k 2,14; G. J 13; D. þvannþoc
 dauþeins Pl. (152) vékþwic das
 Absterben A. k 4,10; éy þauþ-
 þwic in einum in Todesnöten
 k 1,23
 dauþjan sic.V.1 vékroþv tóten
 G 5,5
 af-dauþjan tóten (perfektiv, 291 ff.)

The result



▶ Hyperbaton

▶ WH-movement

Summary

- We have seen that the PROIEL corpus is a small but deeply annotated corpus
 - Morphology
 - Syntax
 - Information structure
 - Discourse (experimental, not shown)

Summary

- We have seen that the PROIEL corpus is a small but deeply annotated corpus
 - Morphology
 - Syntax
 - Information structure
 - Discourse (experimental, not shown)
- The syntax is as theory-neutral as possible

Summary

- We have seen that the PROIEL corpus is a small but deeply annotated corpus
 - Morphology
 - Syntax
 - Information structure
 - Discourse (experimental, not shown)
- The syntax is as theory-neutral as possible
- But conversion is possible and an interesting for hypothesis testing

Summary

- We have seen that the PROIEL corpus is a small but deeply annotated corpus
 - Morphology
 - Syntax
 - Information structure
 - Discourse (experimental, not shown)
- The syntax is as theory-neutral as possible
- But conversion is possible and an interesting for hypothesis testing
- The output could be used as a test suite for a implementing an LFG grammar

Summary

- We have seen that the PROIEL corpus is a small but deeply annotated corpus
 - Morphology
 - Syntax
 - Information structure
 - Discourse (experimental, not shown)
- The syntax is as theory-neutral as possible
- But conversion is possible and an interesting for hypothesis testing
- The output could be used as a test suite for a implementing an LFG grammar
- It can also make the data more widely available to researchers in other frameworks

Outlook

- The New Testament text is available for many low-resources languages



Outlook

- The New Testament text is available for many low-resources languages
- The fine-grained reference system (book, chapter, verse) makes alignment feasible

Outlook

- The New Testament text is available for many low-resources languages
- The fine-grained reference system (book, chapter, verse) makes alignment feasible
- We will experiment with annotation transfer
 - Cooperation with the Linguistic Data Consortium at Penn: alignment, comparison, annotation transfer with phrase structure-based NT corpora

Outlook

- The New Testament text is available for many low-resources languages
- The fine-grained reference system (book, chapter, verse) makes alignment feasible
- We will experiment with annotation transfer
 - Cooperation with the Linguistic Data Consortium at Penn: alignment, comparison, annotation transfer with phrase structure-based NT corpora
 - Cooperation with Iceland and Språkbanken in Gothenburg: alignment and annotation transfer between annotated and unannotated, Nordic bible texts (Old Swedish, Icelandic, possibly Old Finnish)

Availability

- The corpus is available for everyone to use.
- We publish XML files with raw data as well.
- All our data is released under a Creative Commons license.
- Visit <http://www.hf.uio.no/ifikk/proiel/> for details.