INTRODUCING THE INDIANA PARSED CORPUS OF (HISTORICAL) HIGH GERMAN

Christopher Sapp, Elliott Evans, Rex Sprouse, Daniel Dakota Indiana University



### Overview

- Penn family of parsed corpora has not represented High German (HG) until now.
- Consituency parsed; loosely based on Principles and Parameters grammar.
- Existing Germanic corpora: English (PPCHE), Icelandic (IcePaHC), Yiddish (PPCHY), and Low German (CHLG).
- We fill this gap with the Indiana Parsed Corpus of (Historical) High German.

#### **Corpus Structure**

• 165 texts up to 10,000 words each (total 1.4 million words) from 1050-1950. Three subcorpora:

## **Annotation Process**

- 1. Extract each sentence from ReM/ReF/DTA into a single bracketed line, using treetools, C6C, and/or python scripts:
  - (3) (S (PRELS der) (AP (PP (APPR mit) (NA golt)) (ADJV koestlich) (VVPPD belegt)) (VAFIN was)) 'which with gold richly covered was' (1533 Fierrabras, 36)
- 2. Parse with the Berkeley Neural Parser using word+char+dbmdz\_embeddings (predict POS tags as aux task), trained targeting gold ENHG:
  - (4) (WNP (D der)) (IP-SUB (PP (P mit) (NP (N golt))) (ADVP (ADV koestlich)) (VBN belegt) (BEDI was))
- Middle High German (MHG; 1050-1350): 35 texts, selected from the Referenzkorpus Mittelhochdeutsch (ReM). Manually divided into sentences, manually corrected POS and inflection.
- Early New High German (ENHG; 1350-1650): 64 texts, from the Referenzkorpus Frühneuhochdeutsch (ReF). Some manually, some automatically tagged.
- New High German (NHG; 1650-1950): 66 texts, from the Deutsches Textarchiv (DTA). Automatic and often inaccurate sentence divisions and tagging.
- 12 German-speaking regions (map below):
- 10 HG regions (yellow and orange) are represented at every time period.
- 2 Low German regions (blue) represented beginning in the late 16th century.

	Cent West Middle German			West Upper German			East Middle German		East Upper German			(Low German areas)		
	1.	Cologne	Hesse	Alsace	Swabia	Switzerla	Saxony	Thuringia	Nurember	Bavaria	Austria	Northwes	Northeast	
	11.2			Älterer Ph	ysiologus &	Rheinau.			Williram	Otlohs	Gebet			
					Gebete									
	12.1	12.1 Rheinfr. Interlinear &		Alkuins Traktat & Londoner				Wiener Physiologus & Prü		& Prüler				
		Bamb. Arzneibuch		Predigt & Christi Geburt				Steinbuch						
	12.2	Schleizer	Trierer		Trud.	Zuricher		Predigtfra	Schlierbac	Windberg	Spec.eccl.	(		
		Psalm	Interlinear	Lucidarius	Hohelied	Pred.		g.	h Ps.	er Psalter				
	13.1	3 short	Mitteldeut	Millstätter	Hoffmann.	Zweifalten		Mülhäuser	Prager	St. Pauler	Leysersch			
	40.0	texts	sche Pred.	Pred.	Pred.	er		Rechtsbuc	Predigtent	Pred.	e Pred.			
	13.2	Die Lilie	Salomons	Freiburger	Buch der	Hugo V. S.		Jenaer			Admonter	(		
	111	Kaalnar	Haligaplah	Nikolouo	Augoburg		Loinzigor	Porlinor	Engoltholo	Augsburg	Den.			
	14.1	Klosternre	en	Pred	Augsburg	Prediaten	Pred	Evangelist	r	Rechtsh	h			
Northwest Northeast	14 2	Nuwe	Leiden	Büchlein	Rotes	Naturlehre	Altdeutsch	Thüringer	6 Namen	Buch der	Rationale			
Cormon V	14.2	Buch Köln	Christi	d. ew.	Buch Ulm	Mainau	e Pred.	Spiele	Fronleichn	Natur	Rationale			
	15.1	Reimchro	Karrenritte	Nebuchod	Fuchsfalle	Appenweil	Ältestes	Eisenache	Laien-	Sendbrief	Denkwürdi			
		nik	r	onosor		er Chronik	Stadtbuch	r Chronik	doktrinal	von	g-keiten			
	15.2	Koelhoff	Jerusalem	Chirurgie	Verbotene	Edlibach	Tauler	Stolle:	St.	Bairische	Hystoria			
		Chronik			Kunst	Chronik	Sermon	Memoriale	Anselmi	Chronik	Troyana			
	16.1	Junge	Fierrabras	Butzers	Franck:	Olvier und	Bachmann	Vonn	Osiander:	Geistliche	Stiftbriefe			
Saxony		fursten		Predig	Weltbuch	Artus	: Martinus	gehorsam	Grundlich	Mai				
	16.2	Epitome	Wahrhafti	Nachbarn	Beschreib.	Gespenst	Chronica	Thüringisc	Dietriech:	Concordia	Moscouia	Sattler: 3	Wahrhafti	
THE RESERVE AND A STATE OF A STAT		Warhaftig	g historia		der Reise	er	Marsburg	he	Summaria			Predigten	ge	
	17.1	Teutscher	Hessische	Policeij	Lichtkugel	Brun:	Opitz:	Peckenste	Opus	Kurtze vnd	Faber:	Das	Wahrem	
		Nation	Chronik	Ordnung		Schiffarten	Poeterey	in: Theatri	Theatricu	Nothwendi	Probstein	Friede	Christentu	
	17.2	Santa	Becher:	Dannhaue	Zeller:	Heidegger	Weise:	Comoedia	Saar: Ost-	Furttenba	Beer: Der	Hamburg:	Siegemun	
-berg	10 1	Clara:	Psychoso	r: Easamann	Cenuria	Seboueb <del>z</del>	Drey	Vom	Indianisch	Ch: Dockor	verliebte	Statuta	Cottochod	
Margan and a second	10.1	Güldenes	Post-	· Golobrto		or: Natur-	Sector-	Betrachtu	Curiouso	Baumoisto			Gousched	
Alsace/ Bovaria	18.2	Kortum.	Cancrin:	Mesmer	Schiller	Mever	Goeze.	Reichardt	Glück:	Sailer:	Kemnelen <sup>.</sup>	Campe <sup>.</sup>	Forster	
Baden Ourshie Davalla	10.2	Jobsiade	Beschr	Magnetis	Naïve und	Grossen	Zeitvertrei	Land- u	Versuch	Kurzgefas	Maschine	Theophro	Ansichten	
Austria Austria	19.1	Hartwia:	Schopenh	Kerner:	Schmidlin	0.00001	Laube:	Niethamm	Heael:	Steub:	Grillparzer	Wienbarg:	Fouqué:	
Cha Tradition Tradition		Die	auer:	Geschicht	Ueber		Die	er: Streit	Wissensch	Drei	: Ein	Aesthetisc	Frauen in	
And the second s	19.2	Fuhlrott:	Huber:	Laband:	Huber:	Wedekind:	Nietzsche:	Schleicher	Wundt:	Pocci:	Ernst:	Brunn:	Gercke:	
Switer-		Neanderth	Sieben	Das	Geschichtl	Frulings	Also	:	Handbuch	Lustiges	Training	Griechisch	Torpedow	
	20.1	Eggor	Promochoi		Krukonhor	Staigar	Hampo	Kruggor		Wahar	Adlar	Book:	Zotkin:	

- 3. Replace the POS tagging from the parser with ReM/ReF's gold POS tags (if available), but keep the phrasal labels from the parser:
  - (5) (WNP (PRELS der)) (IP-SUB (PP (APPR mit) (NP (NA golt))) (ADVP (ADJV koestlich)) (VVPPD belegt) (VAFIN was))
- 4. Rule-based preprocessing to improve the parse: convert gold POS tags to a version of our tagset (keeping a few distinctions e.g. D-relative), correct or flag obvious errors, insert null elements (here, an erroneous null subject):
  - (6) (WNP (D-relative der))

(IP-SUB (NP-SBJ \*pro\*-CHECK) (PP (P mit) (NP (N golt))) (ADVP (ADV koestlich)) (VBN-adverbial? belegt) (BEDI^3^SG was))

- 5. Hand correct the parse, build higher constituents (here, the CP-REL), remove flags (e.g. changing D-relative to D), add inflection:
  - (7) (CP-REL (WNP-SBJ-2 (D^N^SG der))

```
(C 0)
(IP-SUB (NP-SBJ *T*-2)
(PP (P mit) (NP (N^D^SG golt)))
(ADVP (ADV koestlich))
(VBN belegt)
```



- •When possible, 1 text per region per 50-year time bin (table above).
- Currently 34 ENHG texts gold annotated (shaded green) and available at https://ipchg.iu.edu/
- Balance sociolinguistic factors (genre/register, social class, gender) if possible.

Tagset

- We adapt the Penn tagset to German, mostly following the POS tags of the Old Saxon HeliPaD.
- Source corpora use variants of the Stuttgart-Tübinger Tagset (STTS).
- STTS is often redundant in a parsed corpus (e.g. pre- vs. postpositions).
- STTS also lacks useful distinctions, e.g. collapsing tense & mood as 'finite.'

### **Annotation Issues**

- German-specific adaptations, some following CHLG (Booth et al. 2020). Others are unique to HG, e.g.:
- -Homophonous adverbs/complementizers tagged C and placed in the sub. clause.
- -Exception: *da* following an overt relative pron. is treated as ADV:
  - (1) (CP-REL (WNP-SBJ-1 (WD^N^SG welches))

- (BEDI^3^SG was)))
- 6. Rule-based postprocessing flags remaining errors; final hand correction.
- 7. Inter-annotator reliability 89.01 after step 5; improves to 91.70 after step 6.

# Use Case Study

- Besides filling a gap in the Penn family, the IPCHG is already useful for research in diachronic syntax.
- Adnominal genitives in ENHG precede (8) or follow (9) the head N:
  - (8) [<sub>Gen</sub> meins hrrn] eelicher sun
    'my lord's legitimate son' (1480 *Troyana*, 342)
  - (9) das haubt [<sub>Gen</sub> der heyligen jungfrauwen]
    'the head of the holy virgin' (1486 *Jerusalem*, 24)

1. We discovered two additional structures: 'split' (10) and 'embedded' (11) genitives:

- (10) [<sub>Gen</sub> Josephs] sun [<sub>PP</sub> von aramathia]
  'Joseph of Arimathea's son' (1430 *Karrenritter*, 472)
- (11) eyn besunder [<sub>Gen</sub> Rulands] streitgesel
  'a certain combatant of Ruland' (1533 *Fierrabras*, 196)
- 2. We can determine language-internal effects, e.g. the length of the genitive (Fig. 1):

(C 0) (IP-SUB (NP-SBJ \*T\*-1) (ADVP (ADV da)) (RP an) (VBDI^3^SG kam))) '...[a ship], which arrived (<sup>?</sup>there)' (1557 Staden Historia, 242)

Discourse particles are difficult to distinguish from adverbs, thus tagged ADV.
Superlatives of predicative adjectives (*am kleinsten*, lit. 'on the smallest') look like PPs, but we annotate *am* 'on the' as an unanalyzed particle:



Figure 1: effect of length on order

Figure 2: effect of year on order

3. We can determine external/variationist features; e.g. time (Fig. 2).