

Visualization of Parallel Multilingual Historical Data

Aikaterini-Lida Kalouli, Rebecca Kehlbeck, Rita Sevastjanova,

Katharina Kaiser, Georg A. Kaiser, Miriam Butt

Universität Konstanz

firstname.lastname@uni-konstanz.de

Funded by the German Research Foundation (DFG)
Research Unit 2111 "Questions at the Interfaces"



Motivation & Challenges

Parallel corpora for the analysis of historical linguistic data facilitate:

- **direct comparability** of concrete examples across time periods
- **selective investigation** of passages with potentially relevant structures
- analysis of **languages not spoken** by the researcher, based on the known languages

However, such data are often (cf. [1]):

- too **sparse** for state-of-the-art statistical methods
- too large and **high-dimensional** (time, domain, language) for manual inspection
- unsuitable for learning methods which necessarily **reduce dimensionality**

ParHistVis

Function: interactive **visualization tool** for **parallel**, multilingual data of a) the same time period across languages; b) of different periods of the same language; c) across languages.

Input: tabulated file with aligned data over time or/and language, annotated with features

Output: color-encoded matrix view of the data

URL: typo.uni-konstanz.de/parhistvis/

Parallel and Aggregated Analysis of Linguistic Change

- **preserves dimensionality:** investigation of the data in a parallel manner (Fig. 3)
- **avoids overwhelming:** different features encoded with different colors (Fig. 3)
- allows **detailed view:** user selects subset to investigate, subset gets highlighted (Fig. 3)

Use Case: Romance interrogatives

- **Data:** 3 French and 3 Spanish Bible translations of the 12th, 16th and 20th centuries
- **Features:** a) word order in interrogatives
b) interrogative pronouns and verbs of speaking introducing questions
c) particles used with interrogatives
- **Goal:** investigation of the strict word order in Old Romance vs. the greater word order variation in Modern Romance
- **Observations:** a) emergence of complex inversion in Modern French [2]: the orange stream (whSNPVSCI) first appears in Middle French and increases its frequency in Modern French (Fig. 1)
b) diachronic non-adjacency of the wh-element and the verb when a particle is present: the blue (whPtcVS) stream stays stable over time (Fig. 1)
c) some interrogative pronouns allow more variation in the sentence structure [3]: *why* allows for more frequent use of the particle *pues* and *donc* in Spanish and French, respectively, than other pronouns (Fig. 2)

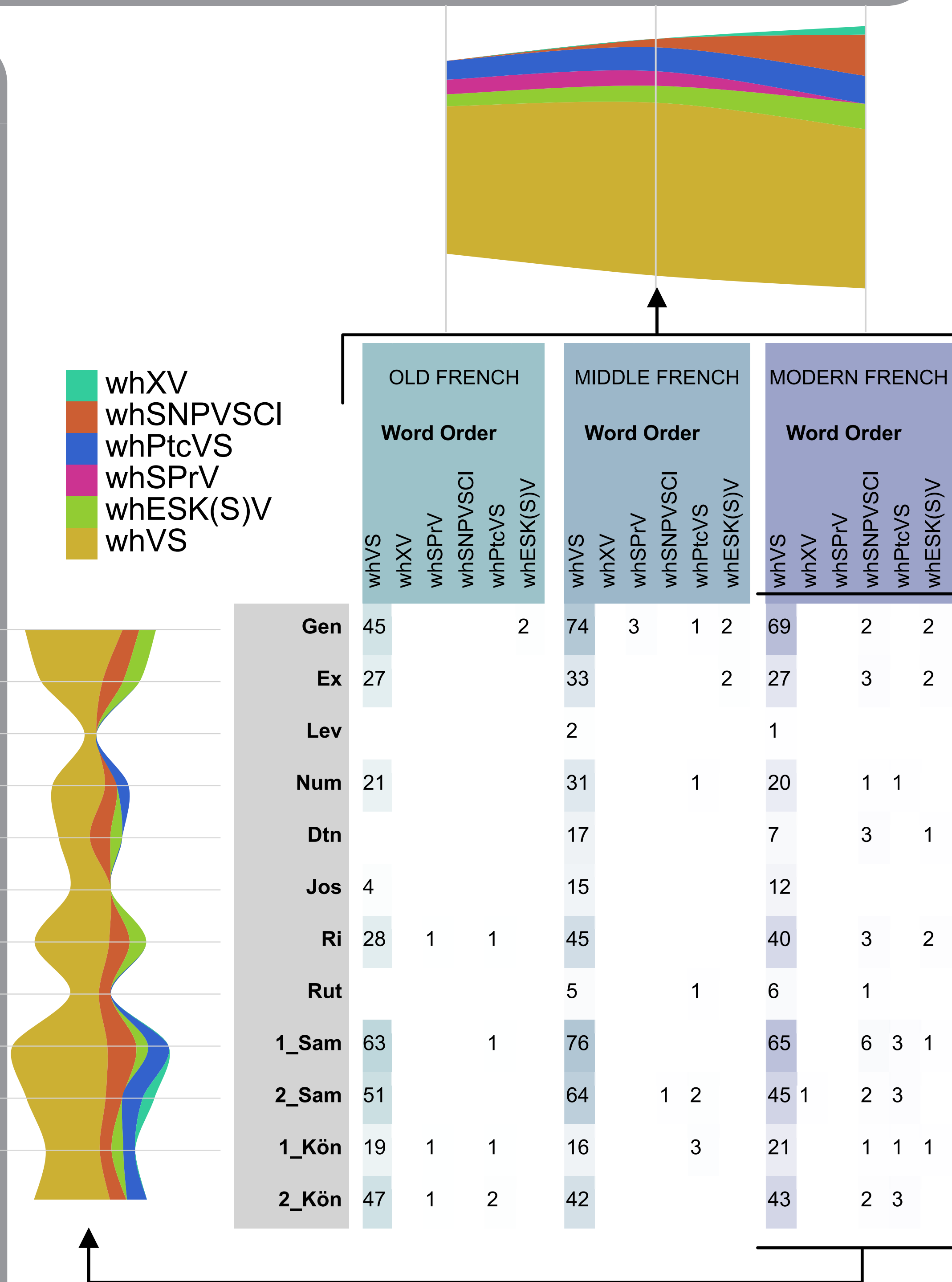


Figure 1: Streamgraphs: Word order in French across time and in Modern French across the aggregated Bible

Pattern Recognition & Interaction

User selects features for comparison:

- **Streamgraphs:** represent selected features as streams (Fig. 1)
- **Sankey diagrams:** represent selected features as nodes and their interaction as a flow between them (Fig. 2)

⇒ **at-a-glance view of patterns and interactions** across time, language and data course



Figure 2: Sankey Diagrams: interactions between particles and interrogative pronouns in French and Spanish

	OLD FRENCH					MIDDLE FRENCH					MODERN FRENCH					OLD SPANISH					MIDDLE SPANISH					MODERN SPANISH														
	Quest. Type	Word Order				Pronoun	Particle	Verb	Word Order	Pronoun	Particle	Verb	Word Order	Pronoun	Particle	Verb	Word Order	Pronoun	Particle	Verb	Word Order	Pronoun	Particle	Verb	Word Order	Pronoun	Particle	Verb												
	ISQ	NISQ	whVS	whXV	whSPrV	whSNPVSCI	whPtcVS	whESK(S)V	whVS	whXV	whSPrV	whSNPVSCI	whPtcVS	whESK(S)V	whVS	whXV	whSPrV	whSNPVSCI	whPtcVS	whESK(S)V	whVS	whXV	whSPrV	whSNPVSCI	whPtcVS	whESK(S)V	whVS	whXV	whSPrV	whSNPVSCI	whPtcVS	whESK(S)V								
Gen	64	45	45			2	47		40	74	3	1	2	80	2		66	69	2	2	70	1		57	40		37		31	78	5	81	1	63	81	4	1	84		65
Ex	35	13	27				26		22	33		2	35				28	27	3	2	31	2		23	32		32		25	37	2	39	2	32	38	1	38		29	
Lev	1		1							2				2			1	1			1			1	2		2		1	2		2		1	2		2		1	
Num	34	6	21				21		14	31		1	31	1			19	20	1	1	21	1		13	33		31		16	29	2	30		15	28	2	28		14	
Dtn	14	15								17			18				12	7	3	1	11			8	18		17		9	21		21		10	21		21		9	
Jos	9	9	4				4		4	15			14	1			12	12			12			10	17		16		14	15	1	17	1	14	16	1	16		15	
Ri	48	20	28	1	1		30	2	23	45			44	2			35	40	3	2	44	2		34	45	2	47		33	50	2	52	1	35	50	3	51		36	
Rut	4		4							5			6	1			5	6	1		7			7			7		7		7		6	7	1	8		7		
1_Sam	96	22	63				64	1	44	76			60	65	6	3	1	73	6		73	6		56	85		83		60	85	6	89	1	65	87	6	90	1	63	
2_Sam	80	14	51				51		32	64		1	2	67	3		52	45	1	2	49	4		34	56		54	1	38	57	1	8	65	8	47	59	1	8	48	
1_Kön	32	6	19	1	1		21	1	16	16		3	19	2			18	21	1	1	24	1		20	22		21		19	23	2	25	2	23	22	2	23		22	
2_Kön	63	4	47	1	2		49	6	36	42			41				38	43	2	3	47	4		41	38	1	38		34	41	3	43	1	34	45	1	45		36	

Figure 3: Aggregated matrix view of the books of the Old Testament across time periods and languages.

References

- [1] Qihong Gan, Min Zhu, Mingzhao Li, Ting Liang, Yu Cao, and Baoyao Zhou. Document visualization: an overview of current research. *Wiley Interdisciplinary Reviews: Computational Statistics*, 6(1):19–36, 2014.
- [2] Ian Roberts. *Verbs and Diachronic Syntax. A Comparative History of English and French*. Kluwer, Dordrecht, 1993.
- [3] Francisco Ordóñez. *Word order and clause structure in Spanish and other Romance languages*. PhD thesis, The City University of New York, 1997.