

LREC 2022, 13th Conference on Language Resources and Evaluation, Marseille, France, 20-25 June 2022

Querying Interaction Structure: Approaches to Overlap in Spoken Language Corpora

Elena Frick¹, Henrike Helmer¹ and Thomas Schmidt² ¹IDS Mannheim (DE), ²University of Basel (CH)



(<annotationBlock/> containing ([word=".*" & !pos="(NGIRR|NGHES|XY)"]

at least one word token that occurs outside of the speaker overlap and is

not a non-word, hesitation, interjection or responsive particle.

!within <speaker-overlap/>)) precededby (<another-speaker/><para/>{0,5})

looks for turn-taking by one of the non-speakers whose contribution contains



Q I ? [pos="NGIRR"] precededby <pause/> Search ✓ ignore punctuation **Examples** Search by individual Speakers <u>XML</u> 🔁 Search in Transcript

	<pre><pause dur="PT0.31S" end="TLI_992" pre="" remains<=""></pause></pre>	id="p301"/>	Results							
	<pre><annotationblock end="TLI_996" start="TLI_99</pre></td><td>2" who="US" xml:id="c667"></annotationblock></pre>		for a line for a line							
	<u></u>			for searching [pos="NG	IRR"] precededby <p< td=""><td>ause/> (in FOLK)</td><td>Group Hits O</td><td>pen Metadata Vi</td><td>iew Dow</td><td>vnload KWIC</td></p<>	ause/> (in FOLK)	Group Hits O	pen Metadata Vi	iew Dow	vnload KWIC
	<pre><accident a="" and="" of="" start="" t<="" td="" the=""><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></accident></pre>									
	<pre><w lemma="ia" norm="ia" nos="NGIRR" type="ol-in" xml:id="w3007">ia</w></pre>			Total: 61054			First Previou	s 1 2 3 4	4 5 Next	Last
	<pre><w type="ol-in" xml:id="w3008">es</w></pre>	pos norm space								
	$\langle w \rangle$ xml:id="w3009" type="ol-in">is									
	<w type="ol-in" xml;id="w3010">ja</w>			1 FOLK_E_00022_S	E_01_T_04 HM (0				DGD	ZuViel
	<w type="ol-in" xml:id="w3011">auch</w>			2 FOLK E 00022 S	E 01 T 04 HM so	aber hh° h° ähm °h (()	äuspert sich)) (0.42) ia	out was halt	DGD	7uViel
	<pre><w type="ol-in" xml:id="w3012">wenn</w> <w type="ol-in" xml:id="w3013">sich</w> <anchor synch="TLI_993" type="ol-end"></anchor></pre>								DCD	
				5 FOLK_E_00022_5	5 FOLK_E_UUU22_SE_UI_I_U4 HMI tatsachlich gucke ahm ((rauspert sich)) (0.78) ja was äh was D					ZuViel
	<w xml:id="w3014">da</w> <w xml:id="w3015">gas</w>			4 FOLK_E_00022_S	E_01_T_04 HM so	o sei ähm hh° (0.32) ja (1.13) wie sieht er		DGD	ZuViel
	<w type="ol-in" xml:id="w3016">au<anchor synch="TLI_994" td="" typ<=""><td>art"/>sbreitet<td>5 FOLK_E_00022_S</td><td>E_01_T_04 SZ ga</td><td>ar net ((räuspert sich))</td><td>okay (0.3) ja aber grad</td><td>e zum</td><td>DGD</td><td>ZuViel</td></td></anchor></w>		art"/>sbreitet <td>5 FOLK_E_00022_S</td> <td>E_01_T_04 SZ ga</td> <td>ar net ((räuspert sich))</td> <td>okay (0.3) ja aber grad</td> <td>e zum</td> <td>DGD</td> <td>ZuViel</td>	5 FOLK_E_00022_S	E_01_T_04 SZ ga	ar net ((räuspert sich))	okay (0.3) ja aber grad	e zum	DGD	ZuViel
	<pre><anchor synch="TLI_995" type="ol-in"></anchor> <w type="ol-in" xml:id="w3017">und</w></pre>			6 FOLK_E_00022_S	E_01_T_04 AW	und denkt hups (.) ja w	o de auch		DGD	ZuViel
	<w type="ol-in" xml:id="w3018">dann</w>	Timo intorvolo	in which char "N	L" is silont but ath	r charlore are	snooking				
	<pre><anchor synch="TL1_996" type="ol-end"></anchor> </pre>		in which speaker in	n is silent, but othe	er speakers are s	speaking				
		<pre>ccnanGrn type</pre>	-"anothor_cnoakor"	subtypo-"timo_bacc	d">					
	<pre><spangrp subtype="time-based" type="speaker-overlap"></spangrp></pre>	- spandip type-	- another-speaker $\frac{1}{2}$	Sublype- line-base	u >					
	LM	<span <="" from="" td=""><td>TLI_992" [0="TLI_992"</td><td>4">US</td>	TLI_992" [0="TLI_992"	4">US		1 [26:45. 2 [26:45.3]		3 [26:46.2] 4 [[26:46.8]	5 [26:47.4]
	<pre>NH</pre>	<span from="</td><td>" ili_992"="" to="ILI_993</td><td>3">LM	US [v]	ja es is ja	auch wenn sich	da gas au st	breitet	und danr		
	AM	<span from="</td><td>'TLI_995" to="TLI_996</td><td>6">US	LM [v]	ia () für	dia	_				
		<spangrp></spangrp>			Dist [4]	ja (.) iui	uic			
			• •					w	ar s auch	
	<pre><annotationblock end="TLI_993" pre="" start="TLI_992" who="LM" xm<=""></annotationblock></pre>	l:id="c668">			AM [v]					
	<u></u>					(0.21)				
	<seg></seg>				[m]	[0.52]				
	<pre><anchor synch="TL1_992" type="01-start"></anchor></pre>									
	<pre><w <="" lemma="ja" norm="ja" pre="" type="ol-in" xml:1d="w3019"></w></pre>	pos="NGIKK">Jak/	w>							
	<pre><pre><pre><pre><pre><pre><pre>cuse rend="(.)" type="micro" xml:id="p302"/></pre></pre></pre></pre></pre></pre></pre>		9 [51:57.5] 10 [51	:58. 11 [51:58:5] 12 [51:5	3.9]		13 [52:00.2]	14 [52:00.5]		15 [52:01.4]
Contact:	<pre><w type="ol-in" xml:id="w3021">die</w></pre>	AM [v]		(flacht))				ia des nein	n nein	es is
Elena Frick M.A.	<pre><anchor synch="TLI_993" type="ol-end"></anchor></pre>		in an () man	de a blinet d'as bli	et istat mislisis	at hart ish maik	niah mia iah a	haaan ooo	مع راء ترسل	
Abteilung Pragmatik Leibniz-Institut für Deutsche Sprache		0.5 [4]	ja so (.) nee				men wie ien s	UCSSCI aus	urueken	SUL
Postfach 10 16 21	<pre><spangrp subtype="time-based" type="speaker-overlap"></spangrp></pre>	NH [v]	((lac	ht)) das jetz vielleic	ht über		trieben			ja
68016 Mannheim, Germany	<pre>US </pre>	[nn]								
Phone: +49 621 1581-137	<pre><annotationblock end="TLT 995" start="TLT 99</pre></td><td>4" who="NH" xml:id="c669"></annotationblock></pre>									
Fax: +49 621 1581-200			1 500 41 51							E 750 44
frick@ids-mannheim.de	<seq></seq>		1 [52:41.5]			2 (52:43, 3 (52:4,	5.9] 4 [52:44.4]			5 [52:43
mercendo maninem.de	<pre><anchor synch="TLI 994" type="ol-start"></anchor></pre>	US [v]	((Lachansatz)) (.)	der muss schal schmecken (.) ((Lachansatz)) °h ((lacht)) also trinken ((lacht))						
	<w type="ol-in" xml:id="w3022">war</w>	NH [v]	also der muss (.) e	eigentlich muss er weg		°h (flacht))				
	<pre><w type="assimilated ol-in" xml:id="w3023">s</w> </pre>	4 B/I [37]				((ler worder im h	ühleebroolz el	leo ie () ie	der kolt
Street Address:	<pre><mail.iu= <anchor="" sauch(="" synch="TLT 995" type="ol-end" wooze="" ws=""></mail.iu=></pre>	2201 [4]				war		unsen ank al	130 13 (.) 15	uei Kält
Leibniz-Institut für Deutsche Sprache		[nn]								

frick@ids-man

Street Address: Leibniz-Institut R 5, 6-13 68161 Mannheim, Germany

Phone: +49 621 1581-0 Fax: +49 621 1581-200 info@ids-mannheim.de www.ids-mannheim.de

© 2022 IDS Mannheim/ÖA

</u>

<spanGrp type="speaker-overlap" subtype="time-based"> US </spanGrp> </annotationBlock>

More about the project:

Fandrych, C., Frick, E., Kaiser, J., Meißner, C., Portmann, A., Schmidt, T., Schwendemann, M., Wallner, F., and Wörner, K. (2022). ZuMult: Neue Zugangswege zu Korpora gesprochener Sprache. In Kämper, H. et al. (Eds.), Sprache in Politik und Gesellschaft: Perspektiven und Zugänge. Jahrbuch des Instituts für Deutsche Sprache 2021. Berlin etc.: de Gruyter.

Frick, E. and Schmidt, T. (2020). Using Full Text Indices for Querying Spoken Language Data. In Proceedings of the LREC 2020 Workshop, Language Resources and Evaluation Conference, 11–16 May 2020, 8th Workshop on Challenges in the Management of Large Corpora (CMLC-8), pages 40–46, Paris: European Language Resources Association (ELRA).

Speaker Overlaps: segment-based vs. contribution-based approach

[word.type=".*ol-in.*"]

looks for word tokens within overlaps; the search pattern containing regular expression characters '.*' from both sides of 'ol-in' is important to match also type-attributes containing multi-word values

[norm="bitte" & word.type=".*ol-in.*"]

looks for any transcribed form of 'bitte' within overlaps

<annotationBlock/> containing [word.type=".*ol-in.*"] looks for all speaker contributions containing overlaps

<speaker-overlap/>

looks for all spans annotated as speaker overlap

<speaker-overlap/> containing [lemma="(Herr|Frau)"]

looks for all spans annotated as speaker overlap and containing any forms of 'Herr' or 'Frau'

<speaker-overlap>[norm="also"]

looks for any transcribed form of 'also' at the beginning of speaker overlaps

<speaker-overlap="SZ"/>

looks for all token sequences overlapping with the contributions of the speaker 'SZ'