

ÚFAL CorPipe at CRAC 2022: Effectivity of Multilingual Models for Coreference Resolution

Milan Straka, Jana Straková



Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks
 - works better than individual models

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks
 - works better than individual models
 - can be used even on unseen languages

CorPipe

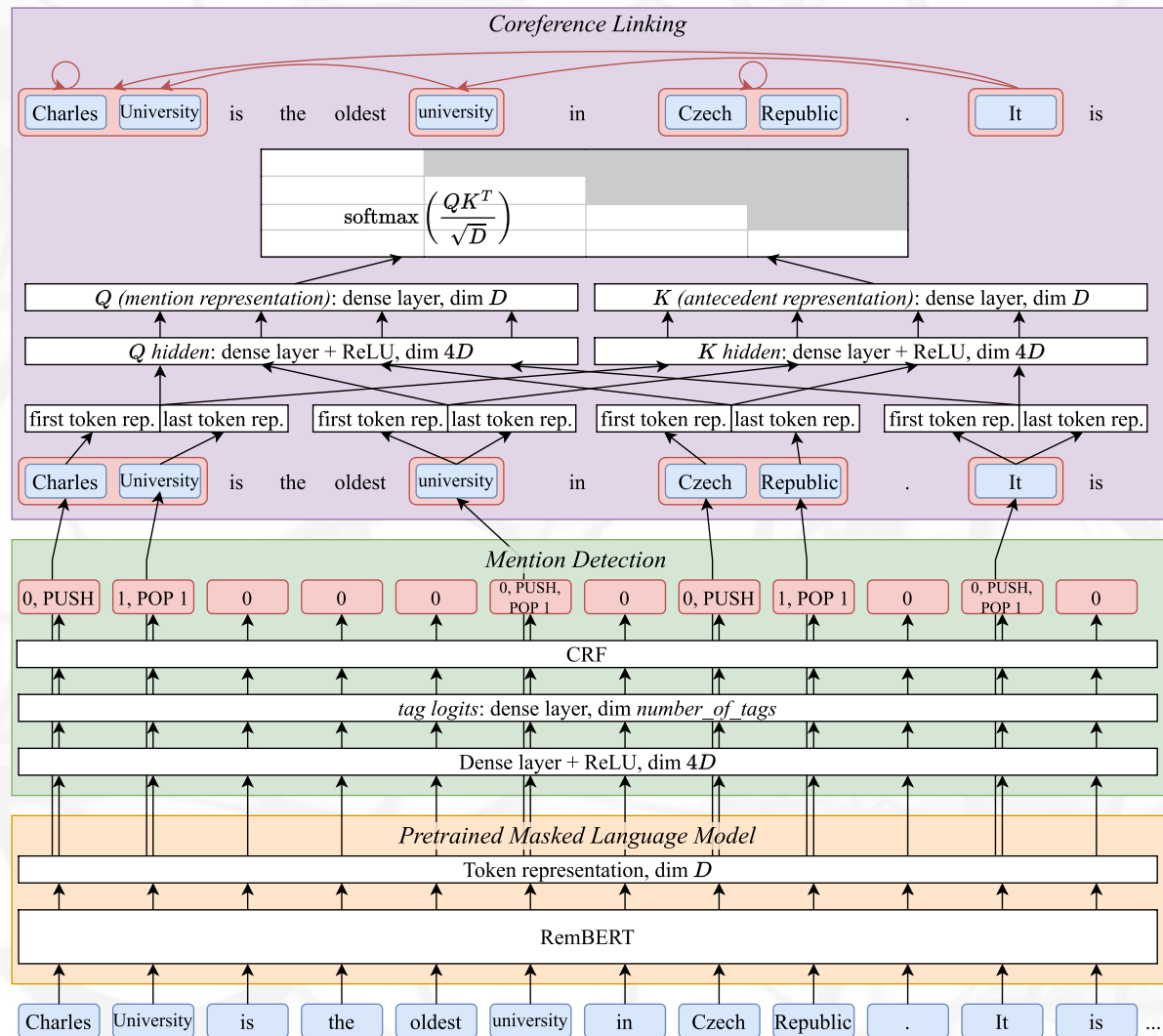
- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks
 - works better than individual models
 - can be used even on unseen languages
 - contrary to the baseline solution:
 - CorPipe first predicts mentions

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks
 - works better than individual models
 - can be used even on unseen languages
 - contrary to the baseline solution:
 - CorPipe first predicts mentions
 - only then it predicts links between the predicted mentions

CorPipe

- winning entry of the CRAC 2022 Shared Task on Multilingual Coreference Resolution
- distinguishing features:
 - a single multilingual model for all 13 treebanks
 - works better than individual models
 - can be used even on unseen languages
 - contrary to the baseline solution:
 - CorPipe first predicts mentions
 - only then it predicts links between the predicted mentions
 - both tasks performed by a single model



Team/Submission	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
ÚFAL CorPipe	70.72	78.18	78.59	77.69	65.52	70.69	72.50	39.00	81.39	65.27	63.15	69.92	78.12	79.34
<i>best dev</i>	1	2	1	1	2	2	2	1	1	1	3	1	2	1
ÚFAL CorPipe	69.56	78.49	78.49	77.57	59.94	71.11	73.20	33.55	80.80	64.35	63.38	67.38	78.32	77.74
<i>multilingual</i>	2	1	2	2	3	1	1	3	2	3	2	3	1	2
UWB	67.64	70.55	74.07	72.42	73.90	68.68	68.31	31.90	72.32	61.39	65.01	68.05	75.20	77.50
<i>ondfa[†]</i>	3	4	4	4	1	3	4	4	4	4	1	2	4	3
ÚFAL CorPipe	64.30	76.34	77.87	76.76	36.50	56.65	70.66	23.48	78.78	64.94	62.94	61.32	73.36	76.26
<i>individual</i>	4	3	3	3	5	5	3	5	3	2	4	6	5	4
Barbora Dohnalová	59.72	64.67	70.56	67.95	38.50	57.70	63.07	36.44	66.61	56.04	55.02	65.67	65.99	68.17
<i>berulasek</i>	5	5	5	5	4	4	5	2	5	5	5	4	6	5
UWB	58.53	63.74	70.00	67.27	33.75	55.44	62.59	36.44	65.98	55.55	52.35	64.81	65.34	67.66
BASELINE [‡]	6	6	6	6	6	6	6	2	6	6	6	5	7	6
Matouš Moravec	55.05	58.25	68.19	64.71	31.86	52.84	59.15	36.44	62.01	54.87	52.00	59.49	63.40	52.49
<i>moravec</i>	7	7	7	7	7	7	7	2	7	7	7	7	8	7

Table 1: Official results of CRAC 2022 Shared Task on the test set (CoNLL score in %). The systems [†] and [‡] are described in [Pražák and Konopik \(2022\)](#) and [Pražák et al. \(2021\)](#), respectively; the rest in [Žabokrtský et al. \(2022\)](#).

Team/Submission	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
ÚFAL CorPipe <i>best dev</i>	70.72	78.18	78.59	77.69	65.52	70.69	72.50	39.00	81.39	65.27	63.15	69.92	78.12	79.34
	1	2	1	1	2	2	2	1	1	1	3	1	2	1
ÚFAL CorPipe <i>multilingual</i>	69.56	78.49	78.49	77.57	59.94	71.11	73.20	33.55	80.80	64.35	63.38	67.38	78.32	77.74
	2	1	2	2	3	1	1	3	2	3	2	3	1	2
UWB <i>ondfa</i> [†]	67.64	70.55	74.07	72.42	73.90	68.68	68.31	31.90	72.32	61.39	65.01	68.05	75.20	77.50
	3	4	4	4	1	3	4	4	4	4	1	2	4	3
ÚFAL CorPipe <i>individual</i>	64.30	76.34	77.87	76.76	36.50	56.65	70.66	23.48	78.78	64.94	62.94	61.32	73.36	76.26
	4	3	3	3	5	5	3	5	3	2	4	6	5	4
Barbora Dohnalová <i>berulasek</i>	59.72	64.67	70.56	67.95	38.50	57.70	63.07	36.44	66.61	56.04	55.02	65.67	65.99	68.17
	5	5	5	5	4	4	5	2	5	5	5	4	6	5
UWB BASELINE [‡]	58.53	63.74	70.00	67.27	33.75	55.44	62.59	36.44	65.98	55.55	52.35	64.81	65.34	67.66
	6	6	6	6	6	6	6	2	6	6	6	5	7	6
Matouš Moravec <i>moravec</i>	55.05	58.25	68.19	64.71	31.86	52.84	59.15	36.44	62.01	54.87	52.00	59.49	63.40	52.49
	7	7	7	7	7	7	7	2	7	7	7	7	8	7

Table 1: Official results of CRAC 2022 Shared Task on the test set (CoNLL score in %). The systems [†] and [‡] are described in [Pražák and Konopik \(2022\)](#) and [Pražák et al. \(2021\)](#), respectively; the rest in [Žabokrtský et al. \(2022\)](#).

Team/Submission	Avg. with singletons
ÚFAL CorPipe, <i>best dev</i>	72.98
ÚFAL CorPipe, <i>multilingual</i>	71.81
ÚFAL CorPipe, <i>individual</i>	67.93
UWB, <i>ondfa</i>	58.06
Barbora Dohnalová, <i>berulasek</i>	50.84
UWB, BASELINE	49.69
Matouš Moravec, <i>moravec</i>	46.79

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
XLM-R base, multilingual	67.8	77.1	75.8	74.3	54.7	66.9	70.1	38.5	77.6	64.2	62.3	69.4	73.3	76.6
Best base model, individual	-5.2	-4.0	+1.5	+2.2	-18.2	-9.8	-3.4	-15.0	-2.4	-2.4	-2.0	-8.1	+0.0	-5.6
Best base model, best dev	+0.4	-0.6	+1.5	+2.2	+2.0	+0.6	-1.0	-0.9	+1.2	-1.1	+0.6	+0.4	+0.0	+0.2
RemBERT, multilingual	+1.8	+1.4	+2.6	+3.3	+5.2	+4.2	+3.1	-4.9	+3.2	+0.1	+1.1	-2.0	+5.0	+1.1
RemBERT, individual	-3.5	-0.7	+2.0	+2.5	-18.2	-10.3	+0.6	-15.0	+1.2	+0.7	+0.7	-8.1	+0.0	-0.4
RemBERT, best dev	+3.0	+1.1	+2.7	+3.4	+10.8	+3.8	+2.4	+0.5	+3.8	+1.0	+0.9	+0.5	+4.8	+2.7

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
C) EFFECT OF MULTILINGUAL DATA AND THE PRETRAINED MODEL														
XLm-R base multilingual	73.3	75.8	76.0	75.0	73.4	74.1	73.1	75.4	78.4	66.1	65.2	78.0	72.1	71.7
XLm-R large multilingual	+1.5	+1.7	+1.8	+2.0	+0.3	+4.1	+2.1	-4.5	+2.2	+1.7	+3.1	-0.0	+2.9	+0.9
RemBERT multilingual	+1.9	+1.6	+3.3	+3.3	+2.9	+2.4	+2.4	-6.1	+2.7	+2.0	+4.0	-1.2	+3.7	+2.9
XLm-R base individual	-4.6	-4.4	-0.3	-1.1	-7.8	-12.1	-1.9	-12.2	-2.8	-3.0	-3.8	-4.6	-2.3	-6.1
XLm-R large individual	-0.6	+0.2	+2.8	+3.0	-7.7	-5.2	-0.9	-4.4	+1.0	+0.3	+3.7	-5.4	+3.5	-1.2
RemBERT individual	-4.7	+0.6	+2.8	+1.9	-23.0	-12.1	+0.7	-30.5	+1.1	+0.7	-0.4	-8.9	+2.7	-1.8
RemBERT 50% additional	+0.3	+1.0	+2.5	+2.4	-1.4	-0.5	+1.7	-8.3	+0.9	+1.3	+1.6	-3.5	+3.6	+1.7

Comparing Pretrained Models on Dev

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
C) EFFECT OF MULTILINGUAL DATA AND THE PRETRAINED MODEL														
XLm-R base multilingual	73.3	75.8	76.0	75.0	73.4	74.1	73.1	75.4	78.4	66.1	65.2	78.0	72.1	71.7
XLm-R large multilingual	+1.5	+1.7	+1.8	+2.0	+0.3	+4.1	+2.1	-4.5	+2.2	+1.7	+3.1	-0.0	+2.9	+0.9
RemBERT multilingual	+1.9	+1.6	+3.3	+3.3	+2.9	+2.4	+2.4	-6.1	+2.7	+2.0	+4.0	-1.2	+3.7	+2.9
XLm-R base individual	-4.6	-4.4	-0.3	-1.1	-7.8	-12.1	-1.9	-12.2	-2.8	-3.0	-3.8	-4.6	-2.3	-6.1
XLm-R large individual	-0.6	+0.2	+2.8	+3.0	-7.7	-5.2	-0.9	-4.4	+1.0	+0.3	+3.7	-5.4	+3.5	-1.2
RemBERT individual	-4.7	+0.6	+2.8	+1.9	-23.0	-12.1	+0.7	-30.5	+1.1	+0.7	-0.4	-8.9	+2.7	-1.8
RemBERT 50% additional	+0.3	+1.0	+2.5	+2.4	-1.4	-0.5	+1.7	-8.3	+0.9	+1.3	+1.6	-3.5	+3.6	+1.7

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
G) EFFECT OF SEVERAL LANGUAGE-SPECIFIC BASE PRETRAINED MODELS														
XLm-R base individual	68.7	71.4	75.7	73.9	65.7	62.0	71.2	63.2	75.6	63.1	61.5	73.4	69.8	65.6
mBERT (Devlin et al., 2019)	-2.8	-1.5	-3.0	-3.4	-3.3	+0.4	-2.8	-1.1	-1.8	-1.1	-2.7	-7.5	-4.4	-3.6
BERTa (Armengol-Estapé et al., 2021)	+1.3													
RobeCzech (Straka et al., 2021)			+2.0	+2.8										
gBERT (Chan et al., 2020)					-9.9	+5.3								
SpanBERT (Joshi et al., 2020)							-0.4	-2.4						
BETO (Cañete et al., 2020)									+0.4					
CamemBERT (Martin et al., 2020)										-0.2				
HuBERT (Nemeskey, 2020)											+3.6			
LitLatBERT (Ulčar and Robnik-Šikonja, 2021)												+2.7		
HerBERT (Mroczkowski et al., 2021)													+1.6	
RuBERT (Kurатов and Arkhipov, 2019)														+0.2
XLm-R large individual	+4.0	+4.6	+3.1	+4.1	+0.0	+6.9	+1.0	+7.8	+3.8	+3.3	+7.4	-0.8	+5.8	+4.8
RemBERT individual	-0.0	+4.9	+3.1	+3.1	-15.2	+0.0	+2.6	-18.3	+3.9	+3.8	+3.3	-4.3	+5.0	+4.3
XLm-R large multilingual	+6.1	+6.1	+2.1	+3.2	+8.0	+16.2	+4.1	+7.7	+5.0	+4.8	+6.9	+4.6	+5.1	+6.9
RemBERT multilingual	+6.6	+6.0	+3.6	+4.4	+10.6	+14.5	+4.3	+6.1	+5.5	+5.1	+7.7	+3.5	+6.0	+9.0

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
D) EFFECT OF MIXING RATIOS USING XLM-R BASE PRETRAINED MODEL														
Logarithmic, w/o corpus id	73.3	75.8	76.0	75.0	73.4	74.1	73.1	75.4	78.4	66.1	65.2	78.0	72.1	71.7
Logarithmic, w/ corpus id	-0.4	-0.5	+0.1	+0.3	-0.8	-0.3	-0.6	-4.6	+0.1	+0.6	+1.3	-0.9	+0.3	-0.7
Uniform, w/o corpus id	-0.8	-0.5	-0.2	-0.9	-1.8	-3.5	-0.2	-1.9	+0.0	-0.0	+0.4	-1.1	+0.0	-1.5
Uniform, w/ corpus id	-1.6	-1.1	-0.5	-0.6	-6.4	-2.1	-0.4	-7.0	+0.1	+0.1	-0.5	-1.1	-0.6	-1.2
Linear, w/o corpus id	-0.3	+0.1	+0.8	+1.1	-1.1	-0.5	-0.5	-3.5	-0.1	+0.3	+1.0	+0.3	-0.1	-1.6
E) EFFECT OF MIXING RATIOS USING REMBERT PRETRAINED MODEL														
Logarithmic, w/o corpus id	75.3	77.4	79.3	78.3	76.3	76.5	75.5	69.3	81.1	68.1	69.2	76.8	75.8	74.6
Logarithmic, w/ corpus id	+0.6	+0.4	+0.1	+0.1	+3.0	+1.2	-0.1	+5.8	+0.3	+0.9	-2.4	-1.3	+0.1	-0.2
Uniform, w/o corpus id	+0.1	+1.2	-0.3	-0.1	+2.4	+0.5	+0.0	-0.9	-0.1	+0.7	-0.6	-0.2	+0.1	-1.2
Uniform, w/ corpus id	-0.1	-0.0	-0.2	-0.3	-4.2	+0.3	-0.1	+4.5	+0.4	+0.6	-1.0	-0.1	+0.1	-1.2
Linear, w/o corpus id	-0.1	+1.3	+0.1	+0.2	-2.3	-0.5	-1.5	+1.9	+0.5	+0.7	-1.0	+0.4	+0.0	-1.3

Experiment	Avg	ca	cs pcedt	cs pdt	de parc	de pots	en gum	en parc	es	fr	hu	lt	pl	ru
F) ZERO-SHOT EVALUATION OF A MULTILINGUAL MODEL														
Multilingual XLM-R base	73.3	75.8	76.0	75.0	73.4	74.1	73.1	75.4	78.4	66.1	65.2	78.0	72.1	71.7
Zero-shot XLM-R base	-17.1	-11.1	-28.6	-23.8	-13.3	-13.8	-19.8	-18.5	-6.8	-7.6	-16.1	-23.8	-24.6	-15.1
Multilingual RemBERT	+1.9	+1.6	+3.3	+3.3	+2.9	+2.4	+2.4	-6.1	+2.7	+2.0	+4.0	-1.2	+3.7	+2.9
Zero-shot RemBERT	-12.5	-6.7	-23.7	-20.6	-11.1	-7.5	-15.6	-9.8	-2.8	-8.3	-10.5	-20.0	-18.3	-7.2

Thank You