

TEHNOLOGII ALE LIMBAJULUI PENTRU LIMBA ROMÂNĂ - PLATFORMA RELATE

Vasile Păiș

Institutul de Cercetare pentru Inteligență Artificială "Mihai Drăgănescu",

Academia Română

vasile@racai.ro

29.09.2021



CUPRINS

- SCURTĂ DESCRIERE
- ARTICOLE
- ARHITECTURA PLATFORMEI
- CARACTERISTICI
- VIZUALIZARE
- CORPORA PROCESATE ÎN RELATE
- DEZVOLTARE

RELATE

O PLATFORMĂ PENTRU TEHNOLOGII ALE LIMBAJULUI ÎN LIMBA ROMÂNĂ

INCLUDE TEHNOLOGII DEZVOLTATE LA ICIA ȘI DE PARTENERI ÎN MULTIPLE PROIECTE:

- CoRoLa, ReTeRom, Robin, Presidency, Marcell, Curlicat

ADAPTATĂ FILOSOFIEI **ELG**:

- WEB SERVICES, REST API, DOCKER
- SERVICIILE POT FI DISTRIBUITE PE NODURI DE REȚEA MULTIPLE
- SERVICIILE POT FI CONSUMATE DIRECT DE LA PARTENERI

RELATE



ARTICOLE ȘTIINȚIFICE

VASILE PĂIȘ, RADU ION, AND DAN TUFIȘ. **“A PROCESSING PLATFORM RELATING DATA AND TOOLS FOR ROMANIAN LANGUAGE”**. ENGLISH. IN:PROCEEDINGS OF THE 1ST INTERNATIONAL WORKSHOP ON LANGUAGE TECHNOLOGY PLATFORMS. MARSEILLE, FRANCE: EUROPEAN LANGUAGE RESOURCES ASSOCIATION, 2020, PP. 81–88. [HTTPS://WWW.ACLWEB.ORG/ANTHOLOGY/2020.IWLTP-1.13](https://www.aclweb.org/anthology/2020.iwltlp-1.13)

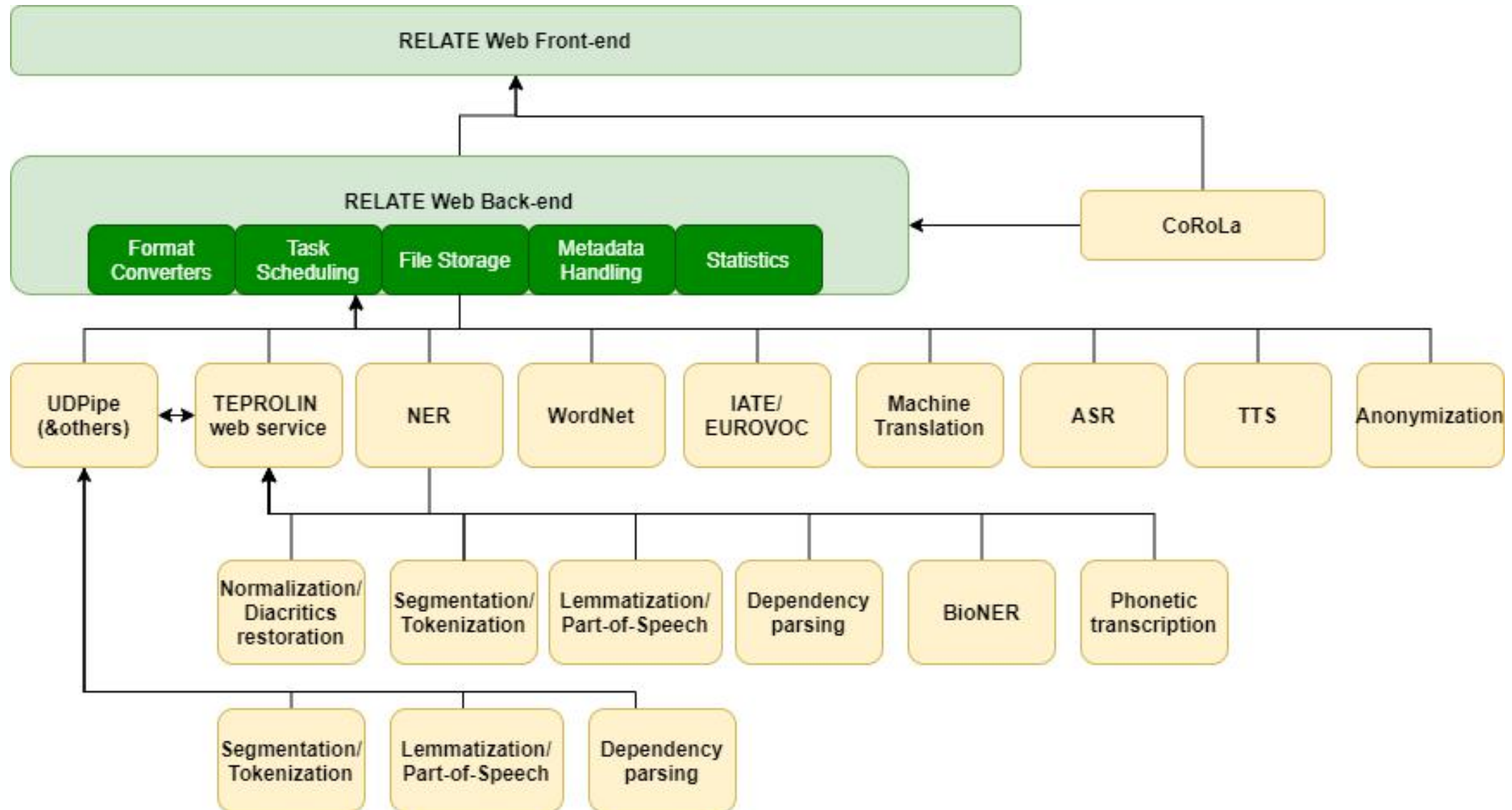
VASILE PĂIȘ. **“MULTIPLE ANNOTATION PIPELINES INSIDE THE RELATE PLATFORM”**. IN:THE 15TH INTERNATIONAL CONFERENCE ON LINGUISTIC RESOURCES AND TOOLS FOR NATURAL LANGUAGE PROCESSING. 2020, PP. 65–75. [HTTPS://PROFS.INFO.UAIC.RO/~CONSILR/WP-CONTENT/UPLOADS/2021/03/VOLUM-CONSILR-V-4-FINAL-REVIZUIT.PDF#PAGE=73](https://profs.info.uaic.ro/~consilr/wp-content/uploads/2021/03/volum-consilr-v-4-final-revizuit.pdf#page=73)

VASILE PĂIȘ, DAN TUFIȘ, AND RADU ION. **“INTEGRATION OF ROMANIAN NLP TOOLS INTO THE RELATE PLATFORM”**. IN:INTERNATIONAL CONFERENCE ON LINGUISTIC RESOURCES AND TOOLS FOR NATURAL LANGUAGE PROCESSING. 2019, PP. 181–192. [HTTPS://PROFS.INFO.UAIC.RO/~CONSILR/2019/WP-CONTENT/UPLOADS/2020/01/CONSILR2019_FINAL_BTT-60-EX-B5.PDF#PAGE=189](https://profs.info.uaic.ro/~consilr/2019/wp-content/uploads/2020/01/consilr2019_final_btt-60-ex-b5.pdf#page=189)

RELATE



ARCHITECTURA PLATFORMEI



RELATE



CARACTERISTICI PENTRU VOLUME MARI DE DATE

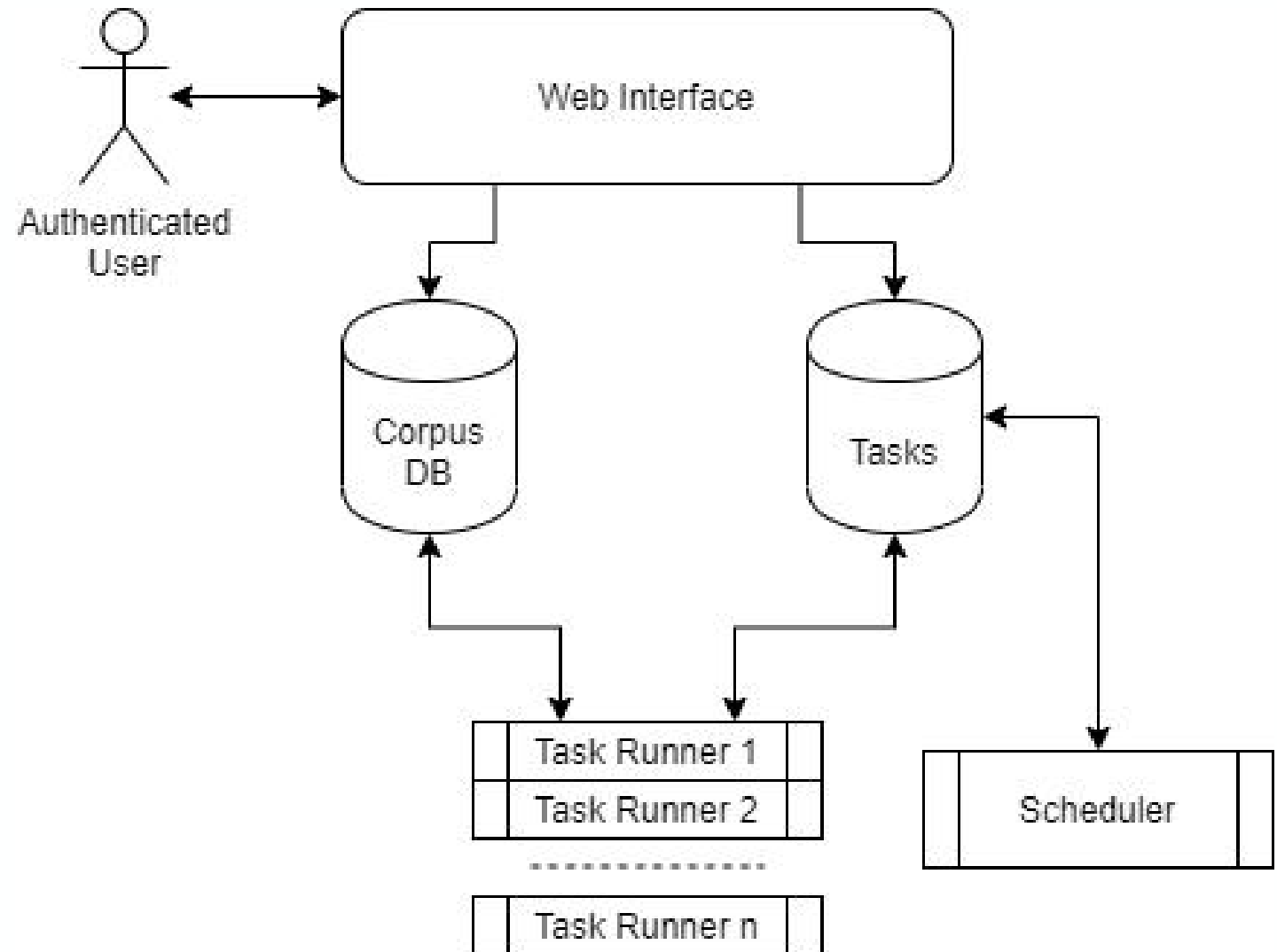
STOCARE ȘI PROCESARE PENTRU CORPORA DE DIMENSIUNI MARI

- TEXT
- TEXT + AUDIO

PROCESARE PARALELĂ:

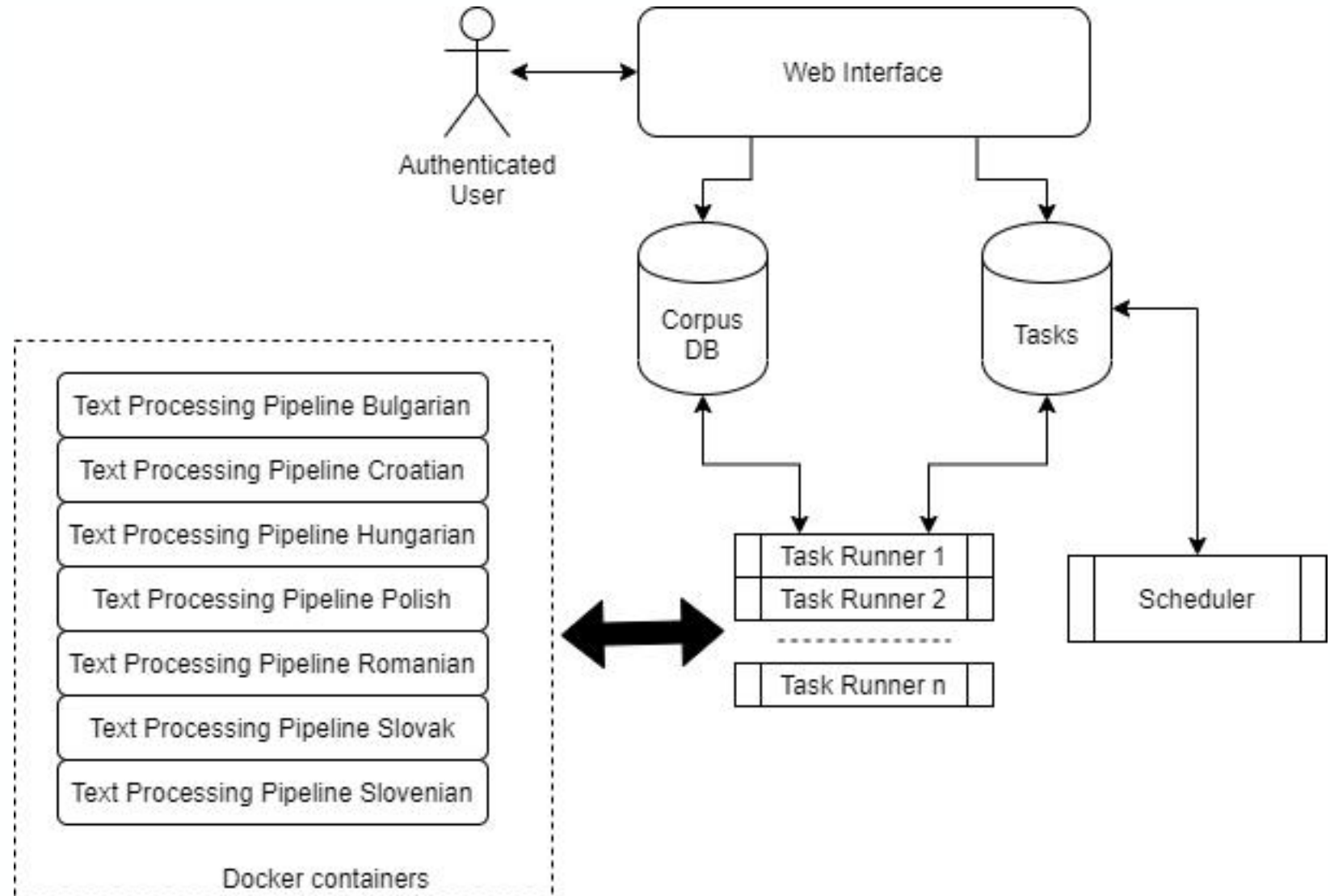
- SISTEM DE TASK-URI
- CONSUM DE SERVICII WEB EXTERNE PLATFORMEI
 - NODURI DE REȚEA DIFERITE
 - DIRECT DE LA PARTENERI

TASK MANAGEMENT



ADAPTABILITATE LA NOI LANȚURI DE PRELUCRARE

MARCELL SUSTAINABILITY PLATFORM



RELATE



CARATTERISTICI GENERALE

CORPUS MANAGEMENT: CREARE, ÎNCĂRCARE, DESCĂRCARE, ARHIVARE, ADNOTARE, STATISTICI, VIZUALIZARE

CREARE CORPUS "GOLD": INTEGREAȚĂ BRAT PENTRU NER, ÎNREGISTRATOR PENTRU CORPUS VOCE-TEXT

ADNOTARE: LEMMA, PART-OF-SPEECH, DEPENDENCY PARSING, SYLLABIFICATION, DIACRITICS RESTORATION, NER, BIONER, IATE, EUROVOC

TRADUCERE: TEXT-TEXT, VOCE-VOCE

VOCE: ASR, TTS, TRADUCERE

WORDNET: ROWORDNET, INTEROGARE ALINIATĂ ENWORDNET

CoRoLA: INTEGRARE KORAP, COMPONENTA AUDIO, WORD EMBEDDINGS

ANONIMIZARE: APLICAȚIE DIN PROIECTUL CURLICAT (ÎN DESFĂȘURARE)

LEGAL NER: PERSOANE, ORGANIZAȚII, LOCAȚII, TIMP, REFERINȚE LEGISLATIVE
- PE BAZA UNUI SUB-CORPUS EXTRAS DIN MARCELL

MODELE DE LIMBĂ: POT FI DESCĂRCATE DIN PLATFORMĂ

FORMATE STANDARD: CoNLL-U, CoNLL-U Plus, XML, JSON, RDF

ARHITECTURĂ MODULARĂ: PLATFORMA SE EXTINDE CU FIECARE NOU PROIECT

RELATE



VIZUALIZĂRI

VIZUALIZARE ADAPTATĂ CONTEXTULUI

- **TABEL** - CORPUS MANAGEMENT
- **ARBORE** - GRAF DEPENDENȚE
- **TEXT** - DOCUMENTE, ALTERNATIVĂ LA TABEL
- **ADNOTĂRI FRAGMENTE TEXT** - NER
- **BUTOANE, FORMULARE** - TASK-URI
- **PLAYER AUDIO** - FIȘIERE AUDIO
- **RECORDER AUDIO** - ÎNREGISTRĂRI, TRADUCERE VOCE

Complete run

JSON CoNLL-U CoNLL-X XML Text Chunks **Tree** Entities

Fiscul va face verificări la firmele indicate de CNSP, iar pe zona de dezvoltare va acorda granturi, precum cele pentru primării.

<

>

Word	acorda
<input type="text" value="Q"/>	
Lemma	acorda
<input type="text" value="Tree"/>	
U-POS	VERB
CTAG	VN
MSD	Vmnp
Chunk	Vp#3
Named Entity	
Phonetic	a.k.o.r.d.a
<input type="text" value="Play"/> <input type="text" value="Refresh"/>	
Syllables	a.cor.d'a
Similar Words	
Similar Lemma	primi solicita acordare beneficia acordat neacordare

Search in Korap

Search in Wordnet

Search in CoRoLa

Text to Speech

Word embeddings from CoRoLa

Word	acorda
	
Lemma	acorda
	
U-POS	VERB
CTAG	VN
MSD	Vmnp
Chunk	Vp#3
Named Entity	
Phonetic	a.k.o.r.d.a
 	
Syllables	a.cor.d'a
Similar Words	acordă acordat acordată primi beneficia acorde aloca acordau acordase acordam
Similar Lemma	primi solicita acordare beneficia acordat

Complete run

JSON

CoNLL-U

CoNLL-X

XML

Text

Chunks

Tree

Entities

1 ORG Fiscul va face verificări la firmele indicate de ORG CNSP , iar pe zona de dezvoltare va acorda granturi , precum cele pentru ORG primării .

2 DISO Diabetul zaharat este un DISO sindrom caracterizat prin valori crescute ale concentrației CHEM glucozei în sânge (DISO hiperglicemie) și dezechilibrarea metabolismului .

Corpus Marcell **Corpora**

- Files
- Tasks
- Basic Tags
- Statistics
- Archives

Files list

+ Add TEXT + Add CSV/TSV + Add ZIP TEXT

	Name	Type	Description	User	Creation Date
1	mj_00000G3W5B04K68KG8E2BJ...	text		pvf	2019-10-23 15:04:43
2	mj_00000G3W5A1UQW3MCZG1U...	text		pvf	2019-10-23 15:04:43
3	mj_00000G3W58159BQ4SSF1FC...	text		pvf	2019-10-23 15:04:43
4	mj_00000G3W547PDTJOU101U6U...	text		pvf	2019-10-23 15:04:43
5	mj_00000G3W544EFEG91D4145P...	text		pvf	2019-10-23 15:04:43
6	mj_00000G3W54086WCLHHF371...	text		pvf	2019-10-23 15:04:43
7	mj_00000G3W53QZCYSPF6P15G...	text		pvf	2019-10-23 15:04:43
8	mj_00000G3W51V24NCDHUZ367...	text		pvf	2019-10-23 15:04:43
9	mj_00000G3W4XPRW2QM1XT38...	text		pvf	2019-10-23 15:04:43
10	mj_00000G3W4XFXSAN4T0R39K...	text		pvf	2019-10-23 15:04:43
11	mj_00000G3W4WNUDI2YL0X04P...	text		pvf	2019-10-23 15:04:43
12	mj_00000G3W4UYQKLSX5JG3EH...	text		pvf	2019-10-23 15:04:43
13	mj_00000G3W4UU6GRLNRJG3E8...	text		pvf	2019-10-23 15:04:43
14	mj_00000G3W4SA4N7P0J0R0CF...	text		pvf	2019-10-23 15:04:43

Page 1 of 7207 20 1 to 20 of 144131

Back

Download

View as Text

File View

	ID	Form	Lemma	UPOS	XPOS	Feats	Head	Deprel	De
1	# sent_id = ro_legal.1								
2	# text = HOTĂRÂRE nr. 1.182 din 4 octombrie 2007								
3	1	HOTĂRÂRE	hotărâre	NOUN	Ncfsrn	Case=Nom Definite=Ind Gender=Fem Number=Sing	0	root	
4	2	nr.	nr.	NOUN	Yn	Abbr=Yes	1	nmod	
5	3	1.182	1.182	NUM	Mc	-	2	nummod	
6	4	din	din	ADP	Spsa	AdpType=Prep Case=Acc	6	case	
7	5	4	4	NUM	Mc	-	6	nummod	
8	6	octombrie	octombrie	NOUN	Ncms-n	Definite=Ind Gender=Masculine Number=Sing	2	nmod	
9	7	2007	2007	NUM	Mc	-	6	nummod	
10									
11	# sent_id = ro_legal.2								

Page 1 of 18

20

1 to 20 of 350

Corpus: **legalnero**

View Corpora

Files

Standoff

Tasks

Annotated

Statistics

Archives

Corpus tasks



- + TEPROLIN
- + UDPipe
- + IATE/EuroVoc
- + Classify EuroVoc
- + Cleanup
- + NER Baseline
- + CoNLLUP2BRAT
- + BRAT2CoNLLUP
- + Gold NE list
- + Statistics
- + ZIP Text
- + ZIP Annotated
- + ZIP Standoff
- + ZIP Gold Standoff
- + ZIP Gold Annotated
- + ZIP Audio
- + Export Marcell
- + Change Terms Marcell
- + Anonymization

	Type	Status	Description	User	Creation Date
1	statistics	DONE		maria@racai.ro	2021-05-24 19:14:35
2	unzip_text	DONE	Unzip TEXT from legalnero.zip	maria@racai.ro	2021-05-24 19:09:59

Speech-to-Speech Translation

File

Recording

Results

Selectați un fișier WAV asociat unui text în limba română. Este indicat să fie înregistrat cât mai clar.

Choose File No file chosen

Lanț de prelucrare:

Traducere

For this implementation we used the following:

- Romanian ASR from the ROBIN Project: Andrei-Marius Avram, Vasile Păiș, Dan Tufiș. 2020. Towards a Romanian end-to-end automatic speech recognition based on DeepSpeech2. Proc. Ro. Acad., Series A, Volume 21, No. 4, pp. 395-402.
- Romanian TTS from <http://romaniantts.com> : Adriana Stan, Junichi Yamagishi, Simon King, Matthew Aylett, "The Romanian speech synthesis (RSS) corpus: Building a high quality HMM-based speech synthesis system using a high sampling rate", In Speech Communication, vol. 53, no. 3, pp. 442-450, 2011.
- Romanian SSLA TTS: Tiberiu Boroș, Ștefan D. Dumitrescu, Vasile Păiș, "Tools and resources for Romanian text-to-speech and speech-to-text applications", CoRR, vol. abs/1802.05583, 2018. <https://arxiv.org/pdf/1802.05583.pdf>
- Translation using the EU Council Presidency Translator developed by Tilde with support from RACAI during the Romanian presidency.
- English DeepSpeech2 ASR from: Amodei et al. 2016. Deep Speech 2 : End-to-End Speech Recognition in English and Mandarin. In Proceedings of The 33rd International Conference on Machine Learning (PMLR), 48:173-182.
- English Mozilla DeepSpeech ASR: Hannun et al. 2016. Deep Speech: Scaling up end-to-end speech recognition. arXiv:1412.5567 [cs.CL]
- English Mozilla TTS: <https://github.com/mozilla/TTS>

RELATE



CORPORA DEZVOLTATE CU AJUTORUL PLATFORMEI

MARCELL-RO

DAN TUFÎȘ, MARIA MITROFAN, VASILE PĂIȘ, RADU ION, AND ANDREI COMAN. **“COLLECTION AND ANNOTATION OF THE ROMANIAN LEGAL CORPUS”**. IN: PROCEEDINGS OF THE 12TH LANGUAGE RESOURCES AND EVALUATION CONFERENCE. MARSEILLE, FRANCE: EUROPEAN LANGUAGE RESOURCES ASSOCIATION, 2020, PP. 2773–2777. URL: [HTTPS://WWW.ACLWEB.ORG/ANTHOLOGY/2020.LREC-1.337](https://www.aclweb.org/anthology/2020.lrec-1.337)

ADNOTARE

STATISTICI

EXPORT ÎN ARHIVE CU FORMAT DIFERIT (TEXT, CONLLU PLUS AND XML)

LEGALNERO

PĂIȘ, VASILE, MITROFAN, MARIA, GASAN, CAROL LUCA, IANOV, ALEXANDRU, GHIȚĂ, CORVIN, CONESCHI, VLAD SILVIU, & ONUȚ, ANDREI. (2021). ***ROMANIAN NAMED ENTITY RECOGNITION IN THE LEGAL DOMAIN (LEGALNERO)*** [DATA SET]. ZENODO. [HTTP://DOI.ORG/10.5281/ZENODO.4772094](http://doi.org/10.5281/zenodo.4772094)

ADNOTARE MANUALĂ NER

ADNOTARE AUTOMATĂ

STATISTICI

FORMAT RDF

LEGARE ENTITĂȚI GEONAMES

ROBIN TECHNICAL ACQUISITION SPEECH CORPUS

PĂIȘ, VASILE, ION, RADU, BARBU MITITELU, VERGINICA, IRIMIA, ELENA, MITROFAN, MARIA, & AVRAM, ANDREI.
(2021). ***ROBIN TECHNICAL ACQUISITION SPEECH CORPUS*** [DATA SET]. ZENODO.

[HTTP://DOI.ORG/10.5281/ZENODO.4626539](http://doi.org/10.5281/zenodo.4626539)

ÎNREGISTRARE

ADNOTARE AUTOMATĂ

STATISTICI

FOLOSIT PENTRU ADAPTAREA ASR-ULUI LA CERINȚELE PROIECTULUI

RELATE



DEZVOLTARE

OPEN SOURCE

[HTTPS://GITHUB.COM/RACAI-AI/RELATE](https://github.com/racai-ai/RELATE)

BAZAT PE COMPONENTE

The screenshot shows a web browser displaying the GitHub repository page for `racai-ai/RELATE`. The address bar shows the URL `github.com/racai-ai/RELATE/tree/master/src/modules/udpipe`. The repository name `racai-ai / RELATE` is visible, along with navigation options: `<> Code`, `! Issues`, `🔗 Pull requests`, `▶ Actions`, and `📁 Projects`. Below the repository name, the current branch is `master`. The file structure for the `RELATE / src / modules / udpiper /` directory is shown, including a pull request by `pvf2005 and vasilie Refactoring (#1)`, a `..` directory entry, and three files: `module.json`, `runner.php`, and `scheduler.php`.

READY? LET'S GIVE IT A TRY!



[HTTPS://RELATE.RACAI.RO](https://relate.racai.ro)

MULȚUMESC PENTRU ATENȚIE!

<https://relate.racai.ro>

