Comparing ML OCR Engines on

Texts from 19th Century Written in the Romanian Transitional Script

Marc Frîncu, Marius E. Penteliuc, Simina Frîncu, Gheorghe Bran, and Manuela Zănescu

Comparing ML OCR Engines on

Texts from 19th Century Written in the Romanian Transitional Script

Talking Points

Context



Proposal



Experiments



Results



Conclusions





Talking Points





Results



Conclusions



Latin >

Lorem ipsum dolor

RTS >

Порет іпсут дол

Cyrillic >

Порем ипсум доло

4500+

RTS Documents from 9 queried Libraries



a lot more...

Documents from 3642 Libraries, Church and Museum Collections



- 90 -

IV.

Inвidie a.*)

Пріп піста dіаволяляї а інтрат шоартеа для ляте пі чеї din партеа ляї для каятъ пре dжися з. Пицел. 2, 24.

Тот пъкател din лъвитря — din a inime товменль Се вргеще de де бань zica бенель тор; Ба е фъкет ди сінцімінте маї паінте de івенль, Маї паінте de че фанта лар веdi ла черкътор. — Кънд ъптей оамен войса фъкъторелей кълкаръ, Мъреціса маї паінте т дитінасе ке пъкат. Маї паінте de че ферел дий дитінде пеагра геаръ, De пъкател лъкоміст кенетет с'а ші стрікат. Маї паінте de че фанта перешінет се плінеще, Сефлетел прін реоа понтъ мънъ дес с'а ші стітт. Маї паінте de че Каін вчісесе тіръпеще,



Talking Points

Context



Experiments



Results



Conclusions





	Tesseract	Transkribus		
Architecture	RNN with LSTM units	VGG with BLSTM units		
Learning Type	Fine-tuning a Cyrillic Script model Russian model			
Learning Rate	0.0001			
Training Epochs	250, 50000	250		



Talking Points

Context



Proposal



Results



Conclusions





Test Set	Training Set	
Training Set		Test Set
Training Set	Test Set	Training Set

Experiments • Dataset

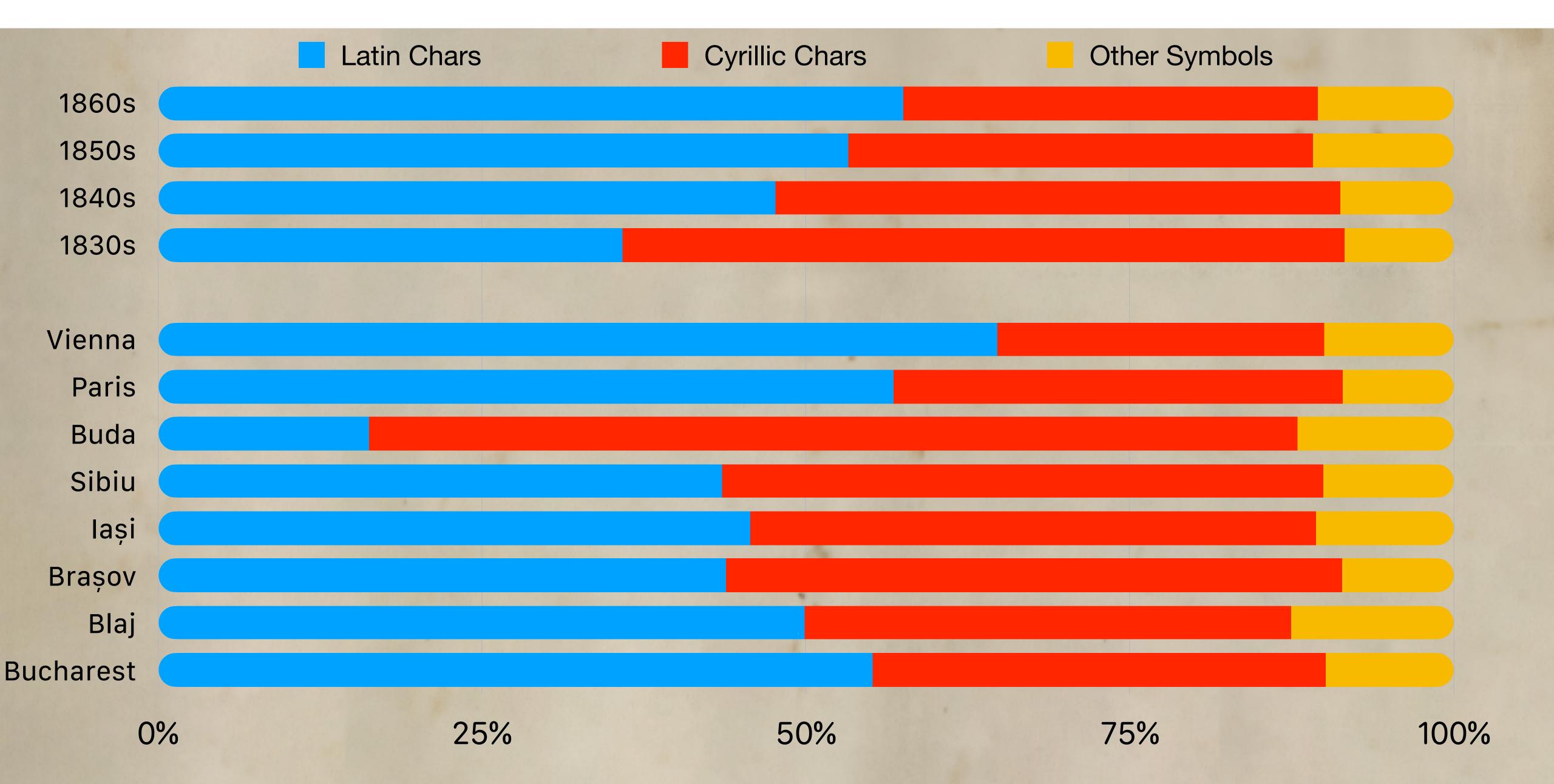
Average Dimension	$300 \times 900 \rightarrow 2000 \times 3000$ pixels			
Average Size	90kB → 10MB			
	29 790 words			
	153 286 characters			

Latin Chars Cyrillic Chars Overlapping Chars Other Symbols

 25837
 59819

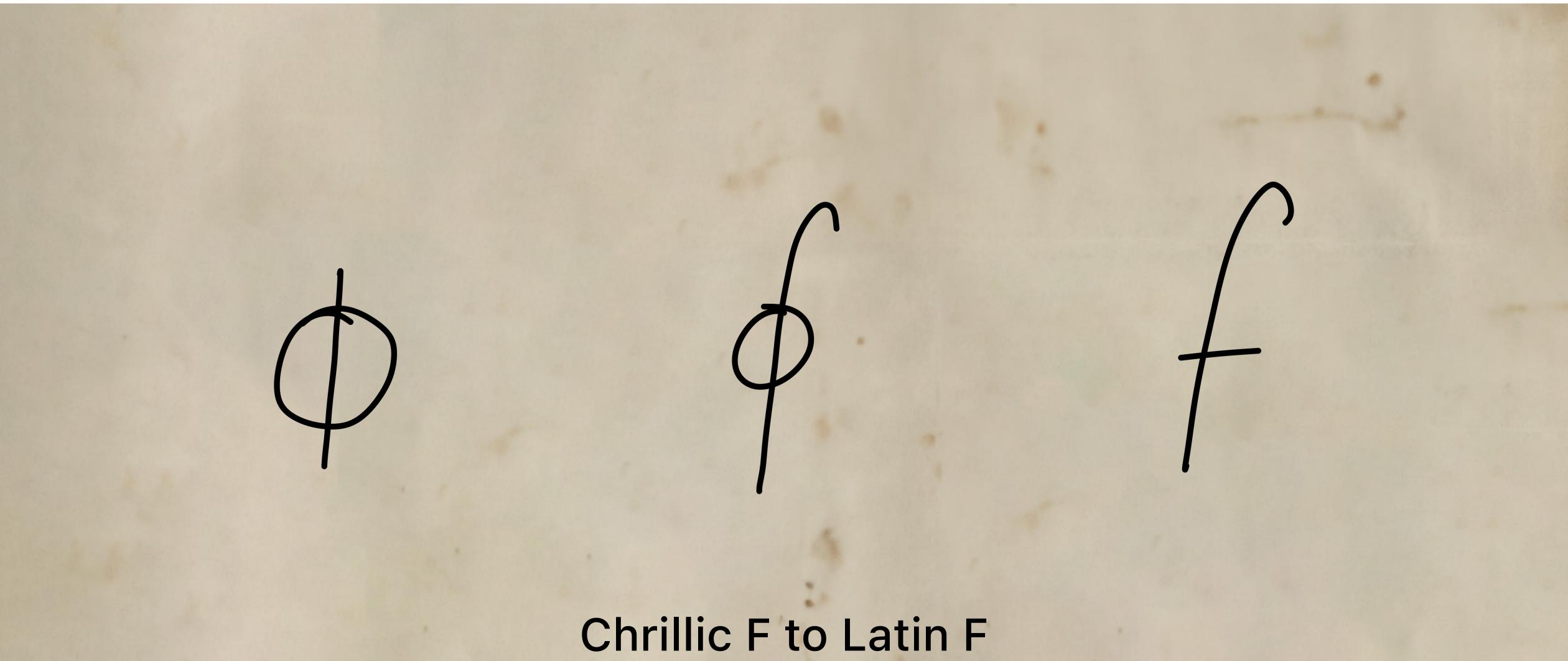
 52163
 15163

Experiments • Dataset





Experiments • Dataset





Experiments • Dataset

Chrillic J to Latin J



Experiments • Challenges

отсерикъ пре чет че пре скіть торії de е филії пре тоці і ай о: Лваці ачестеа de тей касъ de nersві ав дат прілеж ла втеле ляї, тотяші лор, пептря къ щіа а ав авыт ХС копре а dosa паштере пвъцъчей съй дл г аколо кв еї, ші въцъчей ляй Ioann Іпвъцъчеілор деї спре Xc, ка deптвіторівл летеї. в ачеста сав дпi Ioann Anainteкс, къ ав авгіт въцъчеї ape de-Isdea mi car dec врео прітеждіе, pipei care. He

Лицел. 2, 24. а іпітеї токшеаль ъптвітор; nainte de івеаль, і ла черкътор. горуляї кълкаръ, е кв пъкат. nrinde nearpa геаръ, са ші стрікат. пет се плинеще, пъ две с'а щі стіптіт. ста за пі фост къпіт. г, атъръчівпе преа ревчос е, стрікъчівне, че ё вып, жпалт, фримос. осълът дескріре, пе кввжитъ пре фрвтос: се десфатъ порочіе; е алтъл порокос." de въп пічт вреа съ щіе; ръще пе'ппъкат. der tee greens en 3/

emnaці літераці аї повл ачеста: inвіdiea. De фі doapъ maï віпе ал

tisnele lei Innouentis
omanilors mi al Bellors nostri. La annel
X katre Bela Regele

Bela întziš naskutul

n Episkopatal Komase nemeskě Romani, koteskš a fi kreutini, verse riteri wi datine ne. Kzui despretsinds misteriele biseriqii Episkopul Komaniloru dela niute Pserdolors, wi melti Unreri reks din Reratsl Uni, wi ama fakandsse etsesks pe Episkops,

ezisele misterii întro

PI DE PEPITIPE MAI CHEPIOACE! BOIM е ші съ лифіінφZΓZA8mm nAZ= NEZTOP AN ASMITS іжлоква кінврі-Z ПЕ ВАЛ8 РТ Т8 P-CE TPRHAZBETI-PHITOAPE 38 FPZж8р! Тот гра BERRATE PARE AE KINSA SEA MAI 38ME'E WI DE а, ма фачь са ПАЖЧЕРІЛОР ЧЕ COPIE MAHAPE ASI AZKAW 4E THE EPA NEHTP8 A CEMEH'S K8

ші, пеферічіткле шопц! de poame.. Bin' mai e, mi xaideni amundoi асочіаціе пентря ка съ евіповъціа Провіпціаліей фаче къ'ті ещі роб, в аръта песте тот лоn сълватік din Amepidea кът de sine мі фі пътеле de сълбатік, ъ ка ачеаста таї алес, чізілісаці сжитё сокоmi! " — Neamusл meš се е ка о броаскъ ла соаършіт аў пріітіт.. Dea-Domnsляї, па карце ві-'n пасат!. Ns'ї сат впьщігът къте 50 de обе път de гъпъ... Съ тръатапіствл!. Ел ші авdomnitopiĭ веакълъї!. акъ прін ажеторієл леї пре вп от Ancemnat?.. ь фак вр'яп Пріпц...de въхънд пе Клаїне Швабе). ооп фон Клаїне Швапосоторыту і! mommmm

CHENA III.

IAINE-ШВАБЕ (допръ фоарте посоторыт).

ai mai neemuvae d



Experiments • Scenarios

Scenarios	Groups	Tesseract	Transkribus
Region	8		
Decade	4		
Publisher	44		
Transformations	1		
OCR engine (control)	3		

× 3 folds = 180 models

ool Results

Talking Points

Context



Proposal

Experiments



Conclusions



$$CER = \frac{S + D + I}{N}$$

CER Character Error Rate

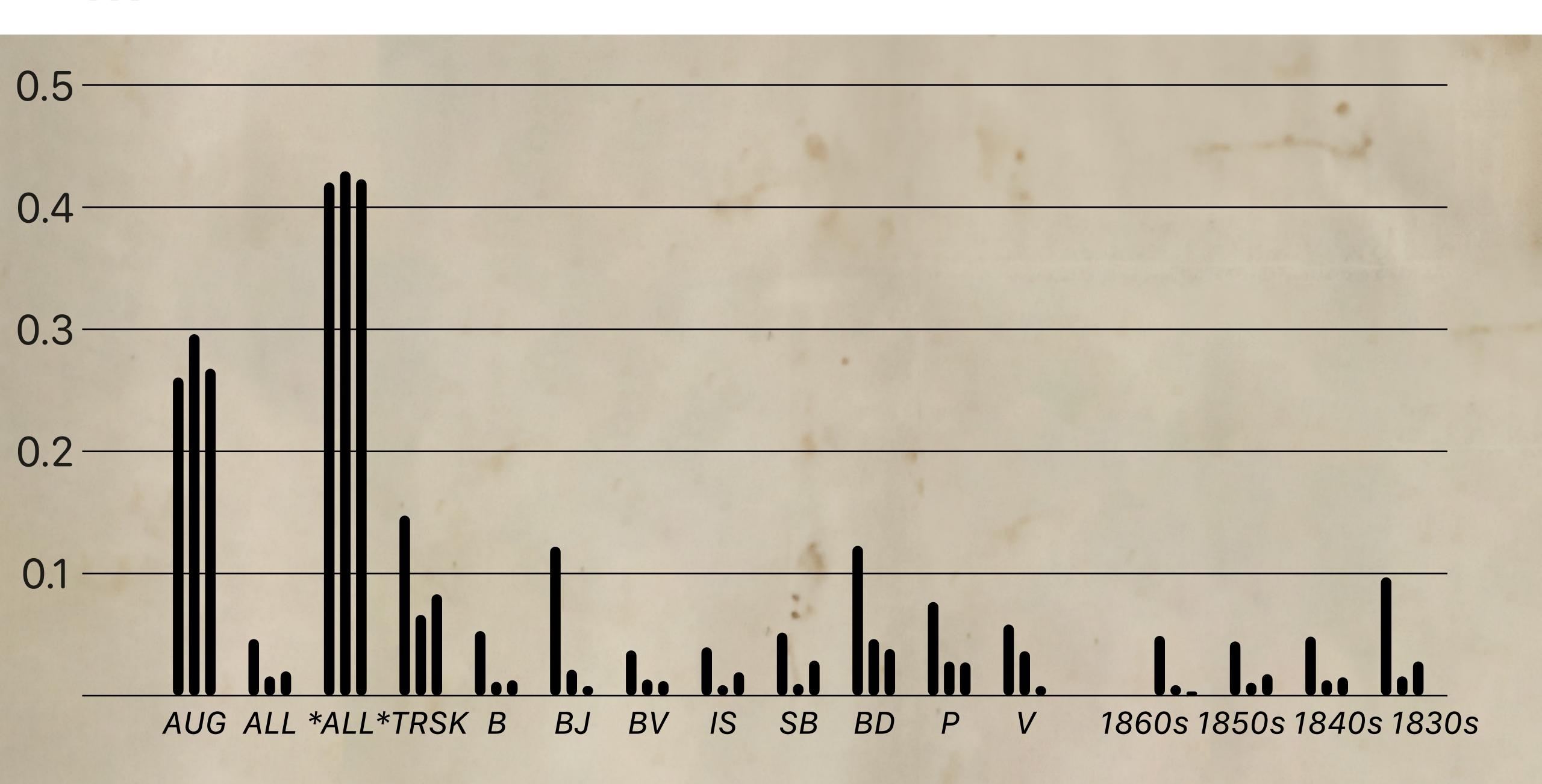
S Substitution

D Deletion

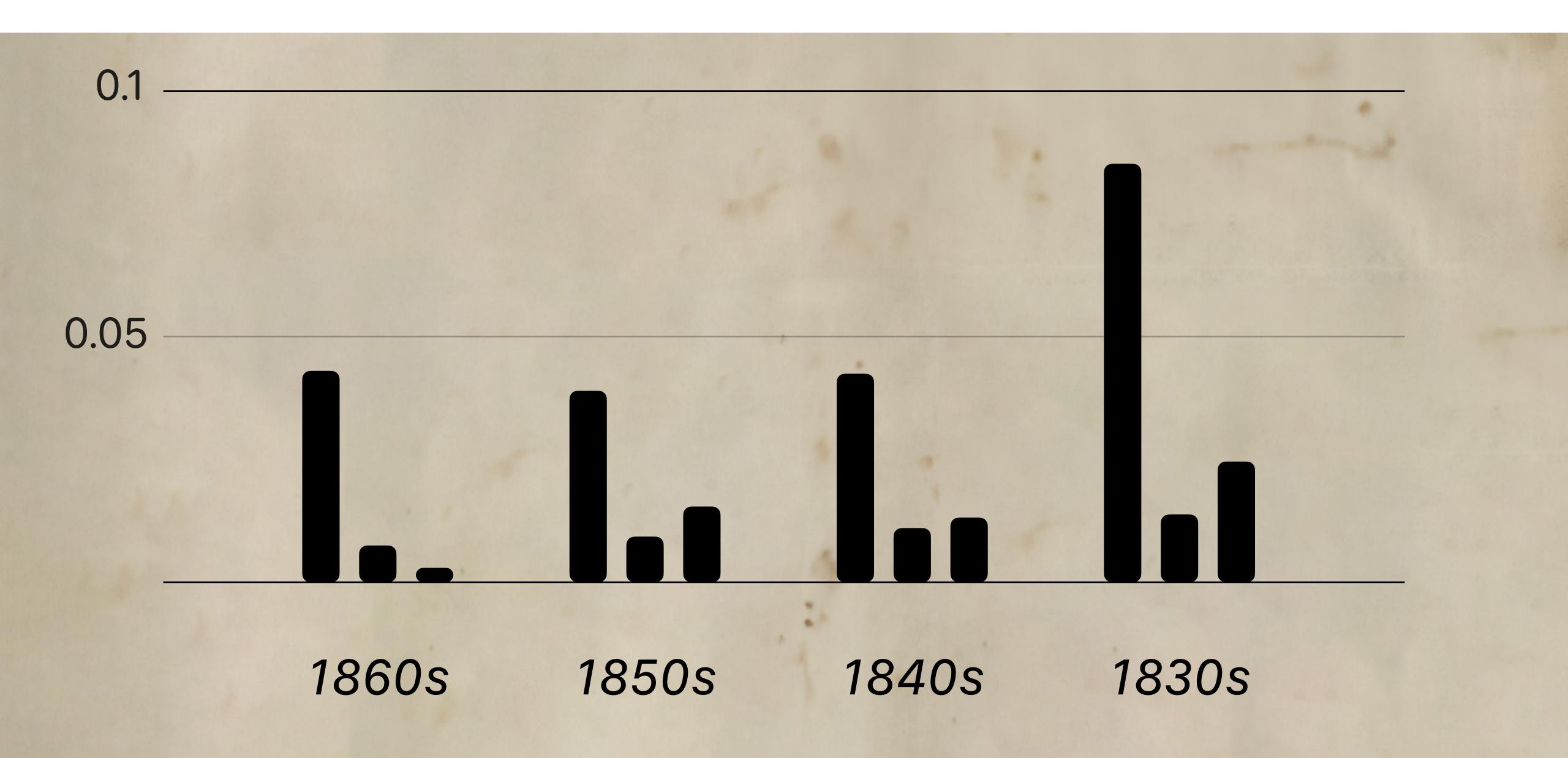
I Insertion

N Number of Characters (baseline)

Results • Performance



Results • Performance



Results • Comparison



Obest

100 worst

Results • Comparison



0 best 100 worst

Results • Comparison



0 best

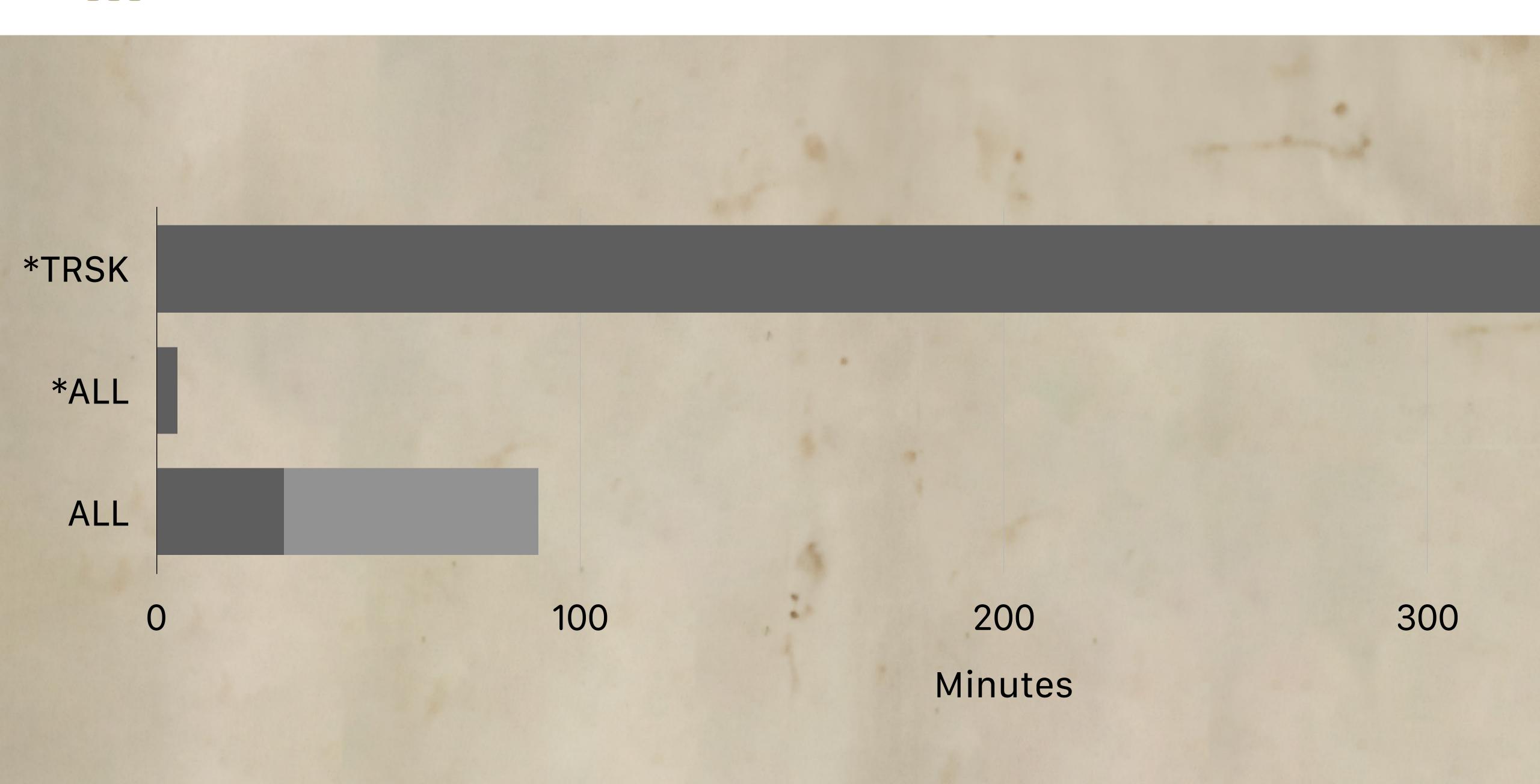
100 worst

		В	BJ	BV	IS	SB	BD	P	V
Bucha	rest	0.04	0.39	0.30	0.13	0.23	0.45	0.27	0.30
	Blaj	0.20	0.08	0.50	0.40	0.50	0.39	0.56	0.52
Bra	ĄSOV	0.04	0.36	0.04	0.03	0.05	0.31	0.20	0.32
	Iași	0.05	0.42	0.15	0.04	0.10	0.36	0.20	0.34
S	ibiu	0.05	0.38	0.10	0.04	0.04	0.29	0.20	0.33
Buda(p	est)	0.15	0.22	0.43	0.41	0.48	0.08	0.58	0.63
P	aris	0.05	0.43	0.25	0.06	0.16	0.47	0.06	0.23
Vi	ena	0.14	0.45	0.44	0.24	0.32	0.58	0.35	0.07

	В	BJ	BV	IS	SB	BD	P	V
Bucharest	0.04	0.39	0.30	0.13	0.23	0.45	0.27	0.30
Blaj	0.20	0.08	0.50	0.40	0.50	0.39	0.56	0.52
Brașov	0.04	0.36	0.04	0.03	0.05	0.31	0.20	0.32
Iași	0.05	0.42	0.15	0.04	0.10	0.36	0.20	0.34
Sibiu	0.05	0.38	0.10	0.04	0.04	0.29	0.20	0.33
Buda(pest)	0.15	0.22	0.43	0.41	0.48	0.08	0.58	0.63
Paris	0.05	0.43	0.25	0.06	0.16	0.47	0.06	0.23
Viena	0.14	0.45	0.44	0.24	0.32	0.58	0.35	0.07

		В	BJ	BV	IS	SB	BD	P	V
Bucha	rest	0.04	0.39	0.30	0.13	0.23	0.45	0.27	0.30
	Blaj	0.20	0.08	0.50	0.40	0.50	0.39	0.56	0.52
Bra	așov	0.04	0.36	0.04	0.03	0.05	0.31	0.20	0.32
	lași	0.05	0.42	0.15	0.04	0.10	0.36	0.20	0.34
S	Sibiu	0.05	0.38	0.10	0.04	0.04	0.29	0.20	0.33
Buda(p	est)	0.15	0.22	0.43	0.41	0.48	0.08	0.58	0.63
P	aris	0.05	0.43	0.25	0.06	0.16	0.47	0.06	0.23
Vi	iena	0.14	0.45	0.44	0.24	0.32	0.58	0.35	0.07

Results • Time





Talking Points

Context



Proposal



Experiments



Results



Conclusions

отсерікъ пре чеї че пре скіть торії de е филії пре тоці і ай о: Лваці ачестеа de тей касъ de nersві ав дат прілеж ла втеле дві, тотвші лор, пептря къ щіа а ав авыт ХС копре а dosa паштере пвъцъчей съй дл г аколо кв еї, ші въцъчей ляй Ioann Іпвъцъчеілор льї спре Xc, ка deптвіторівл летеї. ь ачеста сав лпi Ioann Anainteкс, къ ав aszit въцъчеї ape de-Isdea mi car dec врео прітеждіе, pipei care. He

Лицел. 2, 24. а іпітеї токшеаль ъптвітор; nainte de івеаль, і ла черкътор. горуляї кълкаръ, е кв пъкат. nrinde nearpa геаръ, са ші стрікат. пет се плинеще, пъ две с'а щі стіптіт. ста за пі фост къпіт. г, атъръчівпе, преа ревчос е, стрікъчівне, че ё вып, жпалт, фримос. освляї дескріре, пе кввжитъ пре фрвтос: се десфатъ порочіе; е алтъл порокос." de въп пічі вреа съ щіе; ръще пе'ппъкат. der tee greens en 3/

emnaці літераці аї повл ачеста: inвіdiea. De фі doapъ maї віпе ал tisnele lei Innouentis
omanilors mi al Bellors nostri. La annel
X katre Bela Regele

Bela întziš naskutul

n Episkopatal Komase nemeskě Romani, koteskš a fi kreutini, verse riteri wi datine ne. Kzui despretsinds misteriele biseriqii Episkopul Komaniloru dela niute Pserdolors, wi melti Unreri reks din Reratsl Uni, wi ama fakandsse etsesks pe Episkops, ezisele misterii întro Pi AE PEPITIPE MAI CHEPIOACE! BOIM е ші съ лифіінφZΓZA8mm nAZ= NEZTOP AN ASMITS іжлоква кінврі-Z ПЕ ВАЛ8 РТ Т8 P-CE TPRHAZBETI-PHITOAPE 38 FPZж8р! Тот гра BERRATE PAGE AE KINSA PER MAI 38ME'E WI DE а, ма фачь са ПЛЕЧЕРІЛОР ЧЕ COPIE MAHAPE ASI AZKAW 4E THE EPA NEHTP8 A CEMEH'S K8

ші, пеферічіткле шопц! de poame.. Bin' maï e, mi xaideni amundoi асочіаціе пентря ка съ евіповъціа Провіпціаліей фаче къ'ті ещі ров, в аръта песте тот лоn сълватік din Amepidea кът de sine жі фі пътеле de сълбатік, ъ ка ачеаста таї алес, чізілісаці сжитё сокоmi! " — Neamusл meš се е ка о вроаскъ ла соаършіт аў пріітіт.. Dea-Domnsляї, па карце ві-'n пасат!. Ns'ї сат впьщігът къте 50 de обе път de гъпъ... Съ тръатапіствл!. Ел ші авdomnitopiĭ веакълъї!. акъ прін ажеторівл леї пре вп от Ancemnat?.. ь фак вр'яп Пріпц...de въхънд пе Клаїне Швабе). ооп фон Клаїне Швапосоторыта'і!

CHENA III.

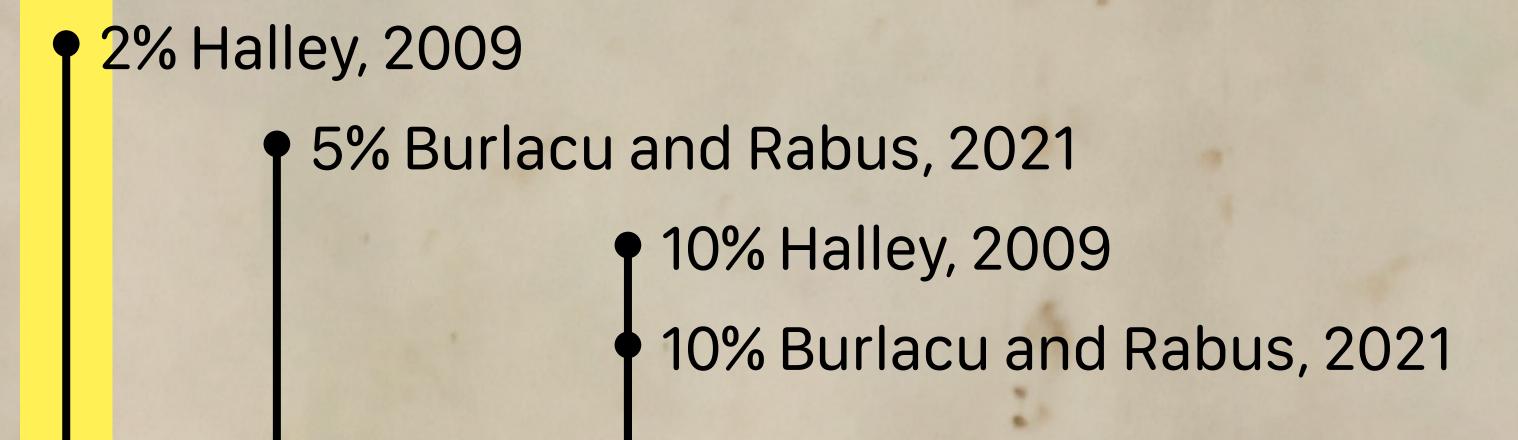
IAINE-ШВАБЕ (допръ фоарте посоторыт).

mommmm

i mai noomana J



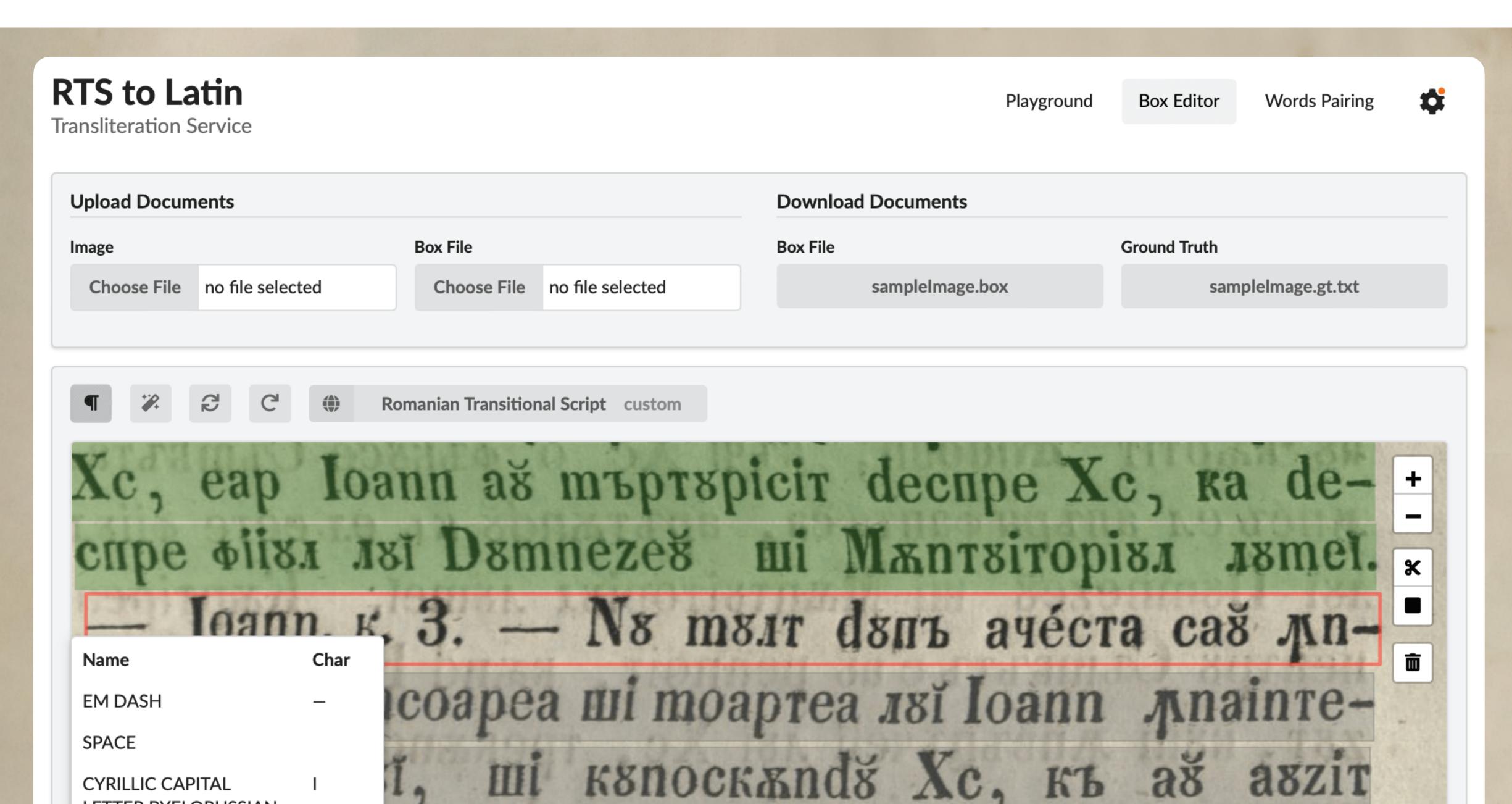
Our models



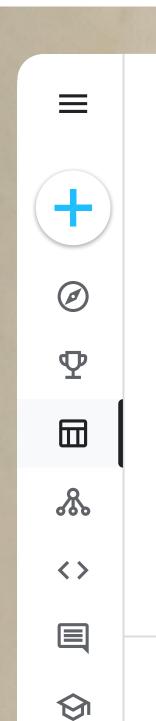
0% CER

30% CER

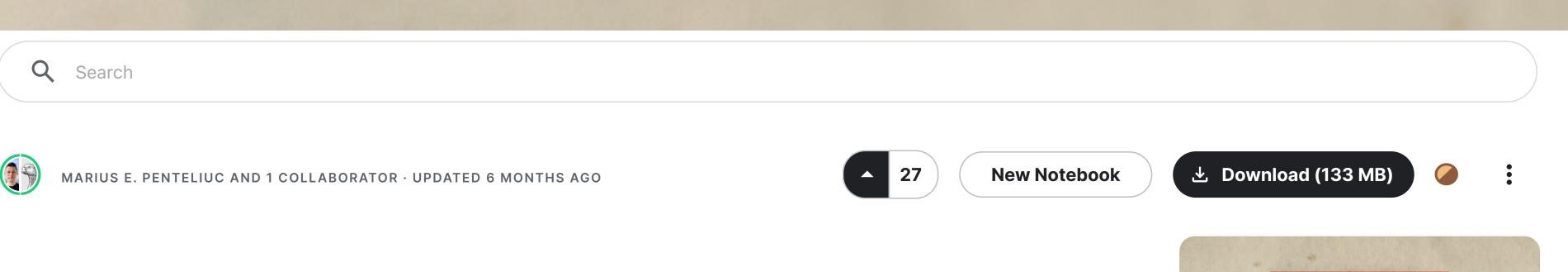








V



19th-Century Romanian Transitional Script

OCR Dataset for the Romanian Transitional Script of 19th Century



Data Card Code (1) Discussion (0) Settings

About Dataset



Usability (i)

10.00

This dataset consists of 147 pages of Romanian texts written in the **Romanian Transitional Script** (RTS). RTS is a mix of Latin and Cyrillic characters that were used in the 19th century in the Romanian provinces to facilitate the transition from the Romanian Cyrillic Script to the modern Latin Script. The images cover the period between 1833 and 1864. The selected texts cover a diverse range of literary genres, including poems, novels, dramas, stories, newspapers, and religious texts.

The dataset was obtained from the Central University Libraries (BCU) of Timişoara, Iaşi, and Cluj-Napoca through their free online platforms or by request. The scanned images are provided in JPEG and PNG formats, with dimensions ranging from approximately 300 by 900 pixels to 2000 by 3000 pixels. The file sizes vary between 70 KB and 10 MB.

To ensure **diversity**, the dataset includes images with *various fonts, styles, regions, publishers, and years*. It *covers all three main Romanian provinces' key publishing regions* (Bucharest - B, lasi - IS, Brasov - BV, Sibiu - SB, Blaj - BJ) as well as *some located outside Romania* that printed texts in RTS (Vienna - V, Budapest - BD, Paris - P). It comprises 4353 lines of text, totaling 29,790 words and 153,286 characters. Among these characters, there are 59,819 Cyrillic characters, 25,837 Latin characters, 52,163 overlapping characters (identical symbols), and 15163 other characters (e.g., punctuation, digits). The images below summarize its content per publisher and decade. More statistics (including

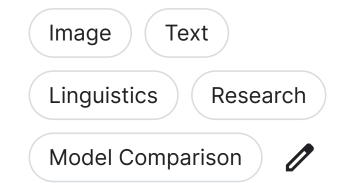
License

Other (specified in descripti... Edit

Expected update frequency

Quarterly Edit

Tags



Comparing ML OCR Engines on

Texts from 19th Century Written in the Romanian Transitional Script

Grazie

Comparing ML OCR Engines on

Texts from 19th Century Written in the Romanian Transitional Script

Backup Slide

RTS	Romanian Transitional Script						
RNN	Recurent Neural Network						
VGG	Very Deep Convolutional Networks						
B/LSTM	Bidirectional / Long Short-Term Memory						
OCR Flavors	Hand-writing, printed text, and case-specific						
Tagging	30 minutes using Custom Designed Box Editor						
Dataset Span	147 pages between 1833 — 1864						
Unique Chars	Not appearing in Unicode						
CER Score	Excluding spaces due to inconsistent use						



