

Developments on Text to Speech Synthesis in Romanian at the Speech Processing Group in Technical University of Cluj-Napoca

Prof. Mircea Giurgiu Speech Processing Research Group (UTCN) <u>http://speech.utcluj.ro</u>



ConsILR 2020

OUTLINE

- **1. Speech processing research laboratory**
- 2. Introduction on Text to Speech (TTS) synthesis
- 3. Natural Language Processing (NLP) tasks related to TTS
- 4. Acoustic modelling using deep neural networks (DNN) for high quality TTS, speaker adaptation, and expressive speech synthesis
- 5 Conclusions and future work





1. Speech processing research laboratory



1. Speech processsing research laboratory (http://speech.utcluj.ro)

Laboratory accreditation (UTCN) in 2003



FP7, Coordinated by University of Cambridge UTCN: ASR in noisy environments, RO-GRID



Simple4ALL (2011-2014)

FP7, Coordinated by University of Edinburgh, UTCN: Supervised and unsupervised tools for TTS



PCCA, Coordinated by UTCN UTCN: TTS assistive tools for laryngectomized patients

SWARA (2014-2017)

PCCDI, Coordinated by RACAI UTCN: Neural based TTS, Expressive TTS

SINTERO (2018-2020)



2. Introduction on Text to Speech Synthesis



2.1. The main components of a TTS system





6

2.2. Representation of the linguistic features



 \rightarrow Based on knowledge about spoken language

- Lexicon, letter-to-sound rules
- Tokenizer, tagger, parser
- Phonology rules





2.3. Representation of the acoustic features







2.4. Concatenative TTS







2.5. Statistical parametric TTS (HMM-based)

HMM: Handle variable length & alignment **Decision tree:** Map linguistic \rightarrow acoustic





$$\begin{split} \boldsymbol{o} \mid \boldsymbol{l}, \lambda) &= \sum_{\forall \boldsymbol{q}} \prod_{t=1}^{T} p(\boldsymbol{o}_t \mid q_t, \lambda) P(\boldsymbol{q} \mid \boldsymbol{l}, \lambda) \quad q_t \text{: hidden state at } t \\ &= \sum_{\forall \boldsymbol{q}} \prod_{t=1}^{T} \mathcal{N}(\boldsymbol{o}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) P(\boldsymbol{q} \mid \boldsymbol{l}, \lambda) \end{split}$$

Regression tree: linguistic features \rightarrow Stats. of acoustic features





2.6. Neural network – based TTS



$$h_t = g \left(\mathbf{W}_{hl} \mathbf{l}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{b}_h \right) \qquad \hat{\lambda} = \arg \min_{\lambda} \sum_t \| \mathbf{o}_t - \hat{\mathbf{o}}_t \|_2$$
$$\hat{\mathbf{o}}_t = \mathbf{W}_{oh} \mathbf{h}_t + \mathbf{b}_o \qquad \lambda = \{ \mathbf{W}_{hl}, \mathbf{W}_{hh}, \mathbf{W}_{oh}, \mathbf{b}_h, \mathbf{b}_o \}$$

FFNN: $\hat{o}_t \approx \mathbb{E}[o_t \mid l_t]$ RNN: $\hat{o}_t \approx \mathbb{E}[o_t \mid l_1, \dots, l_t]$

 $\hat{o}_t pprox \mathbb{E}\left[o_t \mid l_t
ight]
ightarrow \mathsf{Replace decision trees \& Gaussian distributions}$





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.1. Text normalization



Text Normalization – overview

 Input text may contain Non Standard Word (NSW) – examples. Taxonomy (Sproat, 1999). TN is not a trivial task.

Categorie	Codare	Specificație	Exemplu
ALPHA	EXPN	abreviere	Adv., dvs., Dvs., Str., CJ.
	LSEQ	secvență de litere	CD., PP., CIA
	ASWD	citit ca un cuvânt separat	NATO, Nume proprii
	MSPL	ortografiere greșită	Siteză, senal vocal
NUMBERS	NUM NORD NTEL NDIG NADDR NZIP NTIME NDATE NYEAR MONEY PRCT SPLT ROM	număr cardinal număr ordinal număr de telefon sau o parte din el numar ca digiți numar ca adresă stradă npd postal ora data anul nume de bani (1) nume de bani (2) procent secvente mixate sau separate litere romane	12, 45, $\frac{1}{2}$, 0.4, 15.23, 15,23 Al 25-lea, a 7-a zi 0742 551111, 2452 Camera 121, etajul 11 15 Bariţiu str., Bariţiu 26 400568 Cluj-Napoca 3.20, 03:20 2/2/99, 2/2/1999, 2/02/99, 2.02.1999 2015, '90, 12,45RON, €200,09, €200K 12,45 Mil RON, €200,09 mld 74%, 0,34% WS99, x2F, 2Fh, 3-a, II-a IV, IV-a, MMCX2L
MISC	SLNT	nu se pronunță	(…)
	PUNC	nu se pronuntă, sfârșit frază	., (punct)
	FNSP	interjecții	Oaaauuuuuu, ahhh,
	URL	adresa site web	http://com.com
	NONE	ignoră	Caractere nedorite





Text Normalization. Solution (1) – a hybrid approach^(*)



- (a) Rule-based (numerical strings) + (b) Lexicon (@, #, & ..., acronims).
- Cloud-based TTS demonstrator (//speech.utcluj.ro/swara/)

(*) Implemented in SWARA project: <u>http://speech.utcluj.ro/swara/</u>





Text Normalization. Solution (1) - results

Tip NSW	Intrare	leșire
Numere mai mici de 10^12	123	o sută douăzeci și trei
Numere mai mari de 10^12	750580558282384	șapte cinci zero cinci opt zero cinci cinci opt doi opt doi trei opt patru
Numere de telefon	0758123456	zero șapte cinci opt zero unu doi trei patru cinci șase
Numere cu semn	-123	minus o sută douăzeci și trei
Numere cu virgulă	123,5	o sută douăzeci și trei virgulă cinci
Ora ăn format HH:mm	Ședința este la ora 12:30	Ședința este la ora douăsprezece și treizeci de minute
Ora în format hh: <u>mm:ss</u>	11:35:40	Ora unsprezece treizeci și cinci de minute și patruzeci de secunde
Data in formatul <u>d.mon.y</u>	Noul cod fiscal va intra în vigoare la 1.ian.2016	Noul cod fiscal va intra în vigoare de la întâi ianuarie două mii șaisprezece.
Data în formatul d. <u>m.y</u>	Noul cod fiscal va intra în vigoare de la 01.01.2016	Noul cod fiscal va intra în vigoare de la întâi ianuarie două mii șaisprezece
Data în formatul d.mon	Noul an școlar începe în data de <u>15.sept</u> .	Noul an școlar începe în data de cinsprezece septembrie
Data în format <u>m.y</u>	09.1944	septembrie o mie nouă sute patruzeci și patru





Text Normalization. Solution (2) – Language independent number transcription using statistical machine translation (SMT)^(*)



- Text normalization is based on data only, instead of on expert rules.
- (1) a tokenizer, (2) a phrase-based translation, (3) a post-processing

^(*) R. San-Segundo, J.M. Montero, M. Giurgiu, I. Muresan, S. King, "Multilingual Number Transcription for Text-to-Speech Conversion", In Proc. of The 8th Speech Synthesis Workshop, Barcelona, September, 2013.





Text Normalization. Solution (2) – Language independent number transcription using statistical machine translation (SMT)(*)

Tokenization						
EN	BLEU	WER	SER			
1st alternative	97.9	1.6	7.4			
2nd alternative	98.5	0.8	6.8			
ES	BLEU	WER	SER			
1st alternative	97.8	1.9	6.8			
2nd alternative	98.2	0.9	6.1			
RO	BLEU	WER	SER			
1st alternative	98.5	0.9	5.4			
2nd alternative	99.2	0.5	4.6			

Aligment for training the translation model					
EN	BLEU	WER	SER		
grow-diag-final	98.5	0.8	6.8		
srctotgt	99.4	0.4	4.4		
tgttosrc	98.1	0.7	6.4		
ES	BLEU	WER	SER		
grow-diag-final	98.2	0.9	6.1		
srctotgt	98.2	0.9	6.1		
tgttosrc	98.5	0.7	4.5		
RO	BLEU	WER	SER		
grow-diag-final	99.2	0.5	4.6		
srctotgt	98.9	0.5	5.1		
tgttosrc	98.5	0.8	7.8		

Training Set Size						
EN	BLEU	WER	SER			
200 numbers	98.3	0.8	6.4			
400 numbers	99.3	0.4	4.6			
800 numbers	99.4	0.4	4.4			
4000 numbers	99.9	0.2	1.8			
ES	BLEU	WER	SER			
200 numbers	97.3	1.5	8.8			
400 numbers	98.2	0.9	6.0			
800 numbers	98.5	0.7	4.5			
4000 numbers	99.6	0.2	1.0			
RO	BLEU	WER	SER			
200 numbers	95.2	2.3	20.2			
400 numbers	97.7	1.3	11.2			
800 numbers	99.2	0.5	4.6			
4000 numbers	99.7	0.2	2.4			

- Use the SMT, without any rule or language specific interventions
- Small post-corrections (0.1%)
- Small training datasets (200 samples)
- Solution expanded for Twitter messages
- Tested on 3 languages (ES, EN, RO)

(*) "NORMA Toolkit" - http://simple4all.org/product/norma/index.html





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.2. Diacritics restoration





Diacritics Restoration. Overview for Romanian

- [Tufis, 1999] 75% may be <u>deterministically corrected</u>. DIAC system (93% bigram, 98,85% for 8 gram).
- [Burileanu, 2010] n-gram language model, multilevel prediction, 1 M words, (96,93% - 99,63%).
- [Boros, 2013] extends the DIAC with <u>more tokens</u>, errors proper names, editing, (96,89% - 99,31%)
- [Mihalcea, 2002] the method is Instance Base Learning , (96,14% 99,69%)
 - global rate of 98,17%. We started from this approach, incl. CART





Diacritics Restoration. (1) TIMBL – SWARA project^(*)

- Corpora: RomParl (2011-2014), RomLit (2006 2015), RomWiki (2014). RomParl has been used, about 420.000 instances.
- Window size (context of 5, 7, 9, 11 letters), Classification: TIMBL

	Vector (context 3 litere), diacritic						Cuvânt proveniență
İ,	n,	ν,	t,	a,	т,	ă	-înv ă țăm- (învățământului)
n,	ν,	a,	a,	т,	a,	ţ	-nvățămâ- (învățământului)
ν,	a,	t,	т,	a,	n,	ă	-vățămân- (învățământului)
t,	a,	т,	n,	t,	и,	â	-țământu- (învățământului)
t,	<sp></sp>	, C,	r,	и,	İ,	ă	-t c ă rui – (cărui)
и,	I,	t,	t,	İ,	I,	ă	-ult ă țil- (facultățile)
I,	t,	а,	İ,	I,	e,	ţ	-ltă ț ile- (facultățile)



^(*) *M. Giurgiu, A. Stan, Livrabil proiect: D1.2 Sistem preliminar de sinteză text vorbire (2016), Swara project..*





Diacritics Restoration. (1) TIMBL - SWARA project

Diacritic	Dimens	Vectori U / W (Unic / Within features)				Vectori U / W (Unic / Within features)			
	(context)	Nr. vectori antrenare	Entropie	Performanța [<u>%]</u>					
a – ă	5 + 1 + 5	66.000	0,79	94,93					
a - â	5 + 1 + 5	53.000	0,31	99,60					
a – ă - â	5 + 1 + 5	69.000	1,02	94,95					
i - î	5 + 1 + 5	65.000	0,39	99,48					
s - ș	5 + 1 + 5	25.000	0,72	98,20					
t - ț	5 + 1 + 5	42.000	0,69	98.44					
toate	5 + 1 + 5	99.900	2,62	97.63					

Tip eroare	Predicții incorecte
a prezis ca ă	gafă, că, pictă, <u>arăta,attilă</u> , listă, contră, urmă
ă prezis ca a	dusa, ramân, varga, doina, amarui
a prezis ca â	cândida, mânea, până, român, atânasiu, bârbu, ruxândra, jucât, câsnică, rugât
â prezis ca a	varf, dansii, franeze, carpi, stalpilor, coborat
i prezis ca î	înegal, întemperii, îndependent, învers
î prezis ca i	reintregi, indestulat, impreunare, inclin, intreg, indrazneț
s prezis ca ș	clașifica, aușter, școasă, prăpaștia, foloși
ș prezis ca s	sinele, defineste, usurință, desarte, sefii, sapte
t prezis ca ț	paciență, simțe, președințe, oponență, plățindu-l
ț prezis ca t	piata, știti, amenintați, ședintelor, fortelor, finante, prezentă





Diacritics Restoration. (2) CART vs TIMBL - SWARA project

	J 48	TIMBL		
Model pentru predicția	cel mai bun	(nU), cu diacritic inclus	(U), cu diacritic inclus	
perechii	[%]	[%]	[%]	
a-ă	96,04	95,27	92,73	
a - â	99,64	99,71	99,43	
i - î	99,84	99,50	99,29	
s - ș	98,95	98,10	97,20	
t - ț	98,75	98,34	97,09	
a – ă - â	95,10	95,06	92,77	
model global (toate)	98,15	97,71	96,25	





Diacritics Restoration. (3) LSTM & CNN – ReTeRom project^(*)

• CoRoLa, 1 M Tokens, 63.194 words, Input features: OHE, 4 NN structures.



(*) M. Nutu et al, "Deep Learning for Automatic Diacritics Restoration in Romanian", ICCP 2019





Diacritics Restoration. (3) LSTM & CNN – ReTeRom project

Architecture ID	Latent dimension	Batch size	Accuracy		
Architecture ID	Latent dimension	Datch size	3-gram	Word	Character
seq2seq_LSTM	128	512	75.50%	89.98%	71.61%
seq2seq_stacked_1_LSTM	256	128	79%	93 %	78%
seq2seq_stacked_2_LSTM	128	512	84%	94 %	82%
seq2seq_CNN	128	1024	91%	97 %	89%

Architecture ID	Accuracy				
Architecture ID	a-ă-â	i-î	s-ș	t-ț	
seq2seq_CNN	93.51 %	99.44 %	98.39 %	97.94 %	





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.3. Phonetic transcription



Phonetic transcription. Overview for Romanian

- Important task in TTS & ASR. Apart of the 9 phonetic rules in Romanian, there are a number of <u>exceptions</u> (diftongs, hiat).
- Syllabification and lexical stress may help the phonetic transcription.
- Previous work diversity: different text datasets / different classification methods

Paper	Level	Accuracy
Burileanu, 2002	word	98.30 %
Ordean et al, 2009	word	94.80 %
Toma et al, 2009	word	95.00 %
Toma et al, 2013	word	96.68 %
Boroș et al, 2012	word	93.00 %
Boroș et al, 2013	word	96.29 %
Cucu et al, 2014	word	97.24 %
Boroș et al, 2017	word	95.05 %
Domokos et al, 2011	phone	92.83 %
Toma et al, 2017	phone	99.61 %
Stan et al, 2018	phone	99.63 %





Phonetic transcription. (1) Decision trees – NaviRO^(*)

NaviRO – 138.500 words, phonetically transcribed. From the total set of 31 phonemes, prediction have been done for those depending on the context (14ph)

Lungime fereastră (nr. de litere)	Procent date de antrenare [%]	Acuratețe [%]
	25	92.48
4	50	92.38
1	100	92.35
	25	98.35
3	50	98.35
	100	98.37
	25	98.91
5	50	99.15
5	100	99.30
	25	98.72
7	50	99.05
1	100	99.28

^(*) Adriana Stan, Mircea Giurgiu, A Comparison Between Traditional Machine Learning Approaches And Deep Neural Networks For Text Processing In Romanian, in Proc. of the 13th International Conference on Linguistic Resources and Tools for Processing Romanian Language, 22-23 November, Jassy, Romania





Phonetic transcription. (2) DNN approach – s2s (LSTM)(*)

- MaRePhoR datasets: 72.375 words, 591.570 characters (Phonetic)
- DEX: 1.6 M words with their accent (Lexical stress)
- RoSyllabiDict: 507.000 syllabified words (Syllabification)
- → common words: 62.874 words



^(*) Adriana Stan, "Input Encoding for Sequence-to-Sequence Learning of Romanian Grapheme-to-Phoneme Conversion", In Proceedings of the 10th IEEE International Conference on Speech Technology and Human-Computer Dialogue (SpeD), Timisoara, Romania, 2019





Simultaneous predictions (Ph / Syll / Stress), (3) CNN/At, BLSTM^(*)



^(*) Beáta Lőrincz, "Concurrent phonetic transcription, lexical stress assignment and syllabification with deep neural networks", Proceedings of the 24th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems KES2020, 2020





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.4. Syllabification and accent





Syllabification. Results from SWARA project^(*)

- **RoSyllabiDict** reference sylabification dataset
- Simple scheme for feature encoding

Litera	Trasaturi	Clasa
î	<u>01</u> ****înţele	no
n	<u>00</u> ****înţeleg	yes
ţ	<u>00</u> ***înţelegi	no
е	<u>01</u> **înţelegi*	yes
1	<u>00</u> *înţelegi**	no
е	<u>11</u> înţelegi***	no
g	<u>00</u> nțelegi****	no
i	<u>0_1</u> ţelegī****	yes

Setul de antrenare	Clasificator	Acuratete de clasificare	Precizie
	Random Forest	99,46%	99,5%
0-44	SMO	96,27 %	96,3 %
Set 1	Naive Bayes	87,03 %	87,5 %
	Ada Boost	80,92 %	81,3 %
	Random Forest	99,00 %	99.0%
	SMO	96,04%	96,1 %
Set 2	Naive Bayes	87,09 %	87,6%
	Ada Boost	80,92 %	81,3%
	Random Forest	98,93 %	98,9%
Set 3	SMO	96,02 %	96,0%
	Naive Bayes	87,13 %	87,7%
	Ada Boost	80,92 %	81,3%
		00.000/	

^(*) Livrabil D1.15. Sistem de sinteza text vorbire preliminar (2016), Swara Project (Joint research with Prof. Potolea Rodica from the Computer Science Departament)





Syllabification & Lexical stress

Performanța de clasificare obtinuta de RF pentru poziționarea accentului

<u> </u>				
	Test	Nr.	Nr.	Acuratete
		instante	cuvinte	(la nivel litera)
	Set 1	58.434	13.210	96,37 %
	Set 2	58.626	13.209	96,37 %
	Set 3	58.207	13.210	96,26 %
	Set 4	58.293	13.211	96,58 %
	Set 5	58.288	13.210	93,53 %

Tabelul 2.4.7. Performanta de clasificare obtinuta de RF pe sarcina de pozitionare a accentului, dimensiuni variabile ale multimii de evaluare

Nr.	Nr.	Nr.	Acuratete	Cuvinte cu lipsa
Test	cuvinte	instante		predictie accent
1	1.100	4.860	96,04 %	143
2	2.200	9.720	95,72 %	318
3	3.300	14 410	96,12 %	389

^(*) Livrabil D1.15. Sistem de sinteza text vorbire preliminar (2016), Swara Project (Joint research with Prof. Potolea Rodica from the Computer Science Departament)





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.5. Part of speech tagging



Part of Speech tagging (POS). Overview of results, RO

- POS is important for expressive TTS
- POS tags: MSD, C-tagset, root POS (first letter of MSD tag)

Authors	Method	Accuracy	Tagset
Tufis & Mason [2]	Probabilistic	98.39%	MSD
Boros & Dumitrescu [3]	Deep Neural Networks	98.19%	MSD
Simionescu [6]	Probabilistic & Rule-based	97.03%	MSD
Teodorescu et al. [7]	Probabilistic	96.12%	Root POS
Frunza et al. [8]	Machine Learning	95.30%	Root POS





Part of Speech tagging. Results from SWARA project^(*)



NAACL 2003 English – Romanian dataset

^(*) Toth Grigore, "Part of speech tagging for Romanian", Diploma project, 2016 Swara Project, D1.15. Sistem preliminar de sintza text vorbire, 2016

/N /V /I /P /R /C



•

/J /Q /T /D /M /A /U /H /X



Part of Speech tagging. Results from SWARA project^(*)

Tool\Training Corpus (in lines)	1.000 linii	5.000 linii	10.000 linii	20.000 linii	30.000 linii	40.000 linii
Kytea - SVM	76.41%	86.48%	89.30%	91.57%	93.20%	97.14%
Kytea - LR	76.43%	86.55%	89.27%	91.67%	93.25%	97.08%
HMM Tagger	68.99%	84.69%	87.84%	90.00%	92.09%	95.29%
NLTK - TnT	76.23%	86.47%	89.43%	91.54%	93.07%	96.66%

^(*) Toth Grigore, "Part of speech tagging for Romanian", Diploma project, 2016 Swara Project, D1.15. Sistem preliminar de sintza text vorbire, 2016





Part of Speech tagging. Results from ReTeRom project (DNN)^(*)

- Datasets: a) CoRoLa (180Kw), b) DEX (1.9Mw), c) WPT (1.9Mw). •
- Individual words, no context .
- Encoding: OHE + Letter Embedding (Gensim) ۰
- (1) LSTM with dense layers, (2) LSTM sequence to sequence (80%-20%) •

System ID	Dataset	Tag	Network type	Character encoding	Latent dimension	Batch size	Epochs	Accuracy
1	WPT	RPOS	LSTM + Dense (*)	OHE	256	512	50	99.18%
2	WPT	RPOS	LSTM + Dense	OHE	256	512	50	94.85%
3	WPT	RPOS	LSTM + Dense	LE	256	256	25	54.80%
4	WPT	RPOS	seq2seq LSTM	LE	256	256	25	94.99%
5	WPT	RPOS	seq2seq + Embedding layer	OHE	256	256	20	93.88%
6	WPT	MSD	seq2seq LSTM (*)	OHE	512	1024	50	98.25%
7	WPT	MSD	seq2seq LSTM	OHE	512	1024	50	75.28%
8	WPT	MSD	seq2seq + Embedding layer	OHE	256	512	50	76.62%
9	DEX	RPOS	LSTM + Dense	OHE	256	512	50	94%
10	CoRoLa	CTAG	seq2seq LSTM	OHE	256	512	100	97.15%

(*) Beáta Lőrincz, Maria Nutu, Adriana Stan, "Romanian Part of Speech Tagging using LSTM Networks", In Proceedings of the IEEE 15th International Conference on Intelligent Computer Communication and Processing, Clui-Napoca, Romania, 2019,





3. Developments on Natural Language Processing (NLP) tasks related to TTS

3.6. Text style identification





Text Style Identification – ReTeRom (2018)

- Important task to control the expressivity in TTS
- Few research for Romanian (Univ. Bucharest, Romanian Academy).
- We have proposed 2 methods: LDA, CNN classifiers
- Dataset: CoRola
- <u>Text styles:</u> (1) Beletrictic: 251.927 words, (2) Judicial: 337.143, (3) Memorialistic: 315.213, (4) Political: 374.680, (5) Jurnalistic: 325.191





Text Style Identification. LDA (Latent Dirichlet Allocation) (*)



Drawback – the method gives poor result for <u>short texts</u> (sentences, or even paragraphs), so <u>longer context is required</u> to identify the style.

(*) Nicoleta Raiu, "Text style identification for Romanian", Diploma project (2019)





Text Style Identification. CNN approach^(*)

- Inspired by several approaches that use CNN (Kim, 2014)
- (+) Even with a small training data set of 3.000 sentences, > 92% accuracywork
- Future work: status for larger contexts, to explore the integration into a TTS



(*) Adriana Stan, ReTeRom project (2019)





Speech Styles analyses (*)

 Preliminary studies on: Neutral (SWARA), Journalistic (TV news), Narativ (audiobook) – phoneme duration, pitch



(*) Adriana Stan, ReTeRom project (2018)





Adaptation of a statistical TTS for new speaking styles using CSMAPLR (Constrained Maximum Aposteriory Linear Regression) ^(*)



^(*) J. Yamagishi, et al., "Analysis of Speaker Adaptation Algorithms for HMM-Based Speech Synthesis and a Constrained SMAPLR Adaptation Algorithm", IEEE trans on Audio Speech and Language Processing, Feb. 2009.





New speaking styles with CSMAPLR (2018, ReTeRom)

	Neutru	Jurnalistic	Narativ
Voce naturala			
Voce sintetizata			
Adaptare durata			
Adaptare F0			
Adaptare durata si F0			

https://speech.utcluj.ro/sintero/prosody_examples/





4. Acoustic modelling using DNN

4.1. End to end TTS and new voices using Tacotron GST





End to end TTS with Tacotron (*)

- Original Tacotron: only the pair <text, audio> is necessary to train the model
 - Truly end-to-end TTS pipeline by Google
 - Reduces feature engineering
 - Allows conditioning on various attributes
 - Architecture
 - One network based on the sequence-tosequence with attention paradigm
 - Red Encoder
 - Blue Decoder
 - Green Post-processing net



(*) Yuxuan Wang, et al. "Tacotron: Towards End-to-End Speech Synthesis", Interspeech, 2017





Tacotron GST (*). Result on using Mara Dataset (about 11hr speech)

GST – a bank of embeddings jointly trained with Tacotron. In synthesis stage they
can control the speed and speaking styles, independent of text. No explicit text
labels, yet increased expressiveness.



(*) Yuxuan Wang, "Style Tokens: Unsupervised Style Modeling, Control and Transfer in End-to-End Speech Synthesis" (2018)





Adaptation of TTS system to new speakers

- (1) Train Tacotron 2 with a very large dataset: (a) single speaker: Mara, ~11hr; (b) multiple speakers: Swara1, ~21hr. Result → good modelling of speaker(s), hardly to obtain with small datasets
- (2) Adapt the pretrained model(s) to new speaker(s) with a small dataset (10-50min)



Adaptation of TTS system to new speakers – sample results

SPK5

Tok3

https://speech.utcluj.ro/pf_lrec2020/

Trained with Mara, adapted to Swara (10 speakers / 500 sentences) – includes vectorial representation of speakers

Trained with Mara, adapted to Swara (10 speakers / 500 sentences) – does not • include vectorial representation of speakers

Trained with Mara, adapted to Swara (10 speakers / 500 sentences) - freezing • the GST layer Tok3 Tok7



SPK0

Tok₂



ConsILR 2020 National Conference, Bucharest



SPK7

Tacotron – based new voices

Text: "Tezele ar putea fi eliminate în acest an școlar, iar Ministerul Educației analizează mai multe variante."



Note: Only <u>ortographic transcription</u> and <u>audio</u> are used for training. No other NLP annotation is added. (<u>https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/Sisteme-</u>





Speech prosody transfer using Global Style Tokens

- Tacotron GST pretrained expressive model (8.000 utt, Mara)
- Task1: Single speaker prosody transfer (500 utt / speaker)
- Task2: Multiple speakers (10) prosody transfer (500 utt / speaker x 10)
- Prosody is controlled: (a) reference utterance, (b) manual setting of GST weights
- Multispeaker \rightarrow <u>Mara, adaptation, Swara</u>.

System	Nat	ural	Syntetic				
	Mara	neutral	Eigen voice	Imtermediary	D	Different GSTs	
MARA	20102				2010		Sala.
MARA adapted to IPS					2010		
MARA adapted to EME	NIN.				100 M	~n, n, p, p, p, n, n, p, p, n, n, p, n, n, p, n, n, p, n, n, n, n, n, n, n, n, n, n, n, n, n,	N ^I n,





4. Acoustic modelling using DNN

4.2. Multispeaker TTS using DCTTS implementation





End to end TTS with DC-TTS, Deep Convolutional TTS. (*)

An alternative to Tacotron (recurrent layers – time and computing power, non-۲ parallelizable) \rightarrow proposal of a CNN TTS (DC-TTS).



(*) Hideyuki Tachibana,, "EFFICIENTLY TRAINABLE TEXT-TO-SPEECH SYSTEM BASED ON DEEP CONVOLUTIONAL NETWORKS WITH GUIDED ATTENTION", ICASSP 2018.





End to end TTS with DCTTS. Experimental setting

- Audio for training: Mara / Some speakers from Swara
- **Text** <u>automatically annotated with existing tools</u> (some errors may appear), in the following scenarios:

Only text	E însă altceva la mijloc.
Phonetic transcription	e <> a@ n s @ <> alt che v a <> la <> mizh lo k <.>
Phonetic+syllabification	e <> a@ n * s @ <> a l t * ch e * v a <> l a <> m i zh *bl o k <.>
Phonetic+accent	e0 <> a@1 n s @0 <> a1 I t ch e0 v a0 <> I a0 <> m i0 zh I o0 k <.>
Phonetic+accent+sylab	e0 <> a@1 n * s @0 <> a1 l t * ch e0 * v a0 <> l a0 <> m i0 zh * l o0 k <.>





End to end TTS with DCTTS. Analysis of the text embeddings





End to end TTS with DCTTS. Analysis of F0 trajectories



<u>https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/Sistem-DC-TTS-Date-lexicale</u>





End to end TTS with DCTTS. Speech samples

- Reference speech:
- Text to be synthetized:
- "Ce se va întâmpla cu aceasta după ce vom ieși din starea de urgență".
- Only ortographic transcription:
- Only phonetic transcription:
- Phonetic transcription + accent:
- *Phonetic transcription* + *syllabification*:
- Phonetic transcription + accent + syllabification







Speaker adaptation using DCTTS and speaker embeddings

- Speaker adaptation is realized using a <u>speaker dependent cost function</u> (Cosine Similarity (CS), respective EER)
- Systems: Baseline (B), B + CS, B + EER

Sistem	ALL (18sp/21.302)	RND1 (18sp/8.900)	RND1-100 (18sp/1.787)	RND1-SAM (1sp/500)
В	6.94	4.86	8.33	2.43
B+CS	6.25	4.66	6.25	2.43
B+E	4.66	8	6	2.43





Speaker embeddings (t-SNE) natural vs TTS for the Baseline (B)







Speaker embeddings natural vs TTS for the B + EER







Speaker adaptation using DCTTS. Samples

Speaker DDM (Swara)

ALL DATA (1500 utt/spk)

ONLY 100 utt / spk

Baseline

Baseline + CS

Baseline + EER



"La automatele de vândut cartele se poate achita cu numerar."



"De aceea, spune medicul, autoritățile ar trebui să schimbe tonul discuțiilor."





New synthetic voices generated from imperfect data

Imperfect data: (a) recorded speech does not match perfectly with the text, (b) pronunciation errors, (c) possible additional spoken artefacts, (d) noisy speech System: Tacotron 2



"Tezele ar putea fi eliminate în acest an școlar, iar Ministerul Educației analizează mai multe variante."





5. Conclusions and future work





Conclusions and future work

General aspects...

... discussions with the audience

In particular:

... discussions with the audience ... ?



Many thanks to the research TEAM

SI.dr.ing. Adriana Stan drd. Beata Lorincz drd. Maria Nutu

Previous collaborators: Andrei Barbos, Ioana Muresan, Andrei Homodi, Dalia Popescu, Cristian Contan, Jozef Domokos, Mihai Ordean.



THANK YOU

FOR YOUR ATTENTION !



