



## D1.17. Analiza metodelor de control și adaptare automată a expresivității vorbitorilor în sistemele de sinteză text-vorbire

Aceste rezultate au fost obținute prin finanțare în cadrul Programului PN-III Proiecte complexe realizate în consorții CDI, derulat cu sprijinul MEN – UEFISCDI, PN-III-P1-1.2.-PCCDI, nr. 73 PCCDI/2018:

**“SINTERO: Tehnologii de realizare a interfețelor om-mașină pentru sinteza text-vorbire cu expresivitate”**

© 2018-2020 – SINTERO

Acest document este proprietatea organizațiilor participante în proiect și nu poate fi reprodus, distribuit sau diseminat către terți, fără acordul prealabil al autorilor.

Denumirea organizației participante în proiect	Acronim organizație	Tip organizație	Rolul organizației în proiect (Coordonator/partener)
<b>Institutul de Cercetări Pentru Inteligență Artificială “Mihai Drăgănescu”</b>	ICIA	UNI	CO
<b>Universitatea Tehnică din Cluj-Napoca</b>	UTCN	UNI	P1
<b>Universitatea Politehnică din București</b>	UPB	UNI	P2
<b>Universitatea "Alexandru Ioan Cuza" din Iași</b>	UAIC	UNI	P3



### Date de identificare proiect

Număr contract:	PN-III-P1-1.2.-PCCDI, nr. 73 PCCDI/2018
Acronim / titlu:	<b>SINTERO: Tehnologii de realizare a interfețelor om-mașină pentru sinteza text-vorbire cu expresivitate”</b>
Titlu livrabil:	<b>D1.17 Analiza metodelor de control și adaptare automată a expresivității vorbitorilor în sistemele de sinteză text-vorbire</b>
Termen:	<b>Mai 2018</b>
Editor:	<b>Adriana Stan (Universitatea Tehnică din Cluj-Napoca)</b>
Adresa de email editor:	<b>Adriana.Stan@com.utcluj.ro</b>
Autori, in ordine alfabetică:	<b>Mircea Giurgiu, Adriana Stan</b>
Ofițer de proiect:	<b>Cristian STROE</b>

### Rezumat

Acest document prezintă o analiză a metodelor de adaptare și control automat a expresivității în cadrul sistemelor de sinteză text-vorbire. Vor fi prezentate metodele principale disponibile în cadrul sistemelor de sinteză concatenative, a celor statistice bazate pe modele Markov și a sistemelor bazate pe rețele neuronale multistrat.

**Cuprins**

1. Introducere .....	4
2. Controlul expresivității în cadrul sistemelor concatenative .....	4
3. Controlul expresivității în cadrul sistemelor probabilistice bazate pe modele Markov.....	5
4. Controlul expresivității în sisteme de sinteză bazate pe rețele neuronale multistrat .....	6
5. Concluzii .....	7
6. Bibliografie .....	8

## 1. Introducere

Sistemele de sinteză text-vorbire au atins deja un nivel al calității vocii apropiat de cel al vocii naturale. Rămâne însă problema variabilității și a expresivității acestor voci. Această problemă este de actualitate și datorită faptului că expresivitatea sau prozodia nu poate fi evaluată în mod obiectiv printr-un set limitat de parametri și de cele mai multe ori depinde de starea emoțională a persoanei care o evaluează, precum și de fondul cultural, etnic, educațional șamd.

Diferitele tipuri de sisteme de sinteză: concatenative, parametric-statistice sau cele ce modelează direct forma de undă, permit un control al expresivității și prozodiei specific, în funcție de arhitectura sistemului. Secțiunile următoare vor prezenta succint metodele principale de control a expresivității și prozodiei acestor sisteme.

## 2. Controlul expresivității în cadrul sistemelor concatenative

Sistemele de sinteză concatenative au fost până de curând cele mai utilizate sisteme de sinteză pentru aplicațiile comerciale. Principiul de bază al acestora îl reprezintă înregistrarea unui corpus de voce extins, de înaltă calitate, iar apoi concatenarea unor segmente audio pentru a reda informația textuală furnizată la intrare. Aceste segmente au lungimi variabile și pornesc de la nivel de fonem sau silabă și pot ajunge până la nivel de sintagmă.

Problema acestor tipuri de sisteme este faptul că informația audio nu este parametrizată sub nicio formă, astfel că pentru controlul expresivității este necesară manipularea formei de undă. La modul cel mai simplu, controlul expresivității putea fi făcut prin selectarea segmentelor audio de concatenat pe baza unei traiectorii prozodice predefinite sau estimate din text [Hunt96]. Această metodă însă se limitează la utilizarea unor evenimente prozodice existente în corpusul de voce, fără a putea extinde aceste evenimente în funcție de variabilitatea textului de intrare. Utilizarea unor corpusuri audio ce conțin seturi de date cu emoții diferite a fost implementată de [Iida03] și [Johnson02], iar rezultatele au fost considerate satisfăcătoare de către evaluatori la momentul respectiv. Totodată, în cadrul sintezei concatenative bazată pe unități acustice, este necesar ca funcția de selecție a unităților să fie adaptată în conformitate cu anumite etichete atașate informației vocale și care descriu gradul de expresivitate sau tipul emoției [Wang07], [Fernandez07].

Însă înregistrarea aceluiași vorbitor în condiții de expresivitate sau emotivitate variabile este destul de greu de realizat. Drept urmare, pentru a controla rezultatul sintezei, se pot efectua modificări directe ale formei de undă sau parametrizări ale acesteia, în funcție de anumiți parametri prozodici derivați din statistici asupra unor seturi de date audio extinse.

Pentru controlul direct al formei de undă vocale, în sistemele concatenative, cea mai utilizată metodă este PSOLA (Pitch Synchronous Overlap and Add) [Valbret92]. PSOLA utilizează cadre de analiză sincrone cu frecvența fundamentală extrase și prelucrate pentru a obține o altă frecvență fundamentală sau durată a unui fonem sau a unei silabe. De cele mai multe ori această manipulare a formei de undă duce la artefacte -- segmente nenaturale de voce ce se adaugă problemelor apărute la concatenarea diferitelor segmente acustice extrase din contexte diferite.

O altă problemă, o reprezintă faptul că expresivitatea sau emoția prezentă în semnalul vocal nu se limitează doar la modificarea parametrilor prozodici, ci de multe ori se reflectă și în calitatea vocii (tremur, răgușeală, etc.) [Schroder09]. O alternativă la această manipulare a formei de undă pe baza parametrilor prozodici o reprezintă utilizarea unor parametrizări ale semnalului vocal sau vocodere.

[d'Alessandro03] prezintă o metodă de control a rezultatului sintezei prin modificarea spectrului sursei glotale, descris prin formați glotali și panta spectrului. Semnalul vocal este descompus în segmente periodice și aperiodice, modificate individual și apoi re-concatenate. [Matsui03] a propus modificarea înregistrărilor unui vorbitor pentru a reproduce emoția sau expresivitatea unui alt vorbitor folosind transformarea unei parametrizări abstracte a semnalului. [Ye04] utilizează, în schimb, o parametrizare bazată pe modelul sinusoidal pentru a face conversia identității vorbitorului, însă această metodă poate fi aplicată și pentru expresivitate.

Deși au fost dezvoltate și implementate multiple metode de control a emoțiilor și a expresivității în sinteza concatenativă, faptul că această tehnologie se bazează pe forma de undă în sine, cu anumite modificări parametrice ale rezultatului vocal, face ca orice modificare adusă semnalului să introducă artefacte sau erori de sinteză nedorite.

### **3. Controlul expresivității în cadrul sistemelor probabilistice bazate pe modele Markov**

Până în jurul anului 2015, direcțiile de cercetare în cadrul sistemelor de sinteză text-vorbire se bazau în mod preponderent pe modele statistice, dintre care cele mai utilizate erau modelele Markov cu stări ascunse (Hidden Markov Model-based Speech Synthesis System - HTS) [Tokuda13]. Aceste sisteme modelează vorbirea la nivel de fonem, folosind diferite tipuri de parametrizare a formei de undă. Cea mai utilizată metodă de parametrizare sau vocoder este cel STRAIGHT [Kawahara99] ce folosește 4 seturi de parametri: coeficienți Mel-cepstrali, coeficienți de aperiodicitate, frecvența fundamentală (F0) și durata. Pentru controlul expresivității se pot astfel manipula în mod independent modelele pentru F0 și durata. Cu toate acestea, în cadrul sistemelor de tip HTS, problema expresivității a reprezentat o provocare suplimentară deoarece natura vocii sintetizate este condiționată de utilizarea unor înregistrări audio ce conțineau o prozodie cât mai liniară, pentru ca modelele statistice să poată utiliza un număr cât mai mare de exemple fonetice pentru același context.

Există, însă, un număr mare de sisteme bazate pe modele Markov ce înglobează într-o formă sau alta un modul de control al prozodiei, iar cele mai importante dintre acestea vor fi enumerate în paragrafele următoare.

#### **Adnotarea prozodică a informației textuale**

Încă din primele versiuni ale sistemelor, în cadrul etichetelor contextuale utilizate în antrenarea modelelor acustice bazate pe modele Markov, s-a utilizat setul de adnotări prozodice ToBI [Zen09]. Cu ajutorul acestora, se puteau marca atât în setul de date de antrenare, cât și la evaluare evenimentele prozodice existente la nivelul fonemelor individuale. Însă, datorită necesității generării unei traiectorii continue pentru fluxul F0, vocea generată strict pe baza acestei metode era inexpressivă și monotonă [Toda05]. O primă încercare de a permite o mai mare variabilitate a traiectoriei F0 și a duratei a fost cea prin care s-a utilizat așa numita varianță globală (en. global variance) [Toda05], însă rezultatele nu au fost cu mult mai expresive. Iar problema principală a reprezentat-o faptul că sistemele HTS nu foloseau informații suprasegmentale (de ex. la nivel de silabă sau de cuvânt). Iar la evaluarea sistemului (sinteza propriu-zisă), aceste adnotări prozodice trebuiau deduse direct din text, fapt ce este greu de realizat.

#### **Controlul modelelor acustice**

Cel mai simplu mod de a controla modelele acustice este, în mod similar cu sinteza concatenativă, utilizarea unor seturi de date audio ce redau emoțiile sau stilurile de exprimare dorite. [Yamagishi03] prezintă un astfel de sistem de sinteză parametric bazat pe modele Markov și a cărui modele acustice sunt antrenate pe corpusuri audio cu conținut de stil citit,

dur, fericit și supărat. Antrenarea modelelor a fost realizată atât individual, cât și prin interpolarea parametrilor specifici și obținerea unor modele cu stil de exprimare intermediar.

Pe lângă setul de etichete de bază utilizate pentru generarea semnalului vocal, pentru a controla mai ușor expresivitatea vocilor sintetizate, se pot adăuga anumite caracteristici la nivel supra segmental sau metalingvistic sau care să înglobeze elemente de realizare articulatorie. De exemplu, [Miyana04] și [Nose07] utilizează modele Markov ascunse cu regresie multiplă (en. Multiple Regression HMMs), în cadrul cărora mediile modelelor probabilistice sunt controlate la momentul sintezei prin intermediul unor caracteristici auxiliare, cum ar fi stilul de vorbire, emoția, etc. Pe de altă parte, [Ling09] folosește ca parametri de control, adnotări ale mișcărilor organelor fonatoare. În cadrul acestei metode și având la dispoziție un modul ce permite extragerea parametrilor articulatori din semnalul audio, nivelul și tipul expresivității vocii sintetizate poate fi mult mai bine controlat. [Yang10] folosește tot o interpolare a modelelor acustice, însă interpolarea este bazată pe distanța Mahalanobis dintre modelele la nivel de fonem pentru diferite stiluri de exprimare. În [Ohtani15] este prezentată o metodă de control a modelelor acustice pornind de la ideea că emoția este rezultatul însumării caracteristicilor vocii neutre cu un anumit grad de emotivitate descris de un supervector și inclus în antrenarea modelelor expresive.

#### **Adaptarea modelelor acustice**

În cazul în care există un set mai larg de voci înregistrate pentru un sistem de sinteză text-vorbire, există posibilitatea ca din acest set de date să se creeze o voce de tip eigen ce înglobează informația de la toți vorbitorii. Această voce poate fi considerată ca fiind medierea acestor date, iar adaptarea modelelor acustice către un nou vorbitor ar putea fi realizată pornind de la un set de date mai redus, decât prin simpla adaptare a unui vorbitor către un altul. Pornind de la această ipoteză, [Trueba15] include în datele de antrenare caracteristici de emoție și identitate a vorbitorului, caracteristici ce sunt utilizate ulterior în funcția ce controlează gradul de adaptare a modelului rezultat. Emoțiile generate de acest sistem sunt furie, fericire, tristețe și mirare. [Latorre14] realizează mai multe voci de tip eigen, fiecare aparținând unei anumite emoții sau stil de exprimare, prin factorizarea modelelor Markov, similar și cu [Qin06].

#### **4. Controlul expresivității în sisteme de sinteză bazate pe rețele neuronale multistrat**

Rețelele neuronale s-au impus ca algoritm de învățare automată aproape universal începând cu anul 2006, când performanțele mașinilor de calcul au crescut, precum și odată cu prezentarea algoritmului de învățare rapidă a ponderilor rețelelor de către G. Hinton [Hinton06]. Ca urmare, majoritatea problemelor de învățare neliniară și pentru care există suficiente date au fost rezolvate și îmbunătățite cu ajutorul rețelelor neuronale multistrat (en. Deep Neural Networks - DNN). Printre acestea, se numără și sistemele de recunoaștere [Hinton12] și sinteză a vorbirii [Zen13]. În partea de sinteză a vorbirii, sistemele bazate pe rețele neuronale au depășit cu mult naturalețea celor bazate pe modele probabilistice [Oord16]. Au rămas încă deschise problemele de expresivitate a vocii sintetizate, precum și cea de creare a vocilor din seturi de date reduse.

În ceea ce privește expresivitatea sistemelor de sinteză bazate pe rețele neuronale, abordările includ de cele mai multe ori extinderea setului de caracteristici de intrare, cu un set de caracteristici de prozodie sau de stil de vorbire. Problema majoră a acestor sisteme este necesitatea existenței unui corpus de voce de dimensiuni mari pe baza căruia să se realizeze antrenarea rețelei. În cazul în care acest corpus de voce nu este disponibil se poate utiliza învățarea prin transfer. Această metodă presupune pre-antrenarea rețelei cu un set de date

amplu, care însă nu este proiectat specific pentru scopul dat, iar apoi rafinarea acestei rețele cu un set de date specific, dar de dimensiuni reduse [Sawada17].

Dacă aceste date sunt, însă, disponibile, modul în care sunt adnotate datele lingvistice și acustice ce intervin în antrenarea rețelei este esențial. În [Inoue17] sunt analizate 3 arhitecturi ale rețelei neuronale și a caracteristicilor de intrare pentru a modela emoțiile neutru, furie și tristețe. Caracteristicile de intrare sunt augmentate și în [Luong17] prin adăugarea unor codificări suplimentare ale identității vorbitorului. Aceste codificări sau dimensiuni de control a datelor de intrare sunt învățate direct din forma de undă în studiul [Hodari18]. [Fan15] realizează adaptarea la vorbitor prin utilizarea unei rețele comune interioare și cu noduri dependente de vorbitor în stratul de ieșire. Deși metoda este utilizată pentru adaptarea vorbitorilor, o metodă similară poate fi utilizată și în generarea emoțiilor.

În ceea ce privește sistemele de sinteză comerciale, acestea utilizează caracteristici suplimentare, derivate în general din textul de intrare și pe baza cărora, expresivitatea sau stilul semnalului audio rezultat pot fi controlate. Sistemul de sinteză bazat pe rețele neuronale Tacotron [Wang17] de la Google prezintă o extindere a funcțiilor sale de bază pentru a crea voci expresive, prin utilizarea unor caracteristici latente învățate din corpusul de antrenare și utilizate ulterior și la sinteză [Sherry-Ryan18]. Rezultatele lor prezintă și utilizarea unor seturi de parametri de control ai prozodiei dinafara setului de antrenare. Tot în cadrul Tacotron, [Stanton18] a incorporat o reprezentare latentă denumită Global Style Tokens și care poate fi generată automat din textul de intrare.

Cei de la Baidu au introdus un sistem denumit EMPHASIS [Li18] ce modelează dependențele dintre caracteristicile lingvistice și cele acustice folosind o rețea de regresie. Caracteristicile acustice sunt, de asemenea, grupate astfel încât să se maximizeze caracteristicile de emotivitate și prozodie.

Rezultatele ambelor sisteme de sinteză sunt de o calitate foarte bună, rămâne însă deschisă problema generării automate a etichetelor de expresivitate și prozodie din text, precum și transferul stilului de vorbire folosind date cât mai puține.

## 5. Concluzii

Acest document a prezentat succint cele mai importante metode de analiză și control a expresivității în cadrul sistemelor de sinteză text-vorbire. Aceste metode au fost clasificate în funcție de tipul sistemului în cadrul căruia au fost aplicate, în sisteme concatenative, sisteme statistice bazate pe modele Markov și sisteme bazate pe rețele neuronale multistrat.

Este evident faptul că studiul expresivității vocii umane este în continuare un subiect de cercetare important, dat fiind și faptul că evaluarea expresivității este mai degrabă o evaluare subiectivă, fără a fi dependentă în mod clar de anumiți parametri măsurabili. Există, însă, un oarecare nivel de acord între evaluatori în ceea ce privește, de exemplu, realizarea vocală a unor emoții puternice, cum ar fi mânia sau bucuria.

Totodată, nivelul de expresivitate sau prozodia unui vorbitor particular poate fi transpusă unei voci. Din nou, studiile din acest domeniu nu pot să identifice clar un set de parametri sau măsuri obiective a ceea ce reprezintă un anumit stil oratoric. Se pot identifica, însă, un număr redus de modificări de durată sau variații relative ale F0.

Se mai pune problema și de emfază a unor cuvinte ce trebuie accentuate pentru a transmite mesajul în mod clar către ascultător. Această emfază depinde de obicei de vorbitor, precum și de contextul mai larg al discursului din care face parte propoziția sau fraza curentă. Aceste evaluări la nivel de dialog sau discurs amplu țin mai degrabă de analiza textului și identificarea acestei emfaze.

## 6. Bibliografie

- d'Alessandro03 d'Alessandro, C., Doval, B, *Voice quality modification for emotional speech synthesis*. In Proc. Eurospeech 2003, Geneva, Switzerland 1653–1656, 2003
- Fan15 Y. Fan, Y. Qian, F. K. Soong, and L. He, *Multi-speaker modeling and speaker adaptation for DNN-based TTS synthesis*, in Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4475–4479, 2015.
- Fernandez07 R. Fernandez, B. Ramabhadran, Automatic exploration of corpus-specific properties for expressive text-to-speech: A case study in emphasis. In Proc. 6th ISCA Workshop on Speech Synthesis, Bonn, Germany 34–39, 2007
- Hinton06 G. Hinton, S Osindero, YW Teh, *A fast learning algorithm for deep belief nets*, Neural Computation, 2006
- Hinton12 Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury, *Deep Neural Networks for Acoustic Modeling in Speech Recognition*, IEEE Signal Processing Magazine, vol 29, issue 6, pp 82-97, 2012
- Hodari18 Zack Hodari, Oliver Watts, Srikanth Ronanki, Simon King, Learning interpretable control dimensions for speech synthesis by using external data, Proc of Interspeech, 2018
- Hunt96 AJ Hunt, AW Black, Unit selection in a concatenative speech synthesis system using a large speech database, Proc. Of ICASSP, 1996
- Iida03 A. Iida, N. Campbell, Speech database design for a concatenative text-to-speech synthesis system for individuals with communication disorders. International Journal of Speech Technology 6, 379–392, 2003
- Johnson02 W. Johnson, S. Narayanan, R. Whitney, R. Das, M. Bulut, C. LaBore, *Limited domain synthesis of expressive military speech for animated characters*. In Proceedings of the 7th International Conference on Spoken Language Processing, Denver, Colorado, USA, 2002
- Kawahara99 H. Kawahara, I. Masuda-Katsuse, A. de Cheveigné, Restructuring speech representations using a pitch-adaptive time–frequency smoothing and an instantaneous-frequency-based F0 extraction: Possible role of a repetitive structure in sounds. Speech communication, 27(3), 187-20, 1999
- Latorre14 J Latorre, V Wan, J Yanagisawa, *Voice expression conversion with factorised HMM-TTS models*, Proc. Of Interspeech 2014
- Li18 H Li, Y Kang, Z Wang, EMPHASIS: An emotional phoneme-based acoustic model for speech synthesis system, Proc of Interspeech 2018
- Ling09 ZH Ling, K Richmond, J Yamagishi, RH Wang, *Integrating articulatory features into HMM-based parametric speech synthesis*, IEEE Trans Audio Speech Language, vol 17, no 6, pp 1171-1185
- Luong17 Hieu-Thi Luong, Shinji Takaki, Gustav Eje Henter, Junichi Yamagishi, *Adapting and controlling DNN-based speech synthesis using input codes*, IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 4905-4909, 2018



- Matsui03 H. Matsui, H. Kawahara, Investigation of emotionally morphed speech perception and its structure using a high quality speech manipulation system. In Proc. Eurospeech, Geneva, Switzerland, 2113–2116, 2003
- Miyanaga04 K Miyanaga, T Masuo, T Kobayashi, *A style control technique for HMM-based speech synthesis*, Proc of Interspeech 2004
- Nose07 T Nose, J Yamagishi, T Masuko, T Kobayashi, *A style control technique for HMM-based expressive speech synthesis*, IEICE Trans Inf Syst, vol 90D, no 9, pp 1406–1413, 2007
- Ohtani15 Y. Ohtani, Y. Nasu, M. Morita, and M. Akamine, *Emotional transplant in statistical speech synthesis based on emotion additive model*, in Proc. Interspeech, pp. 274–278, 2015
- Oord16 Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu, *WaveNet: A Generative Model for Raw Audio*, arXiv 1609.03499
- Qin06 Long Qin, Zhen-Hua Ling, Yi-Jian Wu, Bu-Fan Zhang, Ren-Hua Wang, *HMM-Based Emotional Speech Synthesis Using Average Emotion Model*, Chinese Spoken Language Processing. ISCSLP 2006. Lecture Notes in Computer Science, vol 4274. Springer, Berlin, Heidelberg, 2006
- Sawada17 Yoshihide Sawada, Yoshikuni Sato, Toru Nakada, Kei Ujimoto, Nobuhiro Hayashi, *All-Transfer Learning for Deep Neural Networks and its Application to Sepsis Classification*, arXiv 1711.04450
- Schroder09 Mark Schroder, *Expressive Speech Synthesis: Past, Present, and Possible Futures*, in Affective Information Processing, Springer, 2009
- SkerryRyan18 RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron J. Weiss, Rob Clark, Rif A. Saurous, *Towards End-to-End Prosody Transfer for Expressive Speech Synthesis with Tacotron*, arXiv 1803.09047
- Stanton18 D Stanton, Y Wang, RJ Skerry-Ryan, Predicting expressive speaking style from text in end-to-end speech synthesis, arXiv 1808.01410
- Toda05 T Toda, K Tokuda, Speech Parameter Generation Algorithm Considering Global Variance for HMM-based Speech Synthesis, Proc. Interspeech 2015
- Trueba15 J Lorenzo-Trueba, R Barra-Chicote, R San-Segundo, J Ferreiros, J Yamagishi, JM Montero, *Emotion transplantation through adaptation in HMM-based speech synthesis*, Computer Speech and Language, vol 34, pp 292-307, 2015
- Valbret92 H.Valbret, E.Moulines, J.P.Tubach, *Voice transformation using PSOLA technique*, Speech Communication, vol 11, issues 2-3, pp 175-187, 1992
- Veaux11 C. Veaux, X Rodet, Prosodic control of unit-selection speech synthesis: A probabilistic approach, Proc. Of ICASSP 2011.
- Wang07 L. Wang, M. Chu, Y. Peng, Y. Zhao, F. Soong, *Perceptual annotation of expressive speech*. In Proc. 6th ISCA Workshop on Speech Synthesis, Bonn, Germany (2007) 46–51
- Wang17 Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, Quoc Le, Yannis Agiomyrgiannakis, Rob Clark, Rif A. Saurous, *Tacotron: Towards End-to-End Speech Synthesis*, arXiv 1703.10135
- Yamagishi03 J. Yamagishi, K. Onishi, T. Masuko, T. Kobayashi, *Modeling of various speaking styles and emotions for HMM-based speech synthesis*. In Proc. Eurospeech,

- Geneva, Switzerland 2461–2464, 2003
- Ye04 Ye, H., Young, S., *High quality voice morphing*. In Proc. ICASSP 2004, Montreal, Canada, 2004
- Zen09 H Zen, K Tokuda, A Black, *Statistical parametric speech synthesis*, Speech Communication, vol 51, no 11, pp 1039-1064, 2009
- Zen13 H Zen, A Senior, M Schuster, *Statistical Parametric Speech Synthesis Using Deep Neural Networks*, Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing, pp. 7962-7966, 2013