



Emir Demirel

Topic: Automatic Lyrics
Transcription and Alignment

Supervisors : Simon Dixon, Sven Ahlback



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 765068



PART 1 :

Computational Pronunciation Analysis and Modelling
of Sung Utterances

&

PART 2 :

Low Resource Audio-to-Lyrics Alignment

PART 1 :

Computational Pronunciation Analysis and Modelling of Sung Utterances

Contents

- Introduction
- Background Info
- Pronunciation Analysis
- Lyrics Transcription Experiments
- Conclusion & Future Work

- Introduction - Background Info - Pronunciation Analysis - Experiments - Conclusion & Future Work

The word recognition rates of industry-level ASR systems can be higher.



- Introduction - Background Info - Pronunciation Analysis - Experiments - Conclusion & Future Work

The word recognition rates of industry-level ASR systems can be higher.

There are a number of factors for lower rates in singing performances:

Background music



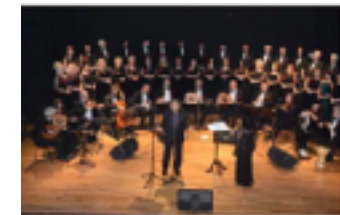
The utterance of words



vs.



Simultaneous utterance



Training data

WELCOME TO THE DALI DATASET: a large Dataset of synchronised Audio, Lyrics and vocal notes.

You can find a detailed explanation of how DALI has been created at: [\(Meseguer-Brocal, 2018\)](#).
Meseguer-Brocal, A., Cohen-Hadria and G. Peeters. DALI: a large Dataset of synchronised Audio, Lyrics and notes, automatically created using teacher-student machine learning paradigms in ISMIR Paris, France, 2018.

On this paper:

(preprint) Meseguer-Brocal, 2018, Author = Meseguer-Brocal, A. and Cohen-Hadria, G. and Peeters, G. (2018), DALI: a large Dataset of synchronised Audio, Lyrics and notes, automatically created using teacher-student machine learning paradigms in ISMIR Paris, France, 2018.

See it on YouTube

DAMP-MVP: Digital Archive of Music Performances - Smule Multilingual Vocal Performance 500x50x2

Smule Inc.

The Smule 500x50x2 dataset contains recordings of sung karaoke tracks, lyrics text files, and video metadata describing the songs being performed by each singer. This dataset has collected from performances on Smule by selecting the most popular singers, female and male, of the 100 most popular

Uploaded on May 11, 2018

Sound effects

Synthesized Vocals



- Introduction - Background Info - Pronunciation Analysis - Experiments - Conclusion & Future Work

The word recognition rates of industry-level ASR systems can be higher.

There are a number of factors for lower rates in singing performances:

Background music



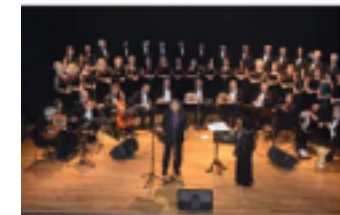
The utterance of words



vs.



Simultaneous utterance



Training data

WELCOME TO THE DALI DATASET: a large Dataset of synchronised Audio, Lyrics and vocal notes.

You can find a detailed explanation of how DALI has been created at: [\(Meseguer-Brocal, 2018\)](#).
Meseguer-Brocal, A., Cohen-Hadria and G. Peeters. DALI: a large Dataset of synchronised Audio, Lyrics and notes, automatically created using teacher-student machine learning paradigms in ISMIR Paris, France, 2018.

On this paper:

(preproceedings) Meseguer-Brocal, 2018, Author = Meseguer-Brocal, Gabriel and Cohen-Hadria, Alize and Peeters, Geoffrey, Booktitle = 79th International Society for Music Information Retrieval Conference, Editor = (ISMIR), Month = (September), Title = DALI: a large Dataset of synchronised audio, lyrics and notes, automatically created using teacher-student machine learning paradigms.

See it on YouTube

DAMP-MVP: Digital Archive of Music Performances - Smule Multilingual Vocal Performance 500x50x2

smule: 194

The Smule 500x50x2 dataset contains recordings of sung karaoke tracks, lyrics text files, and video metadata describing the songs being performed by each singer. This dataset has collected from performances on Smule by selecting the most popular singers, female and male, of the 100 most popular

Uploaded on May 11, 2018

Sound effects

Synthesized Vocals



- Introduction - Background Info - Pronunciation Analysis - Experiments - Conclusion & Future Work

The word recognition rates of industry-level ASR systems can be higher.

There are a number of factors for lower rates in singing performances:

Background music



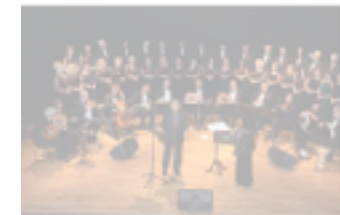
The utterance of words



vs.



Simultaneous utterance



Training data

WELCOME TO THE DALI DATASET: a large Dataset of synchronised Audio, Lyrics and vocal notes.

You can find a detailed explanation of how DALI has been created at [@Moreau-Brosal_2018](#).
Moreau-Brosal, A., Cohen-Hadria and G. Heuser, DALI: a large Dataset of synchronised Audio, Lyrics and notes, automatically created using teacher-student machine learning paradigms in ISMIR Paris, France, 2018.

On this paper:

Representings@Moreau-Brosal_2018, Author + @Moreau-Brosal, Gabriel and Cohen-Hadria, Alice and Heuser, Geoffroy, Available + 70th International Society for Music Information Retrieval Conference, Editor + @ISIRI, Month + September, Title + DALI: a large Dataset of synchronised audio, lyrics and notes, automatically created using teacher-student machine learning paradigms.

Full & Demo (4.3)

Download

Download Source

View

DRMP-MVP: Digital Archive of Music Performances - Smule Multilingual Vocal Performance 500x50x2

Smule Inc.

The Smule 500x50x2 dataset contains recordings of sung karaoke tracks, lyrics text files, and video metadata describing the songs being performed by each singer. This dataset has collected 1000 performances on Smule by selecting the most popular singers, female and male, of the 500 most popular

Uploaded on May 10, 2018

Pronunciation in singing



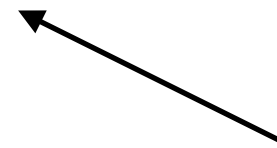
- Introduction - **Background Info** - Pronunciation Analysis
 - Experiments - Conclusion & Future Work

- (Automatic) Lyrics Transcription is the process of recognising the most likely sequence of words from sung utterances.

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

- (Automatic) Lyrics Transcription is the process of recognising the most likely sequence of words from sung utterances.

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}}(P(\mathbf{w}|\mathbf{X}))$$



Fundamental equation of (ASR)

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

- (Automatic) Lyrics Transcription is the process of recognising the most likely sequence of words from sung utterances.

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}}(P(\mathbf{w}|\mathbf{X}))$$

OR applying the Bayes' rule

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}}(P(\mathbf{X}|\mathbf{w})P(\mathbf{w})).$$

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

- (Automatic) Lyrics Transcription is the process of recognising the most likely sequence of words from sung utterances.

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}}(P(\mathbf{w}|\mathbf{X}))$$

OR applying the Bayes' rule

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}}(P(\mathbf{X}|\mathbf{w})P(\mathbf{w})).$$

Acoustic Model

Language Model

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

- In Large Vocabulary Continuous Speech Recognition (LVCSR), the search space of words is very large.

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} (P(\mathbf{w}|\mathbf{X}))$$

- We need to take into account pronunciation variations

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} (P(\mathbf{X}|\mathbf{w})P(\mathbf{w})).$$

Thus, we build the acoustic model for **phonemes (Q)**;

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) \sum_{\mathbf{Q}} P(\mathbf{X}|\mathbf{Q})P(\mathbf{Q}|\mathbf{w}).$$

Language Model
↗

Acoustic Model
↑

Pronunciation Model
↖

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work


$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) \sum_{\mathbf{Q}} P(\mathbf{X}|\mathbf{Q})P(\mathbf{Q}|\mathbf{w}).$$

Through Viterbi decoding, we find the most likely word sequence;

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) \max_{\mathbf{Q} \in Q_w} P(\mathbf{X}|\mathbf{Q})P(\mathbf{Q}|\mathbf{w}),$$

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

Valid pronunciation of words defined by the lexicon (*)

$$P(\mathbf{Q}|\mathbf{w}) = \prod_{l=1}^L P(\mathbf{q}^{w_l}|w_l)$$


Phonetic Lexicon

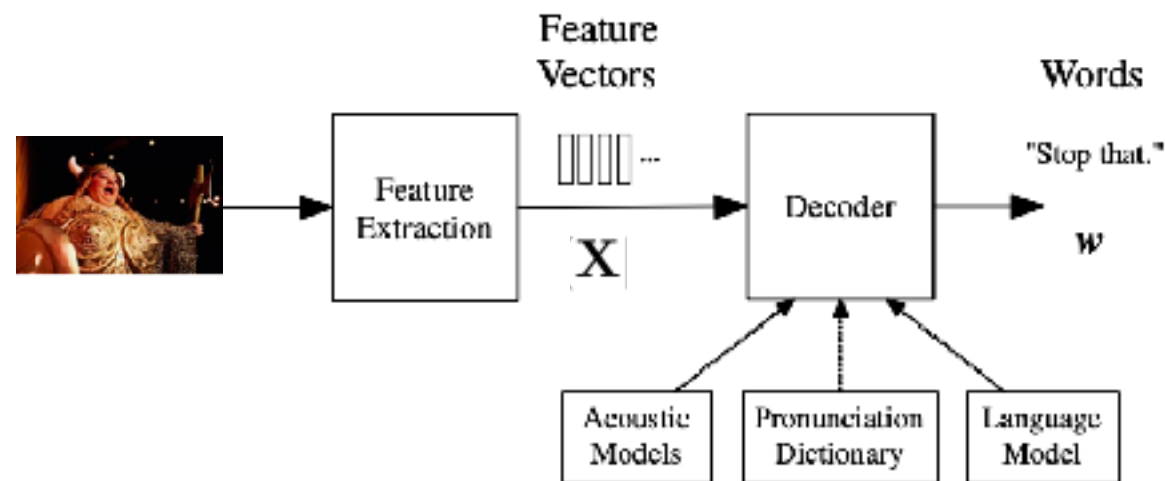
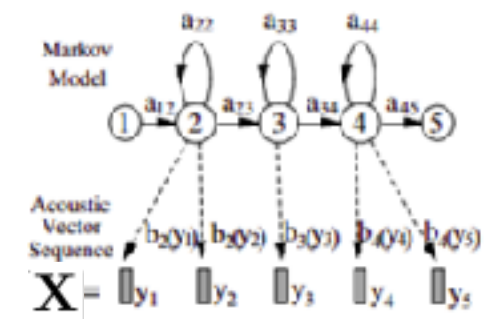
GIMME	G IH IH M IY
GIMME	G IH M IY IY
GIMME	G IH M IY
A	AA
MAN	M AE N
AFTER	AE F T ER
MIDNIGHT	M IH D N AY AY T

- Through vectorizing words with phoneme sequences, we allow multiple pronunciation variants of a word possible.

(*) *Gales & Young (2008)*

- Introduction - **Background Info** - Pronunciation Analysis
- Experiments - Conclusion & Future Work

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) \max_{\mathbf{Q} \in Q_{\mathbf{w}}} P(\mathbf{X}|\mathbf{Q})P(\mathbf{Q}|\mathbf{w}),$$



1 Architecture of an HMM-based Recogniser.

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

Analysis Data: Native ENG speakers within the NUS Sung on Spoken Lyrics Corpus

Duan, Zhiyan, et al. "The NUS Sung and Spoken Lyrics Corpus: A Quantitative Comparison of Singing and Speech."

TABLE III
SUBJECTS IN THE NUS CORPUS

Code	Gender	Voice Type	Sung Accent	Spoken Accent
01	F	Soprano	North American	North American
02	F	Soprano	North American	North American
03	F	Soprano	North American	Mild Local Singaporean
04	F	Alto	Mild Malay	Mild Malay
05	F	Alto	Malay	Malay
06	F	Alto	Mild Malay	Mild Malay
07	M	Tenor	Mild Local Singaporean	Mild Local Singaporean
08	M	Tenor	Northern Chinese	Northern Chinese
09	M	Baritone	North American	North American
10	M	Baritone	North American	Standard Singaporean
11	M	Baritone	North American	North American
12	M	Bass	Local Singaporean	Local Singaporean

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

Example 4.1 Substitutions in the utterance 'AND THE WONDER OF IT ALL'. w and \hat{w} are the human annotated ground truth and predicted word sequences. Q and \hat{Q} are the corresponding phonemic representations. In the bottom, the pronunciations obtained from the CMU English dictionary are provided. The word & pronunciation errors are highlighted with bold font.

	w	AND	THE	WONDER	OF	IT	ALL
	\hat{w}	AND	THOUGH	ONE DARE	OUR	FEET	ALL
Ground Truth	Q	HH EH N <i>eps</i>	D OW	W AH N D EH R	AO F	IH T	AO AH
Prediction	\hat{Q}	<i>eps</i> AE N D	DH OW	W AH N D EH R	AA R	F IY T	AO L

- Introduction - Background Info - Pronunciation Analysis
- Experiments - Conclusion & Future Work

We compute the alignment score matrix, \mathbf{D} , by performing Levenshtein alignment, lev between the phoneme sequences of the predictions, $\hat{\mathbf{Q}}_M$ and the ground truth \mathbf{Q}_N ,

$$\mathbf{D}_{M \times N} = lev(\hat{\mathbf{Q}}_M, \mathbf{Q}_N) \quad (7)$$

and find the best alignment path, $\mathbf{A}_{2 \times K}$ through reverse tracing to find the path with the lowest pairwise gap cost:

$$\mathbf{A}_{2 \times K} = \begin{pmatrix} \cdots & \widehat{\phi_{k-1}^*} & \widehat{\phi_k^*} & \widehat{\phi_{k+1}^*} & \cdots \\ \cdots & \phi_{k-1}^* & \phi_k^* & \phi_{k+1}^* & \cdots \end{pmatrix}$$

.....

Example 4.1 Substitutions in the utterance ‘AND THE WONDER OF IT ALL’. w and \hat{w} are the human annotated ground truth and predicted word sequences. Q and \hat{Q} are the corresponding phonemic representations. In the bottom, the pronunciations obtained from the CMU English dictionary are provided. The word & pronunciation errors are highlighted with bold font.

w	AND	THE	WONDER	OF	IT	ALL
\hat{w}	AND	THOUGH	ONE DARE	OUR	FEET	ALL
Q	HH EH N eps	D OW	W AH N D EH R	AO F	IH T	AO AH
\hat{Q}	eps AE N D	DH OW	W AH N D EH R	AA R	F IY T	AO L

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

Example 4.1 Substitutions in the utterance 'AND THE WONDER OF IT ALL'. w and \hat{w} are the human annotated ground truth and predicted word sequences. Q and \hat{Q} are the corresponding phonemic representations. In the bottom, the pronunciations obtained from the CMU English dictionary are provided. The word & pronunciation errors are highlighted with bold font.

	w	AND	THE	WONDER	OF	IT	ALL
	\hat{w}	AND	THOUGH	ONE DARE	OUR	FEET	ALL
Ground Truth	Q	HH EH N eps	D OW	W AH N D EH R	AO F	IH T	AO AH
Prediction	\hat{Q}	eps AE N D	DH OW	W AH N D EH R	AA R	F IY T	AO L

Insertion

Deletion

Substitution

There are three operations defined on these phoneme pairs to match \hat{Q}_M to Q_N : insertions (I), substitutions (S) and deletions (D). These operations are represented in \mathbf{A} with the symbol ϵ . An alignment instance $a_k = \begin{pmatrix} \epsilon \\ \hat{\phi}_k^* \end{pmatrix}$ is a deletion and the opposite case would be an insertion.

- Introduction - Background Info - Pronunciation Analysis

- Experiments - Conclusion & Future Work

Example 4.1 Substitutions in the utterance ‘AND THE WONDER OF IT ALL’. w and \hat{w} are the human annotated ground truth and predicted word sequences. Q and \hat{Q} are the corresponding phonemic representations. In the bottom, the pronunciations obtained from the CMU English dictionary are provided. The word & pronunciation errors are highlighted with bold font.

	w	AND	THE	WONDER	OF	IT	ALL
	\hat{w}	AND	THOUGH	ONE DARE	OUR	FEET	ALL
Ground Truth	Q	HH EH N <i>eps</i>	D OW	W AH N D EH R	AO F	IH T	AO AH
Prediction	\hat{Q}	<i>eps</i> AE N D	DH OW	W AH N D EH R	AA R	F IY T	AO L

Let the number of correctly matching pairs in \mathbf{A} be C , then the confidence score per phoneme type, c_ϕ , can be retrieved as:

$$c_\phi = \frac{\sum_i^T C_{\phi,i} - (S_{\phi,i} + I_{\phi,i} + D_{\phi,i})}{\sum_i^T C_{\phi,i} + S_{\phi,i} + I_{\phi,i} + D_{\phi,i}}, \quad \phi \in \Omega_E \quad (8)$$

where T is the number of utterances in the analysis set and Ω_E is the English phoneme set used in our analysis. The denominator is necessary to normalize with respect to the total number of pairs in \mathbf{A} , since the phonemes in Ω_E are not necessarily represented equally in the analysis dataset.

- Introduction - Background Info - Pronunciation Analysis
- Experiments - Conclusion & Future Work

Vowels	ϕ	$c_\phi(R)$	Φ'_N
Short Vowels	AE	-0.42 (38)	AH, EH, AA
	AH	0.17 (33)	AA, EH, OW
	EH	0.3 (32)	AH, AE, IH
	IH	0.48 (25)	IY, AH, EY
	UH	0 (36)	AO, UW, AH
Long Vowels	AA	0.5 (24)	AO, AW, AE
	AO	0.06 (35)	AA, AH, OW
	ER	0.36 (31)	AH, OW, EH
	IY	0.87 (6)	EY, IH, EH
	UW	0.88 (4)	OW, AH, UH
Diphthongs	AY	0.86 (8)	AA, AH, EH
	AW	0.71 (13)	AA, AH
	EY	0.87 (7)	IY, AY, EH
	OW	0.76 (17)	AO, AA, AH
	OY	0.4 (28)	OW, AO, AY

Table 1: Results of the phonetic analysis (vowels)

Consonants	ϕ	$c_\phi(P)$	Φ'_N
Plosives	B	0.77 (16)	D, P, W
	D	0.16 (34)	T, N, JH
	G	0.77 (15)	NG, K
	K	0.85 (15)	G, HH
	P	0.78 (14)	B, M, F
Affricates	T	0.37 (29)	D, S, CH
	CH	0.79 (13)	JH, SH, T
	JH	0.88 (5)	CH, S, Y
Nasals	M	0.93 (2)	N, NG
	N	0.85 (12)	M, NG, D
	NG	0.85 (9)	N, M, T
Fricatives	DH	0.36 (30)	TH, D, N
	F	0.91 (3)	V, P, TH
	HH	0.70 (19)	DH, W, Y
	S	0.95 (1)	Z, TH, T
	SH	0.85 (10)	CH, S, Z
	TH	0.57 (21)	S, T, DH
	V	0.56 (22)	F, R, DH
	Z	-0.05 (37)	S, T
	ZH	N/A	N/A
Approximants [†]	L	0.44 (27)	AA, OW, AH
	R	0.48 (25)	AA, AH, JH
	W	0.66 (20)	AA, OW, V
	Y	0.55 (23)	IH, AH, IY

Table 2: Results of the phonetic analysis (consonants)

$$-1 \leq c_\phi \leq 1$$

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

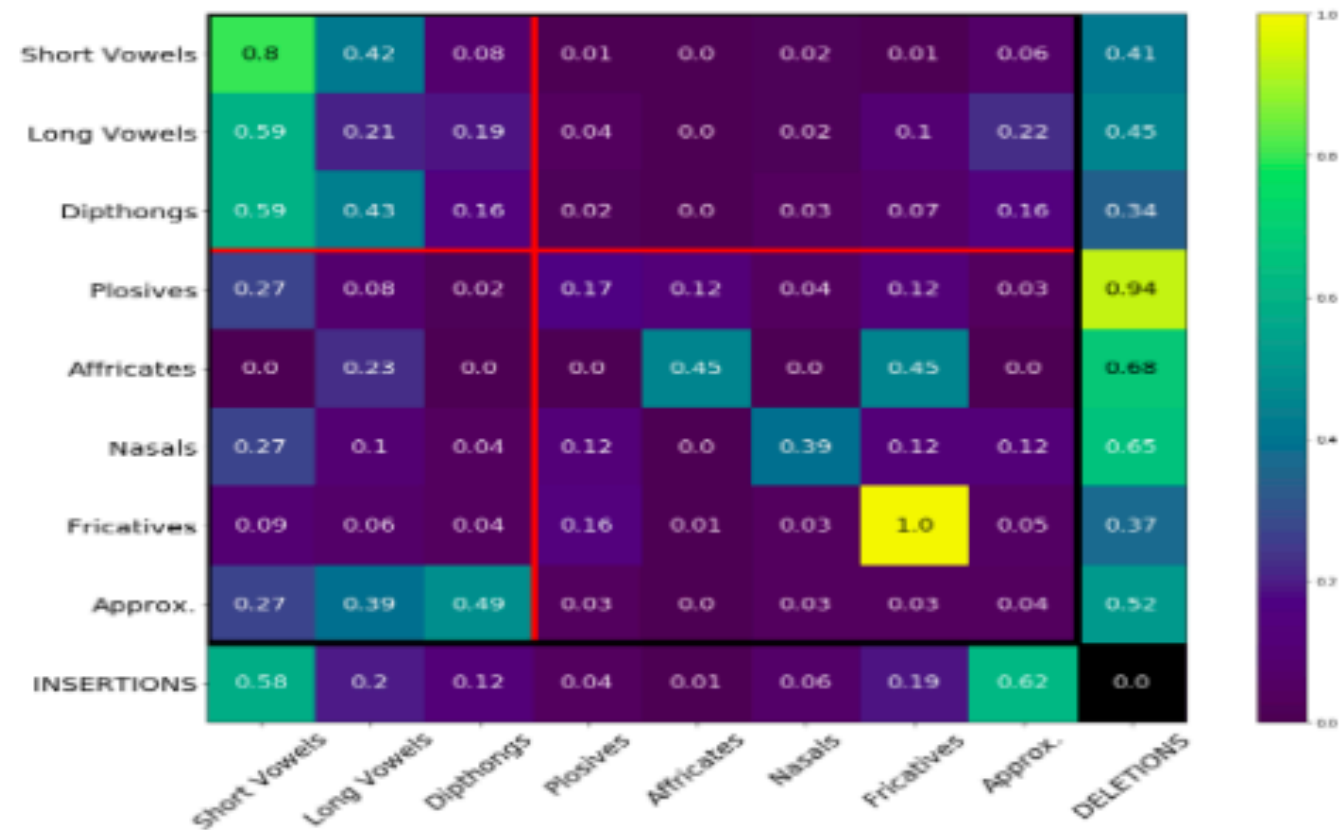


Fig. 1: Confusion matrix w.r.t. phoneme categories in Tables 1 and 2. The red lines highlight the within-class regions for vowels and consonants. The numbers in cells are normalized values. The labels on the horizontal and the vertical axes represent the ground-truth and predictions respectively.

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

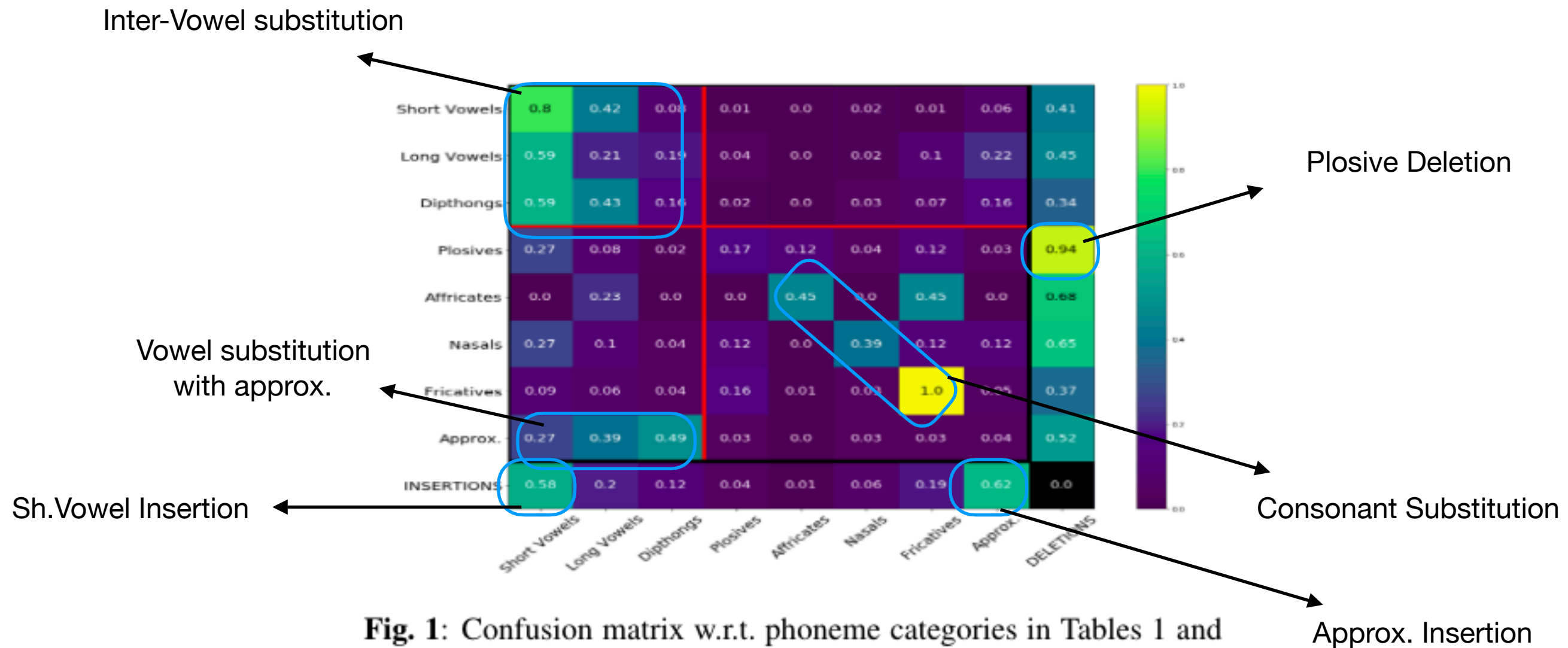


Fig. 1: Confusion matrix w.r.t. phoneme categories in Tables 1 and 2. The red lines highlight the within-class regions for vowels and consonants. The numbers in cells are normalized values. The labels on the horizontal and the vertical axes represent the ground-truth and predictions respectively.

- Introduction - Background Info - **Pronunciation Analysis**
- Experiments - Conclusion & Future Work

Three observations when uttering words in singing:

- Longer vowels
- Omitted plosives & approximants
- Triphone substitutions

- Introduction - Background Info - **Pronunciation Analysis**
 - Experiments - Conclusion & Future Work

Three observations when uttering words in singing:

- **Longer vowels**
- Omitted plosives
- Triphone substitutions

	Speech	Singing
<i>Duration (min)</i>	0.03	0.18
<i>Duration (max)</i>	0.77	3.86
<i>Duration (avg.)</i>	0.10	0.34
<i>Articulation Rate (per min)</i>	266.25	172.50

Table 4: Mean articulations rates (syllables per minute) and duration stats (in seconds)

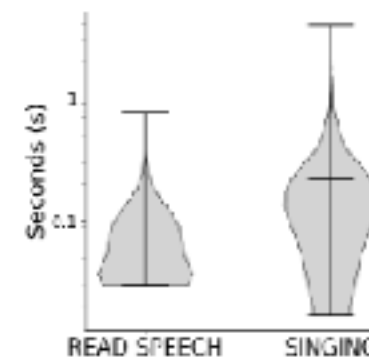


Figure 2: Duration distributions of vowels

- Introduction - Background Info - **Pronunciation Analysis**
 - Experiments - Conclusion & Future Work

Three observations when uttering words in singing:

- Longer vowels
- **Omitted plosives**
- Triphone substitutions

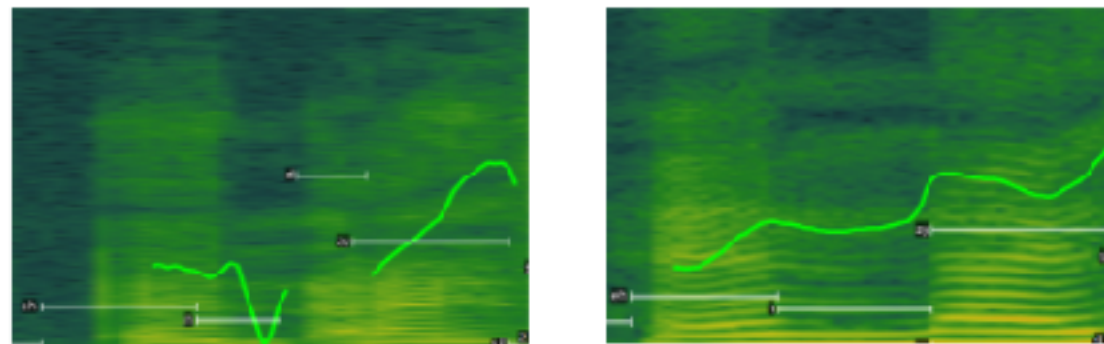


Fig. 2: An example of an omitted plosive in singing. $W = \text{'AND I'}$; $Q^{read} = \text{'\text{AE N D AY'}}$ (left); $Q^{sing} = \text{'EH N AY'}$. The gray horizontal lines show the temporal phoneme regions and the bright green curves are the pitch tracks extracted using pYIN [16].

- Introduction - Background Info - Pronunciation Analysis

- Experiments - Conclusion & Future Work

Three observations when uttering words in singing:

- Longer vowels
- Omitted plosives
- **Triphone substitutions**

Vowels	ψ	$c_{\psi}(R)$	Φ'_N
Short Vowels	AE	-0.42 (38)	AH, EH, AA
	AH	0.17 (33)	AA, EH, OW
	EH	0.3 (32)	AH, AE, IH
	IH	0.48 (26)	IY, AH, EY
	UH	0 (36)	AO, UW, AH
Long Vowels	AA	0.5 (24)	AO, AW, AE
	AO	0.06 (35)	AA, AH, OW
	ER	0.36 (31)	AH, OW, EH
	IY	0.87 (6)	EY, IH, EH
	UW	0.88 (4)	OW, AH, UH
Diphthongs	AY	0.86 (8)	AA, AH, EH
	AW	0.71 (18)	AA, AH
	EY	0.87 (7)	IY, AY, EH
	OW	0.76 (17)	AO, AA, AH
	OY	0.4 (28)	OW, AO, AY

Table 1: Results of the phonetic analysis (vowels)

Consonants	ϕ	$c_{\phi}(R)$	Φ'_N
Plosives	B	0.77 (16)	D, P, W
	D	0.16 (34)	TN, JH
	G	0.77 (15)	NG, K
	K	0.85 (15)	G, HH
	P	0.78 (14)	R, M, F
	T	0.37 (29)	D, S, CH
Affricates	CH	0.79 (13)	JH, SH, T
	JH	0.88 (5)	CH, S, Y
Nasals	M	0.93 (2)	N, NG
	N	0.84 (13)	M, NG, D
	NG	0.83 (9)	N, M, T
Fricatives	DH	0.36 (30)	TH, D, N
	F	0.91 (3)	V, P, TH
	HH	0.70 (19)	DH, W, Y
	S	0.95 (1)	Z, TH, T
	SH	0.83 (10)	CH, S, Z
	TH	0.57 (21)	S, DH, DH
	V	0.56 (22)	F, R, DH
	Z	0.05 (37)	S, T
	ZH	N/A	N/A
	L	0.44 (27)	AA, OW, AH
Approximants [†]	R	0.48 (25)	AA, AH, IH
	W	0.66 (20)	AA, OW, V
	Y	0.51 (23)	IH, AH, IY

Table 2: Results of the phonetic analysis (consonants)

- Introduction - Background Info - Pronunciation Analysis
- **Experiments** - Conclusion & Future Work

We extend the pronunciations CMU dictionary with longer vowels and omitted plosives...

FLOAT	F	L	OW	T	
FLOAT	F	L	OW	OW	T
FLOAT	F	L	OW		
FLOAT	F	L	OW	OW	T

FLOATING	F	L	OW	T	IH	NG	
FLOATING	F	L	OW	OW	T	IH	NG
FLOATING	F	L	OW	T	IH	IH	NG
FLOATING	F	L	OW	T	IH	N	
FLOATING	F	L	OW	OW	T	IH	N
FLOATING	F	L	OW	T	IH	IH	N

And test the effectiveness of these pronunciation extensions in the context of word recognition.

- Introduction - Background Info - Pronunciation Analysis

- Experiments - Conclusion & Future Work

In the experiments, the s.o.t.a lyrics transcription framework is employed (Demirel, 2020) - open source toolkit!.

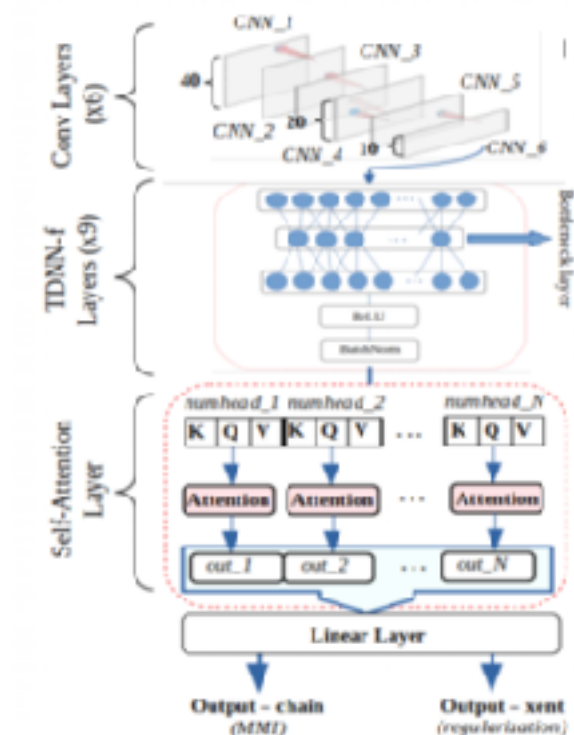
ALTA - (A)utomatic (L)yrics (T)ranscription & (A)lignment

A kaldi recipe for automatic lyrics transcription and audio-to-lyrics alignment tasks.

If you use this repository, please cite it as follows:

```
@inproceedings{demirel2020,  
  title={Automatic lyrics transcription using dilated convolutional neural networks with self-attention},  
  author={Demirel, Enir and Ahlback, Sven and Dixon, Simon},  
  booktitle={International Joint Conference on Neural Networks},  
  publisher={IEEE},  
  year={2020}  
}
```

Link to paper : <https://arxiv.org/abs/2007.06486>



- Introduction - Background Info - Pronunciation Analysis
- **Experiments** - Conclusion & Future Work

Evaluation sets:

- A capella recordings
- Both sung and read
- Gender-wise balanced

Mixed (predominantly non-native)

	Char.	Words	Sentences	Recordings
NUS_read	21935	5788	781	32
NUS_sing	21935	5788	1029	32
DAMP_test	17609	4626	479	70

Native from UK

Table 3: Statistics of evaluation sets

- Introduction - Background Info - Pronunciation Analysis
- **Experiments - Conclusion & Future Work**

Speech lexicon

Singing adapted lexicon

	<i>L_{CMU}</i>				<i>L_{sing}</i>			
	<i>ER</i>	<i>S</i>	<i>I</i>	<i>D</i>	<i>ER</i>	<i>S</i>	<i>I</i>	<i>D</i>
word								
DAMP_test	17.21	10.67	1.43	5.66	15.61	10.83	1.53	3.24
NUS_read	10.51	7.52	1.07	1.91	9.26	6.41	1.01	1.80
NUS_sing	13.19	8.60	1.63	2.95	10.06	7.16	1.25	1.65
character								
DAMP_test	11.41	4.78	1.55	4.79	9.88	4.55	1.63	3.50
NUS_read	5.57	2.73	1.38	1.47	5.10	2.62	1.11	1.37
NUS_sing	7.07	3.13	1.60	2.33	6.14	3.03	1.36	1.75

Table 4: Word and character error rates using standard (*L_{CMU}*) and singing-adapted (*L_{sing}*) pronunciation dictionaries.

	<i>L_{CMU}</i>			<i>L_{sing}</i>		
	<i>ER</i>	<i>S</i>	<i>D</i>	<i>ER</i>	<i>S</i>	<i>D</i>
word (ending with plosives)						
DAMP_test	22.84	13.06	7.78	17.67	10.15	7.21
NUS_read	9.74	8.82	0.91	9.01	7.90	1.10
NUS_sing	14.01	7.76	5.73	7.94	5.73	2.21
vowel						
DAMP_test	13.20	6.47	6.72	9.80	5.59	4.21
NUS_read	4.02	2.44	1.58	3.99	2.55	1.44
NUS_sing	7.23	2.98	4.26	6.71	3.03	3.68

Table 5: Error analysis for plosives and vowels

- Introduction - Background Info - Pronunciation Analysis
- **Experiments - Conclusion & Future Work**

Speech lexicon

Singing adapted lexicon

	<i>L_{CMU}</i>				<i>L_{sing}</i>			
	<i>ER</i>	<i>S</i>	<i>I</i>	<i>D</i>	<i>ER</i>	<i>S</i>	<i>I</i>	<i>D</i>
word								
DAMP_test	17.21	10.67	1.43	5.66	15.61	10.83	1.53	3.24
NUS_read	10.51	7.52	1.07	1.91	9.26	6.41	1.01	1.80
NUS_sing	13.19	8.60	1.63	2.95	10.06	7.16	1.25	1.65
character								
DAMP_test	11.41	4.78	1.55	4.79	9.88	4.55	1.63	3.50
NUS_read	5.57	2.73	1.38	1.47	5.10	2.62	1.11	1.37
NUS_sing	7.07	3.13	1.60	2.33	6.14	3.03	1.36	1.75

Table 4: Word and character error rates using standard (*L_{CMU}*) and singing-adapted (*L_{sing}*) pronunciation dictionaries.

	<i>L_{CMU}</i>			<i>L_{sing}</i>		
	<i>ER</i>	<i>S</i>	<i>D</i>	<i>ER</i>	<i>S</i>	<i>D</i>
word (ending with plosives)						
DAMP_test	22.84	13.06	7.78	17.67	10.15	7.21
NUS_read	9.74	8.82	0.91	9.01	7.90	1.10
NUS_sing	14.01	7.76	5.73	7.94	5.73	2.21
vowel						
DAMP_test	13.20	6.47	6.72	9.80	5.59	4.21
NUS_read	4.02	2.44	1.58	3.99	2.55	1.44
NUS_sing	7.23	2.98	4.26	6.71	3.03	3.68

Table 5: Error analysis for plosives and vowels

- Consistent WER improvement
- Marginal improvement for *read* samples

- Introduction - Background Info - Pronunciation Analysis
 - Experiments - **Conclusion & Future Work**

- A number of frequent pronunciation variances during singing are identified using a novel computational method that combines human annotations and an AI-based lyrics transcription system.
- Using a singing-adapted lexicon can yield to improvement in word recognition rates.
- Sentence-level annotations are provided for NUS Corpus, which can be leveraged for both training and evaluation in the context of automatic lyrics transcription.

- Introduction - Background Info - Pronunciation Analysis
 - Experiments - **Conclusion & Future Work**

- Add new pronunciations based on 'triphone substitutions'.
- Obtain pronunciation probabilities.



Emir Demirel

Topic: Automatic Lyrics
Transcription and Alignment
Supervisor: Simon Dixon



e.demirel@qmul.ac.uk

Q&A?