



Animate the tongue and jaw from only speech signal to add realism to facial animations.

Accurately animating the tongue is difficult since:

- Performance capture is not reliable as tongue and teeth are partially visible.
- Manually animating the tongue is nearly impossible.

# EMA Tongue Motion Dataset

15:40:29:05

EMA Sensor	Placement
TD	Tongue Dorsum
ТВ	Tongue Blade
BR	Tongue Blade Right
BL	Tongue Blade Left
TT	Tongue Tip
UL	Upper Lip
LC	Right Lip Corner
LL	Lower Lip
LI	Jaw, medial incisors
LJ	Jaw, canine & first premolar

We captured the first large scale electromagnetic articulography (EMA) tongue dataset with parasagittal sensors for animation purposes.

### Carstens AG501<sup>1</sup>

- Sample Rate: 250 Hz
- Capture Error < 1mm
- 10 sensors on tongue and lips
- 3 sensors for bite plane

Total samples: 2160

### Sensor Placement



We placed 5 sensors on the tongue, 2 on the jaw and 3 on the lips.

() Data available for download at https://salmedina.github.io/tongue-anim









# **Speech Driven Tongue Animation**

Salvador Medina<sup>1,2</sup>, Denis Tome<sup>2</sup>, Carsten Stoll<sup>2</sup>, Mark Tiede<sup>3</sup>, Kevin Munhall<sup>4</sup>, Alex Hauptmann<sup>1</sup>, Iain Matthews<sup>2</sup>

Our Approach \_\_\_\_\_

# References

[1] Carstens Medizinelektronik GmbH. 3D electromagnetic articulograph. https://www.articulograph.de/ ] Rothauser, E. H. "IEEE recommended practice for speech quality measurements." IEEE Trans. on Audio and Electroacoustics 17 (1969): 225-246. [3] Zue, Victor, Stephanie Seneff, and James Glass. "Speech database development at MIT: TIMIT and beyond." Speech communication 9.4 (1990): 351-356. [4] Amodei, Dario, et al. "Deep speech 2: End-to-end speech recognition in english and mandarin." International conference on machine learning. PMLR, 2016. [5] Schneider, Steffen, et al. "wav2vec: Unsupervised pre-training for speech recognition." arXiv preprint arXiv:1904.05862 (2019).



# Evaluation

Our method produces realistic tongue animations due to low error inner-mouth pose estimation.



### Tongue motions with more complexity are not modeled as accurately.



### Our best results combine Wav2Vec-C features with a bidirectional 5-layered GRU.

Decoder $\setminus$ Feature	Phone	MFCC	DS2	W2V-C	W2V-Z	Num. Parameters	Inference [ms]	Latency [ms]					~ ~							- 2.6
MLP 15:5	2.445	2.075	2.393	1.959	1.937	$6.62 \times 10^{7}$	0.232	300	Phone	- 2.3	2.6	2.2	2.2	2.7	1.3	1.8	2.6	2.1	2	
LSTM-1L	4.207	2.344	2.269	2.047	2.140	$3.17 \times 10^{6}$	1.150	20												- 2.4
LSTM-2L	4.209	2.178	4.206	1.990	4.212	$5.27 \times 10^{6}$	2.238	20	MECC	- 1.8	2	1.7	1.8	2.1	1.1	1.5	1.8	1.5	1.4	- 2 2
LSTM-5L	2.656	2.037	2.264	1.999	1.960	$1.16 \times 10^{7}$	5.432	20		2.10			110					_10	2	2.2
Bi-LSTM-1L	3.664	2.346	2.375	2.373	3.481	$6.33 \times 10^{6}$	2.229	300												- 2.0
Bi-LSTM-2L	4.577	2.109	2.844	2.188	3.874	$1.26 \times 10^{7}$	4.512	300	DeepSpeech2 -	ch2 2	23	2	2	23	12	17	21	1 0	1 8	2.0
Bi-LSTM-5L	4.365	1.912	2.218	1.927	2.929	$3.15 \times 10^{7}$	11.000	300			2.5			2.5	1.2	1.7	2.1	1.9	1.0	- 1.8
GRU-1L	4.150	2.290	2.250	1.949	2.071	$2.38 \times 10^{6}$	1.144	20												
GRU-2L	2.623	2.117	2.179	1.897	1.980	$3.95 \times 10^{6}$	2.193	20	Wav2Vec-C	- 1.9	2.1	1.8	1.8	2.1	1.1	1.5	1.8	1.5	1.4	- 1.6
GRU-5L	2.661	2.006	2.184	1.916	1.954	$8.68 \times 10^{6}$	5.339	20								1.0		110		
Bi-GRU-1L	4.405	2.368	2.529	2.055	2.613	$4.76 \times 10^{6}$	2.290	300												- 1.4
Bi-GRU-2L	3.143	1.953	2.947	1.932	2.513	$9.48 \times 10^{6}$	4.439	300	Wav2Vec-Z	1 0	2	1 7	10	2 1	1 1	1 /	10	1 /	1 /	
Bi-GRU-5L	2.341	1.973	2.058	1.757	1.784	$2.37 \times 10^{7}$	10.955	300		- 1.c	2	1.7	1.0	2.1	1.1	1.4	1.0	1.4	1.4	- 1.2
Transformer	2.368	2.283	2.168	1.935	1.942	$5.045 \times 10^{7}$	3.515	300		TD	т ТВ	, BR	BL	π	UL	LC	Ļ	Ļ	Ļ	
Experimental Results. Error is temporal mean MSE [mm].											nd	ma	rk	Pre	edio	ctic	on l	Erro	or [	mm]

# Conclusions

- Ο
- Ο
- Limited lip animation due to the sparsity of the sensors

Animations produced from our method are preferred over a no tongue or mismatched animation, and confused with the GT.



Complex motion

Our inner-mouth mocap dataset enables the training of data-driven models Deep Learning audio representations outperform traditional methods for speech-animation Simple-RNN based articulation decoders generalize across gender, age, and prosody