





### Introduction:

- A 3-dimensional convolutional neural network is trained on unlabeled ultrasound videos to predict an upcoming tongue image from previous ones.

- The network obtains results superior to those of simpler predictors
- A starting point for exploiting the higher-level representation of the tongue learned by the system in a variety of applications in speech research.

## Why doing this:

- To use Unsupervised Learning to extract features from tongue images for speech recognition. - To prove simple time regression problem can be handled well without RNN.

### What did we do:

- Prediction of ultrasound video frame (WSJO & TJU data)
- Prediction of utrasound tongue contour (Cross data)
- Evaluation metric design:
- calculating the loss (MSE) between different predictors and the target image
- overlapping different prediction results and the target image
- validating motion detection method to detect the contour change in the predicted frame

# **Different training datasets**

# Input images



8 consecutive images to predict the 9th image,

### **TJU Data**









### **Cross Data**

8 consecutive images to predict the snake contour of the 9th image











Wechat



Github

WSJ0 data

n) and 9th image (red) and 3DCNN predicted 9th image

- The red tongue contour is **above** the green one in both images.
- This means the predicted image actually moved rather than just blurring.

TJU data

Left: comparison of the contours of the 8 ge (green) and Right: comparison of the contours of the 8th i

nage (green) and the predicted 9th image (red).

- Initialize the video with the same 12 points.
- Extract the contours from the real images and the predicted
- If the predicted image's snake has the same position as the real one, then we learnt the 'movement'.

Cross data

Contour of the 9th image (left, in green); Predicted contour of the 9th image (center, in red); Overlay of the two contours (right).

The technique works extremely well, as can be seen in the videos available (scan the QR code)

video for the two left columns



right column

- The Mean Sum of Distances, MSD, between the 3DCNN prediction and the 9th image snake was only **1.1 pixels**, corresponding to **0.4 mm**.