# An Overview about
# Lip Reading
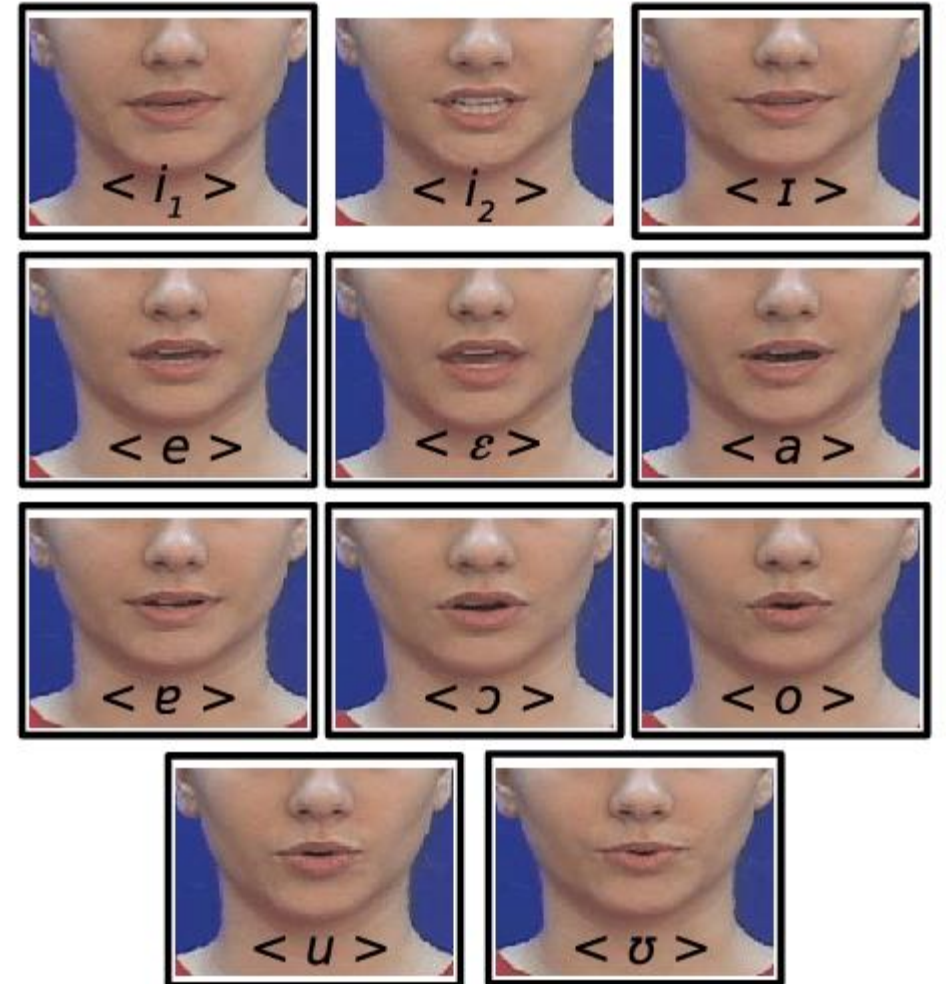# &
# Audio-visual Speech Recognition

Chen Chen

2022/03/18

# Catalog

- Introduction
- Pipeline
- Datasets & Performance Evaluation
- Methods

# Definition

- Lip Reading
  - recognize what is being said from visual information alone
- Audio-visual Speech Recognition
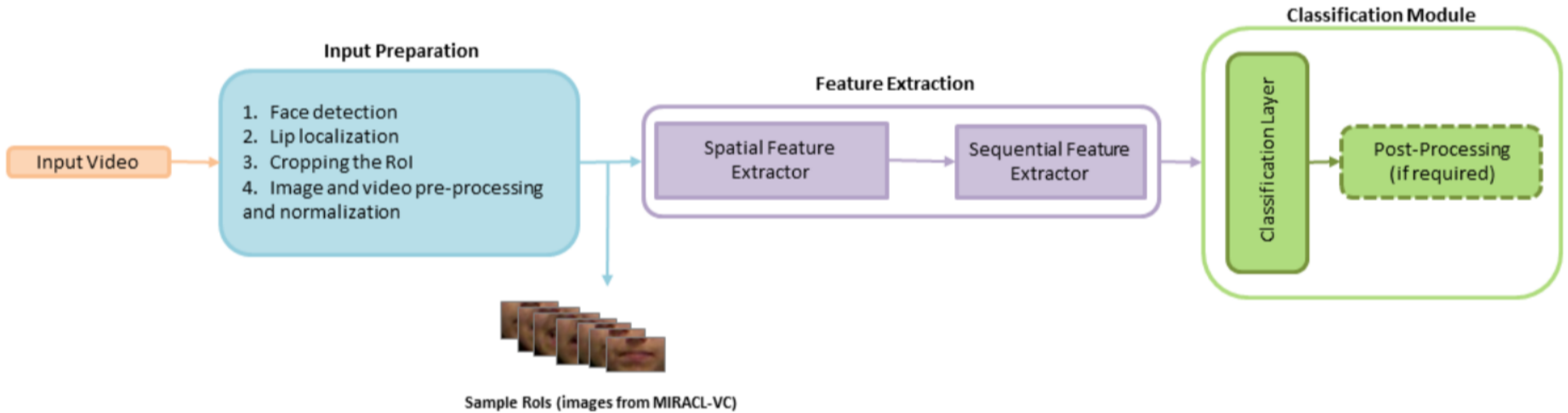  - recognize what is being said from both audio and visual information



Viseme

# Challenges

- Subject dependent factors
  - speaker variation

- Video quality factors
  - pose variation
  - unsynchronized audio & video

- Content-based factors
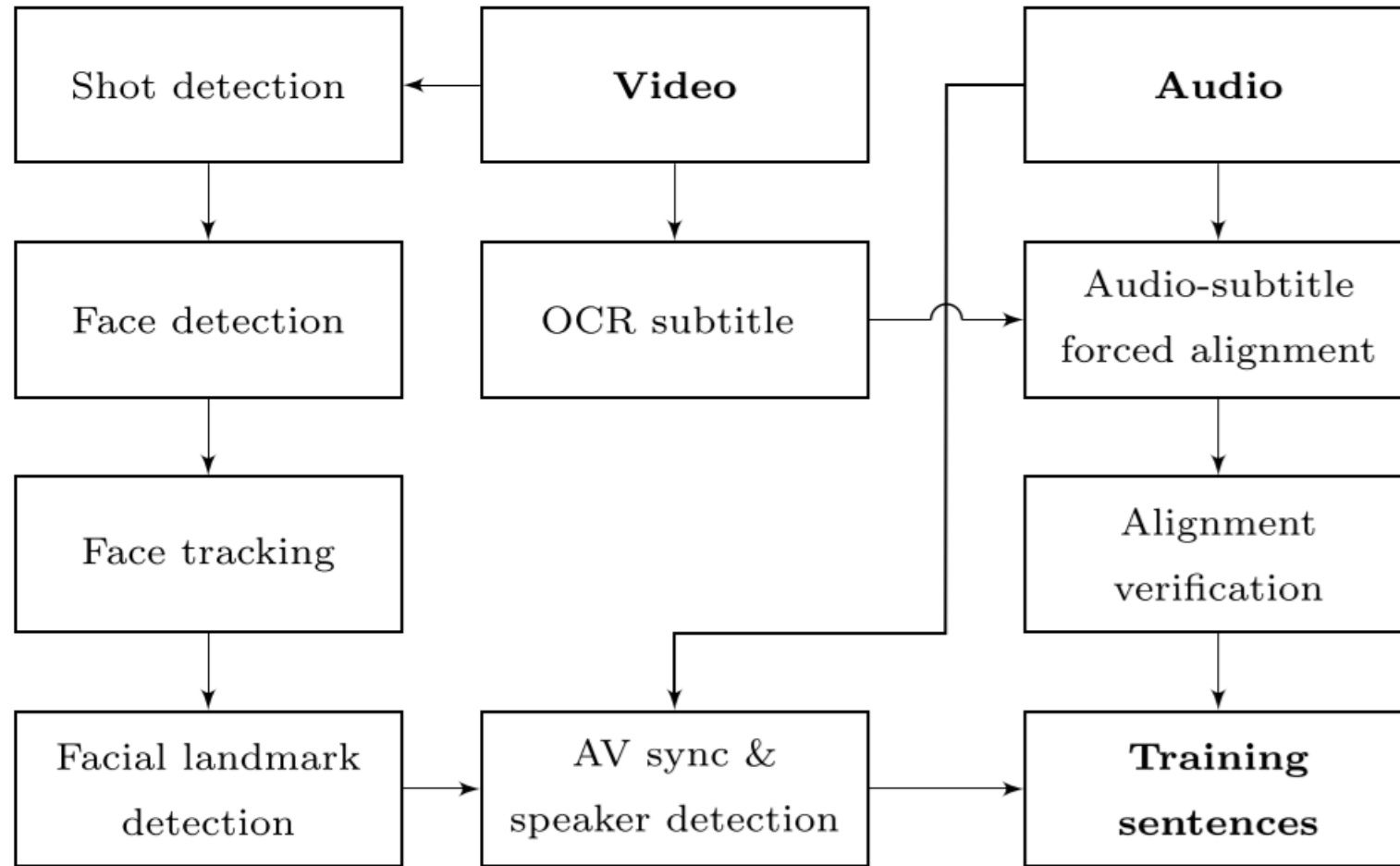  - homophones (same viseme but different phoneme)

- Single modality dominate

- …

# Pipeline of Lip Reading



**Input Preparation**

Input Video →

1. Face detection
2. Lip localization
3. Cropping the RoI
4. Image and video pre-processing and normalization

Sample RoIs (images from MIRACL-VC)

**Feature Extraction**

Spatial Feature Extractor → Sequential Feature Extractor

**Classification Module**

Classification Layer → Post-Processing (if required)

# Datasets

| Dataset | Source | Split | Dates | # Spk. | # Utt. | Word inst. | Vocab | # hours |
|---|---|---|---|---|---|---|---|---|
| GRID [16] | - | - | - | 51 | 33,000 | 165k | 51 | 27.5 |
| MODALITY [17] | - | - | - | 35 | 5,880 | 8,085 | 182 | 31 |
| LRW [12] | BBC | Train-val | 01/2010 - 12/2015 | - | 514k | 514k | 500 | 165 |
| | | Test | 01/2016 - 09/2016 | - | 25k | 25k | 500 | 8 |
| LRS [11] † | BBC | Train-val | 01/2010 - 02/2016 | - | 106k | 705k | 17k | 68 |
| | | Test | 03/2016 - 09/2016 | - | 12k | 77k | 6,882 | 7.5 |
| MV-LRS [14] † | BBC | Pre-train | 01/2010 - 12/2015 | - | 430k | 5M | 30k | 730 |
| | | Train-val | 01/2010 - 12/2015 | - | 70k | 470k | 15k | 44.4 |
| | | Test | 01/2016 - 09/2016 | - | 4,305 | 30k | 4,311 | 2.8 |
| **LRS2-BBC** | BBC | Pre-train | 01/2010 - 02/2016 | - | 96k | 2M | 41k | 195 |
| | | Train-val | 01/2010 - 02/2016 | - | 47k | 337k | 18k | 29 |
| | | Test | 03/2016 - 09/2016 | - | 1,243 | 6,663 | 1,693 | 0.5 |
| | | Text-only | 01/2016 - 02/2016 | - | 8M | 26M | 60k | - |
| **LRS3-TED** | TED & TEDx (YouTube) | Pre-train | - | 5,543 | 132k | 4.2M | 52k | 444 |
| | | Train-val | - | 4,004 | 32k | 358k | 17k | 30 |
| | | Test | - | 451 | 1,452 | 11k | 2,136 | 1 |
| | | Text-only | - | 5,543 | 1.2M | 7.2M | 57k | - |

# LRS dataset pipeline

# Evaluation Metrics

- Error Rate
  - PER (phoneme error rate)
  - CER (character error rate)
  - WER (word error rate)

$$ER = \frac{S + D + I}{N}$$

- BLEU (BiLingual Evaluation Understudy)
  - a Method for Automatic Evaluation of Machine Translation

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}.$$

Then,

$$\text{BLEU} = BP \cdot \exp \left( \sum_{n=1}^{N} w_n \log p_n \right).$$

c is the length of the candidate translation
r is the effective reference corpus length

$$p_n = \frac{\sum\limits_{C \in \{Candidates\}} \sum\limits_{n\text{-}gram \in C} Count_{clip}(n\text{-}gram)}{\sum\limits_{C' \in \{Candidates\}} \sum\limits_{n\text{-}gram' \in C'} Count(n\text{-}gram')}.$$

# Methods

| Method | Year | WER on LRS3 | Create Dataset | Train on |
|---|---|---|---|---|
| LipNet | 2016 | | | GRID |
| WAS | 2017 | | LRS | LRS |
| TM-seq2seq | 2018 | 58.9 | LRS2-BBC, LRS3-TED | MV-LRS, LRS2, LRS3 |
| CTC-V2P | 2018 | 55.1 | LSVSR | LSVSR |
| RNN-T | 2019 | 33.6 | YT | YT |
| VTP | 2021 | 30.7 | TEDx_ext | LRS2, LRS3, MV-LRS, TEDX_ext |
| AVHuBERT | 2022 | 26.9 | | VoxCeleb2, LRS3 |

# LipNet



t frames     STCNN + Spatial Pooling     Bi-GRU     Linear     CTC loss
(x3)           (x2)

Figure 1: LipNet architecture. A sequence of $T$ frames is used as input, and is processed by 3 layers of STCNN, each followed by a spatial max-pooling layer. The features extracted are processed by 2 Bi-GRUs; each time-step of the GRU output is processed by a linear layer and a softmax. This end-to-end model is trained with CTC.

# WLAS

A B C D E F G H I J K L M N O P Q
R S T U V W X Y Z 0 1 2 3 4 5 6 7
8 9 , . ! ? : ' [sos] [eos]
[pad]

Table 10. The output characters



Figure 1. *Watch, Listen, Attend and Spell* architecture. At each time step, the decoder outputs a character $y_i$, as well as two attention vectors. The attention vectors are used to select the appropriate period of the input visual and audio sequences.

$$f_i^v = \text{CNN}(x_i^v) \tag{5}$$

$$h_i^v, o_i^v = \text{LSTM}(f_i^v, h_{i+1}^v) \tag{6}$$

$$s^v = h_1^v \tag{7}$$

$$h_j^a, o_j^a = \text{LSTM}(x_j^a, h_{j+1}^a) \tag{8}$$

$$s^a = h_1^a \tag{9}$$

$$h_k^d, o_k^d = \text{LSTM}(h_{k-1}^d, y_{k-1}, c_{k-1}^v, c_{k-1}^a) \tag{10}$$

$$c_k^v = \mathbf{o}^v \cdot \text{Attention}^{\mathbf{v}}(h_k^d, \mathbf{o}^v) \tag{11}$$
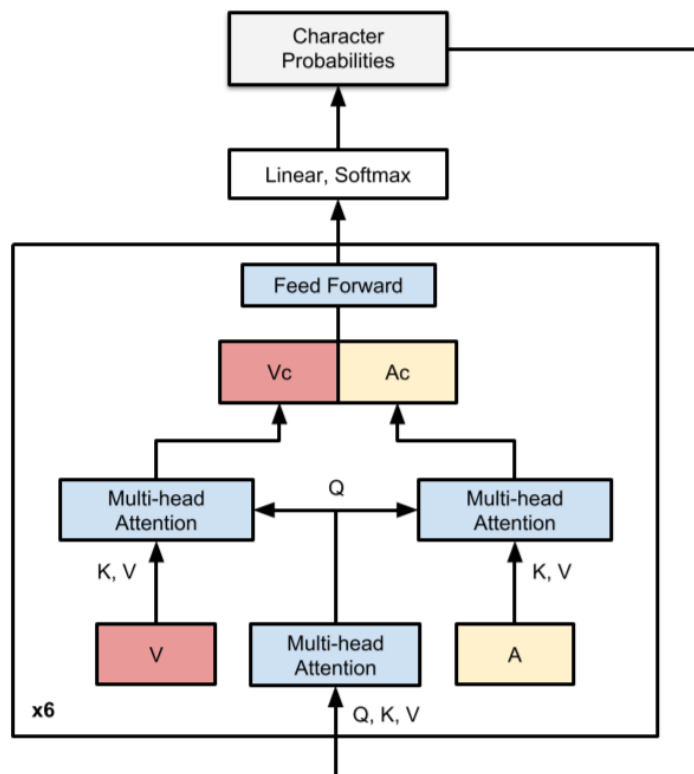
$$c_k^a = \mathbf{o}^a \cdot \text{Attention}^{\mathbf{a}}(h_k^d, \mathbf{o}^a) \tag{12}$$

$$P(y_i | \mathbf{x}^v, \mathbf{x}^a, y_{<i}) = \text{softmax}(\text{MLP}(o_k^d, c_k^v, c_k^a)) \tag{13}$$
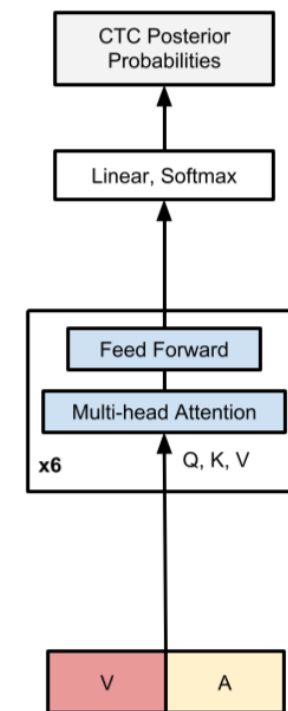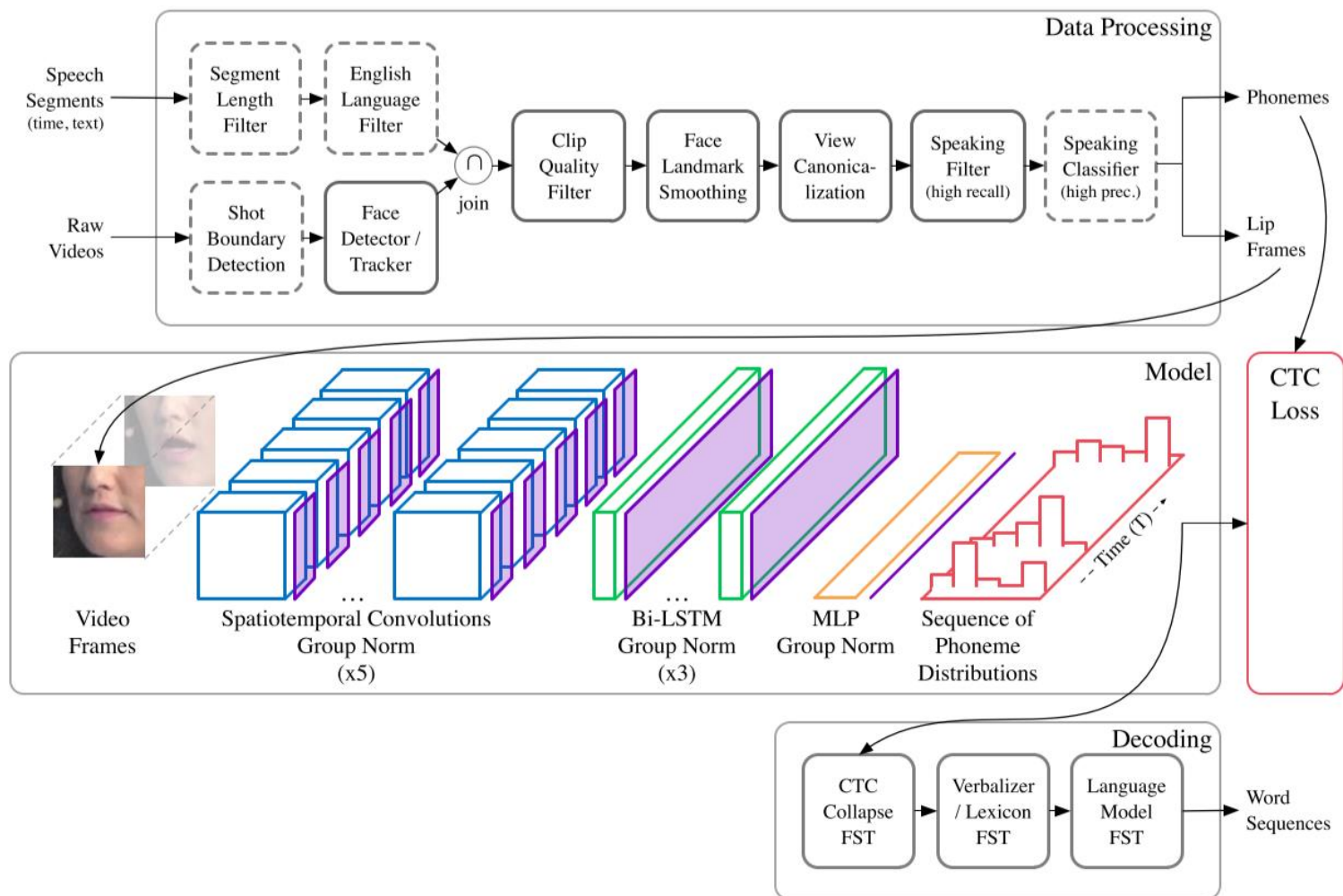
# TM-seq2seq



a. Common Encoder
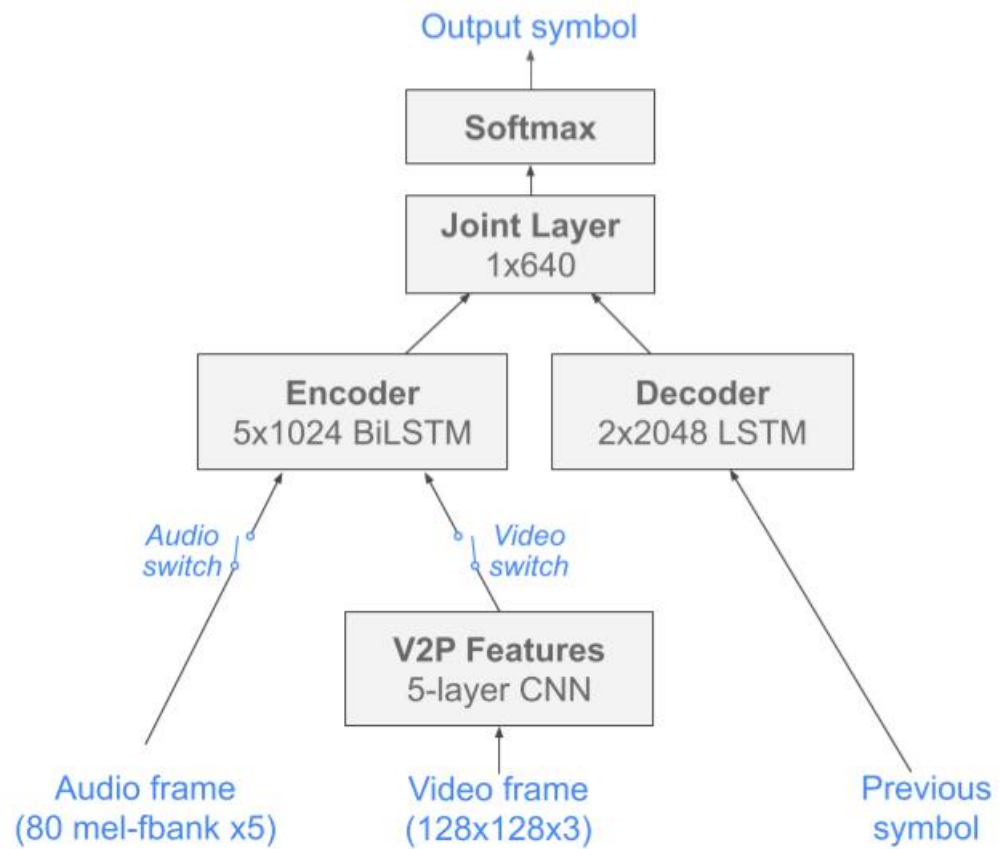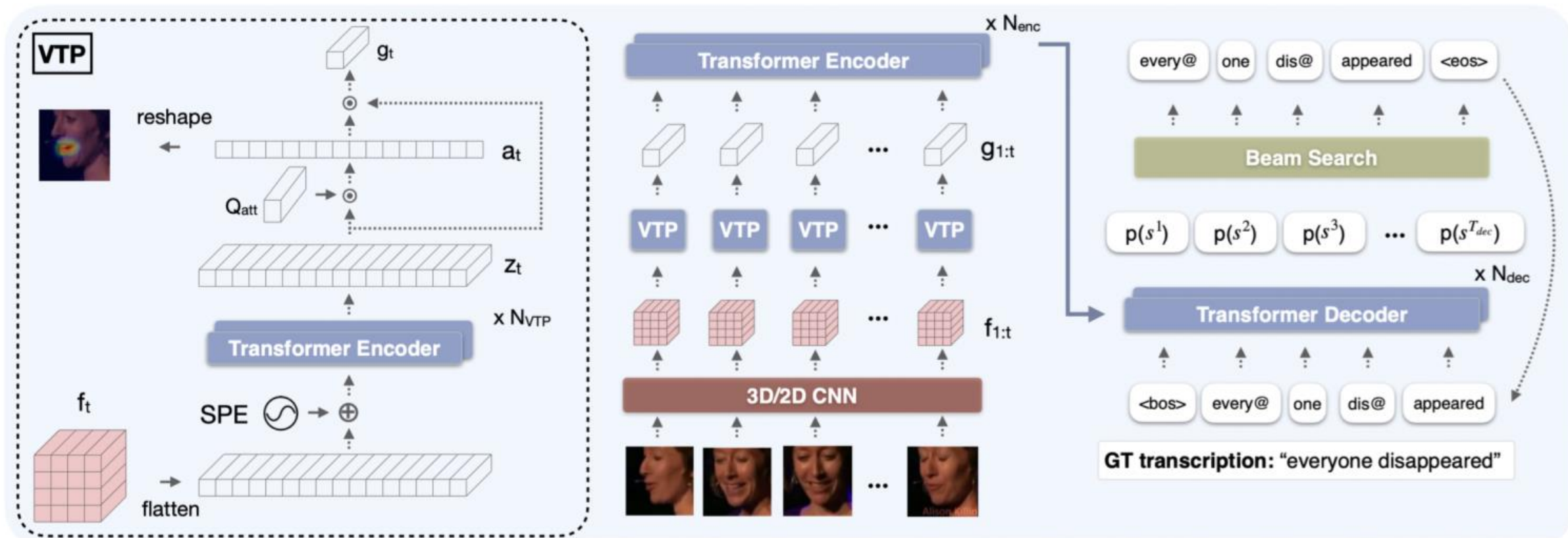
b. Transformer seq2seq
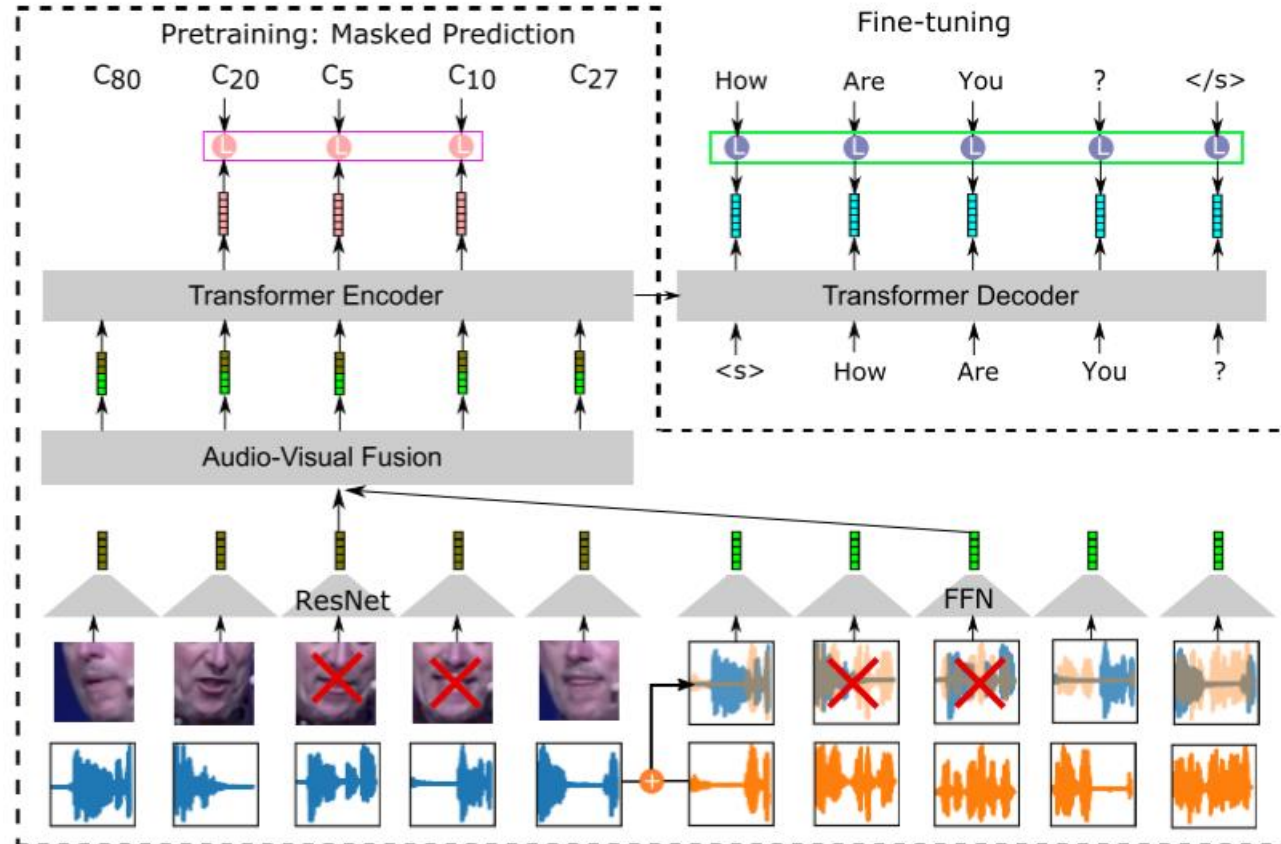
c. Transformer CTC

# CTC-V2P

# RNN-T

# VTP

# AV-HuBERT

# Conclusion

- It is critical to develop a better way to capture both spatial and temporal information
- It is (used to be) important to build a dataset
  - a pipeline for pre-process is needed
  - language variation, vocabulary variation, ...
  - different video quality
  - bigger dataset
- Training protocol is important
  - curriculum learning
  - pre-training
  - avoid single modality dominates
  - ...

# Reference

| Method | Year | Paper Title | PDF |
|---|---|---|---|
| LipNet | 2016 | LIPNET: END-TO-END SENTENCE-LEVEL LIPREADING | https://arxiv.org/abs/1611.05358v2 |
| WAS | 2017 | Lip Reading Sentences in the Wild | https://arxiv.org/abs/1611.05358v2 |
| TM-seq2seq | 2018 | Deep Audio-visual Speech Recognition | https://arxiv.org/abs/1809.02108 |
| CTC-V2P | 2018 | LARGE-SCALE VISUAL SPEECH RECOGNITION | https://arxiv.org/abs/1807.05162 |
| RNN-T | 2019 | RECURRENT NEURAL NETWORK TRANSDUCER FOR AUDIO-VISUAL SPEECH RECOGNITION | https://arxiv.org/abs/1911.04890 |
| VTP | 2021 | Sub-word Level Lip Reading With Visual Attention | https://arxiv.org/abs/2110.07603v2 |
| AVHuBERT | 2022 | Robust Self-Supervised Audio-Visual Speech Recognition | https://arxiv.org/abs/2201.01763 |