

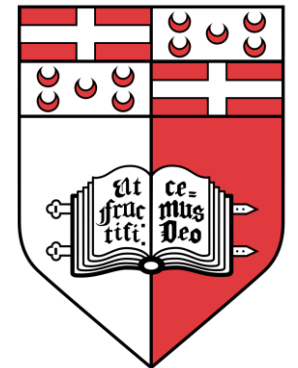
# Sign Language Detection “in the Wild” with Recurrent Neural Networks

**Mark Borg**  
**Kenneth P. Camilleri**

Systems and Control Engineering,  
**University of Malta**

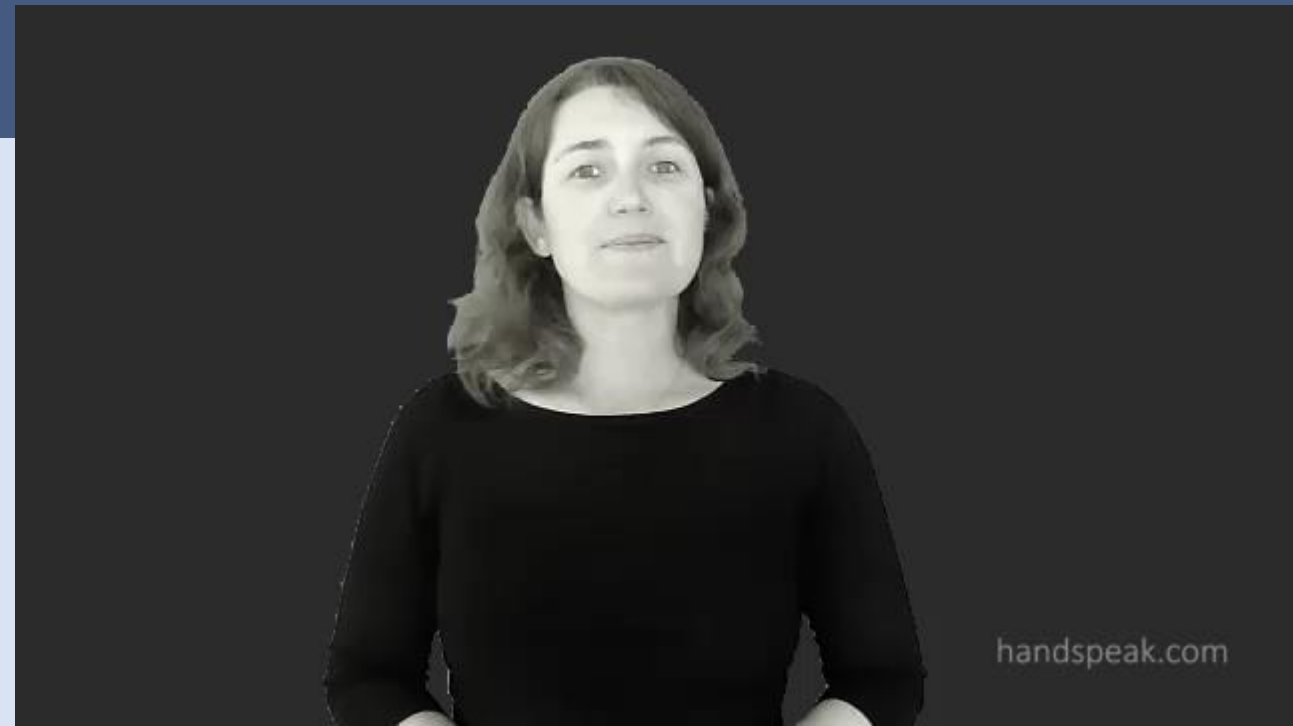
**IEEE 2019 ICASSP**

International Conference on Acoustics, Speech  
and Signal Processing



# Sign Languages

- **Visual languages**
- Multi-modal
- Concurrent modalities
- **Articulators:**
  - Manual
    - Hand motion
    - Hand shapes
    - Place of articulation
  - Non-Manual
    - Mouth patterns
    - Facial Expressions
    - Body posture



**American Sign Language (ASL): PERSON WHATEVER IX-they JUDGE IX1 REALLY WASTE, MINUTE [shook-head] INSTEAD-OF MINUTE IX WHAT-conj LOVE ACCEPT WHO**

**English equivalent: For every minute we judge, we have squandered a minute we could have used to accept and love someone.**

Source: HandSpeak

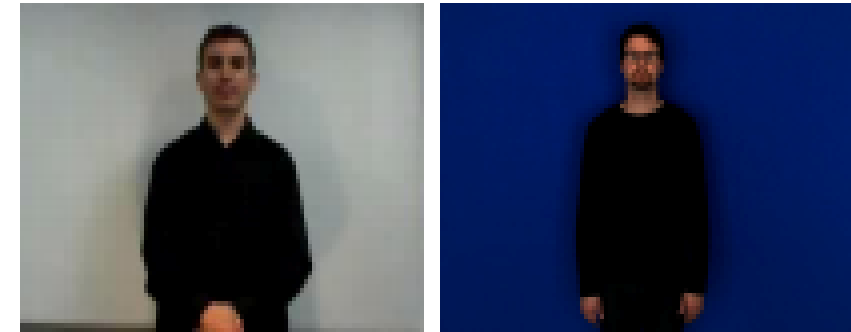


# Sign Language Detection – state of the art

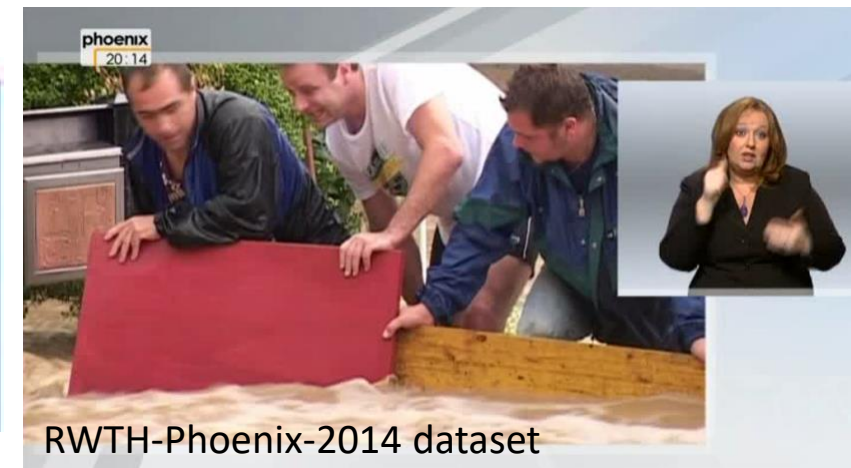
- Monteiro et al. (2012 SIGACCESS)
  - Face detection, background subtraction
  - Hand-crafted visual features: velocity-based
  - SVM
- Shipman et al. (2015 JCDL, 2017 SIGACCESS)
  - Face detection, background subtraction
  - Hand-crafted visual features: polar motion profiles
  - SVM
- Gebre et al. (2013 ICIP)
  - Face detection, skin detection
  - hand-crafted visual features
  - random forests
- Yanovich (2016 LREC)
  - Identification of major sign language constructs: fingerspellings, classifiers, ...
  - Hand-crafted visual features
  - k-NN classifier
- Gebre et al. (2014 Comp. Ling.)
  - Identification of particular sign languages: BSL, DSL, FBSL, FSL, GSL, NGT
  - Sparse auto-encoder and 3D CNN

# Datasets – the need for sign language detection datasets

- No signing in generic video action recognition datasets, like AVA, THUMOS, ...
- Previous work in SL detection
  - Datasets not made publicly available
  - Small size (~200 videos)
- Sign Language Recognition (ASLR) datasets, Phoenix, SIGNUM, VGG BBC pose, ...
  - Trimmed
  - Captured under constrained conditions



Src: Dreuw et al. (2010)



RWTH-Phoenix-2014 dataset

# “Signing in the Wild” dataset

- Untrimmed videos
- Each video can include multiple signing and non-signing events
- Harvested from YouTube

3 categories:

Signing

Speaking

Other

- **1120** video segments
- Each video segment:
  - Up to 6.6 minutes (sampled at 5 Hz)
  - Up to 2000 frames long
- **1.45 million frames** in total

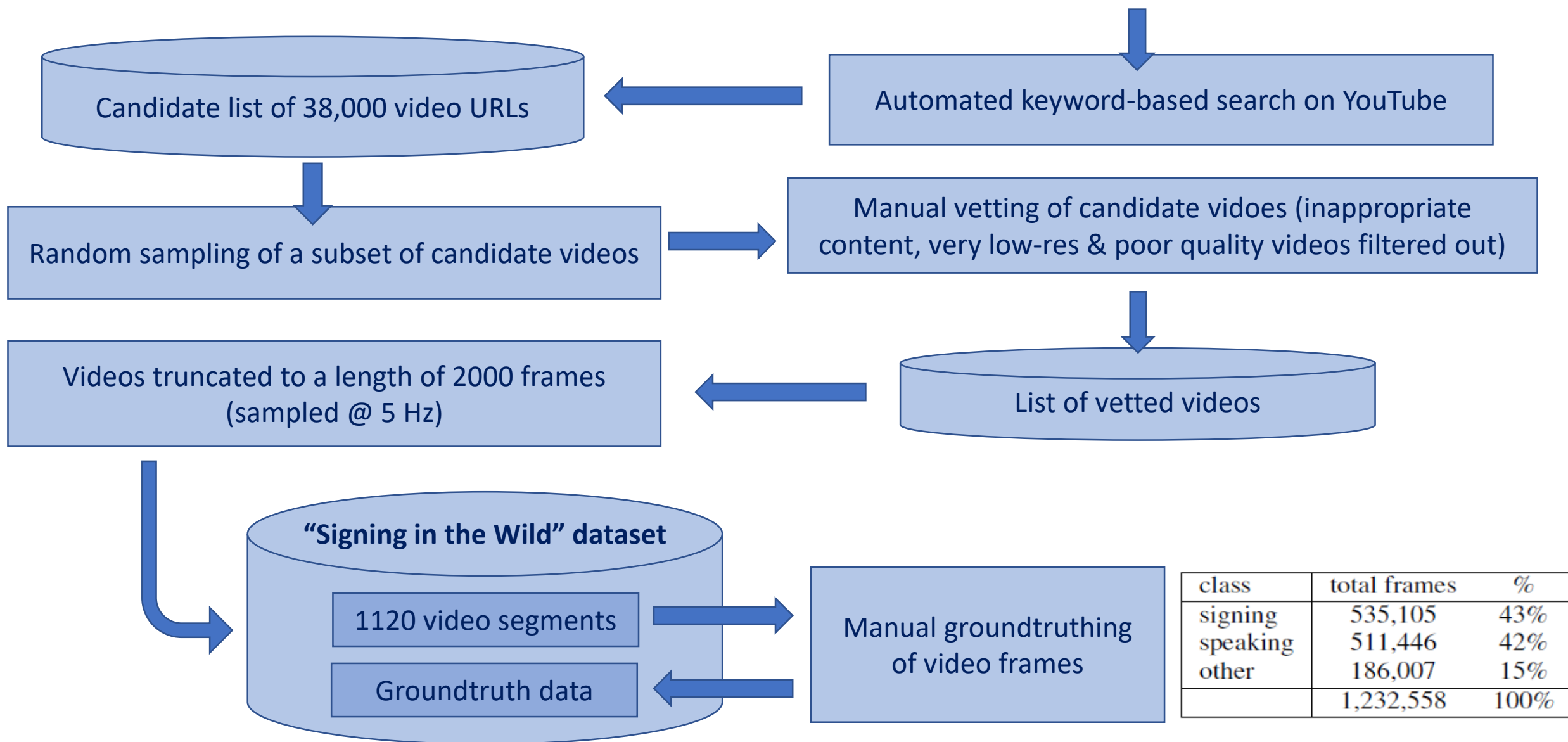
Groundtruthing:

- Frame-level
- 10-frame temporal context
- **1.23 million frames**

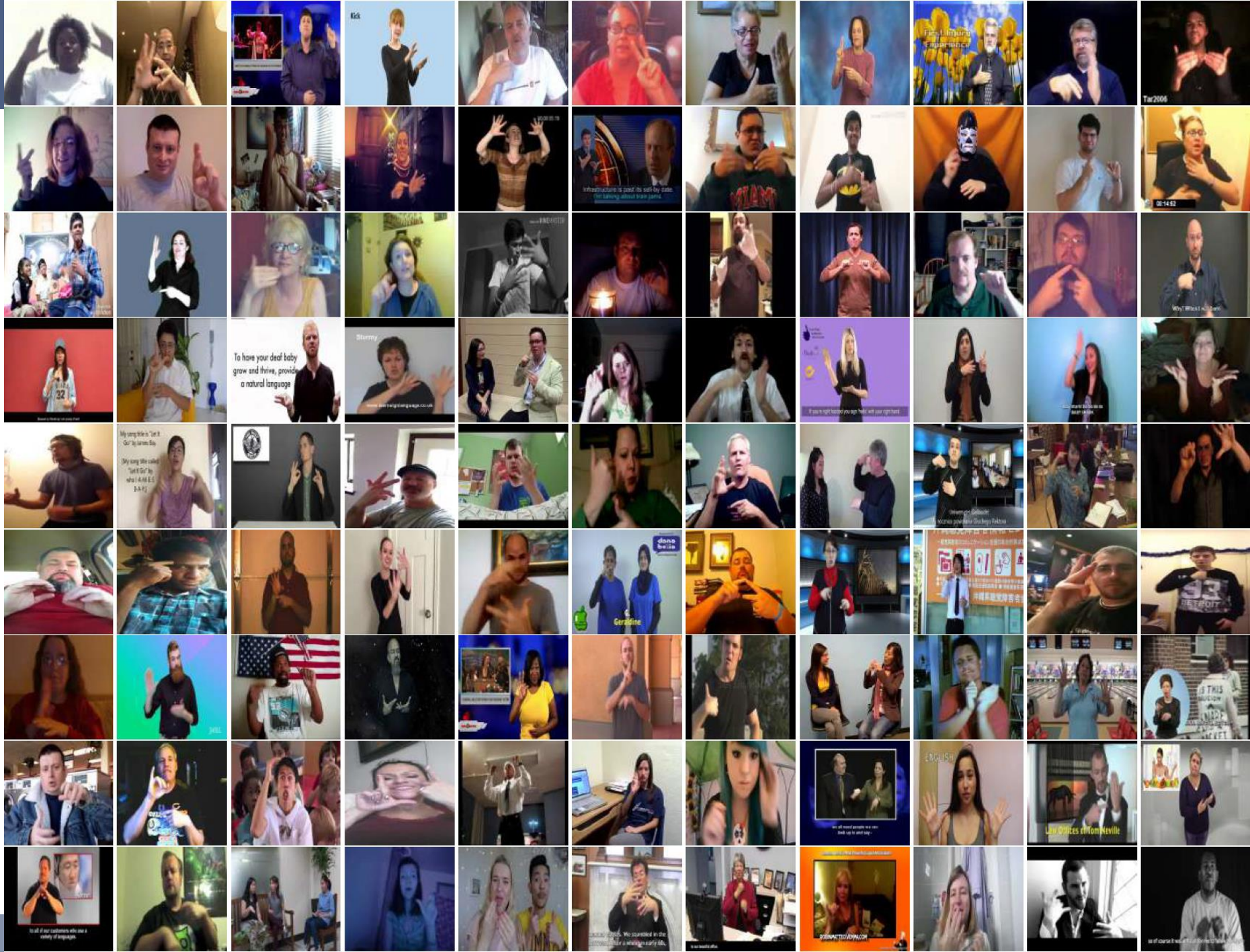
Publicly available:

- **IEEE DataPort**
- <https://github.com/mark-borg/Signing-in-the-Wild-dataset>

# “Signing in the Wild” dataset



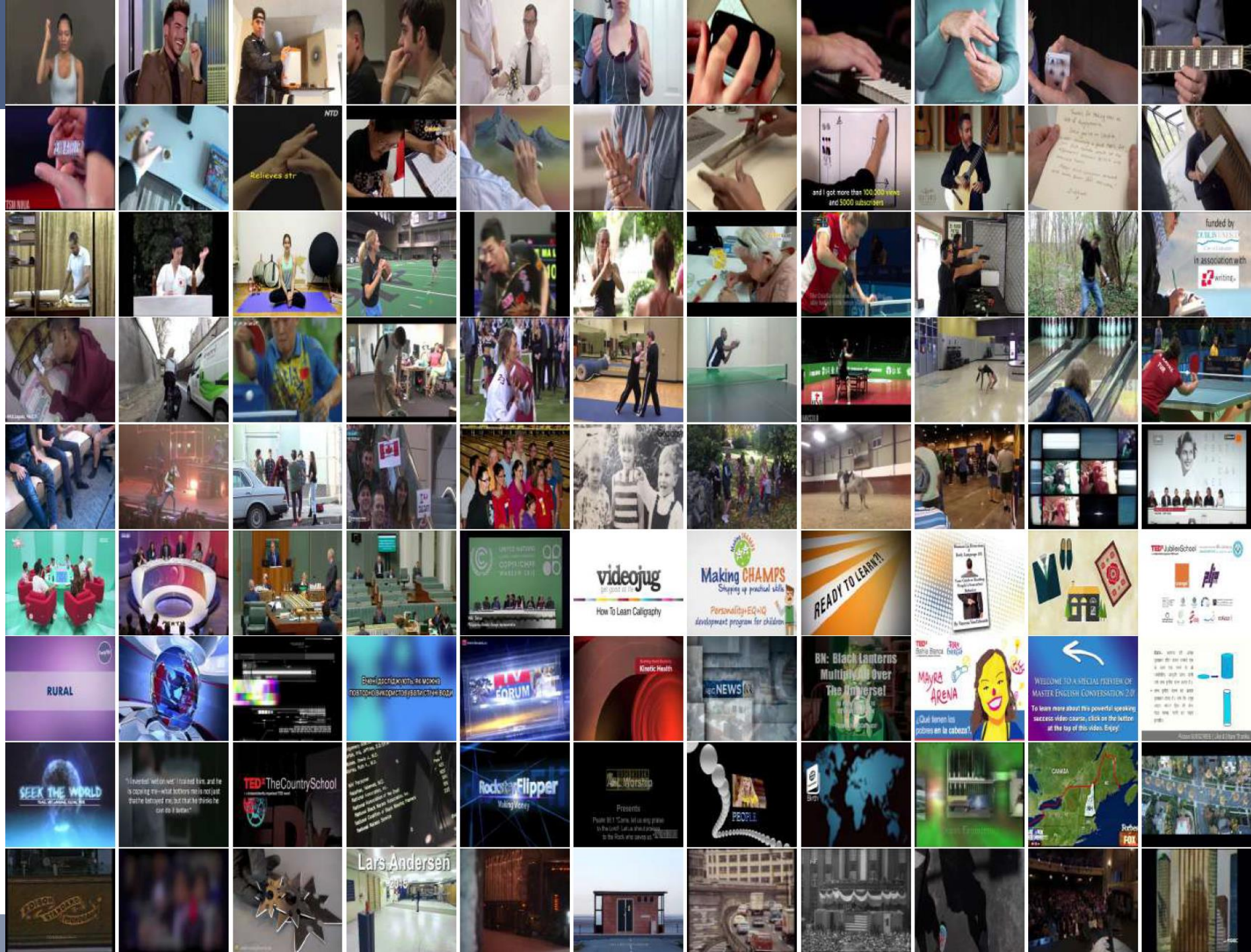
# Example frames from class signing







Example frames  
from class  
other



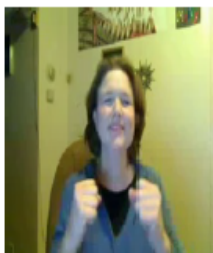
# Sign Language Detection – proposed approach

- **Automated extraction of features** using a **Convolutional Neural Network (CNN)**
- Combining both **visual features** and **motion features**
- Use of a **Recurrent Neural Network (RNN)** to handle the dynamic temporal patterns present in sign languages

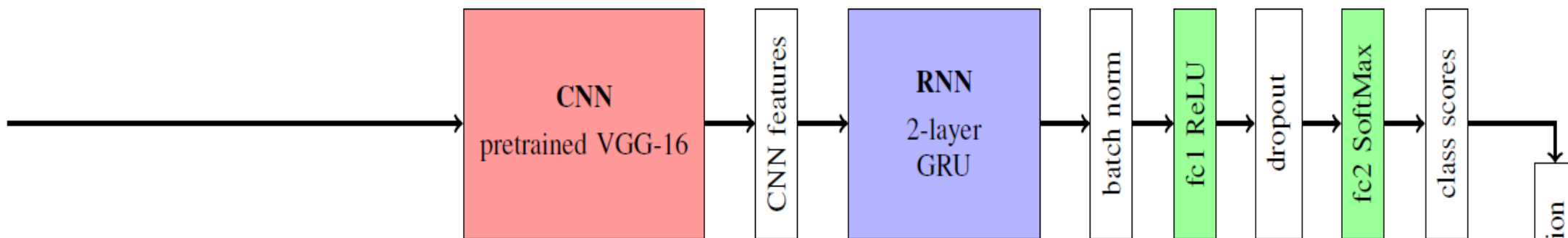
# Proposed architecture

- Two-stream approach (Simonyan, 2014)

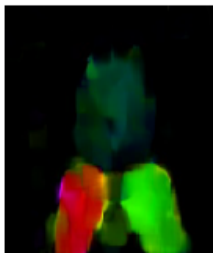
video frame data



224, 224, 3



stacked flow data



224, 224, 15

or

motion history



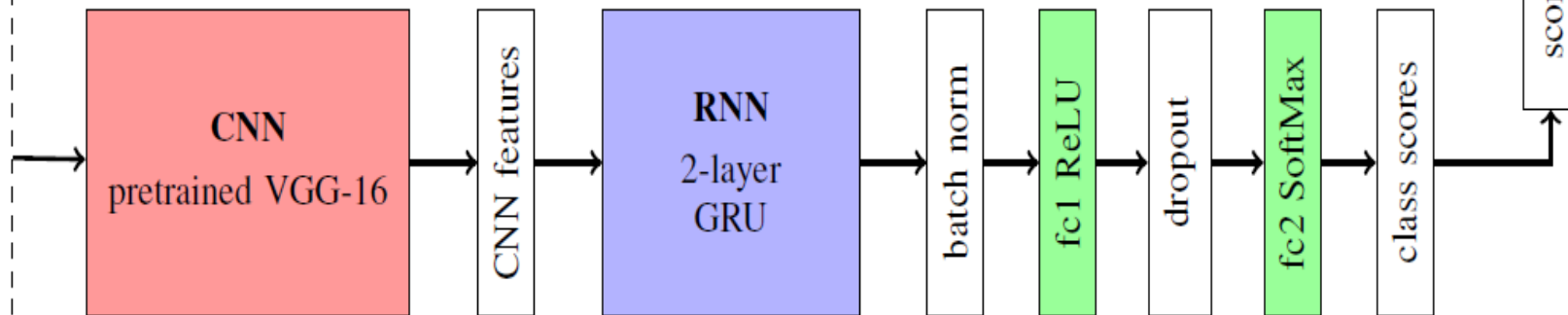
224, 224, 3

or

frame difference



224, 224, 5

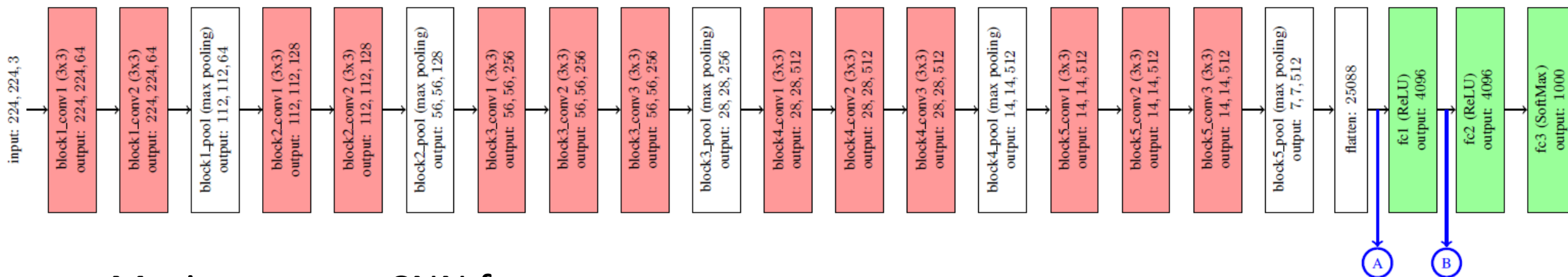


# Proposed approach

- Motion stream:
  - Performance vs. computational efficiency
- Investigated:
  - Optical Flow
  - Motion History Images (MHI)
  - Multi-frame differencing

# CNN features

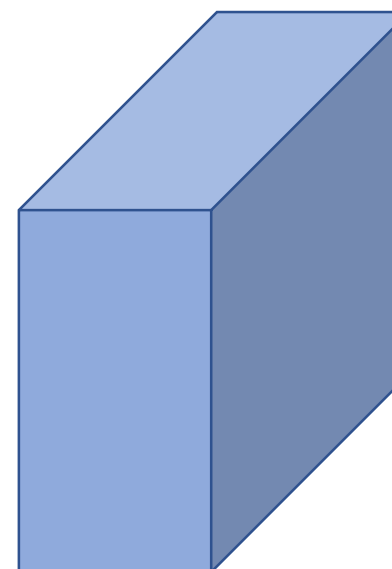
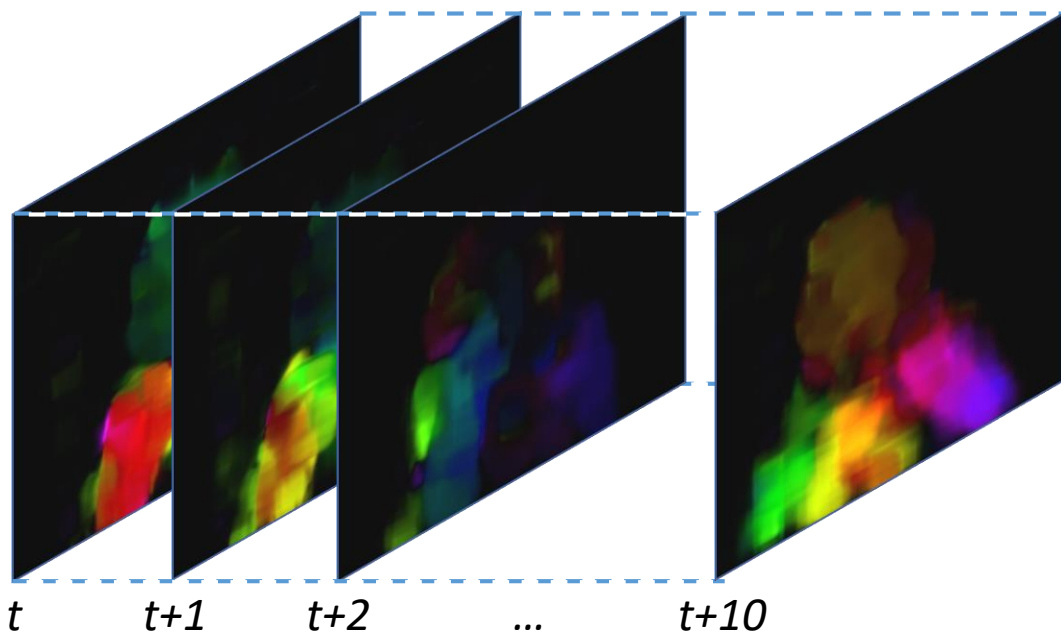
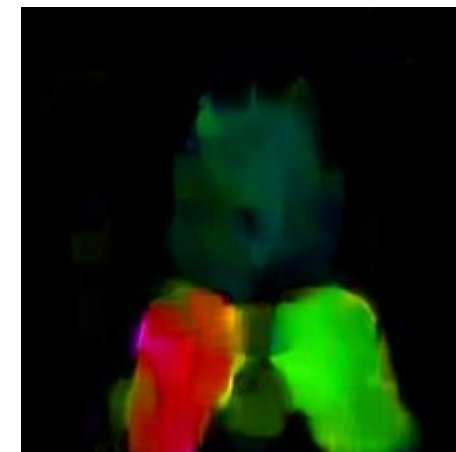
- CNN streams
  - Pre-trained VGG-16 (Simonyan 2014)
  - CNN features:
    - Ⓐ  $7 \times 7 \times 512 = 25088$  feature map from 'block5\_conv3' layer
    - Ⓑ 4096 feature map from 'fc1' layer



- Motion stream CNN features:
  - We use transfer learning from a distant task (unrelated data) vs. Training from scratch (Yosinski et al., 2014)
- No fine-tuning of VGG-16 layers

# Motion data

- Optical flow
  - Dense optical flow (Farneback's algorithm)
  - Encoded as RGB
    - Flow vector magnitude  $\rightarrow$  luminance channel
    - Flow vector angle  $\rightarrow$  chrominance channels



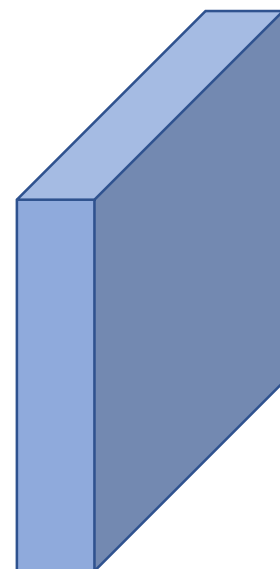
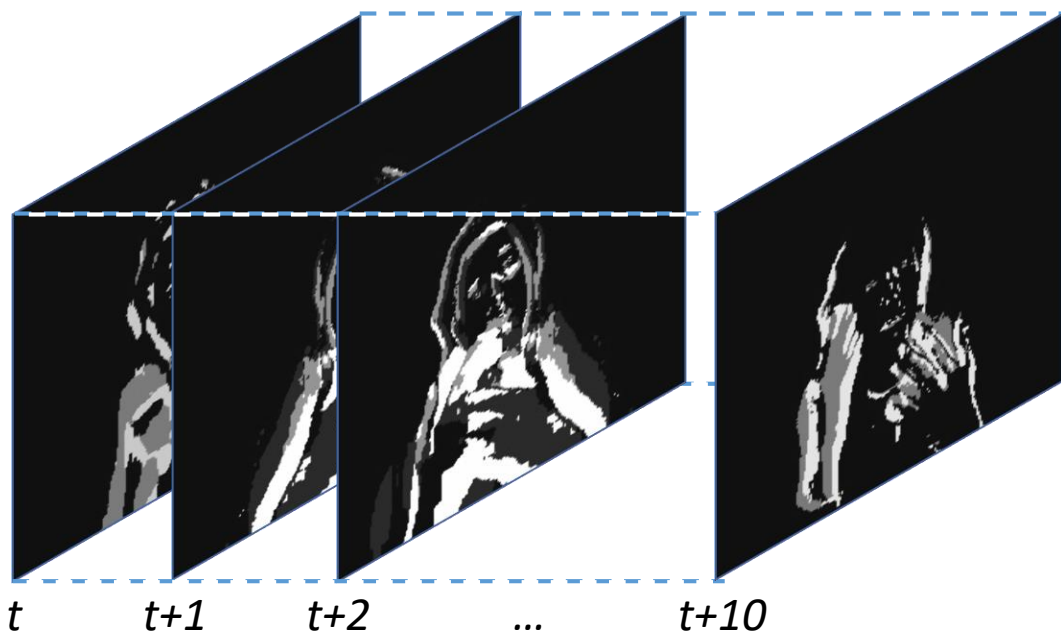
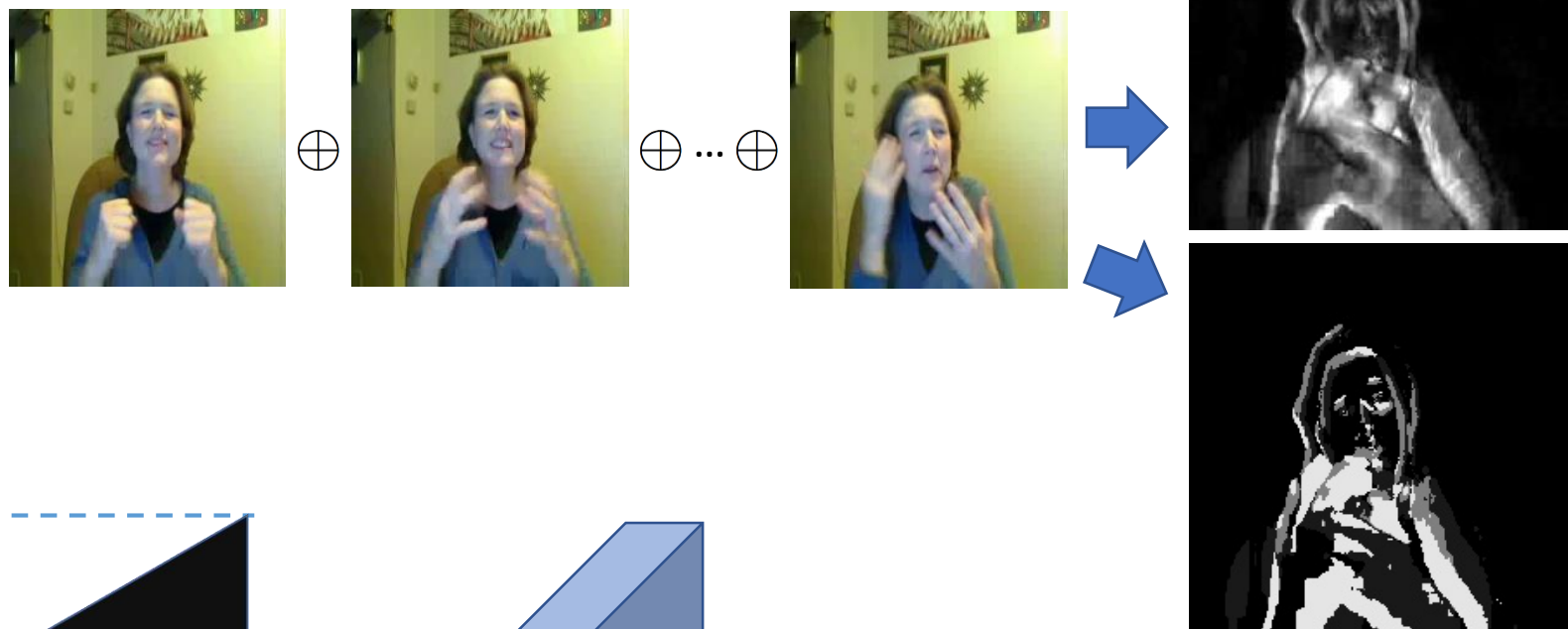
224 x 224 x 30



Filter weights of first layer replicated

# Motion data

- Multi-frame differencing
- Motion History Images
  - 5 frame temporal window



$224 \times 224 \times 10$



Filter weights of first layer  
averaged to 1-channel, then  
replicated 10 times



# RNN

- Various RNN options: LSTMs and GRUs

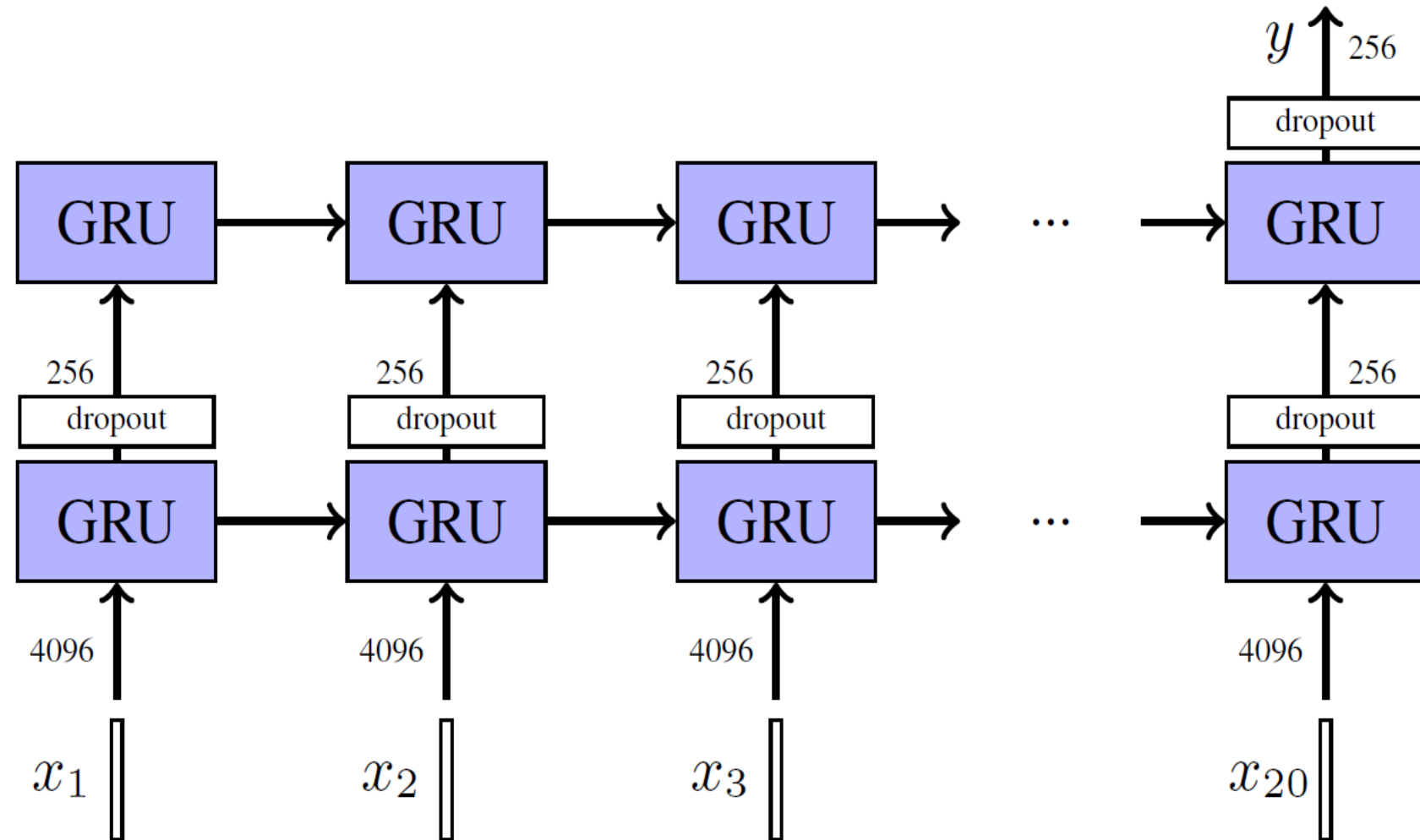
- Stacked RNNs

- 2-layer GRU

- 256 hidden units

- 20 timesteps

- (2.5 seconds with a 5Hz sampling rate)



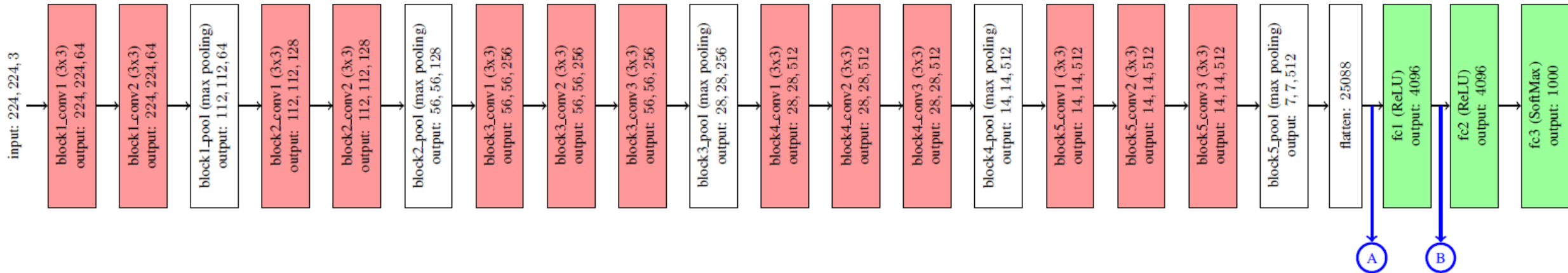
# RNN training

- Stratified partitioning of the dataset
- video frames from a single video appear in only one partition
- 5 fold cross-validation
- Mini-batch stochastic gradient descent (SGD)
- Adam optimizer
- Training for 500 epochs, with early stopping (validation cross-entropy loss)
- Training strategy:
  - Initial mini-batch size of 32, learning rate of 0.001
  - Reduce learning rate when validation loss stops improving for the current combination of mini-batch size and learning rate
  - Increase mini-batch size when no more change in validation loss is observed for the given mini-batch size despite the changes to the learning rate

# Results

- Evaluation of different feature maps from the CNN network

CNN layer	Feature size	Loss (validation set) ↓
VGG-16 block5_conv3 (A)	25088	0.6681
VGG-16 fc1 (B)	4096	<b>0.5037</b>



# Results

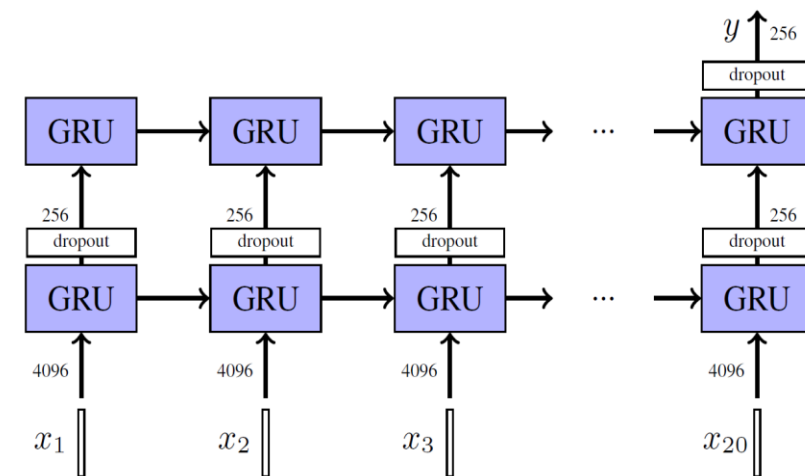
- Evaluation of the individual performance of the different streams, and when fusing both the motion and RGB streams together

Modality	Loss ↓	Accuracy ↑	Time (ms) ↓	
RGB stream only	0.5128	85.01%	—	
optical flow only	<b>0.5387</b>	83.12%	57.8	
MHI only	0.5445	83.67%	17.1	
multi-frame diff. only	0.5738	84.08%	<b>9.7</b>	
RGB stream + optical flow stream	—	<b>87.67%</b>	—	
RGB stream + MHI stream	—	87.60%	—	
RGB stream + frame diff. stream	—	85.61%	—	

# Results

- Evaluation of different RNN architectures

RNN	layers	trainable parameters	Loss (valid. set) ↓
LSTM	1	4,474,627	0.5144
LSTM	2	4,999,939	0.6413
LSTM	3	5,525,251	0.5714
GRU	1	3,360,259	0.5267
GRU	2	3,754,243	<b>0.5037</b>
GRU	3	4,148,227	0.6028



# Results

- Ablation studies on the proposed RNN network

Model settings	Cross-entropy loss on validation set ↓						
proposed model	<b>0.504</b>						
no batch normalisation	0.609	(≈ 20% increase in loss)					
no dropout layer	0.715	(≈ 42% increase in loss)					
no GRU dropout	0.693	(≈ 38% increase in loss)					
no classifier fc1 layer	0.649	(≈ 29% increase in loss)					
with dropout layer	rate:	0.1	0.2	0.3	<b>0.4</b>	0.5	0.6
	loss:	0.605	0.577	0.575	<b>0.504</b>	0.511	0.602
with GRU dropout	rate:	0.1	0.2	<b>0.3</b>	0.4	0.5	0.6
	loss:	0.628	0.601	<b>0.548</b>	0.649	0.554	0.552

# Results

- Comparison with the state-of-the-art in sign language detection
  - $\approx$  18% improvement over baseline
  - $\approx$  9% improvement when using an RNN versus SVM

Method	Feature type & Classifier	Loss ↓	Precision ↑
baseline method †	hand-crafted features + SVM	1.114	69.23%
baseline+RNN	hand-crafted features + RNN	0.841	78.02%
CNN+SVM	2-stream CNN features + SVM	–	79.15%
our method	2-stream CNN features + RNN	<b>0.573</b>	<b>87.67%</b>

† Shipman et al. (JCDL 2015, ACM SigAccess 2017)

# Results

- Confusion matrix

		Actual label		
		sign	speak	other
Predicted label	sign	4992	481	143
	speak	681	4297	227
	other	65	219	1290

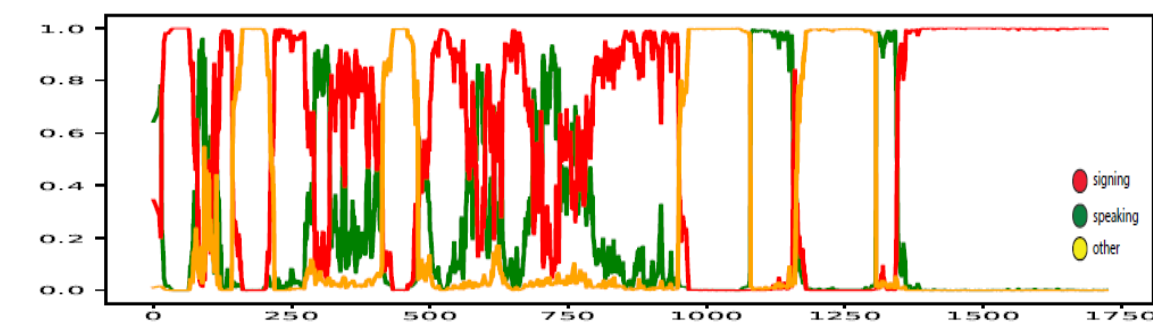
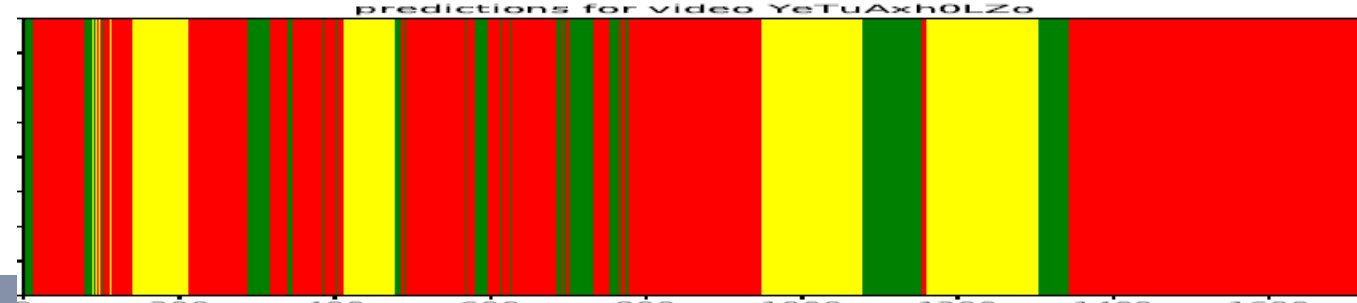
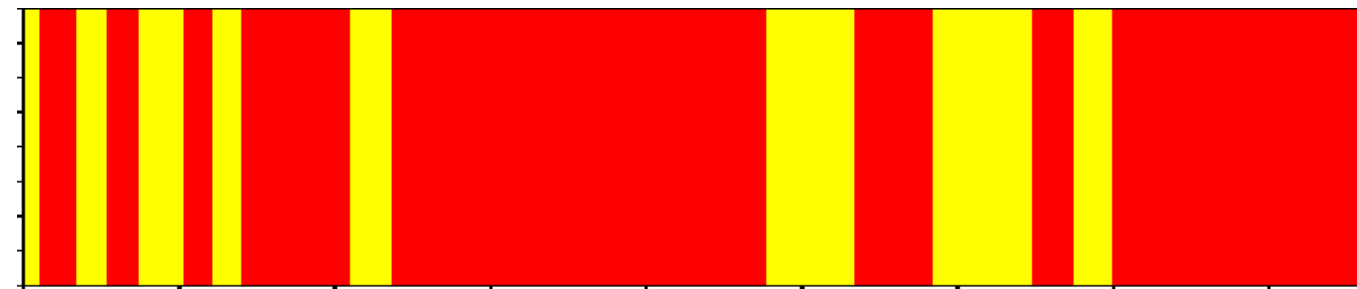
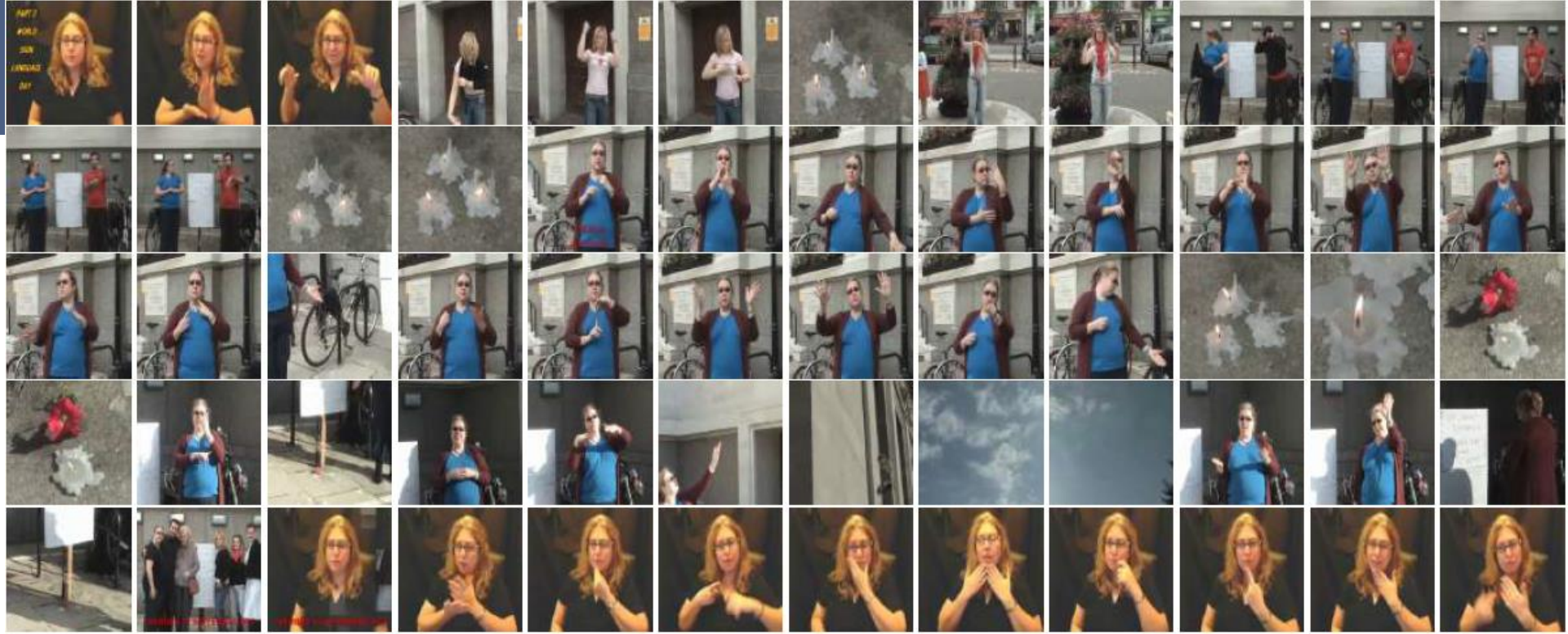


# Results

Video "YeTuAxb0LZo"

Sample video frames  
(every 100)

Groundtruth + prediction  
strips (signing, speaking,  
other)

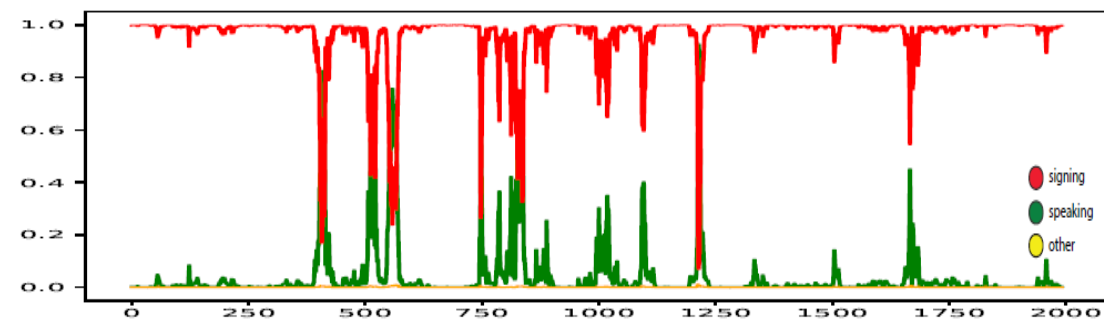
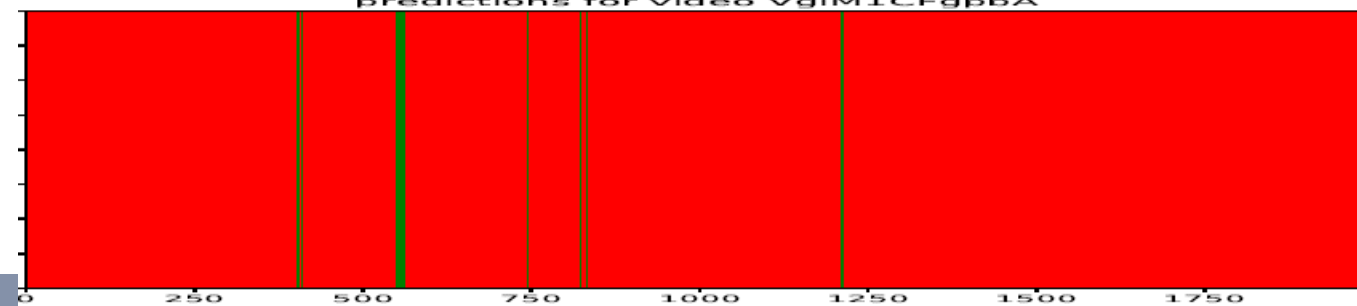


# Results

Video "VgiM1CFgpbA"



predictions for video VgiM1CFgpbA

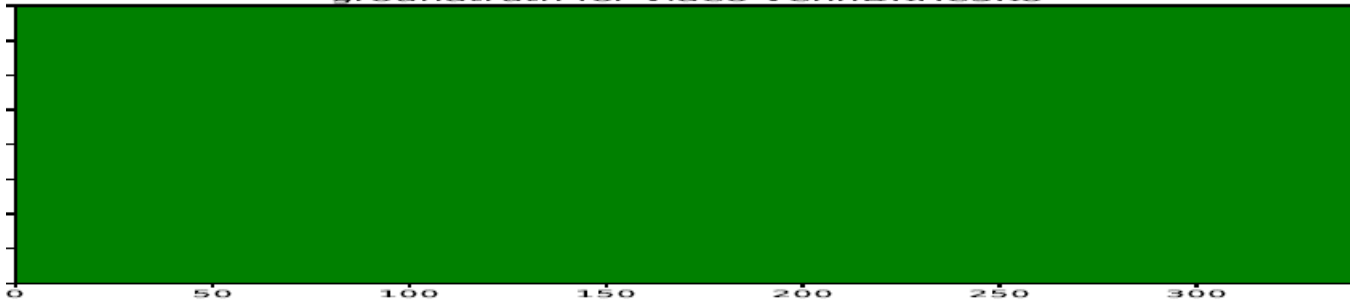


# Results

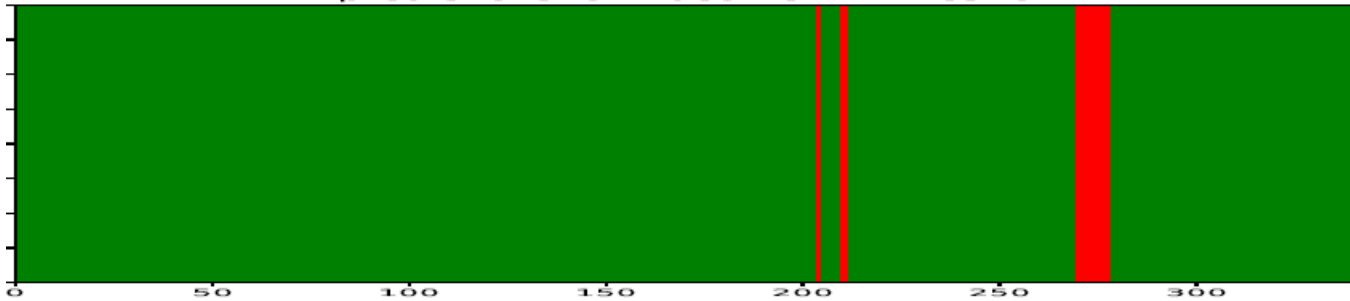
Video "VehhDKHe5Ko"



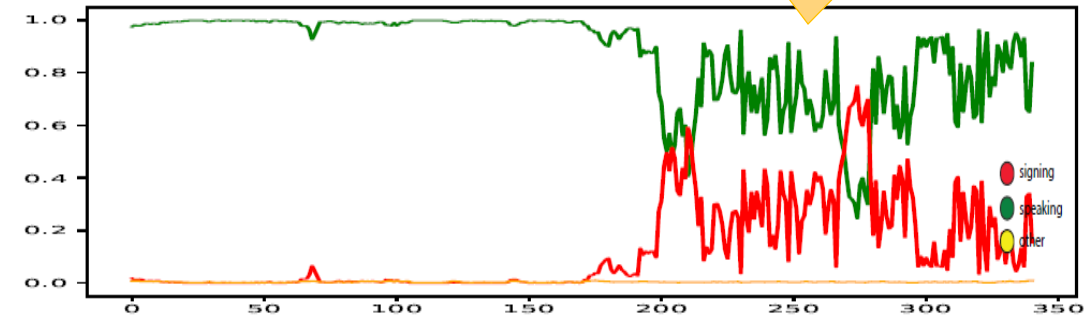
groundtruth for video VehhDKHe5Ko



predictions for video VehhDKHe5Ko



Second part of this video contains hand motion, clapping and singing

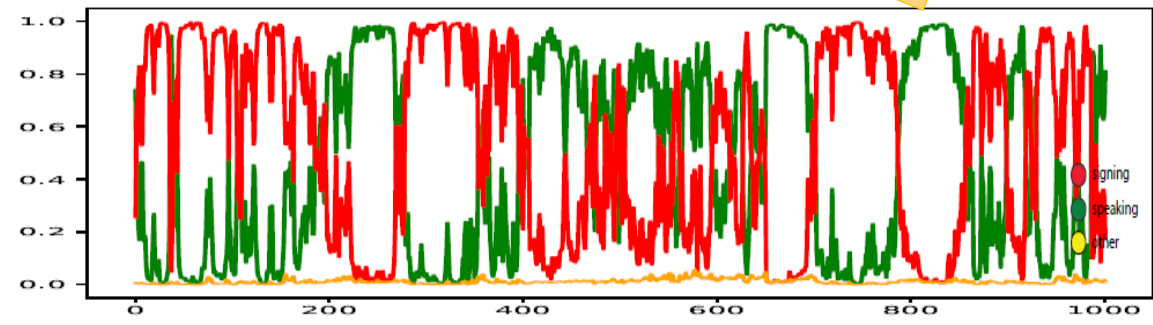
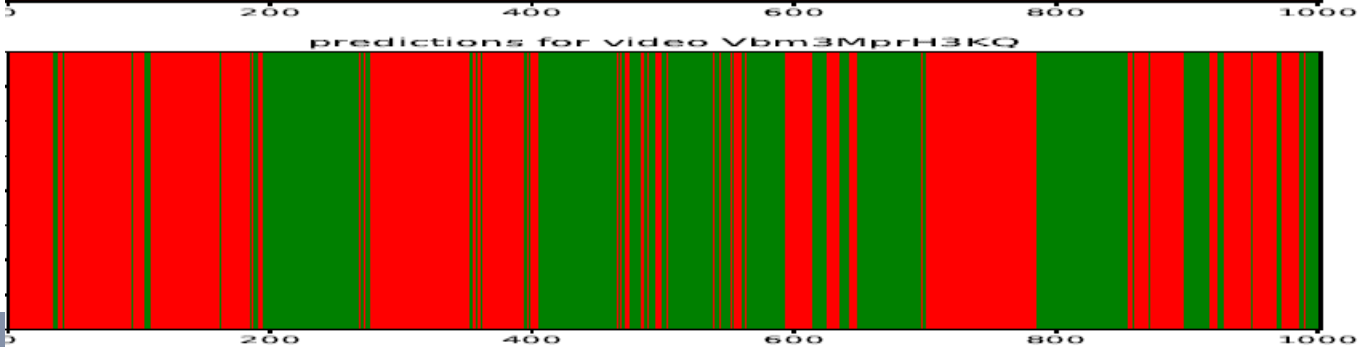


# Results

Video "Vbm3MprH3KQ"



Two signers around a table. Several segments mislabeled as speech

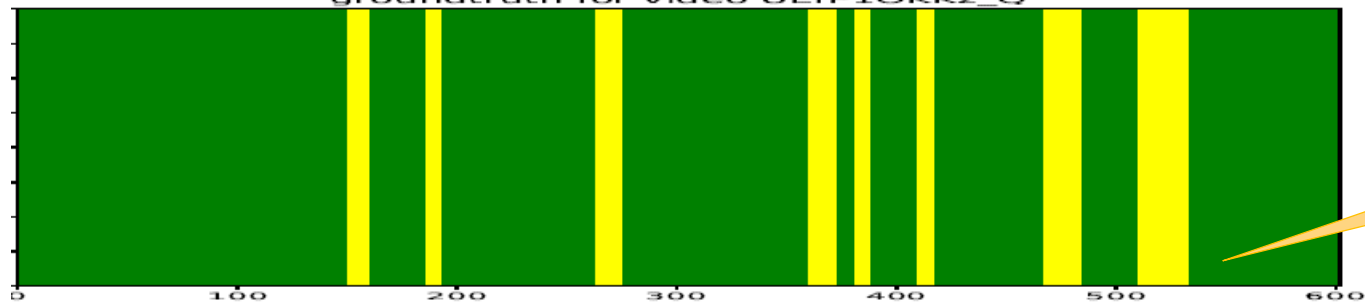


# Results

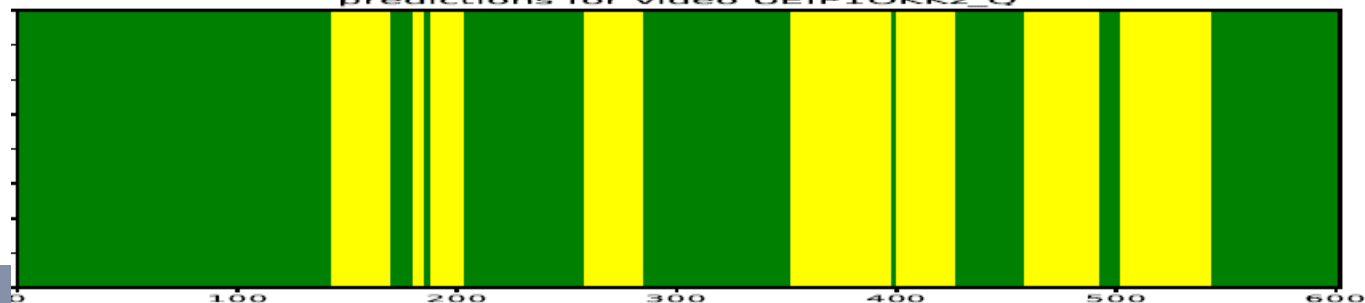
Video "UEfP1OKKz\_Q"



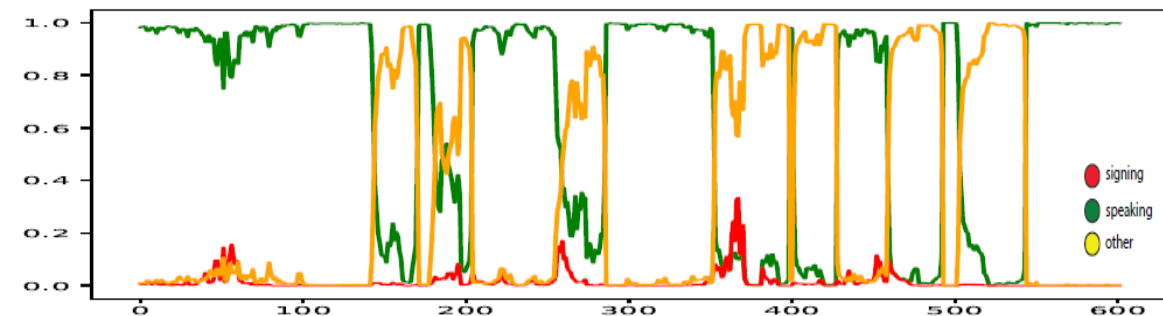
groundtruth for video UEfP1OKKz\_Q



predictions for video UEfP1OKKz\_Q

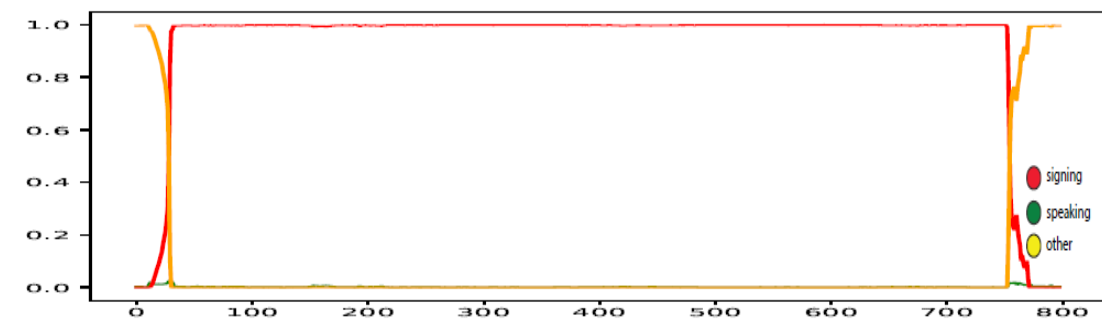
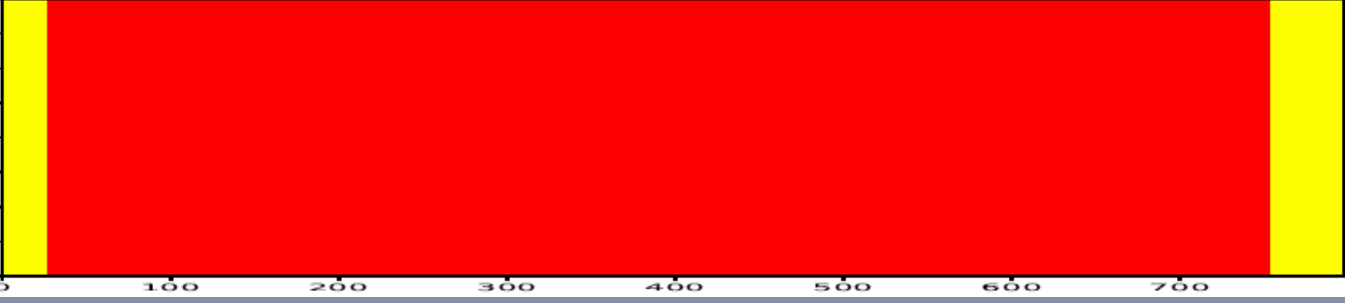
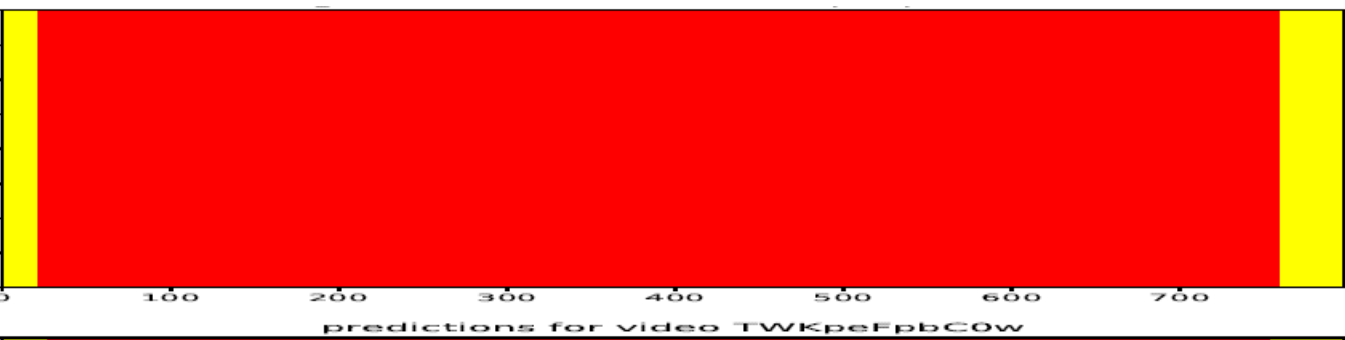


Note the boundary errors between speech and other categories



# Results

Video "TWKpeFpbC0w"

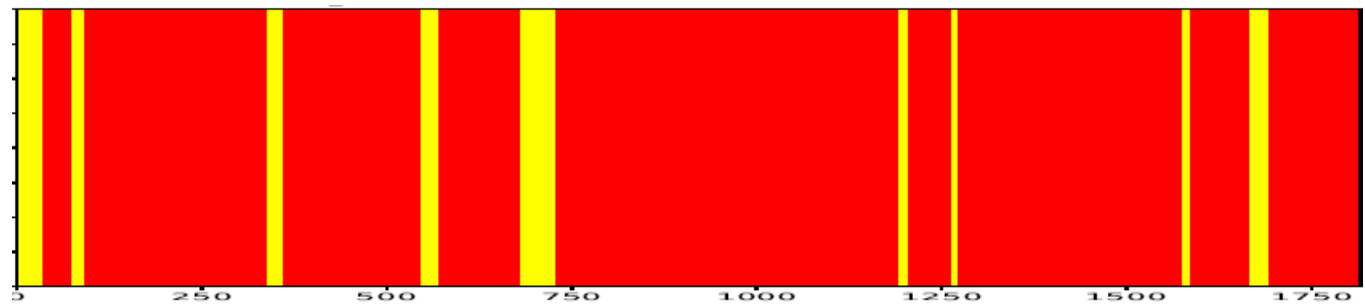


# Results

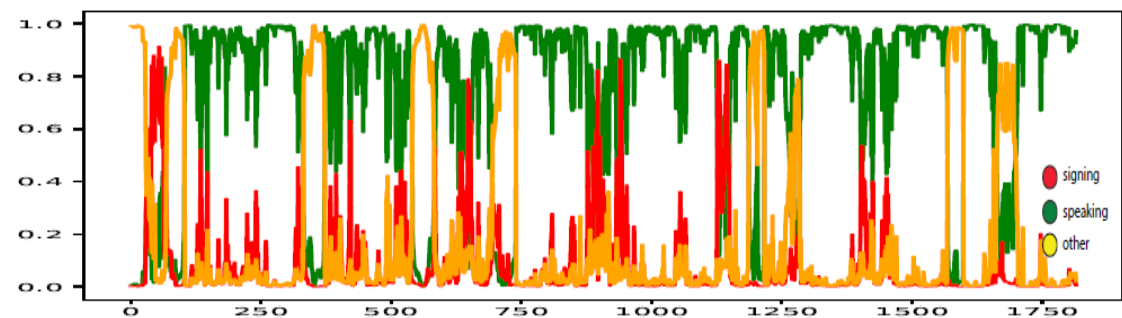
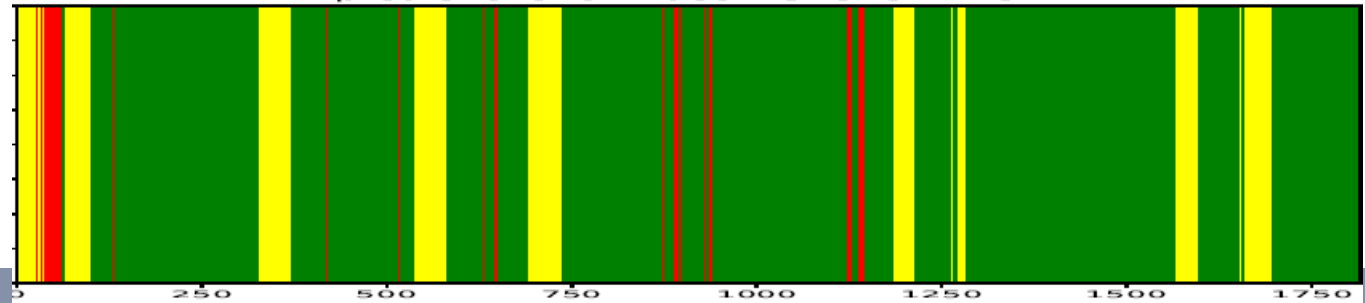
Video "wSA84cXmvCA"



A failure case

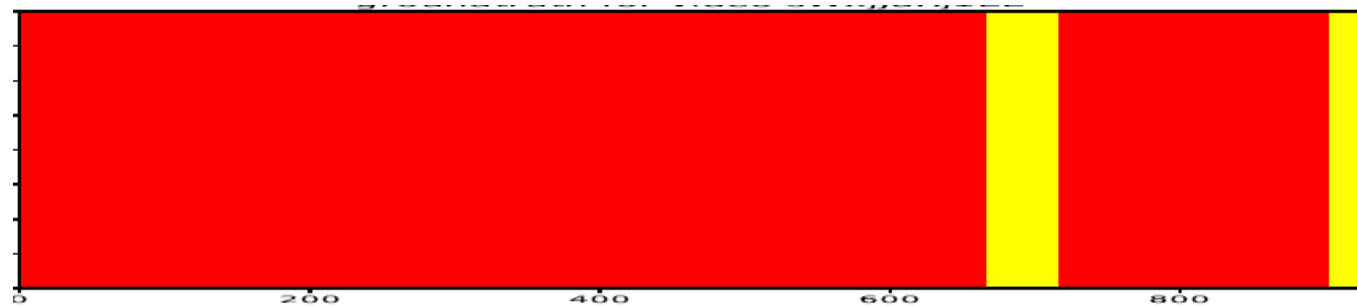


predictions for video wSA84cXmvCA

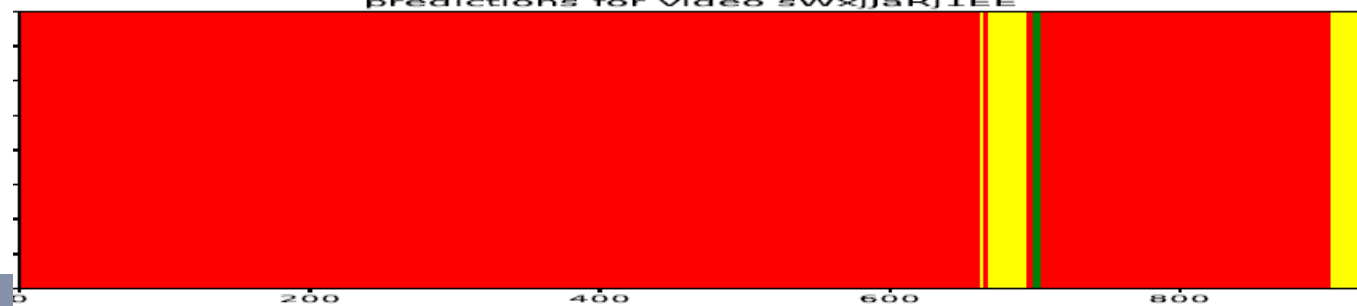


# Results

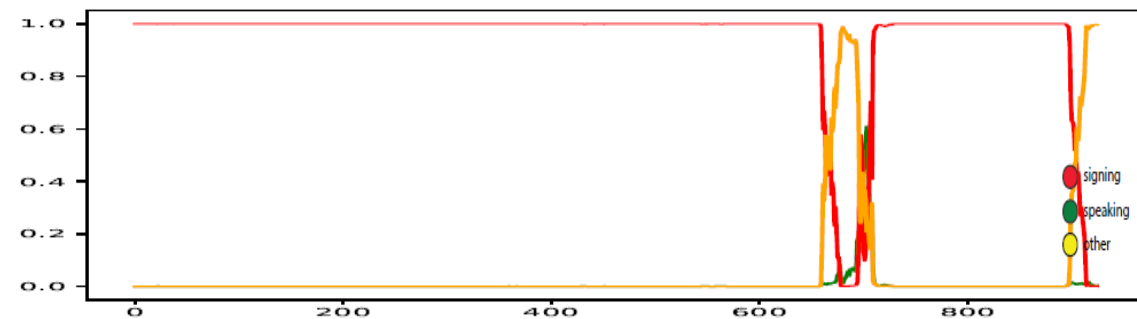
Video "sWxjJaRj1EE"



predictions for video sWxjJaRj1EE



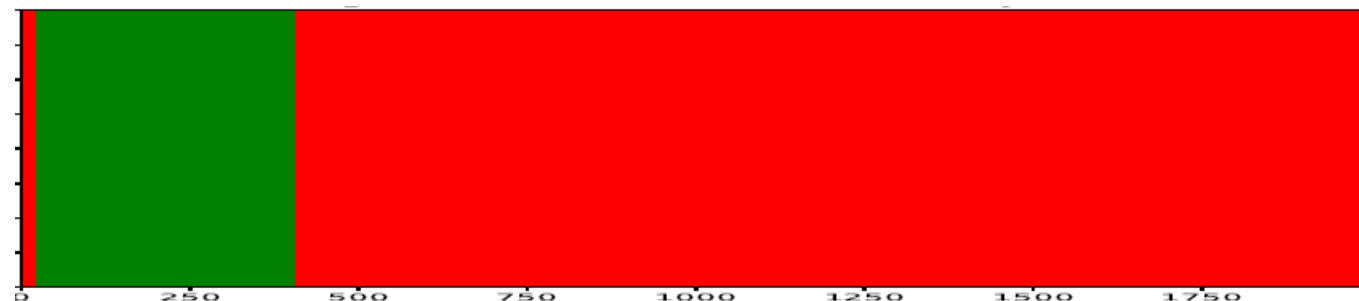
Correct detection. Signer is wearing a mask



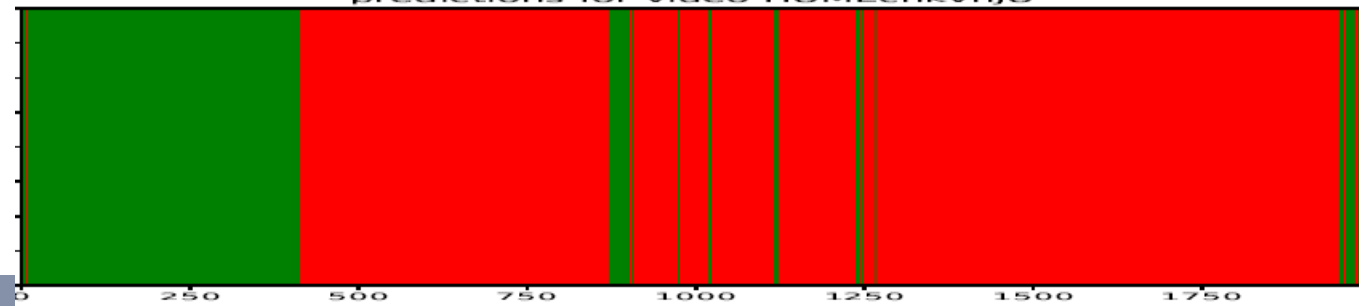


# Results

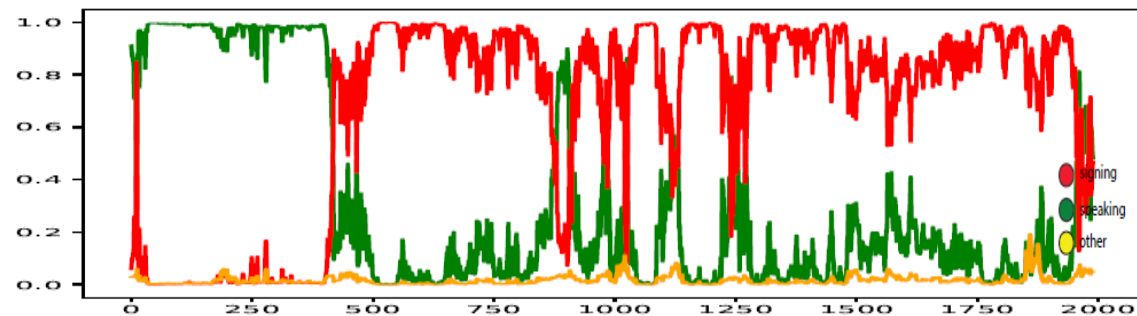
Video "HUMEcnkvhJU"



predictions for video HUMEcnkvhJU



Same person first speaks, then starts signing

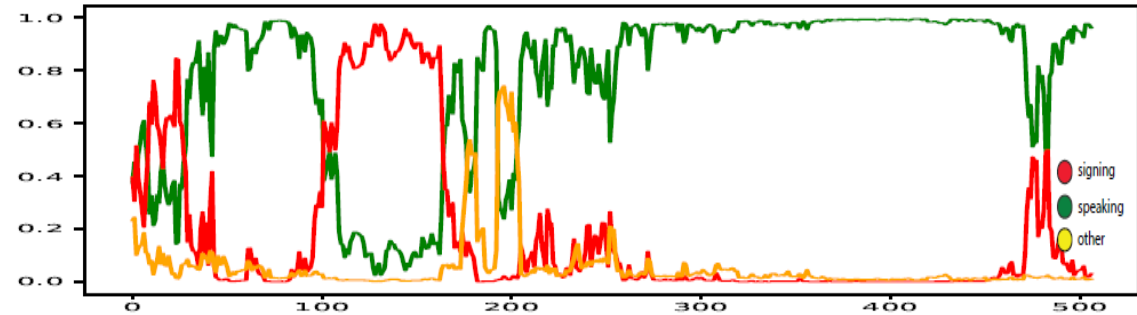
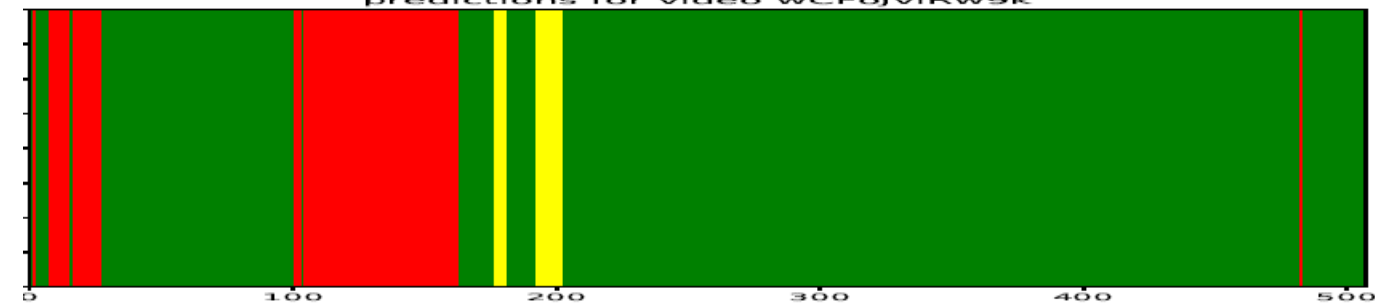
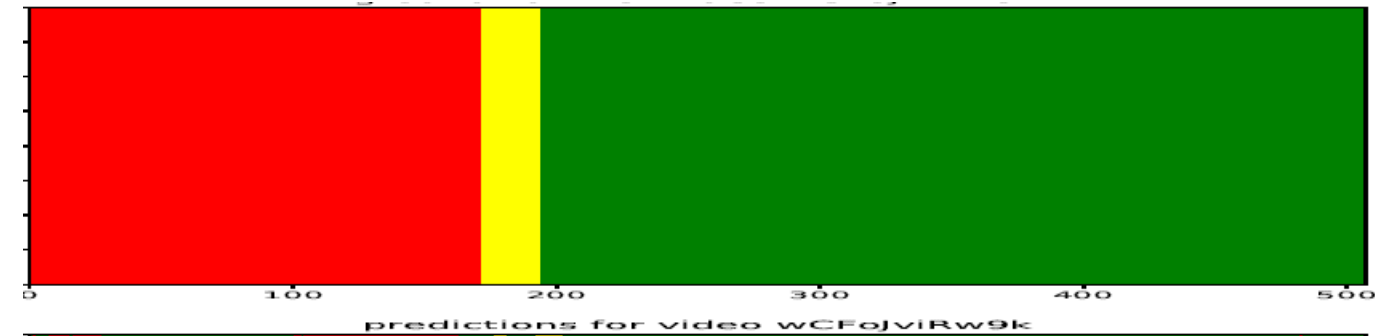


# Results

Video "wCFoJviRw9k"



Video contains a person doing both signing and speaking



# Conclusion & Future Work

- A new dataset “Signing in the Wild”
- Successful sign language detection via a two-stream CNN+RNN
  - $\approx 18\%$  improvement over state-of-the-art
- Visual + Motion features are both important for signing
  
- Future Work:
  - Signer localisation
  - Identification of particular sign languages & sign language constructs
  - Investigating what the CNN+RNN is basing recognition on

- Thank you for your attention