

## Introduction

- We proposed the **Correlational Recurrent Neural Network (CorrRNN)**, a novel temporal fusion model for fusing multiple input modalities that are inherently temporal in nature.

- Characteristics of multimodal temporal data

- **Temporal consistency**  
→ Multimodal GRU Module (GRU)
- **High correlation**  
→ Correlation Module (Corr)
- **Dynamic noise-level changes**  
→ Dynamic Weighting Module (DW)

- Multimodal learning setup (unsupervised)

Tasks	Feature Learning	Supervised Training	Testing
Multimodal Fusion	$X + Y$	$X + Y$	$X + Y$
Cross Modality Learning	$X + Y$	$X$	$X$
Shared Representation Learning	$X + Y$	$Y$	$Y$
Shared Representation Learning	$X + Y$	$X$	$Y$
Shared Representation Learning	$X + Y$	$Y$	$X$

## Proposed Model

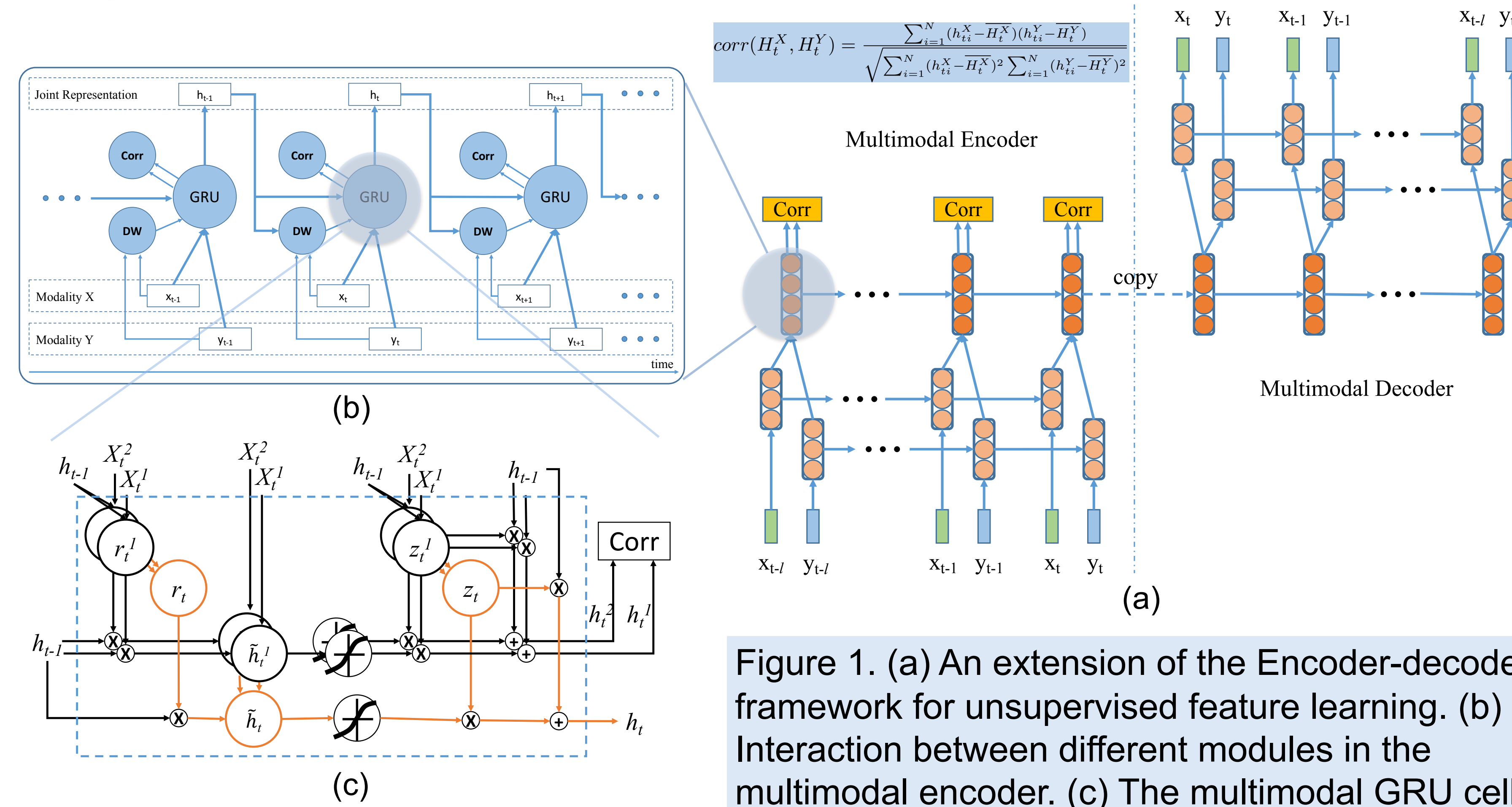


Figure 1. (a) An extension of the Encoder-decoder framework for unsupervised feature learning. (b) Interaction between different modules in the multimodal encoder. (c) The multimodal GRU cell.

## Experiments on Audio-Video Data

- We test our model on two classic datasets.
  - **AVLetters** includes audio and video of 10 speakers uttering the English alphabet three times each.
  - **CUAVE** consists of videos of 36 speakers pronouncing the digits 0-9.
- For fair comparison, we perform the same data pre-processing approaches as those for RTMRBM (2016).

Method (years)	Accuracy	
	AVLetters	CUAVE
MDAE (2011)	62.04	66.70
MDBN (2012)	63.2	67.20
MDBM (2012)	64.7	69.00
RTMRBM (2016)	66.04	-
CRBM (2014)	67.10	69.10
<b>CorrRNN (ours)</b>	<b>83.40</b>	<b>95.9</b>

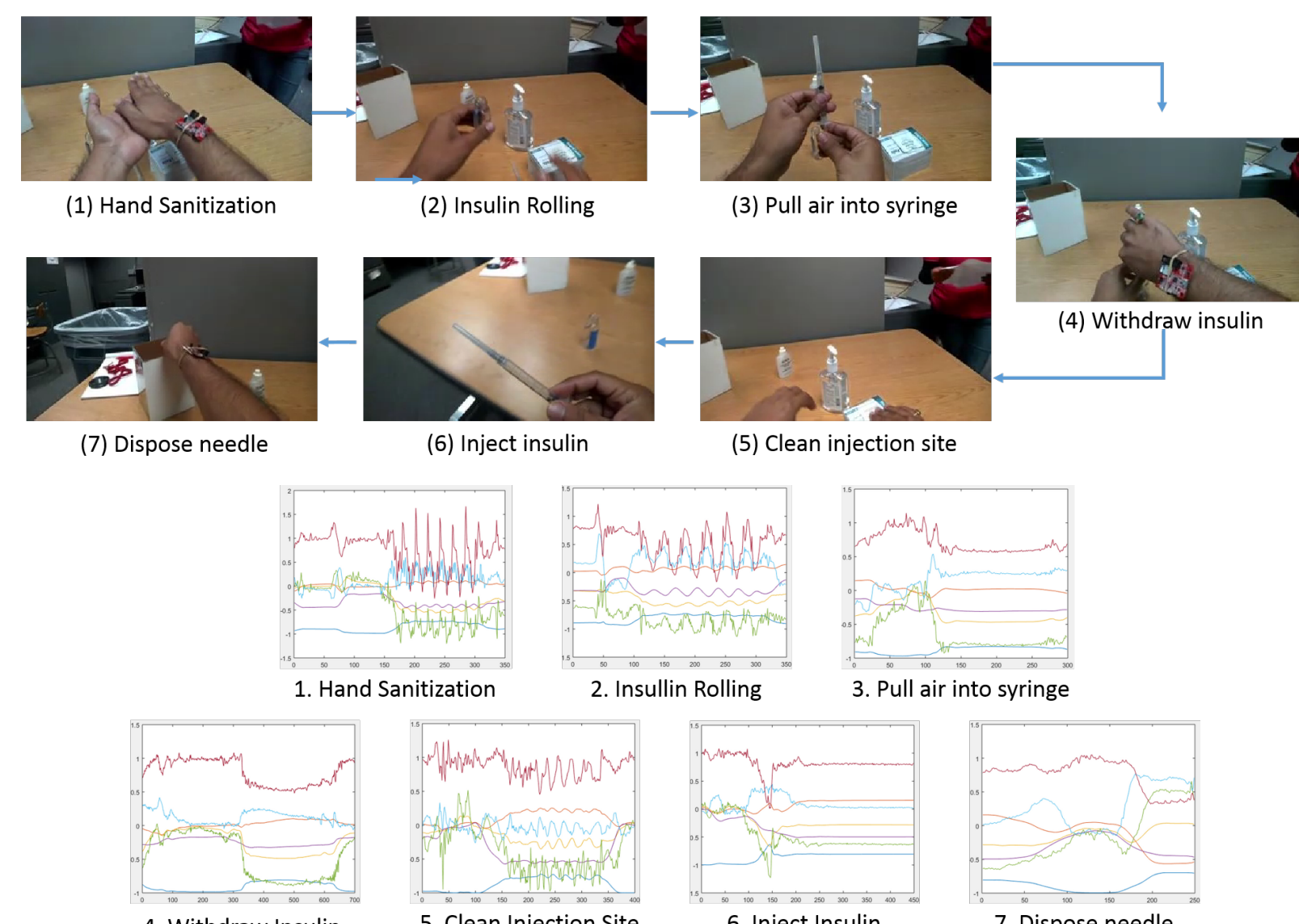
Table 1. Classification performance for audio-visual speech recognition compared to the best published results in literature, using the fused representation of the two modalities.

Method	Accuracy	
	Clean Audio	Noisy Audio
MDAE	94.4	77.3
Audio RBM	95.8	75.8
MDAE + Audio RBM	94.4	82.2
CorrRNN	<b>96.11</b>	<b>90.88</b>

Table 2. Our proposed model is more robust to noise, compared with other multimodal learning methods.

## Experiments on Video-Sensor Data

- The Insulin Self-Injection (ISI) dataset
- 11 subjects perform seven actions related to an insulin self-injection activity.
- Egocentric video data acquired using a Google Glass wearable camera, and motion data acquired using an Invensense motion wrist sensor.
- Sensor data includes Quaternion rotation, Linear acceleration and Gravity compensated acceleration.



- We perform extensive experiments for different model configurations.
- Each loss component contributes to better performance, especially in the settings of cross-modality learning and shared representation learning.

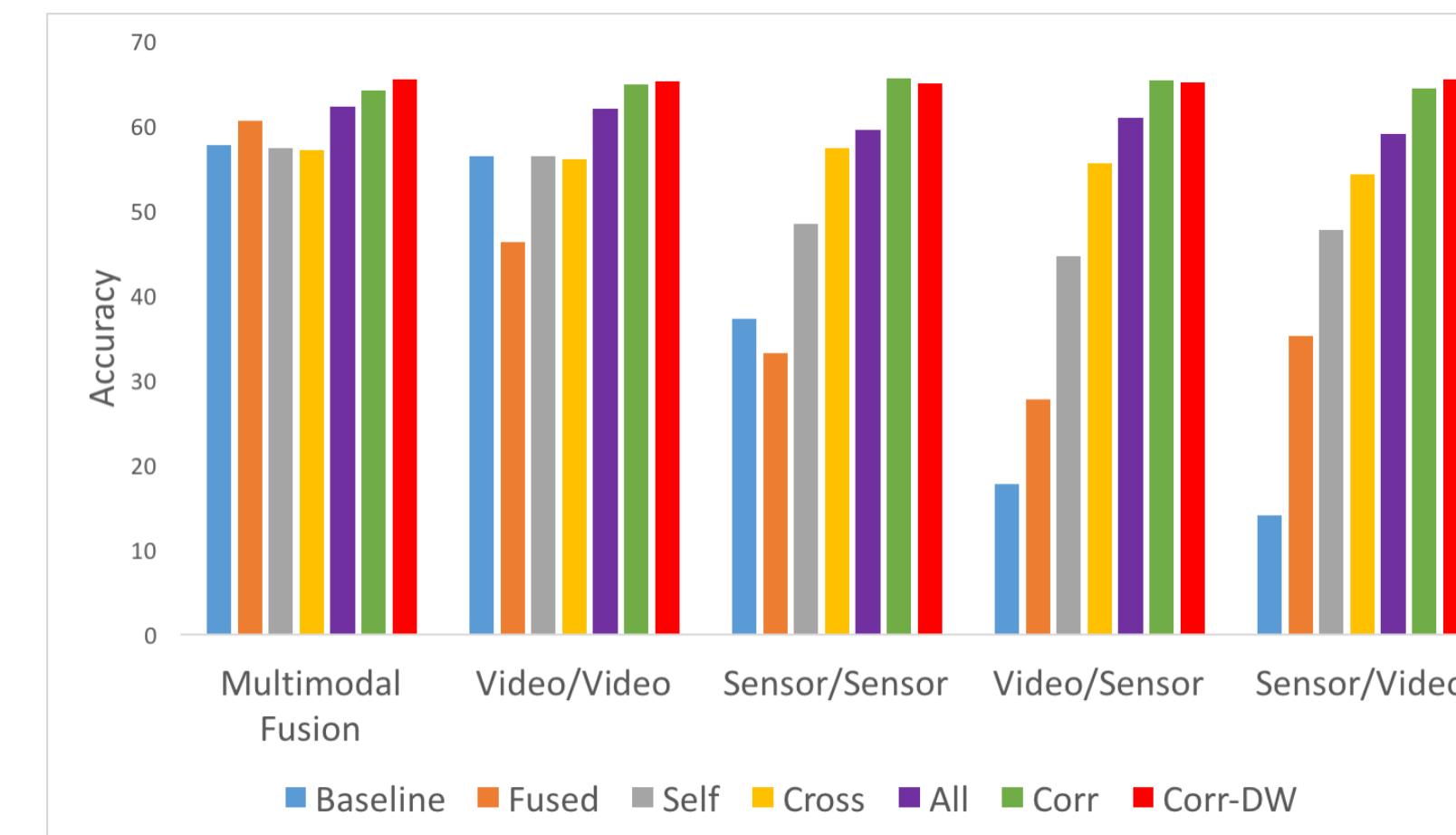


Figure 2. Classification accuracy on the ISI dataset for different model configurations.