



University of  
**Nottingham**

UK | CHINA | MALAYSIA

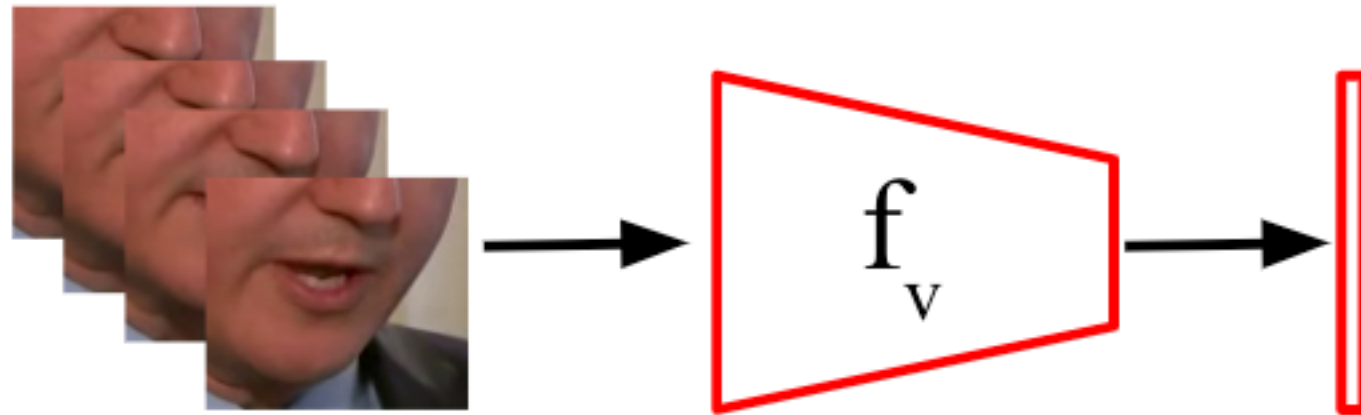


# Audio-Visual Predictive Coding for Self-Supervised Visual Representation Learning

Mani Kumar<sup>1</sup>, Michel Valstar<sup>1</sup>, Michael Pound<sup>1</sup>, Timo Giesbrecht<sup>2</sup>

<sup>1</sup>University of Nottingham, <sup>2</sup>Unilever R&D Port Sunlight, UK

- **Problem Statement:** To learn a visual representation function ( $f_v$ ) from unlabeled video data



### Directly Supervised Representation Learning

- Labeled Data:  $\{X, Y\}$

$$X \xrightarrow{f} Y$$

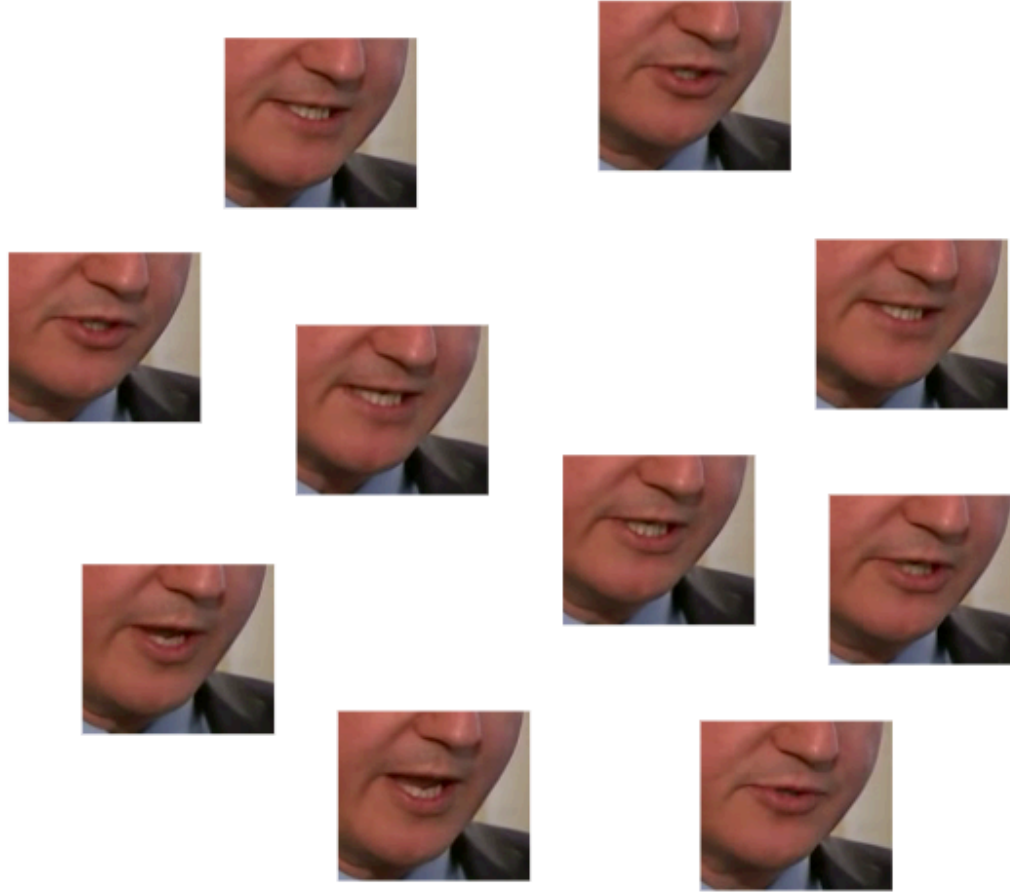
### Self-Supervised Representation Learning

- Unlabeled Data:  $\{X\}$

$\Rightarrow$  Proxy learning task:  $\{X, \hat{Y}\}$

$$X \xrightarrow{\hat{f}} \hat{Y}$$

# Unlabeled Data Points: Intrinsic Correspondences



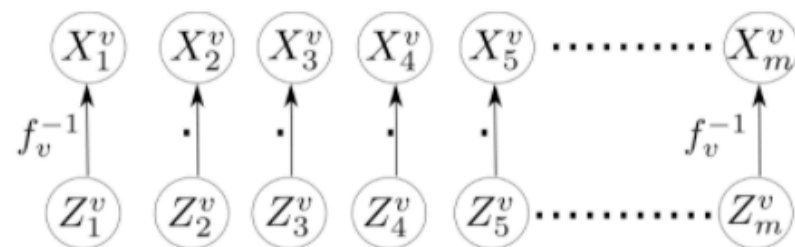
Data-points as i.i.d  
(independent and  
identically distributed)  
samples

Time (t)

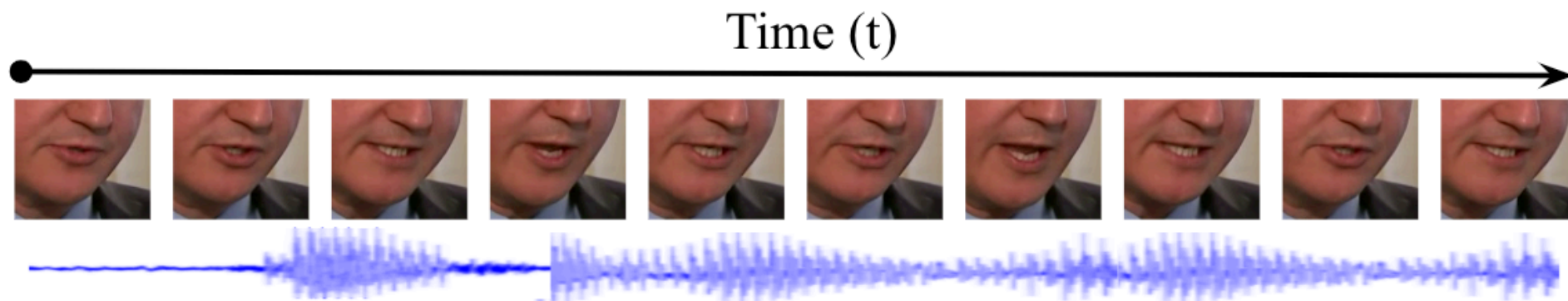


Intrinsic Data-point Correspondences

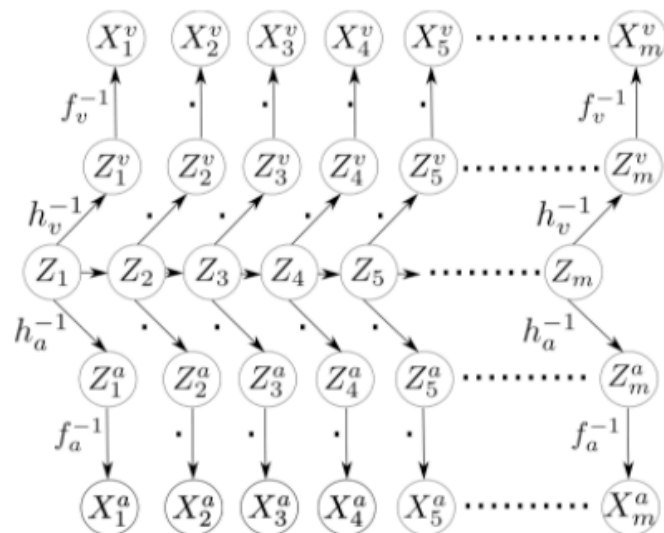
- Intramodal (Temporal correlations)



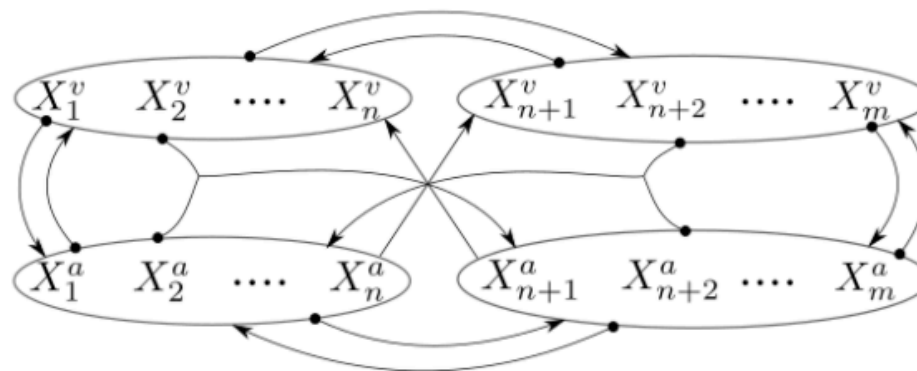
# Audio-Visual Predictive Coding



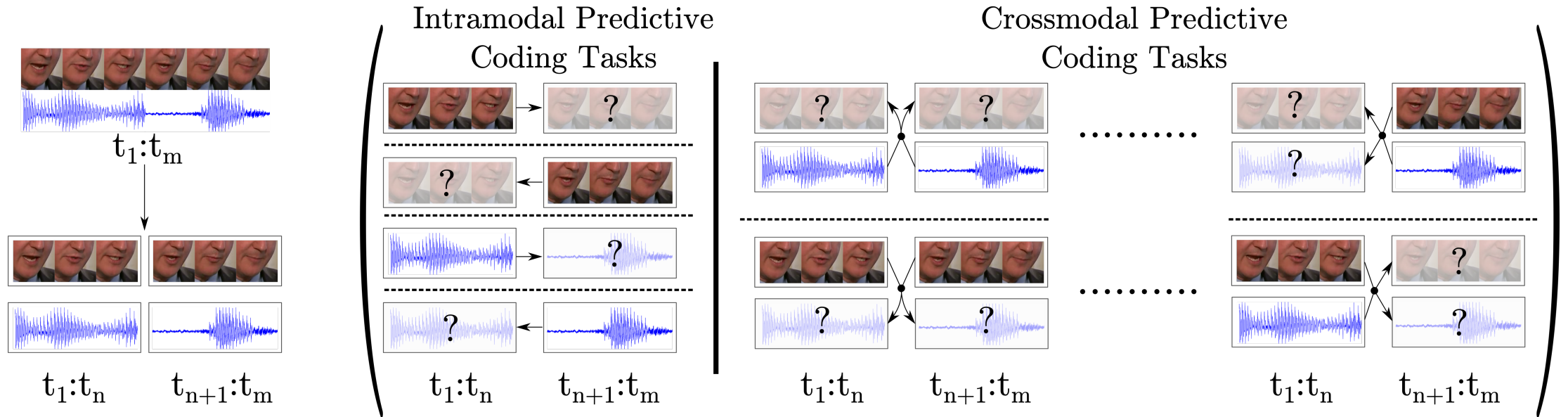
Exploiting temporal and crossmodal correspondences jointly

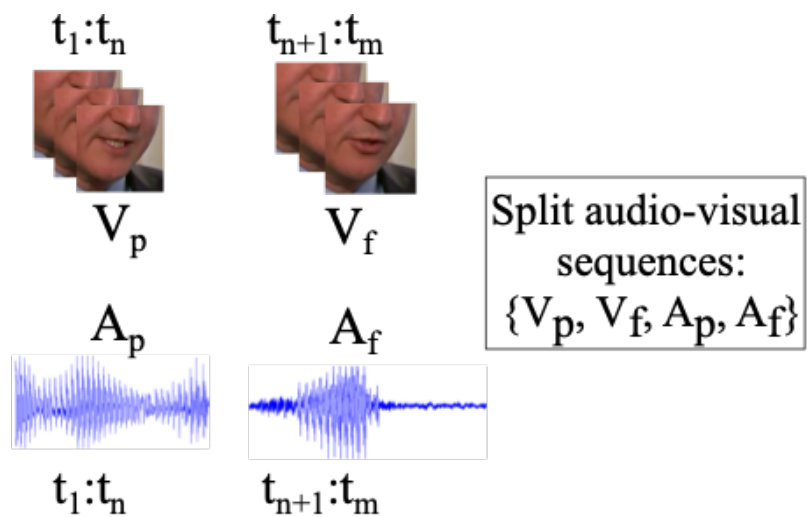


Audio-Visual Permutative Predictive Coding (AV-PPC)  
(Ours)



# Audio-Visual Permutative Predictive Coding

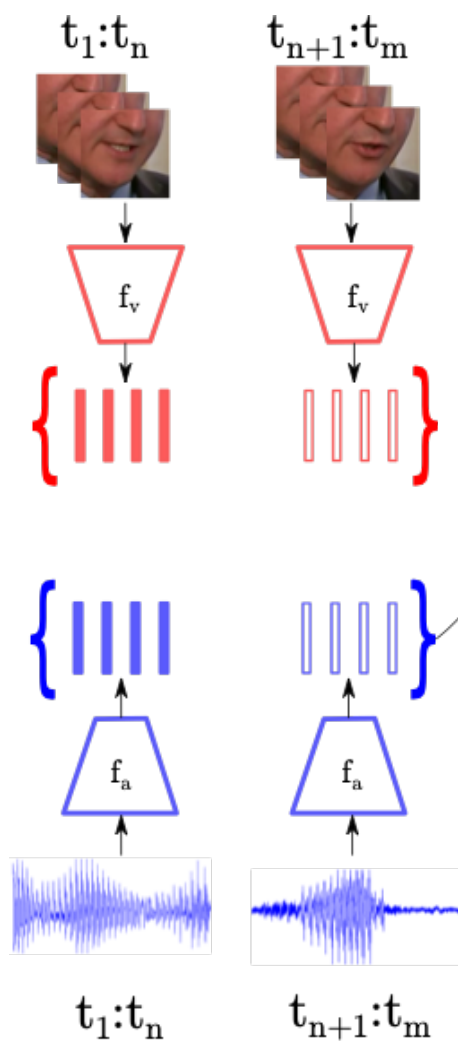


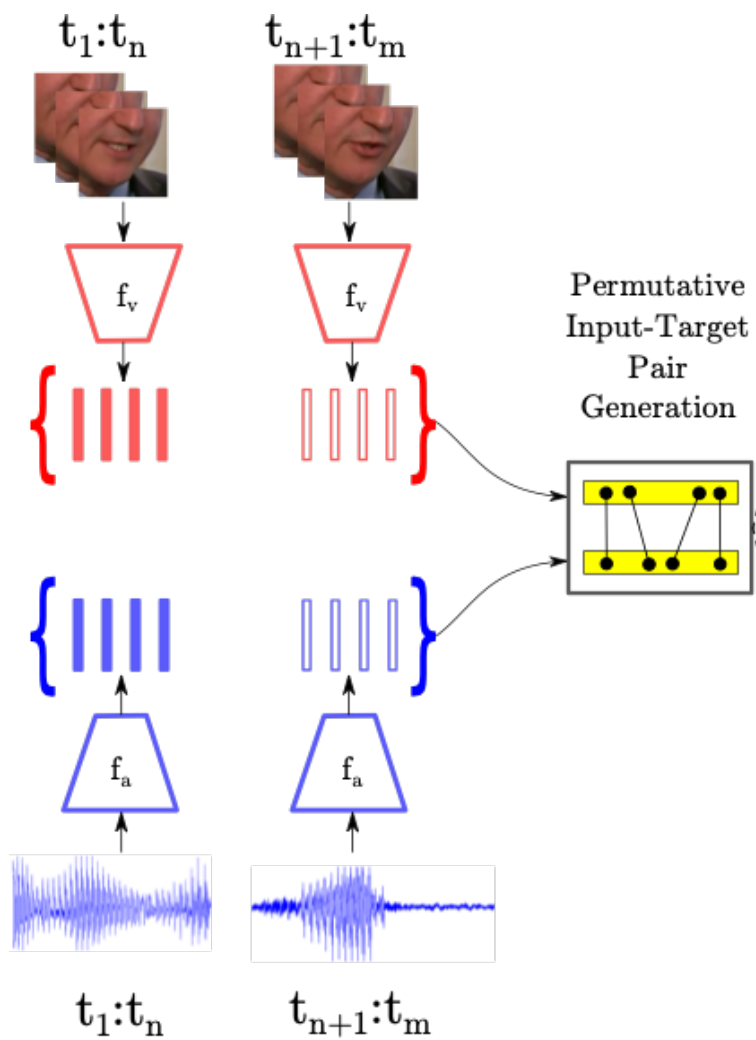


Permutative predictive coding sub-tasks (Input  $\rightarrow$  Target)

One-to-One (#12 tasks)	One-to-Two (#12 tasks)
$V_p \rightarrow V_f$ $A_p \rightarrow V_p$ $A_p \rightarrow A_f$ $\dots$ $\dots$ $V_f \rightarrow V_p$	$V_p \rightarrow (V_f, A_f)$ $V_f \rightarrow (V_p, A_p)$ $A_p \rightarrow (V_f, A_f)$ $\dots$ $\dots$ $V_f \rightarrow (A_p, A_f)$
Two-to-One (#12 tasks)	Two-to-Two (#6 tasks)
$(V_p, A_p) \rightarrow V_f$ $(A_p, A_f) \rightarrow V_p$ $(A_p, V_f) \rightarrow A_f$ $\dots$ $\dots$ $(V_f, A_f) \rightarrow V_p$	$(V_p, V_f) \rightarrow (A_p, A_f)$ $(A_p, A_f) \rightarrow (V_p, V_f)$ $(V_p, A_p) \rightarrow (V_f, A_f)$ $(V_f, A_f) \rightarrow (V_p, A_p)$ $(V_p, A_f) \rightarrow (A_p, V_f)$ $(A_p, V_f) \rightarrow (V_p, A_f)$

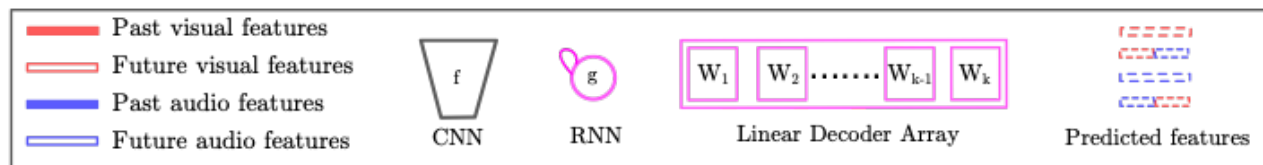
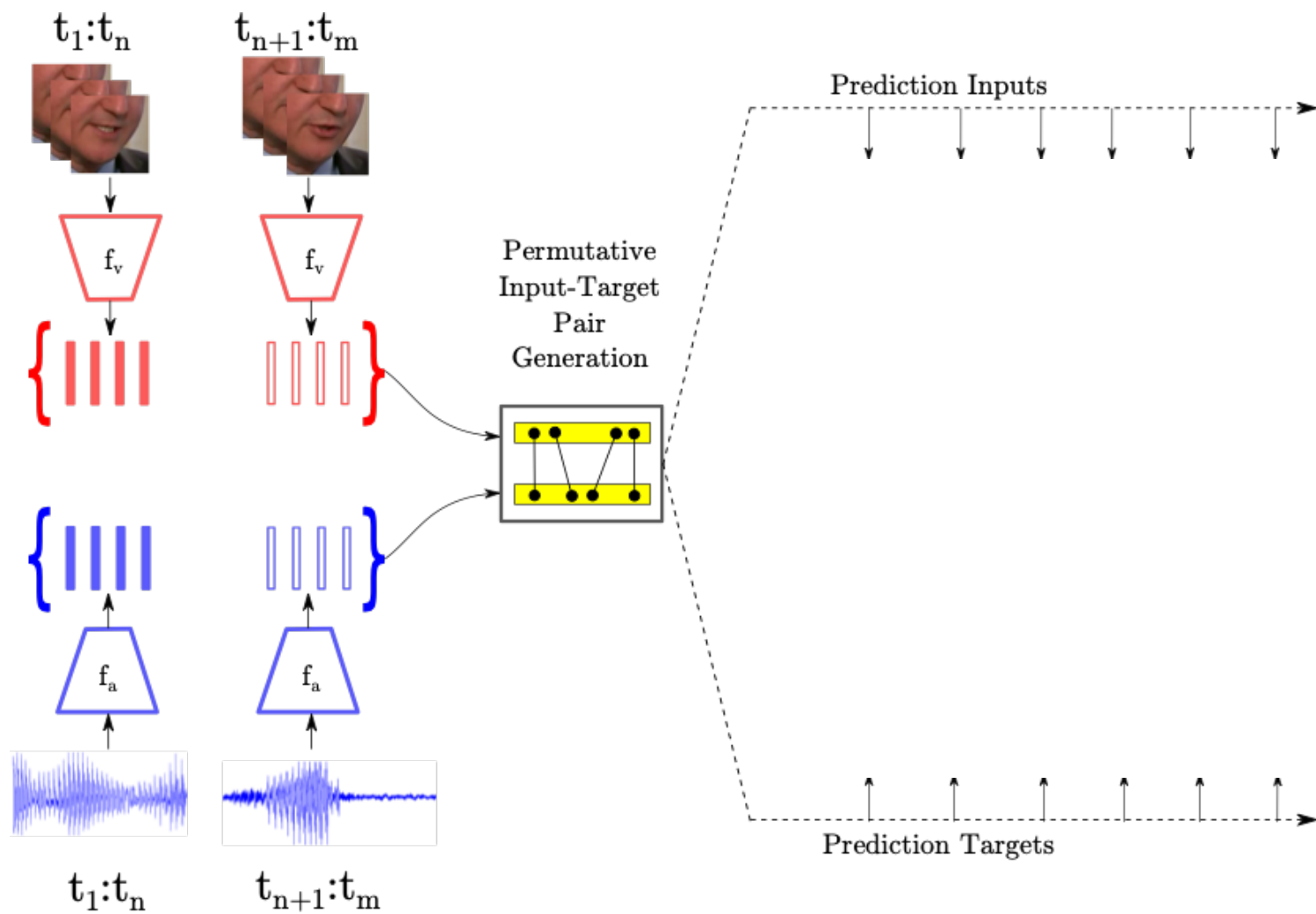


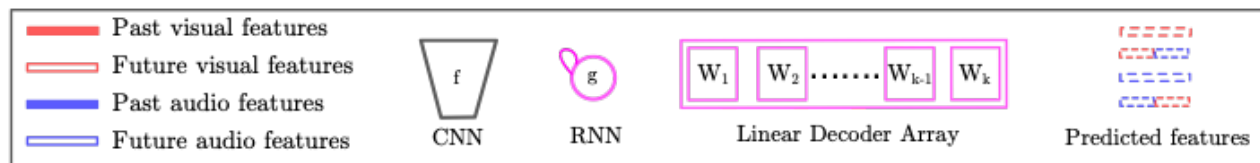
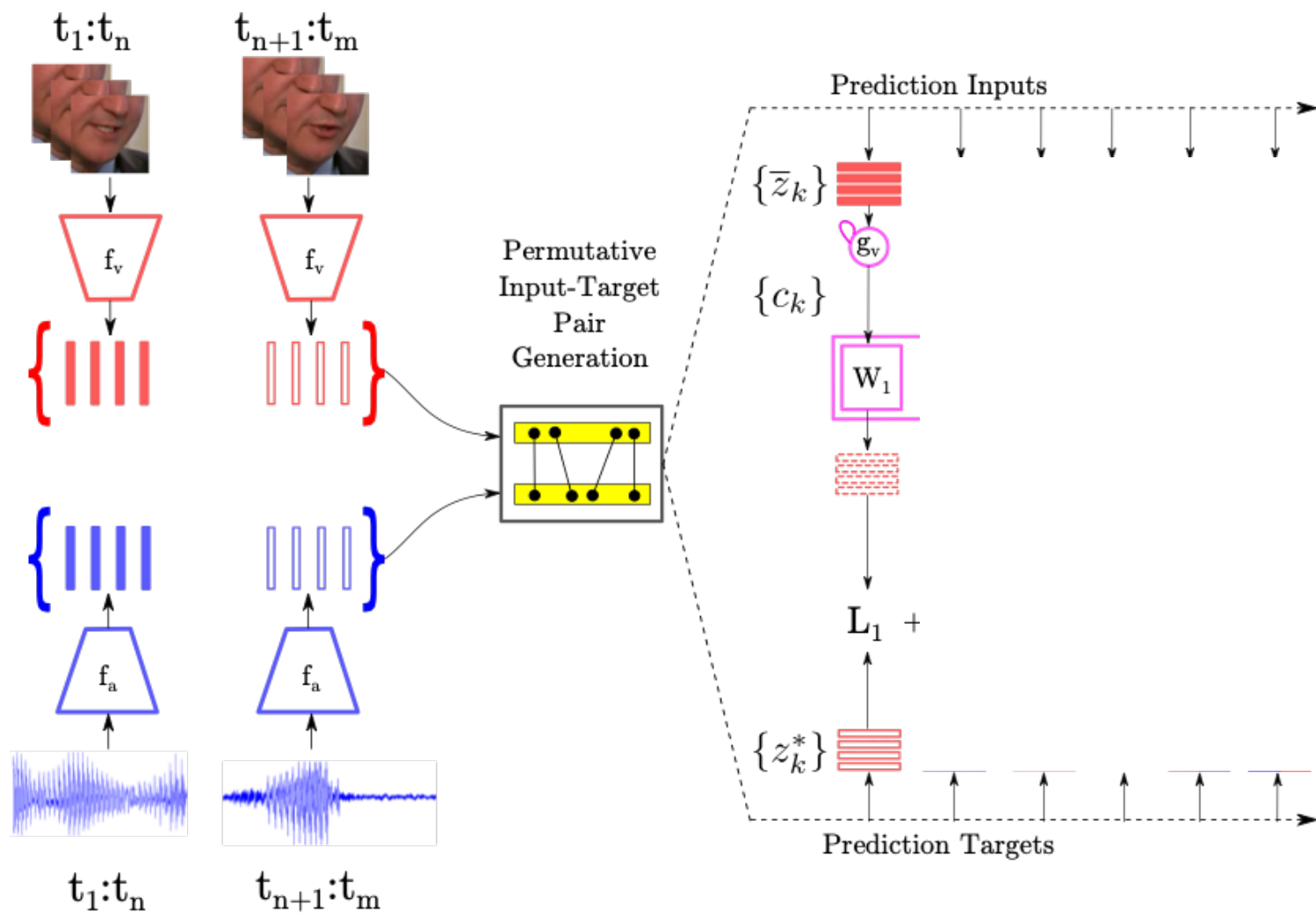


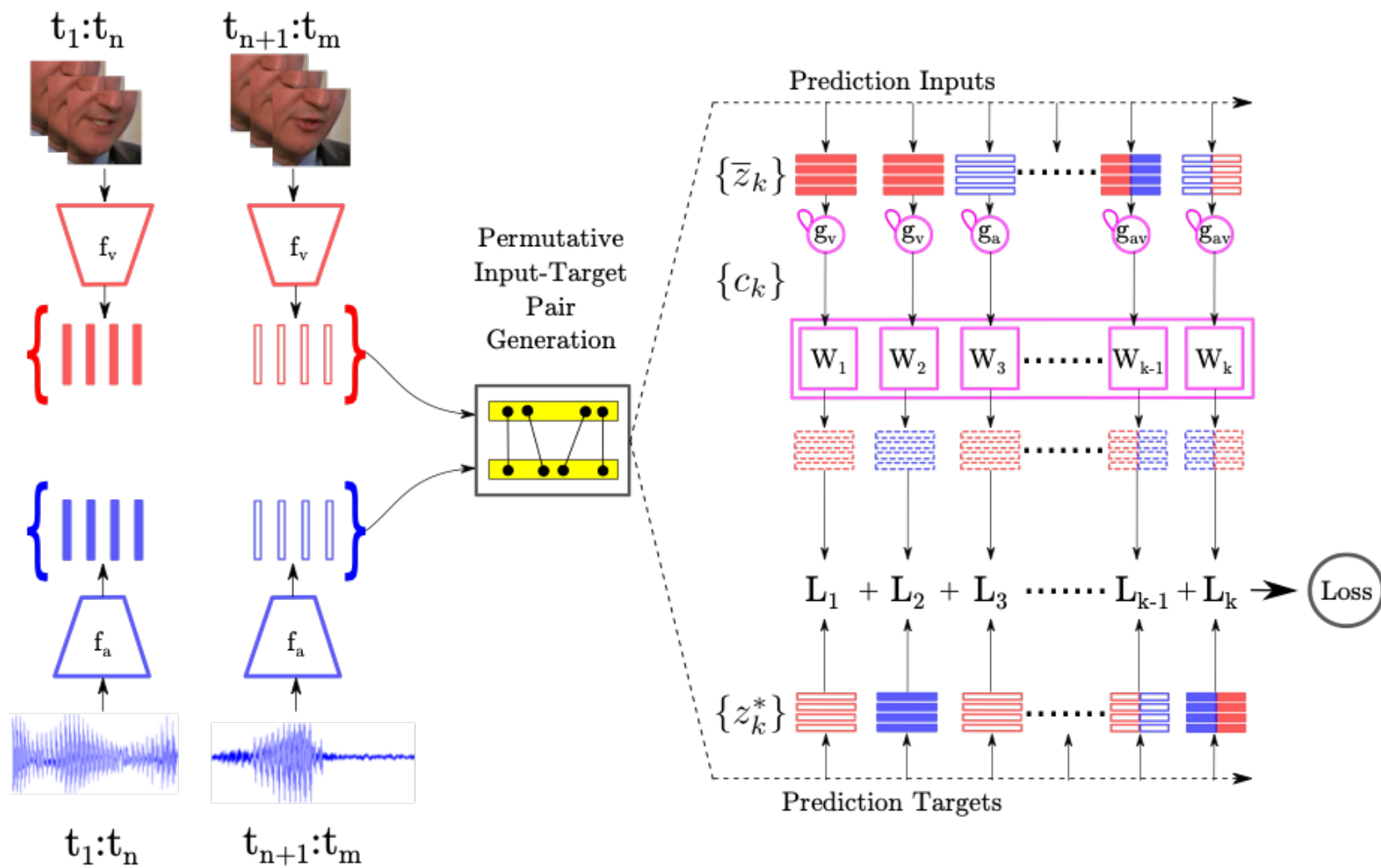


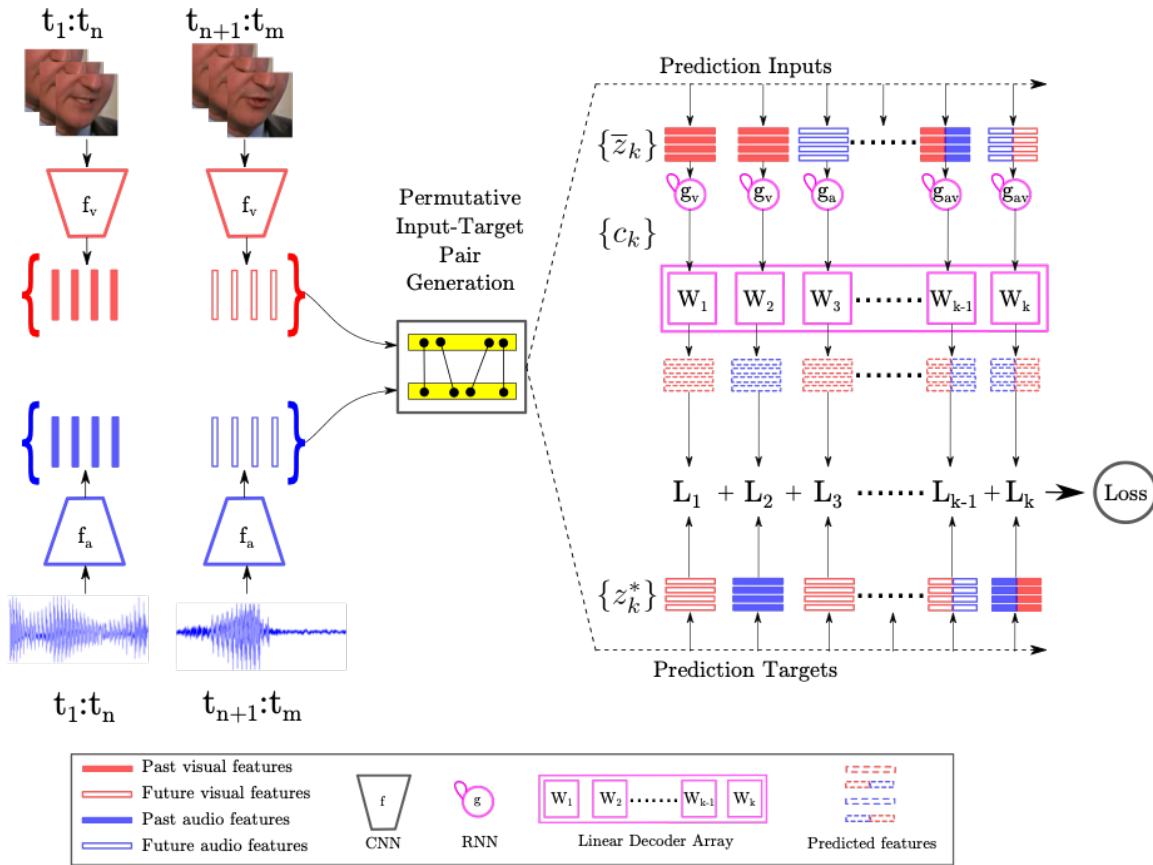
Permutative predictive coding sub-tasks (Input  $\rightarrow$  Target)

One-to-One (#12 tasks)	One-to-Two (#12 tasks)
$V_p \rightarrow V_f$ $A_p \rightarrow V_p$ $A_p \rightarrow A_f$ $\dots$ $\dots$ $V_f \rightarrow V_p$	$V_p \rightarrow (V_f, A_f)$ $V_f \rightarrow (V_p, A_p)$ $A_p \rightarrow (V_f, A_f)$ $\dots$ $\dots$ $V_f \rightarrow (A_p, A_f)$
Two-to-One (#12 tasks)	Two-to-Two (#6 tasks)
$(V_p, A_p) \rightarrow V_f$ $(A_p, A_f) \rightarrow V_p$ $(A_p, V_f) \rightarrow A_f$ $\dots$ $\dots$ $(V_f, A_f) \rightarrow V_p$	$(V_p, V_f) \rightarrow (A_p, A_f)$ $(A_p, A_f) \rightarrow (V_p, V_f)$ $(V_p, A_p) \rightarrow (V_f, A_f)$ $(V_f, A_f) \rightarrow (V_p, A_p)$ $(V_p, A_f) \rightarrow (A_p, V_f)$ $(A_p, V_f) \rightarrow (V_p, A_f)$









## Contrastive Learning: InfoNCE Loss

(Noise Contrastive Estimation)

$$I(z_k^*; c_k) = \sum_{z_k^*, c_k} p(z_k^*, c_k) \log \left( \frac{p(z_k^* | c_k)}{p(z_k^*)} \right)$$

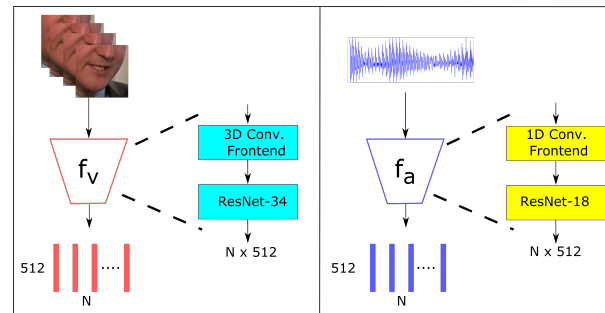
$$f_k(z_k^*, c_k) \propto \frac{p(z_k^* | c_k)}{p(z_k^*)}$$

$$f_k(z_k^*, c_k) = \exp(z_k^{*T} \cdot W \cdot c_k)$$

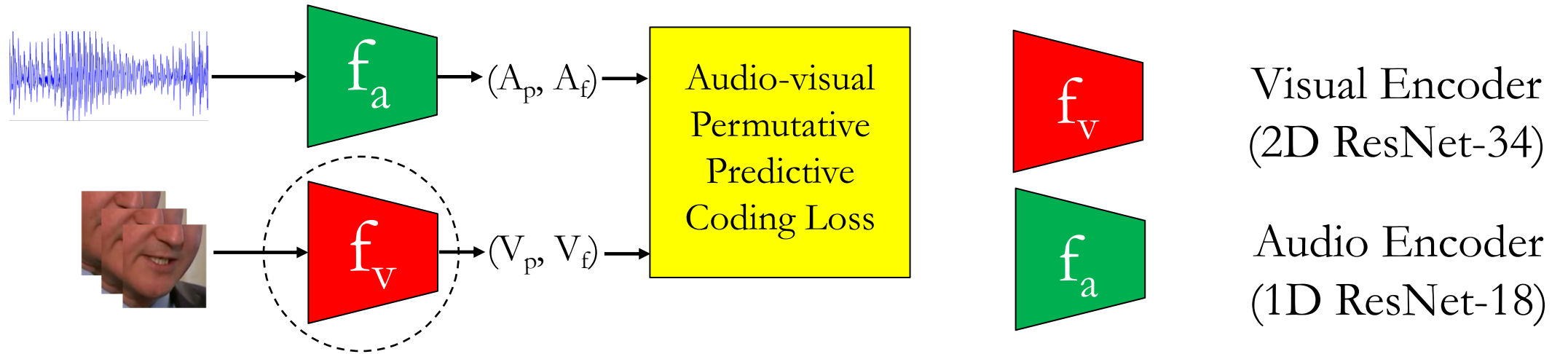
$$L_i = -E_B \left[ \log \frac{f_k(z_k^*, c_k)}{\sum_{z_j \in B} f_k(z_j, c_k)} \right]$$

where B is a mini batch of N samples,

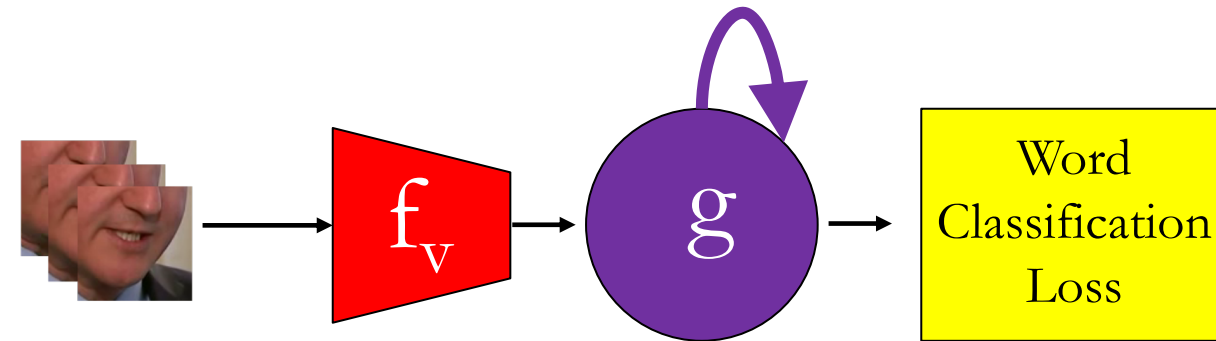
- with 1 positive sample and N-1 negative samples



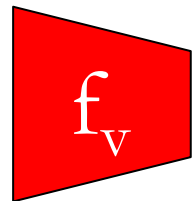
# Summary of Self-Supervised Learning Stage



# Downstream Evaluation Task: Lip-Reading



**Task:** Predict the word uttered in a video  
**Dataset:** LRW with 500-word class labels  
**Metric:** Word Classification Rate (WCR)



Visual Encoder  
(2D ResNet-34)



Temporal Model  
(GRU/Temporal Conv)

**Evaluation Protocol:** Measure WCR

A. before finetuning the visual encoder

B. after finetuning the visual encoder

A. using the entire train data and

B. using small amounts of train data.



# Performance of Different Proxy Tasks on the Lip-Reading Task (Word Classification Rates)

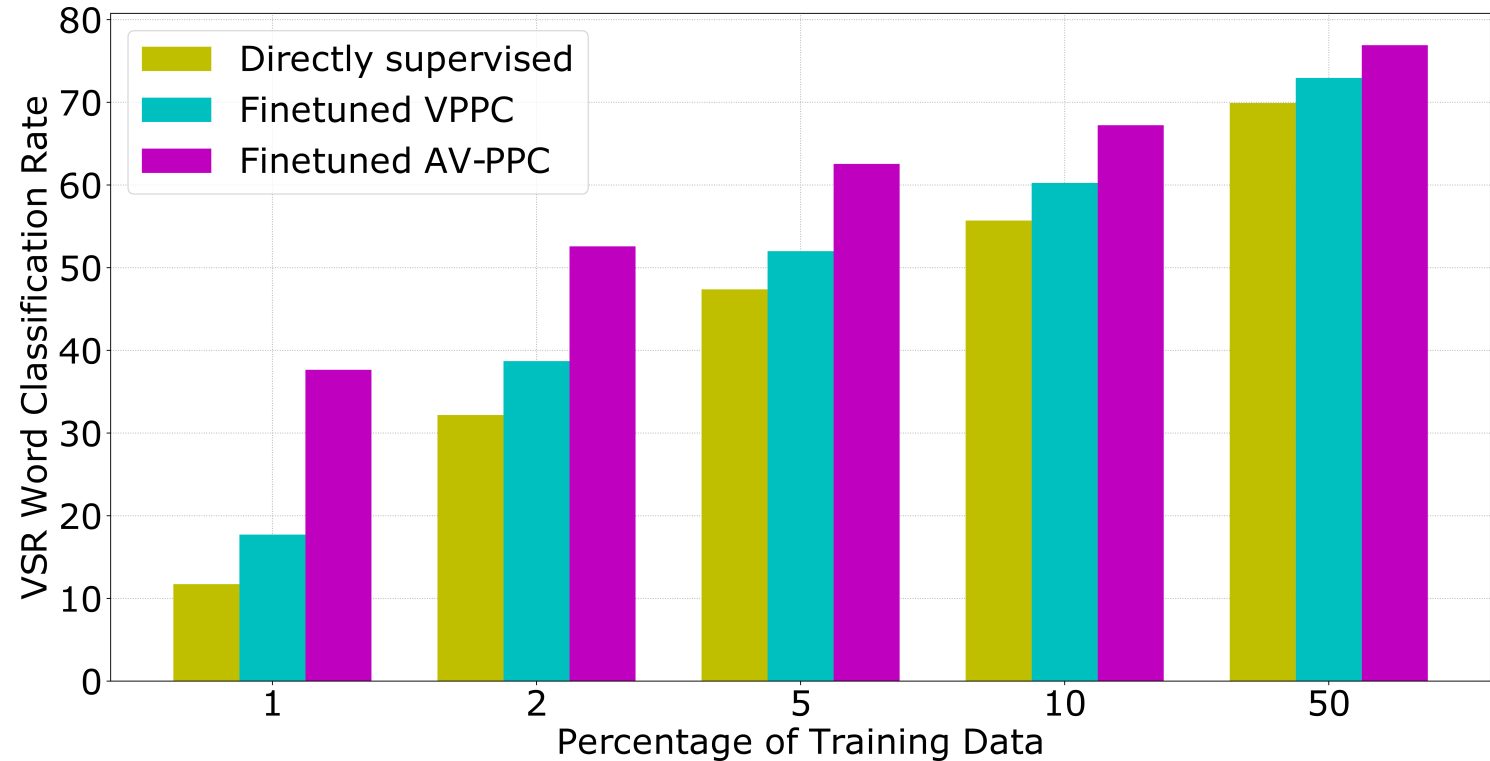
Proxy Task	Using Temporal Conv	Using GRU
AV Synchronization	50.70 (74.17)	55.26 (76.92)
Time-Arrow	52.42 (75.80)	59.88 (78.26)
AV Correspondence	56.22 (74.23)	61.90 (77.90)
Vis. Permutative Pred. Coding	60.77 (77.95)	67.62 (81.76)
AudVis. Permutative Pred. Coding (ours)	<b>76.47 (80.44)</b>	<b>80.30 (83.16)</b>

# Data-Efficiency Evaluation

- Number of labeled instances required to learn lip-reading task

- With 1% of train data (10 instances per word class),

- Our method: 38% WCR
- Fully-supervised: 11% WCR



# Take-home Idea

A potential approach to unsupervised representation learning:

**Leveraging rich intrinsic data-point correspondences**  
**- temporal and cross-modal semantic correlations -**  
**as natural supervision signals in the self-supervised setting.**

Thank You