



# SPATIO-TEMPORAL ACTION DETECTION WITH A MOTION SENSE AND SEMANTIC CORRECTION FRAMEWORK



湖州师范学院  
Huzhou University



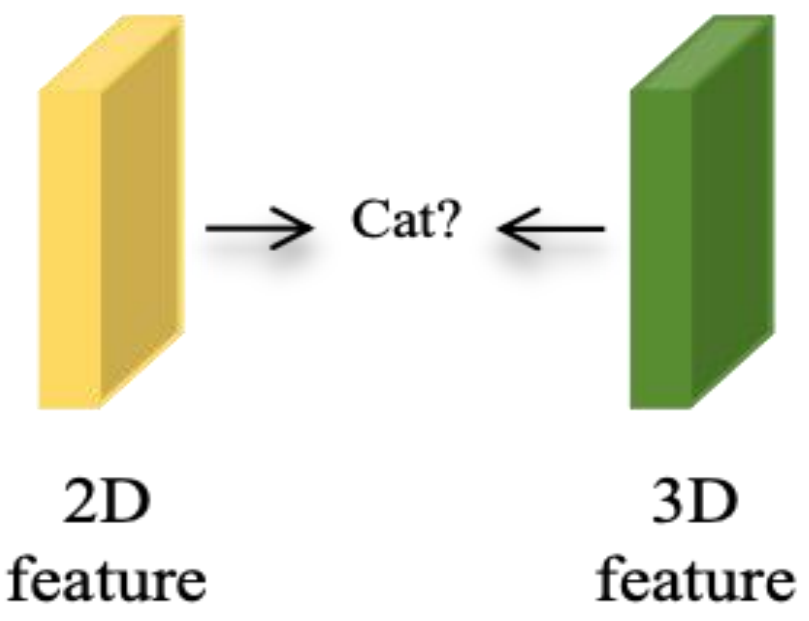
遼寧師範大學  
Liaoning Normal University

Yong Zhang<sup>1,2</sup>, Chunan Yu<sup>1</sup>, Chenglong Fu<sup>1</sup>, Yuanqi Hu<sup>1</sup>, Ying Zang<sup>1,†</sup>  
<sup>1</sup>Huzhou University, <sup>2</sup>Liaoning Normal University

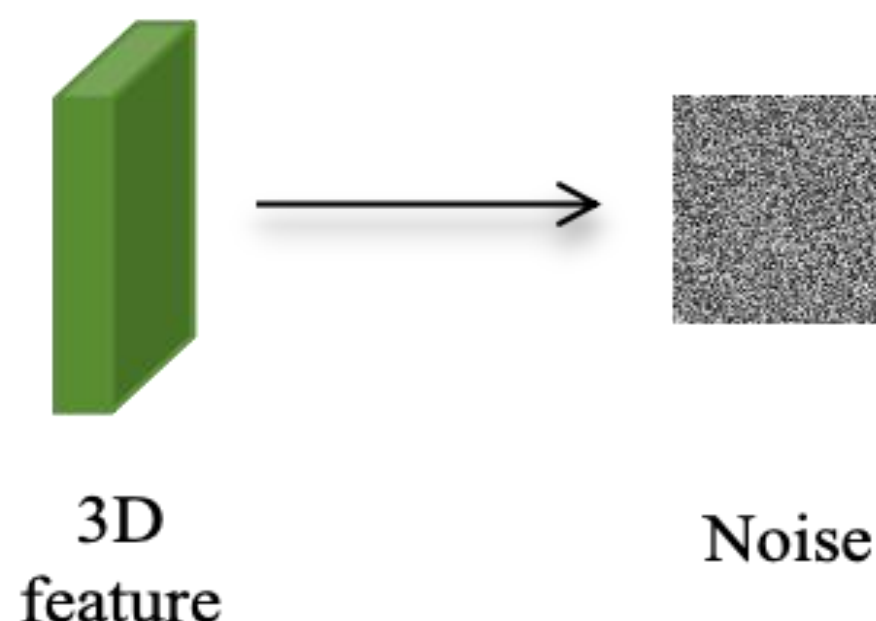
ICASSP  
2024 KOREA

## Motivation

How efficient fusion?



How filter noise?



### Action Detection:

Aimed at utilizing model frameworks to recognize actions in videos or sequences of images.



## Introduction

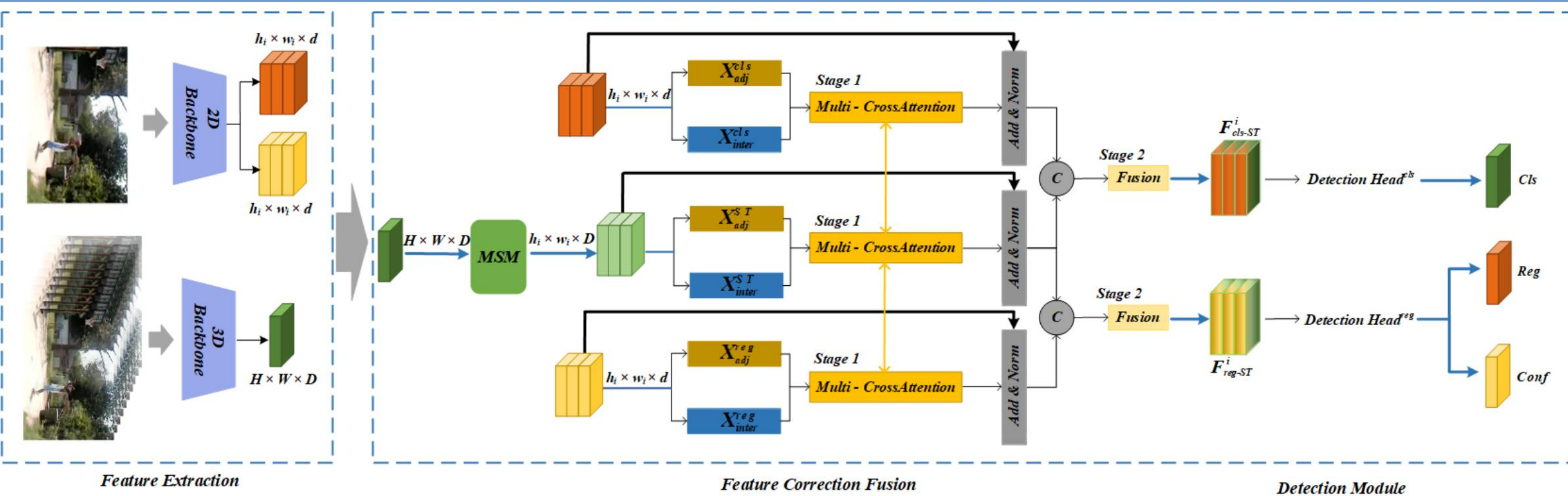
**Q1:** How to address high feature similarity, making it difficult to differentiate between action and background features?

**Contribution 1:** Motion Sense Module (MSM) to significantly increase the feature distance between action and non-action features in the semantic space, enhancing feature discriminability.

**Q2:** How to achieve multimodal feature fusion that  $1+1>2$ , fulfilling fusion between different modalities?

**Contribution 2:** Considering the complementary nature of information across different modalities, an efficient Semantic Correction Fusion Module (SFM) is introduced to facilitate interaction between features of distinct modalities and maximize their complementary information integration.

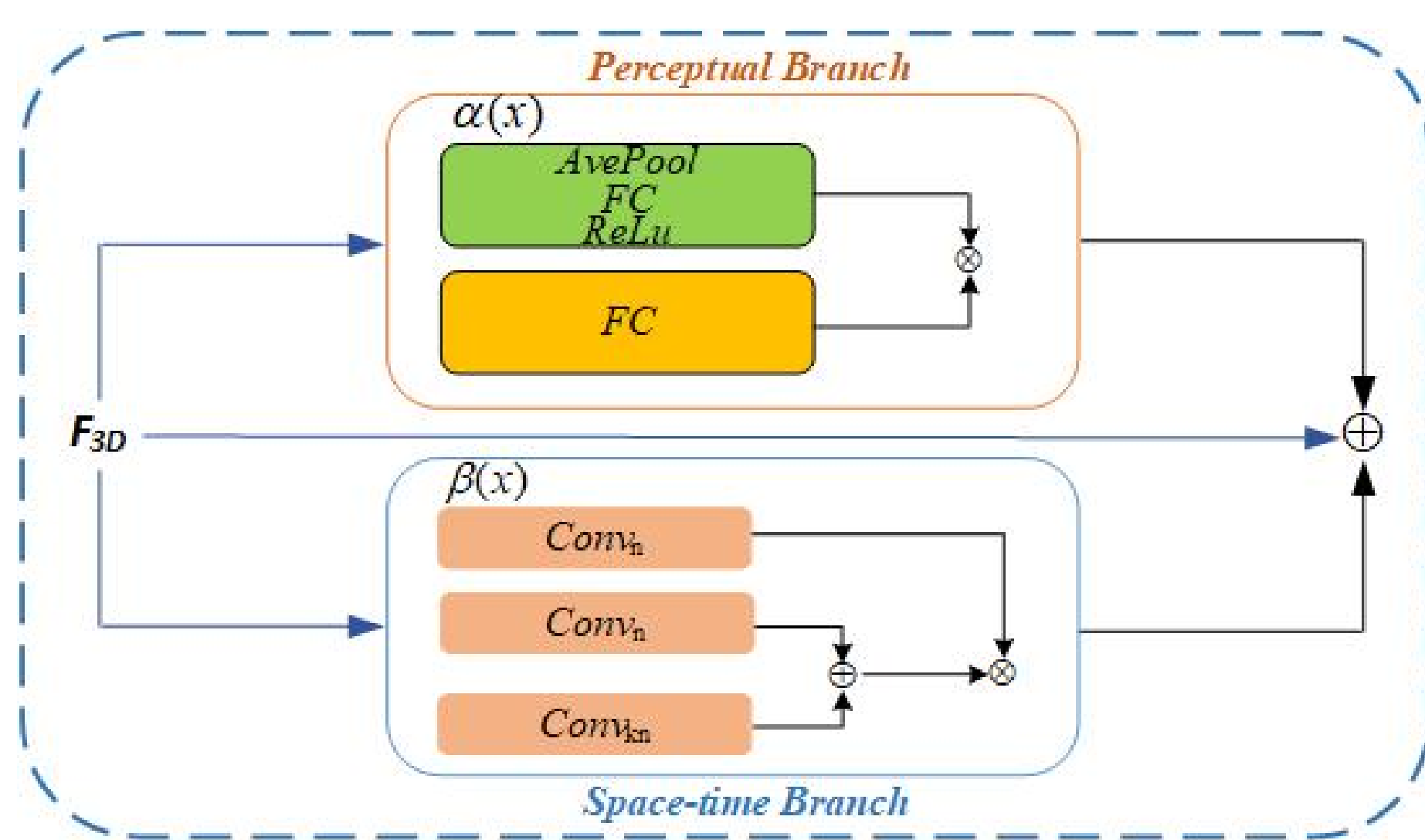
## Method Overview



● This paper proposes a novel Motion Sense and Semantic Correction framework (MS-SC).

● The MS-SC framework achieves accurate detection by fusing features from images (spatial dimension) and videos (spatio-temporal dimension).

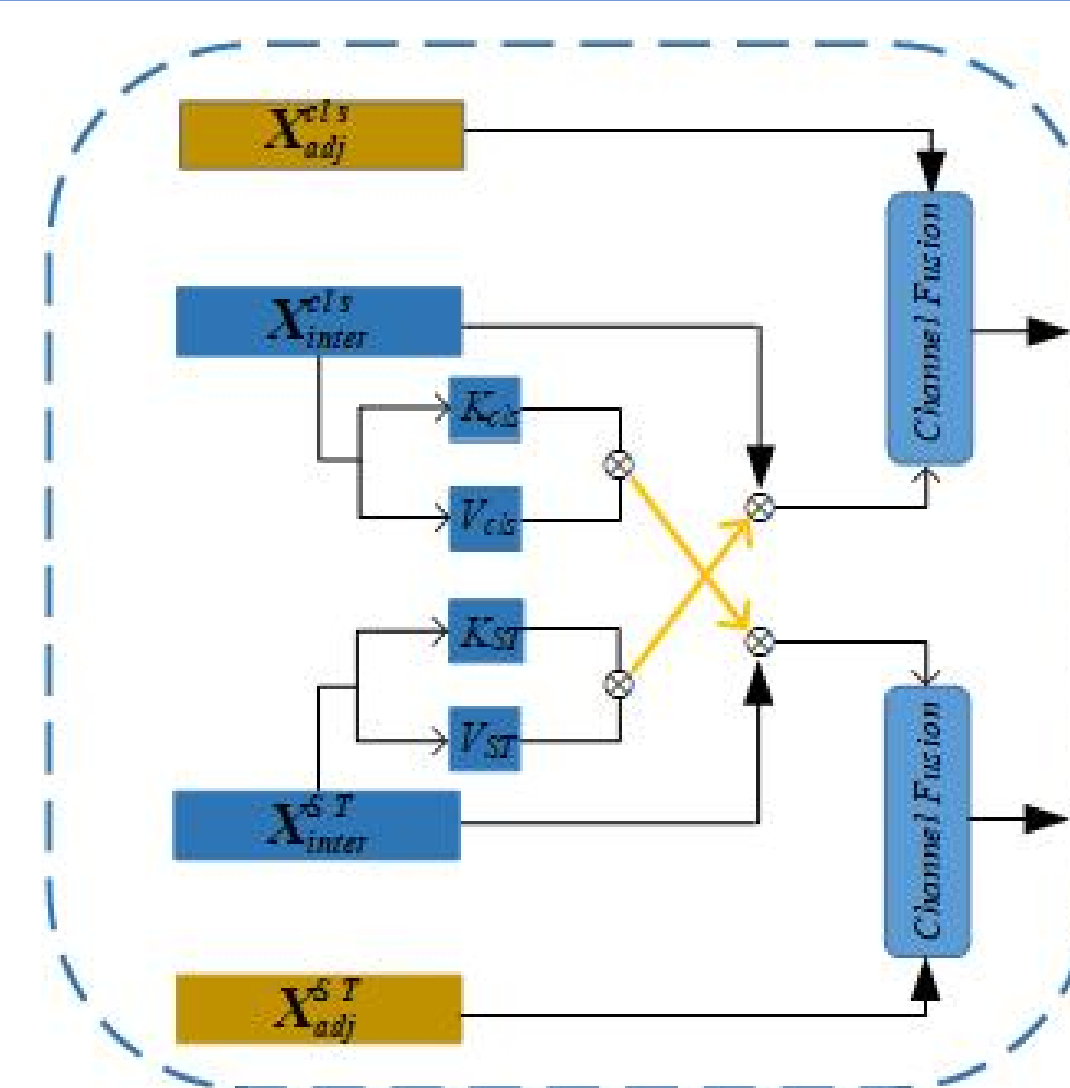
## Motion Sense Module



● **Perceptual Branch:** The objective is to enhance the distinguishability between action features and background features by increasing their feature distance from the average features at the video level.

● **Spacetime Branch:** The objective is to incorporate semantic content from a wider receptive field, which helps in dynamically focusing on features of a specific scale.

## Semantic Correction Fusion Module



● **Stage 1: Information exchange stage.**

In order to correct each other between spatiotemporal features and spatial features, cross-attention method is used for feature guidance.

● **Stage 2: Fusion stage.**

We use a simple channel embedding to merge the features of the two paths, which is achieved by a  $1 \times 1$  convolutional layer.

## Experiments

