# Deep Active Learning from Multispectral Data through Cross-Modality Prediction Inconsistency

Heng ZHANG, Elisa FROMONT, Sébastien LEFEVRE, Bruno AVIGNON

IRISA Laboratory, ATERMES Company
{heng.zhang, elisa.fromont, sebastien.lefevre}@irisa.fr   bavignon@atermes.fr

2021 IEEE International Conference on Image Processing

# Multispectral scene analysis


**Day**


**Night**

❖ The performance of scene analysis applications using only RGB images may be compromised in many real life situations (such as nighttime or shaded areas).

❖ Multispectral systems use two types of camera sensors (**RGB** and **Thermal**) to provide complementary information under various illumination conditions.

❖ However, collecting labelled multi-sensor data is **expensive and time-consuming**.

# Proposed solution: Active learning



Consistent detections        Inconsistent detections

**Multi-sensor redundancy**: detection results from the two modalities are similar in most cases

**Multi-sensor complementarity**: at least one modality is wrong when the detections are contradictory

We rely on the **cross-modal predictions' inconsistency** to adaptively select the multispectral samples to be annotated.
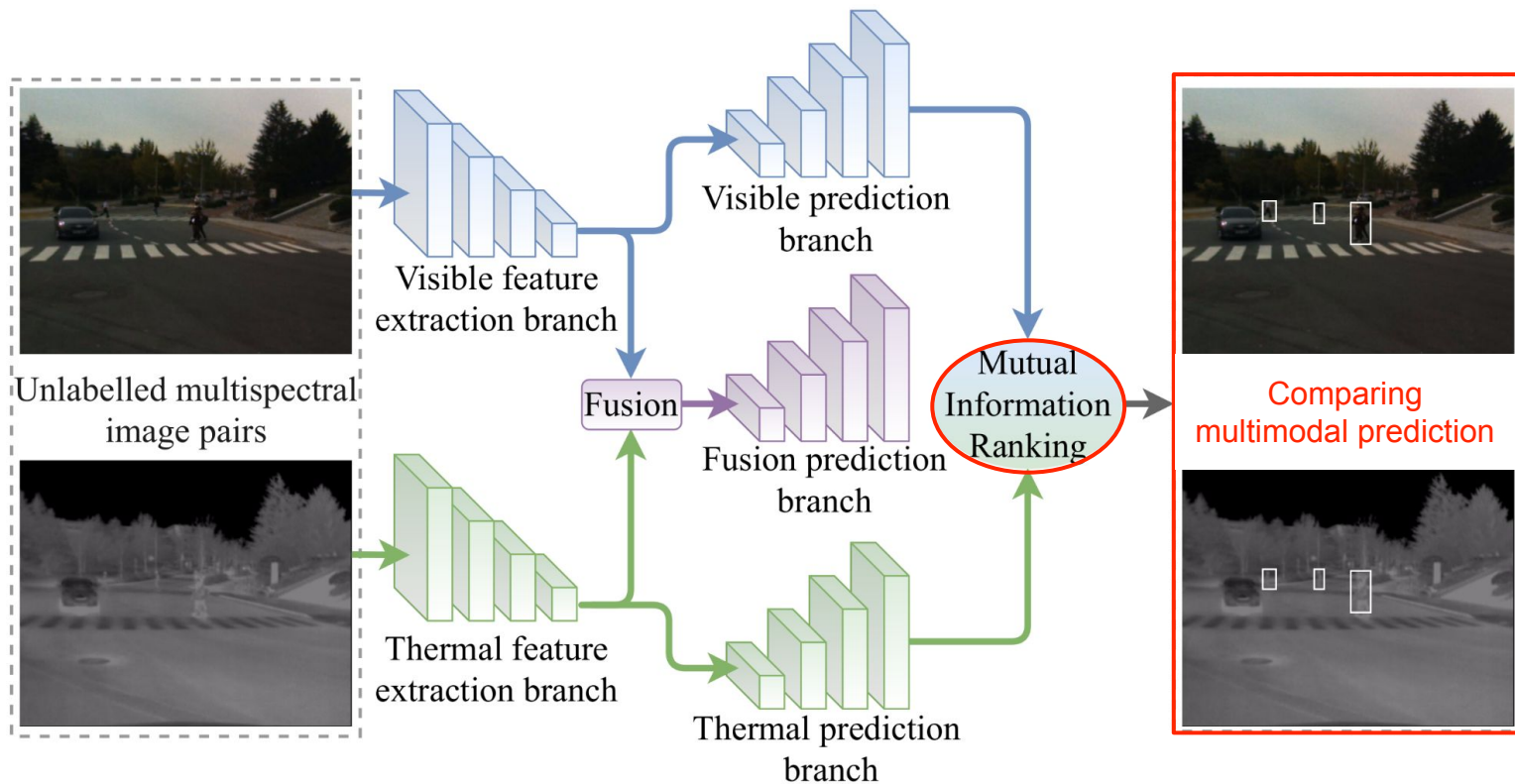
3

# An introduction to Active learning



Active learning loop diagram

❖ The model inference is performed on the unlabelled dataset to select the most informative samples (i.e., multispectral image pairs in our work).

❖ These selected samples are then sent to an external oracle for annotation and appended to the labelled dataset.

❖ The model is consequently fine-tuned on the labelled dataset.

# Overview of the proposed model

# Cross-modality prediction inconsistency

For each prediction p, its **inconsistency score** is defined as:

$$\mathcal{I} = \mathcal{H}\left(\overline{p}\right) - \frac{1}{2}\sum_{m \in \{v,t\}} \mathcal{H}\left(p_m\right)$$



where $p_v$ and $p_t$ denote the prediction from visible and thermal prediction branches; $\overline{p}$ is the average of both predictions; $\mathcal{H}$ is the set entropy function calculated as:

$$\mathcal{H}\left(p\right) = -p\log p - \left(1 - p\right)\log\left(1 - p\right)$$

# Experiments (Active VS Random)



KAIST dataset for
multispectral pedestrian detection

FLIR dataset for
multispectral object detection

TOKYO dataset for
multispectral semantic segmentation

**Observation**: our **active strategy** achieves statistically significant better performance than the **random strategy** for all multispectral scene analysis tasks.

# Inconsistency visualization (I)



KAIST Dataset

FLIR Dataset

Visible camera     Thermal camera     Inconsistency map     Visible camera     Thermal camera     Inconsistency map

# Inconsistency visualization (II)



TOKYO Dataset

background | car | person | bike | curve | car stop | guardrail | color_cone | bump

Visible camera          Thermal camera          Inconsistency map

# Thanks for your attention

Heng ZHANG, Elisa FROMONT, Sébastien LEFEVRE, Bruno AVIGNON

IRISA Laboratory, ATERMES Company
{heng.zhang, elisa.fromont, sebastien.lefevre}@irisa.fr  bavignon@atermes.fr

2021 IEEE
International Conference on Image Processing