# YMX2L at the NTCIR-18 Transfer-2 Task

## Riku Mizuguchi<sup>†1</sup>, Takeshi Yamazaki<sup>†1</sup>, and Shuhei Yamamoto<sup>†2</sup>

<sup>†1</sup> College of Knowledge and Library Sciences, School of Informatics, University of Tsukuba, Japan <sup>†2</sup> Institute of Library, Information and Media Science, University of Tsukuba, Japan

### Background

The advancements in LLMs have paved the way for cross-modal information processing, enabling interactions across various modalities, including text, images, and audio.







#### While large text and image corpora exist (e.g., Microsoft COCO and Flickr30k), they often lack corresponding data from other modalities, limiting their utility in developing robust cross-modal information access technologies.

### Our Task

- We introduce a dense retrieval method for retrieving relevant data across multiple modalities.
- In our approach, image and sensor data are projected into this shared vector space using modality-specific encoders.

![](_page_0_Figure_12.jpeg)

Approach 2: time-consuming, costly, and error prone

### **Contrastive learning for dense retrieval**

![](_page_0_Picture_15.jpeg)

- Dense retrieval training tends to benefit from a larger set of hard negatives and inbatch negatives.
- In our task, positive (relevant) and hard negative (irrelevant) samples are not directly available.
- Therefore, it is necessary to extract an appropriate set of hard negative samples and positive sample.

#### **Object detection-based approach**

#### **Data augmentation-based approach**

Image data augmentation

![](_page_0_Figure_23.jpeg)

### Results

- The experimental evaluation utilizes LSC'24 dataset, which is a multi-modal dataset capturing users' daily activities.
- The dataset consists of images captured by egocentric camera and sensor data by a smart tracker.

#### **Details of the experimental dataset**

Split	Period	# Data Points	# Days
Training	Jan. 01, 2020 ~ Feb. 29, 2020	27,907	60
Validation	Mar. 01, 2020 ~ Mar. 16, 2020	21,198	16

#### **Official Results (evaluated by Mean Reciprocal Rank)**

![](_page_0_Figure_30.jpeg)

#### **Examples of created topics**

Top	ic ID	Timestamp	Query	Relevant Data
\$en2ir	ug_0316	2020-03-16 16:50:00	hr: 0.015 lat: -0.121 lng: 0.055	20200316_165048.jpg
ing2s	m_0319	2020-03-19 09:03:00	20200319_090305.jpg	hr: 1.891 lat: -0.121 lng: 0.055

![](_page_0_Picture_33.jpeg)

![](_page_0_Picture_34.jpeg)

 $20200316\_165048.jpg$ 

![](_page_0_Picture_36.jpeg)

20200319\_090305.jpg

Experimental setup				
lmage Encoder	<ul> <li>Linearly transformed d</li> <li>3 x 224 x 224.</li> <li>Two Resnet50 pretrained by ImageNet/Places365</li> </ul>			
Sensor Encoder	<ul> <li>Normalized to avg. 0 and s.d. 1.</li> <li>Two MLP layers.</li> </ul>			
# dense dimentions	512			
Mini-batch	256			
Optimizer	Adam			
# Iterations	100, 200, 500			
Evaluation metrics	$MRR = \frac{1}{ Q } \sum_{i=1}^{Q} \frac{1}{rank_i}$			