Overview of NTCIR-18 Transfer-2 Task

Hideo Joho¹, Atsushi Keyaki², Yuuki Tachioka³, Shuhei Yamamoto¹

¹University of Tsukuba, Japan ²Hitotsubashi University, Japan ³Denso IT Laboratory, Japan

Introduction

The first round of the NTCIR Transfer Task, called Transfer-1, was a pilot task held at NTCIR-17. The outcomes of Transfer-1 demonstrated the potential of dense retrieval methods for Japanese, but also underscored the need for more diverse and comprehensive datasets to address the full range of challenges in the domain. Based on the insight from the first round, Transfer-2 at NTCIR-18 sets two subtasks: Retrieval Augmented Generation (RAG) and Dense Multimodal Retrieval (DMR). A total of 14 runs was submitted from two teams and organiser.



Figure 2: Overview of the DMR subtask.

To effectively advance this subtask, we reuse Lifelog Search Challenge 2024 (LSC'24) dataset [?], one of the largest multi-modal datasets derived from users' daily activi-

Retrieval Augmented Generation (RAG)

RAG subtask aims to develop a retrieval module suitable for Retrieval-Augmented Generation (RAG). RAG utilizes external knowledge retrieved by a retrieval module during response generation by an LLM to produce high-quality responses. At Transfer-2, we focused on the open-domain factoid question answering (QA) task for the RAG subtask because facts are presented in response to questions in the open-domain factoid QA task, making evaluation easier.



Figure 1: Overview of the RAG subtask.

We used the AIO Official Dataset Version 2.0 which contains passage-level relevance assessments for the retriever and answer sets for the reader in the question-answering

ties. The dataset was generated by one active lifelogger and is 18 months in length. It includes non-linguistic modalities such as ego-centric images, heart rate, and location information (e.g., latitude/longitude), providing a robust foundation for testing our dense retrieval approach in a real-world, multi-modal context. Thereby, deemed conducive to the realization of the subtask's objectives.

We created 140 retrieval topics for the formal run. The topics were extracted from the period from March 1, 2020, to June 30, 2020. Each topic is generated on a daily basis, and sensors or images recorded at the same time as the query are extracted as relevant (correct) data. Thus, the retrieval target for each topic is the set of data recorded on each day, with only one relevant data set and 199 irrelevant data sets for each topic. Additionally, we designed two types of retrieval tasks: one that retrieves images from sensors (reffered to as *sen2img*) and another that retrieves sensors from images (reffered to as *img2sen*). Examples of created topics are shown in Figure 3. Note that the sensor data were normalized in all topics.

Topic ID	Timestamp	Query	Relevant Data
3*img2sen_0309	3*2020-03-09 08:39:24	3*20200309_083924.jpg	hr: 1.891
			lat: -0.121
			Ing: 0.055
3*sen2img_0626	3*2020-06-26 20:47:05	hr: 0.015	3*20200626_204705.jpg
		lat: -0.121	

stage, which aligns with the goals of our subtask. The dataset includes 22,335 QA pairs for training, 1,000 QA pairs for development, and 1,000 questions for testing. The target corpus for RAG consists of Wikipedia articles

Additionally, we have adopted a two-stage retrieval model for the open-domain factoid QA task, specifically the retriever-reader model used in open-domain QA systems. For the first stage, the input is a natural language question, and the output is the top 100 ranked passage IDs corresponding to the natural language question. We use training data and development and evaluation data from the AIO Official Dataset Version 2.0 for training and evaluation, respectively. The evaluation metrics are meanAveragePrecision(mAP), HitRate(HR)@k (k = 1, 5, 10, 50, 100) and nDCG@k (k = 1, 5, 10, 50, 100).

For the second stage, the input is a natural language question and the k passages retrieved in the first stage, and the output is the answer to the natural language question. The evaluation metric is the accuracy. The baseline model provided by the organizers adopts the Fusion-in-Decoder (FiD), a model specialized for open-domain QA tasks.

Dense Multimodal Retrieval (DMR)

The DMR subtask is deinfed as follows: given a non-linguistic modality query q and

Ing: 0.055



G@k 20200309_083924.jpg 20200626_204705.jpg Figure 3: Examples of created topics. Timestamps denote the time when the data was recorded. "hr", "lat", and "Ing" represent heart rate, latitude, and longitude, respecs re- tively.

The organizing team developed two baseline runs and conducted evaluations with 10 and 100 training iterations (ORG:baseline-10 and ORG:baseline-100).

Participated Systems and Results

We have four runs for RAG subtask and ten runs for DMR subtask. Please visit individual poster presentations for the details of participated systems and their results!

a corpus $\{d_1, d_2, \cdots, d_n\}$ consisiting of n data points in another modality query q and of cross-modal retrieval is to find the k data points that are most relevant to the query q, where $k \ll n$. The task is to efficiently generate the top-k candidates that are relevant to the query based on the similarity metric $sim(q, d) \in \mathbb{R}$. In this context, the query q and data d can refer to sensor data $\{q^{sen}, d^{sen}\} \in \mathbb{R}^S$ or image data $\{q^{img}, d^{img}\} \in \mathbb{R}^{3 \times W \times H}$, where S denotes the number of dimensions of the sensor data and W and H represent the width and height of the image data. The evaluation metric is the mean reciprocal rank (MRR).

Conclusions & Acknowledgements

In conclusion, the Transfer-2 task provided valuable benchmarks, resources, and insights for future research on retrieval augmented applications and cross-modal information access. We thank the orgnizers of AIO and NTCIR-18 Lifelog-6 Task for making their datasets available to us.

NTCIR-18 Conference, 10-13 June 2025, NII, Tokyo, Japan