# IMU2CLIP: MULTIMODAL CONTRASTIVE LEARNING FOR IMU MOTION SENSORS FROM EGOCENTRIC VIDEOS AND TEXT NARRATIONS

*Seungwhan Moon\*, Andrea Madotto\*, Zhaojiang Lin, Alireza Dirafzoon, Aparajita Saraf*
*Amy Bearman, Babak Damavandi*

Meta Reality Labs

## ABSTRACT

We present IMU2CLIP, a novel pre-training approach to align Inertial Measurement Unit (IMU) motion sensor recordings with video and text, by projecting them into the joint representation space of Contrastive Language-Image Pre-training (CLIP). The proposed approach allows IMU2CLIP to translate human motions (as measured by IMU sensors) into their corresponding textual descriptions and videos – while preserving the *transitivity* across these modalities.

We explore several *new* IMU-based applications that IMU2CLIP enables, such as motion-based media retrieval and natural language reasoning tasks with motion data. In addition, we show that IMU2CLIP can significantly improve the downstream performance when fine-tuned for each application (*e.g.* activity recognition), demonstrating the universal usage of IMU2CLIP as a new pre-trained resource. Our code will be made publicly available.

***Index Terms*—** IMU modeling, Multimodal learning

## 1. INTRODUCTION

With the growing popularity of smart glasses or new-generation wearable devices, *first-person* or *egocentric* videos have recently become much more prevalent than ever before [1, 2, 3]. These egocentric videos are often accompanied by the parallel head-mounted IMU sensor readings, which record devices' linear and rotational movements and accelerations.

Given its low power consumption and low privacy implications, IMU is regarded as an important modality for powering various on-device models that require understanding of device wearer's movement patterns (*e.g.* exercise / activity recognition for health applications). The previous works on IMU modeling typically focus on the purpose-built datasets with manual annotations [4, 5], which are limited in their scale. Consequently, the utilization of IMU models in real-world scenarios has been confined to a relatively small number of use cases.

On the contrary, for the modalities that are widely studied (*e.g.* text, video), there are vast large-scale resources such as
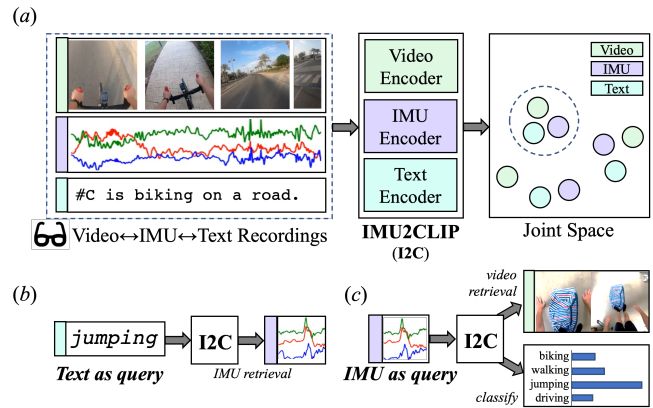


**Fig. 1**: Illustration of IMU2CLIP (I2C): (*a*) The model aligns the parallel video↔IMU↔text data in the joint space. Once trained, IMU2CLIP is used as a retriever for both (*b*) IMU and (*c*) videos, or as a classifier for downstream applications.

BERT [6] and GPT [7] for text, or CLIP4Clip [8] for videos. These powerful pre-trained resources have driven the development of many application-oriented models, showing significant improvements when fine-tuned for each respective task [9]. To the best of our knowledge, however, the study on the equivalent resources for encoding IMU signals has been lacking.

Inspired by the recent works that leverage large pre-trained models for other modalities, we present IMU2CLIP, a new approach to pre-train an IMU encoder by aligning the parallel Video ↔ IMU ↔ Text data in an un-supervised manner via multimodal contrastive training. Specifically, we propose to use CLIP [10], which contains the video encoder and the language model pre-trained on the large parallel image-text data, from which the IMU encoder can learn a semantic representation of various scenes transferred from other modalities.

To show the efficacy of the proposed approach, we evaluate our models on several benchmark tasks as well as new applications that IMU2CLIP enables, such as IMU-based media retrieval, leveraging the modality-transitivity that IMU2CLIP exhibits (Fig. 1). Most importantly, we show that the fine-tuned IMU2CLIP can significantly improve the performance of several downstream tasks, when compared to the identical IMU model trained from scratch.
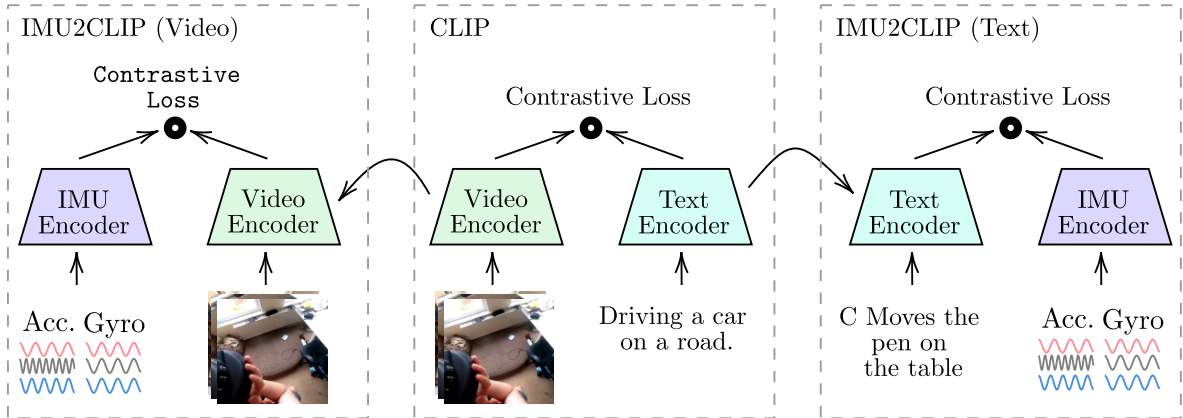
---

\*: Joint First Authors.

**Fig. 2**: Illustration of the proposed multimodal contrastive learning for IMU2CLIP. CLIP [14] is used to align IMU↔Video (left), and IMU↔Text (right). ❄: the parameters of the encoder are frozen during training.

**Our contributions** are as follows: (1) We propose a novel large-scale pre-training approach for IMU sensors, and release the resulting large pre-trained IMU encoders for future research. (2) We provide an in-depth empirical analysis evaluating the pre-trained models, for both upstream and downstream fine-tuning tasks. (3) Lastly, we present novel applications that show the feasibility of a wider usage for IMU sensor signals.

## 2. RELATED WORK

**Contrastive Learning** is as an efficient self-supervised framework applied across multiple domains, which learns similar/dissimilar representations from data that are organized into similar/dissimilar pairs. For instance, SimCLR [11] is a unimodal application of contrastive learning in the data augmentation setting, in which the authors propose to learn a vision encoder given a set of perturbed images. As an example in multimodal settings, Contrastive Language–Image Pre-training (CLIP) [10, 12] learns visual representations from natural language supervision using image and text pairs, achieving competitive results in *e.g.* zero-shot image classification, image retrieval via text, and image/caption generation. Similarly, WAV2CLIP [13] proposes to learn audio representation by distilling it from CLIP. We extend this line of work on contrastive learning to a unique multimodal setting that utilizes IMU signals, which is specific to a new generation of devices (such as smart glasses) that are equipped with such sensors.

**Pre-training Resources**: There are numerous pre-trained resources for well-studied modalities such as images or text. Many popular computer vision models (*e.g.* ResNet [15]) are typically trained on large supervised datasets such as ImageNet [16], *etc*. For language processing, the most popular language models (LM) include BERT [6, 17], GPT-2 [18], and GPT-3[19], which typically use self-superivsion techniques such as next-word predictions or masked token predictions, thus

without any explicit task labels. Studies report that these pre-trained resources achieve competitive zero-shot performance [10], and when fine-tuned, often outperform fully supervised models on several downstream tasks [9].

To our knowledge, the equivalent resource for encoding IMU signals is not made publicly available. Inspired by this line of work, we propose to perform large-scale pre-training for the unique sensor (IMU) signals dataset, and show that such pre-training significantly improves the performance for the downstream applications as well.

**Egocentric Datasets**: We are particularly interested in egocentric (first-person) datasets, for understanding of users' activities from head-mounted devices. Several data collection efforts have been made for building egocentric datasets, including Ego4D [1], Epic-Kitchens [2], and Aria [3] datasets.

Using these datasets, we propose various sub-tasks that can effectively evaluate diverse capabilities of IMU2CLIP, and demonstrate the feasibility of future applications. In addition, we implement a universal multimodal data loader to allow for easy cross-modality and cross-domain (dataset) studies.

**IMU Modeling**: IMU signals have been widely used in various motion recognition tasks, such as pose estimation [20], walking speed estimation [20], foot placement prediction [5]. Various deep learning architectures have been explored for modeling IMU in downstream tasks, including Transformer-CNN based IMU models [4] for gesture recognition, 1D-CNN and GRU ensemble IMU models [21] for clinical balance assessment, and Bi-LSTM IMU models [22] for human activity recognition. Our work proposes a new IMU model architecture, and conducts ablation studies over other models above. Different from prior work modeling IMU in a specific task, however, our work focuses on learning general IMU representations by aligning IMU with other modalities (*e.g.* images and text), which can enable wider downstream applications.

## 3. METHODS

For IMU pre-training, we propose to align the parallel portions of the IMU $\leftrightarrow$ Video $\leftrightarrow$ (optionally) Text data, using multimodal cross-modal contrastive learning schemes [10, 8]. In a nutshell, we train an IMU encoder such that the IMU representation of a given clip resembles the representation of its corresponding video frames, and optionally, corresponding textual descriptions or narrations (Section 3.1). We perform ablation studies over multiple IMU encoder architectures, as detailed in Section 3.2. Fig. 2 illustrates the overall approach.

### 3.1. Cross-modal Contrastive Learning for IMU

We consider a batch of $B$ ground-truth $\underline{\text{IMU}} \leftrightarrow \underline{\text{V}}\text{ideo} \leftrightarrow \underline{\text{T}}\text{ext}$ parallel windows: $\{(\mathbf{i}_1, \mathbf{v}_1, \mathbf{t}_1), ..., (\mathbf{i}_B, \mathbf{v}_B, \mathbf{t}_B)\}$, where the embeddings of each modality lies on the unit hypersphere $S^D$. Since the embeddings are unit-normalized, the similarity can be simply calculated as their inner product:

$$\text{sim}(\mathbf{i}_i, \mathbf{v}_j) = \langle \mathbf{i}_i, \mathbf{v}_j \rangle \tag{1}$$

$$\text{sim}(\mathbf{i}_i, \mathbf{t}_j) = \langle \mathbf{i}_i, \mathbf{t}_j \rangle \tag{2}$$

We can then define the IMU-to-Video ($\mathbf{i2v}$) and IMU-to-Text ($\mathbf{i2t}$) retrieval distributions based on the cross-modal similarities across the parallel signals:

$$P_{\mathbf{i2v}}(\mathbf{v}_j|\mathbf{i}_i) = \frac{\exp(\text{sim}(\mathbf{i}_i, \mathbf{v}_j))^{1/\gamma}}{\sum_{k=1}^{B} \exp(\text{sim}(\mathbf{i}_i, \mathbf{v}_k))^{1/\gamma}} \tag{3}$$

$$P_{\mathbf{i2t}}(\mathbf{t}_j|\mathbf{i}_i) = \frac{\exp(\text{sim}(\mathbf{i}_i, \mathbf{t}_j))^{1/\gamma}}{\sum_{k=1}^{B} \exp(\text{sim}(\mathbf{i}_i, \mathbf{t}_k))^{1/\gamma}} \tag{4}$$

where $\gamma$ is a temperature parameter that controls the concentration of the distributions.

We then train three flavors of IMU2CLIP: (a) aligning IMU$\leftrightarrow$Video, (b) IMU$\leftrightarrow$Text (when the text narration data are available), and (c) IMU$\leftrightarrow$Video$\leftrightarrow$Text. Specifically, we propose to project the IMU representations into the joint CLIP space [10] to leverage the visual and textual knowledge already encoded in CLIP. Similar to [8, 10], we propose to minimize the symmetric cross-modal contrastive loss. For the IMU$\leftrightarrow$Video alignment, we use the sum of IMU-to-Video and Video-to-IMU cross-entropy losses:

$$\mathcal{L}_{\mathbf{i2v}} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(\text{sim}(\mathbf{i}_i, \mathbf{v}_i))^{1/\gamma}}{\sum_{k=1}^{B} \exp(\text{sim}(\mathbf{i}_i, \mathbf{v}_k))^{1/\gamma}}$$

$$\mathcal{L}_{\mathbf{v2i}} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(\text{sim}(\mathbf{i}_i, \mathbf{v}_i))^{1/\gamma}}{\sum_{k=1}^{B} \exp(\text{sim}(\mathbf{i}_k, \mathbf{v}_i))^{1/\gamma}}$$

$$\mathcal{L}_{\mathbf{i}\leftrightarrow\mathbf{v}} = \frac{1}{2}(\mathcal{L}_{\mathbf{i2v}} + \mathcal{L}_{\mathbf{v2i}}) \tag{5}$$

The loss for IMU$\leftrightarrow$Text alignment ($\mathcal{L}_{\mathbf{i}\leftrightarrow\mathbf{t}}$) can be defined similarly, and consequently $\mathcal{L}_{\mathbf{i}\leftrightarrow\mathbf{v}\leftrightarrow\mathbf{t}} = \mathcal{L}_{\mathbf{i}\leftrightarrow\mathbf{v}} + \mathcal{L}_{\mathbf{i}\leftrightarrow\mathbf{t}}$. To



**Fig. 3**: Illustration of the Stacked RNN architecture for the IMU encoder used in IMU2CLIP.

| Ego4d | Tra. | Val. | Tst. |
|---|---|---|---|
| # of Media files | 1444 | 161 | 688 |
| Total Media Durations | 540h | 60h | 265h |
| # of IMU$\leftrightarrow$Text/Video windows (5s) | 528K | 68K | 266K |
| # of IMU$\rightarrow$4 classes windows (5s) | 1552 | 760 | 241 |
| **Aria** | **Tra.** | **Val.** | **Tst.** |
| # of Media files | 747 | 259 | 277 |
| Total Media Durations | 138h | 43h | 51h |
| # of IMU$\leftrightarrow$Video windows (1s) | 496K | 157K | 184K |
| # of IMU$\rightarrow$5 classes windows (1s) | 25K | 138K | 162K |

**Table 1**: Dataset Statistics for Ego4D and Aria.

preserve the text-vision alignment that CLIP already exhibits, we freeze the parameters of the image and text CLIP encoders.

**Implementation details**. To expedite the training, we pre-process each media to have equal-sized parallel windows (IMU $\leftrightarrow$ Video $\leftrightarrow$ Text). The data module retrieves the parallel data of a requested window size at a given timestamp, and caches them for faster training. In addition, to accommodate the memory constraints, we pool the negative samples within the same batch (randomly shuffled), reducing the load on each GPU. We optimize the parameters with Adagrad [23] with batch size 16, learning rate 0.01, epsilon $10^{-8}$, and decay 0.1.

### 3.2. IMU Encoder Architectures

For the IMU encoder, we propose a stack of 1D-CNNs and RNN-based architecture (Figure 3), which performed the best in our ablation studies[1]. First, we perform a GroupNorm operation to normalize the Accelerometer (3D) and the Gyroscope (3D) signals independently. We then perform a stack of $N$ 1D-CNN, a Max Pooling with kernel size 5, and then another GroupNorm to normalize the output features. Finally, we use an RNN (*i.e.* GRU in our experiments) to combine the CNN output, and thus generating the final embedded representation.

## 4. EXPERIMENTS

### 4.1. Dataset

We use Ego4D [1] and Aria [3] as the main datasets for the experiments below, both of which feature parallel video and

---

[1] Ablation results can be found: `tinyurl.com/imu2clip-experiments`

| Train Modalities | | | IMU → Text | | | | Text → IMU | | | | IMU → Video | | | | Video → IMU | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IMU | Video | Text | R@1 | R@10 | R@50 | MRR | R@1 | R@10 | R@50 | MRR | R@1 | R@10 | R@50 | MRR | R@1 | R@10 | R@50 | MRR |
| ✓ | ✓ | | 4.86 | 18.75 | 48.26 | 0.104 | 4.17 | 15.62 | 43.06 | 0.084 | 9.06 | 43.13 | 78.75 | 0.2011 | 12.19 | 45.31 | 80.00 | 0.226 |
| ✓ | | ✓ | 5.21 | 25.00 | 60.42 | 0.123 | 7.29 | 28.82 | 60.07 | 0.143 | 3.75 | 25.94 | 62.81 | 0.105 | 3.75 | 24.06 | 56.88 | 0.098 |
| ✓ | ✓ | ✓ | 4.52 | 22.91 | 56.60 | 0.118 | 5.90 | 22.92 | 56.60 | 0.139 | 8.75 | 40.63 | 73.44 | 0.183 | 11.56 | 42.19 | 75.94 | 0.213 |
| | | | **(Video → Text)** | | | | **(Text → Video)** | | | | | | | | | | | |
| (CLIP) | ✓ | ✓ | 6.94 | 32.29 | 64.24 | 0.150 | 8.33 | 33.68 | 65.28 | 0.168 | | | | | | | | |

**Table 2**: Text↔IMU and Video↔IMU retrieval performances of the pre-trained IMU2CLIP models on Ego4D, with different modalities used for training. The last row shows the video retrieval performance of OpenAI's CLIP model on the same test set.

IMU signals. For Ego4D, a subset of the clips are also annotated with their corresponding narrations. We split the data into train, validation, and test sets (split by video IDs). The statistics of the datasets are provided in Table 1.

## 4.2. Tasks

Note that the proposed pre-training approach enforces the alignment of the IMU, video, and text representations, which allows for new and unique cross-modal applications. We propose the following real-world downstream applications as the novel tasks to evaluate the performance of the IMU encoders.

**Task 1. IMU Retrieval via Textual Queries (Text→IMU)**, where the goal is to retrieve a window of IMU signals given free-form textual queries. Once the IMU signals are retrieved, we can also retrieve the corresponding videos, thus allowing for a new and power-efficient way of performing media retrieval or online action detection. The retrieval performance is measured on the held-out test set (Recall@$k$ and Mean Reciprocal Rank (MRR)), using the text narrations as queries and the IMU signals as the retrieval pool.

**Task 2. Video Retrieval based on IMU (IMU→Video)**, where the goal is to retrieve videos based on IMU signals, allowing for an intuitive way of analyzing motion signals data. We measure the performance on the held-out test set, using the IMU signals as queries and the videos as the retrieval target.

**Task 3. IMU-based Activity Recognition**, where the goal is to predict a discrete activity label given a window of IMU signals. We use the manual motion annotations for Aria (*e.g.* hiking, running, biking) and soft annotations for Ego4D via text matching of the narrations provided.

## 4.3. Results

Table 2 shows the IMU↔Text and IMU↔Video retrieval performance on the Ego4D test set, of IMU2CLIP trained via different combinations of modalities.

**Results 1: IMU-based media search with textual queries**: The Text→IMU column in Table 2 shows performances on Task 1. Note that the CLIP embeddings already exhibit the Video ↔ Text transitivity, and thus IMU2CLIP trained using IMU ↔ Video achieves a competitive zeroshot performance
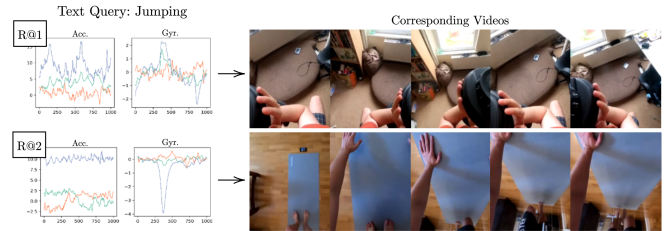


**Fig. 4**: Illustration of IMU-based media retrieval. Given a free-form textual query (*e.g.* "*jumping*"), (Left): IMU2CLIP's predictions of the semantically closest IMU signals from the Ego4D test set (top-2). (Right): the gold-parallel videos corresponding to the retrieved IMU signals (as a reference). The retrieved media match the semantics of the input query.

for IMU ↔ Text retrieval as well. When text narrations are used for pre-training, the model achieves an even higher recall performance. See Fig. 4 for visualizations[2].

To help contextualize the recall performances in Table 2, we also show (as a reference) the Text ↔ Video retrieval performance of the near-SOTA video encoder (CLIP) (bottom). The narrow margin in the performances (*e.g.* MRR=0.143 for Text→IMU *vs.* MRR=0.168 for Text→Video) shows that the IMU encoder could serve as a power-efficient alternative for a video encoder in many applications.

**Results 2: We can search for videos, given IMU recordings**. The IMU→Video column in Table 2 shows the Ego4D performances on Task 2 (See Fig. 5 for examples). We observe higher recall performances in general, showing that the IMU signals and videos have a higher compatibility. When the model is trained on all three modalities, we observe competitive results across all tasks, while the best performances are from the bi-modal models aligned with each respective task.

We observe similar patterns on the Aria data as well: for the IMU→Video retrieval, IMU2CLIP achieves MRR=0.182, and R@$\{1, 10, 50\}$ of $\{8.48, 38.83, 77.67\}$, respectively. For the Video→IMU retrieval, IMU2CLIP achieves MRR=0.190, and R@$\{1, 10, 50\}$ of $\{8.48, 44.19, 78.57\}$. Note that the Aria dataset does not have text narrations annotated.

---

[2]For better readability, we also provide the animated GIF visualizations for all experiments at: tinyurl.com/imu2clip-visualizations

**Fig. 5**: Illustration of the IMU to Video Retrieval. (Top): IMU signals and their corresponding ground-truth video. (Bottom): IMU2CLIP 's model predictions of their corresponding videos from the Ego4D test set (top-5), given the IMU signals. It can be seen that the videos retrieved based on IMU are visually and semantically similar to the gold video.

| Models | | Ego4D | | Aria | |
|---|---|---|---|---|---|
| | | F1 | Acc. | F1 | Acc. |
| Random Init. IMU Encoder | | 23.23 | 49.92 | 56.35 | 76.11 |
| IMU2CLIP $(\mathbf{i} \leftrightarrow \mathbf{v})$ | + Zeroshot | 19.39 | 23.08 | 18.46 | 21.52 |
| | + Probing | 40.55 | 61.46 | **62.52** | **83.54** |
| | + Fine tuning | 43.07 | **65.87** | 61.77 | 82.31 |
| IMU2CLIP $(\mathbf{i} \leftrightarrow \mathbf{t})$ | + Zeroshot | 31.89 | 36.38 | - | - |
| | + Probing | 45.12 | 58.01 | - | - |
| | + Fine tuning | **45.15** | 63.14 | - | - |
| IMU2CLIP $(\mathbf{i} \leftrightarrow \mathbf{v} \leftrightarrow \mathbf{t})$ | + Zeroshot | 27.24 | 24.26 | - | - |
| | + Probing | 42.16 | 55.21 | - | - |
| | + Fine tuning | 44.17 | 62.66 | - | - |

**Table 3**: IMU-based activity recognition on Ego4D and Aria datasets, comparing the randomly initialized model and the pre-trained IMU2CLIP models, with IMU↔Video, IMU↔Text and IMU↔Video↔Text pre-training. **Bold** denotes the best performance for each metric: F1 and Accuracy (Acc).

**Results 3: Fine-tuned IMU2CLIP significantly outperforms the vanilla model with the same architecture, on downstream tasks.** Table 3 shows the activity recognition results on Ego4D and Aria datasets. For all experiments, we use the same IMU architecture (Stacked RNN). For zeroshot experiments, we encode the surface names of each activity (*e.g.* hiking) with the CLIP text encoder, and use the nearest-neighbor classifier on the projected IMU embeddings (thus not using any supervision labels). Probing adds a linear layer on top of the IMU encoder while keeping the IMU encoder frozen, and for fine-tuning we allow all parameters of the encoder to be trainable. The consistent improvements in the fine-tuning performances (*e.g.* ∼16 points absolute improvement in accuracy for Ego4D, comparing the randomly initialized vanilla model *vs.* fine-tuned model) show that IMU2CLIP can learn high quality representations for IMU signals.

Comparing the pre-trained models trained via various combinations of modalities again shows that IMU2CLIP preserves the transitivity among modalities (video ↔ text ↔ IMU).



```
[15:44] Motion: no activity, Audio: no activity
[15:45] Motion: walking, Audio: TV noise
[15:50] Motion: sits down, Audio: watching TV
[15:52] Motion: look around, Audio: kitchen sounds
[16:03] Motion: walking, Audio: no activity
Question: how long did I sit down and focus on watching TV?
Answer: From 15:50 to 16:03, for a total of 13 minutes.
```

**Fig. 6**: Demonstration of an LM-based multimodal reasoning model, using two ambient sensors: IMU and audio. Given the sensor logs and the question, LM generates a response grounded on the multimodal context.

## 4.4. Qualitative Analysis: Multimodal Reasoning with Ambient Sensors

Further exploring the benefit of IMU2CLIP that translates sensor signals into text, we demonstrate a multimodal reasoning model that operates only on the ambient sensor logs (Fig. 6). Specifically, we run IMU2CLIP as a zeroshot tagging model on the clips from Ego4D, to obtain textual descriptions of the IMU and the Audio sensor readings. We then use a large LM (*i.e.* GPT-3 [19]) to use the sensor logs as conditioning context to answer summarizing questions such as: "*What can you tell me about the user activity using the motion logs?*", or memory recall prompts such as: "*What time did I start biking?*" The LM then generates a response via a causal language inference step, which completes the process for zeroshot multimodal reasoning. Unlike the similar approach such as Socratic Models [24], the proposed approach does not rely on the video signals at all – which would incur much higher power consumption – thus operating better under real-world constraints.

## 5. CONCLUSIONS

With the growing popularity of wearable devices of diverse form factors (*e.g.* smart glasses), it is important to study the capability of the ambient sensors such as IMU motion signals. To this end, we propose a new multimodal contrastive training approach for representing IMU signals, and release the pre-trained IMU encoders to be used for future research. Our empirical analysis highlights the efficacy of the proposed approach on many existing and new IMU-based applications. In addition, we show that IMU2CLIP can significantly improve the downstream performance when fine-tuned, demonstrating the universal usage of IMU2CLIP.

## 6. REFERENCES

[1] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al., "Ego4d: Around the world in 3,000 hours of egocentric video," in *CVPR*, 2022.

[2] Dima Damen, Adriano Fragomeni, Jonathan Munro, Toby Perrett, Daniel Whettam, Michael Wray, Antonino Furnari, Giovanni Maria Farinella, and Davide Moltisanti, "Epic-kitchens-100-2021 challenges report," 2021.

[3] Zhaoyang Lv, Edward Miller, Jeff Meissner, Luis Pesqueira, Chris Sweeney, Jing Dong, Lingni Ma, Pratik Patel, Pierre Moulon, Kiran Somasundaram, Omkar Parkhi, Yuyang Zou, Nikhil Raina, Steve Saarinen, Yusuf M Mansour, Po-Kang Huang, Zijian Wang, Anton Troynikov, Raul Mur Artal, Daniel DeTone, Daniel Barnes, Elizabeth Argall, Andrey Lobanovskiy, David Jaeyun Kim, Philippe Bouttefroy, Julian Straub, Jakob Julian Engel, Prince Gupta, Mingfei Yan, Renzo De Nardi, and Richard Newcombe, "Aria pilot dataset," 2022.

[4] Yujian Jiang, Lin Song, Junming Zhang, Yang Song, and Ming Yan, "Multi-category gesture recognition modeling based on semg and imu signals," *Sensors*, 2022.

[5] Xinxing Chen, Kuangen Zhang, Haiyuan Liu, Yuquan Leng, and Chenglong Fu, "A probability distribution model-based approach for foot placement prediction in the early swing phase with a wearable imu sensor," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 29, pp. 2595–2604, 2021.

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[7] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al., "Improving language understanding by generative pre-training," 2018.

[8] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li, "CLIP4Clip: An empirical study of clip for end to end video clip retrieval," *arXiv preprint arXiv:2104.08860*, 2021.

[9] Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith, "Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping," *arXiv preprint arXiv:2002.06305*, 2020.

[10] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, 2020.

[12] Kevin Qinghong Lin, Alex Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Zhongcong Xu, Difei Gao, Rongcheng Tu, Wenzhe Zhao, Weijie Kong, et al., "Egocentric video-language pretraining," *arXiv preprint arXiv:2206.01670*, 2022.

[13] Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello, "Wav2clip: Learning robust audio representations from clip," in *ICASSP*, 2022.

[14] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009.

[17] Nils Reimers and Iryna Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *EMNLP*, 2019.

[18] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al., "Language models are unsupervised multitask learners," 2019.

[19] Luciano Floridi and Massimo Chiriatti, "Gpt-3: Its nature, scope, limits, and consequences," *Minds and Machines*, vol. 30, no. 4, pp. 681–694, 2020.

[20] Shaghayegh Zihajehzadeh and Edward J Park, "A gaussian process regression model for walking speed estimation using a head-worn imu," in *EMBC*, 2017.

[21] Yeon-Wook Kim, Kyung-Lim Joa, Han-Young Jeong, and Sangmin Lee, "Wearable imu-based human activity recognition algorithm for clinical balance assessment using 1d-cnn and gru ensemble model," *Sensors*, 2021.

[22] Sara Ashry, Tetsuji Ogawa, and Walid Gomaa, "Charmdeep: Continuous human activity recognition model based on deep neural network using imu sensors of smartwatch," *IEEE Sensors Journal*, vol. 20, no. 15, pp. 8757–8770, 2020.

[23] John Duchi, Elad Hazan, and Yoram Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *JMLR*, 2011.

[24] Andy Zeng, Adrian Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, et al., "Socratic models: Composing zero-shot multimodal reasoning with language," *arXiv preprint arXiv:2204.00598*, 2022.