YOUR ACTIONS TALK: DUET – A MULTIMODAL DATASET FOR CONTEXTUALIZABLE DYADIC ACTIVI TIES

Anonymous authors

005 006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

026

028

030

031

032

034

040

041

Paper under double-blind review

ABSTRACT

Human activity recognition (HAR) has advanced significantly with the availability of diverse datasets, yet the field remains limited by a scarcity of datasets focused on two-person, or "dyadic," interactions. Existing datasets primarily cater to single-person activities, overlooking the complex dynamics and contextual dependencies present in interactions between two individuals. Failing to extend HAR to dyadic settings limits opportunities to advance areas like collaborative learning, healthcare, robotics, augmented reality, and psychological assessments, which require an understanding of interpersonal dynamics. To address this gap, we introduce the Dyadic User Engagement dataseT (DUET), a comprehensive dataset designed to enhance the understanding and recognition of dyadic activities. DUET comprises 14,400 video samples across 12 interaction classes, capturing the highest sample-to-class ratio of dyadic datasets known to date. Each sample is recorded using RGB, depth, infrared, and 3D skeleton joints, ensuring a robust dataset for multimodal analysis. Critically, DUET features a taxonomization of interactions based on five fundamental communication functions: emblems, illustrators, affect displays, regulators, and adaptors. This classification, rooted in psychology, supports dyadic human activity contextualization by extracting the embedded semantics of bodily movements. Data collection was conducted at three locations using a novel technique that captures interactions from multiple views with a single camera, thereby improving model resilience against background noise and view variations. We benchmark six state-of-the-art, open-source HAR algorithms on DUET, demonstrating the dataset's complexity and current HAR models' limitations in recognizing dyadic interactions. Our results highlight the need for further research into multimodal and context-aware HAR for dyadic interactions, and provide a dataset to support this advancement. DUET is publicly available at "Anonymized DUET Repository", providing a valuable resource for the research community to advance HAR in dyadic settings.

- 038 1 INTRODUCTION
 - 1.1 MOTIVATION

Human activity recognition (HAR) is a field within artificial intelligence focused on identifying and
analyzing human actions from sensor data, and it has achieved significant success across various
domains. The success of HAR can be attributed to many factors, including the commitment of
the field to producing publicly available datasets that can be used to help refine data-driven deep
learning algorithms across various contexts. While there is an abundance of HAR datasets already
available, the majority pertain to single-person—or *monadic*—activities. A better understanding of
two-person—or *dyadic*—interactions is essential for enhancing the accuracy, responsiveness, and
overall capabilities of systems where human interaction plays a central role.

Dyadic interactions, which involve the interplay between two individuals, convey deeper communicative and cultural significance. Despite their complexity, HAR for dyadic interactions offers several advantages. The inclusion of a second subject improves the performance of many HAR tasks
by introducing an additional distinguishing factor (Adeli et al., 2020). For instance, consider the actions "waving in" and "thumbs up." These two movements appear similar at first glance, as both

054 involve extending one's arm. However, their small-scale hand movements differ only slightly, mak-055 ing them difficult to distinguish in isolation. The distinction becomes much clearer when another 056 subject is involved—specifically, by observing the initiating action of one subject and the reaction 057 of the other. The reacting subject may physically approach if the initiating subject waves them in, 058 while they may simply nod in acknowledgment if the initiating subject gives a "thumbs up". These differing responses provide valuable contextual cues that improve the accuracy of recognizing and differentiating between the two activities. The study of dyadic interactions allows for a more ac-060 curate understanding of human behaviors that are absent in monadic activities, enabling systems to 061 better interpret and respond to social dynamics. For example, telepresence avatars in augmented 062 reality provide digital representations of participants during remote conferences. By analyzing the 063 subtleties of interactions between individuals, dyadic activity analysis enhances user experience 064 and increases the authenticity of digital environments (Ahuja et al., 2019). Similarly, social robots 065 designed to provide companionship for children leverage dyadic analysis to recognize dangerous 066 situations and intervene in a timely, adaptive manner. These robots also utilize two-person datasets 067 to deliver more natural and engaging conversational interactions, supporting the social, cognitive, 068 and emotional development of children (Chen et al., 2022). Additionally, in public infrastructure, 069 accurately recognizing dyadic social activities enhances safety by detecting potential dangers and enables the provision of more personalized services in public spaces (Coppola et al., 2020). 070

071

073

072 1.2 REVIEW OF EXISTING DATASETS

Despite these advantages, the availability of dyadic datasets remains limited, particularly in compari-074 son to the abundance of monadic datasets. This scarcity poses an increasing challenge as interactions 075 between humans and technical systems grow more complex. The research community's uneven em-076 phasis on these two activities types is reflected in their differing recognition performance levels. 077 Lin et al. (2024b) showed that monadic algorithms, which have achieved outstanding benchmarking records for monadic activities, do not perform nearly as well for dyadic interactions. This highlights 079 the disparity between monadic and dyadic activities, which stems from the greater variety of expressive and cultural signals, as well as the increased complexity of spatial and temporal coordination 081 between two or more subjects. To reconcile this discrepancy and improve dyadic HAR, there is a need for more datasets tailored to dyadic interactions. As highlighted in the IEEE Control Systems 083 Society's report on control for societal-scale challenges, traditional boundaries between humans and technology are blurring, and emerging fields like cyber-physical-human systems (CPHS) face chal-084 lenges in designing robust interactions between humans and control systems (Annaswamy et al., 085 2023). One of the central CPHS research challenges is characterizing how humans adapt during interactions. Dyadic datasets are critical for developing models that can enhance system adaptability, 087 safety, and trustworthiness in these complex environments (Annaswamy et al., 2023). 088

- Besides increasing the number, diversity, and quality of dyadic datasets, contextualizing activities 089 has proven effective in improving the performance of HAR tasks (Niemann et al., 2021). Contextualization distills meanings embedded in body language, such as emotional and cultural significance, 091 adding another layer of comprehension to the tracking of bodily movements. For instance, a "thumbs 092 up" signifies approval in most Western cultures but represents a profanity in Greece and several Middle Eastern countries. Contextualization enables the interpretation of the cultural significance 094 of gestures, such as recognizing the nuanced meanings of a "thumbs up." In addition to enhancing 095 HAR accuracy, contextualization supports the development of various downstream applications. For 096 example, certain branches of CPHS investigate how humans interact with and benefit from the built 097 environment (Doctorarastoo et al., 2023a;b). A critical aspect of this framework is understanding the 098 embedded semantics of human behaviors through bodily movements. This understanding provides stakeholders with deeper insights into system use, improving infrastructure design, maintenance, 099 and operation. Contextualization also paves the way for automating psychological and sociological 100 assessments—such as sociometric tests (Moreno, 1941)—that currently rely on self-reported data. 101 These manual evaluations are labor-intensive and also prone to attribution bias. By integrating con-102 textualization with dyadic HAR, these processes can be automated, extracting user preferences from 103 bodily movements (Lin et al., 2024a) and addressing these limitations. For instance, contextualiza-104 tion enhances telepresence avatars by capturing nonverbal cues and paralinguistic signals, improving 105 the quality and authenticity of remote communication (Ahuja et al., 2019). 106
- 107 Despite these recognized benefits, the few available dyadic datasets—listed in Table 1—are inadequate for extracting the underlying semantics of bodily movements. While some datasets focus

on healthcare activities, others are restricted to tracking bodily movements within specific action categories. No existing dataset selects activity classes using scientifically grounded methods that prioritize semantic cohesion to capture the social embeddings of activities. This lack of structured selection limits the ability to understand functional relationships between actions, hindering models from generalizing effectively to new, unlabeled behaviors. A dyadic dataset that fully supports contextualization is still absent in the research community.

114

115 1.3 OBJECTIVES AND NOVELTY OF THIS PAPER

To enhance HAR performance for dyadic activities through contextualization, we introduce the
Dyadic User Engagement dataseT (DUET). Featuring 12 taxonomized interactions, DUET helps
to bridge monadic and dyadic HAR while connecting HAR to other disciplines. It is publicly available under an MIT License at "Anonymized DUET Repository" (Authors, 2024).

121 Instead of repeating previous approaches that arbitrarily select activity categories, our dataset is built 122 on a psychology-based classification that identifies five core communication functions in human in-123 teractions: emblems, illustrators, affect displays, regulators, and adaptors. This taxonomy provides 124 a scientifically grounded framework for integrating HAR with interdisciplinary applications. For 125 instance, lie detection often relies on emblematic slips—unconscious, fragmented gestures that deceivers attempt to suppress while lying. Similarly, emotion detection heavily depends on adaptors, 126 which reveal physical or emotional discomfort. By capturing interactions from all five categories, 127 DUET addresses critical gaps left by existing datasets. This stands in stark contrast to existing 128 datasets, as shown in Table 1, particularly the largest dyadic dataset to date, NTU RGB+D 120 129 (Liu et al., 2019). While NTU RGB+D 120 includes dyadic interactions, it represents only three 130 of the five categories—illustrators, affect displays, and regulators. This imbalance prevents it from 131 supporting applications that require a comprehensive understanding of the taxonomy. In contrast, 132 DUET deliberately incorporates interactions from all five categories, ensuring semantic cohesion 133 and preserving the functional relationships between actions. This design enables HAR to generalize 134 more effectively, recognizing both labeled and unlabeled actions by aligning them with shared traits 135 of existing categories. As a result, DUET facilitates connections between HAR and fields like psy-136 chology, sociology, and behavioral sciences, paving the way for applications like automated emotion 137 recognition and the analysis of social behaviors in complex, real-world scenarios. By bridging this gap, DUET stands as a critical step toward advancing both HAR and its wider applications. 138

The dataset was collected using the Microsoft Azure Kinect v2 (Microsoft, 2024a) (hereafter referred to as the Azure Kinect), a high-quality, multimodal camera capable of capturing RGB, infrared (IR), depth, and 3D skeletal joint data. Over the span of one year, 23 participants contributed to the dataset, generating 14,400 video samples, with 1,200 samples recorded for each interaction category. To our knowledge, this dataset features the highest sample-to-class ratio published to date.

144 The testbeds consist of three locations across a university campus in the United States (US): an open 145 indoor space, a confined indoor space, and an outdoor area. These settings were chosen to represent 146 a range of environments where human activities commonly take place. For example, the compan-147 ionship and support provided by social robots may take place in a small bedroom (confined indoor space). A sociological evaluation of a group of students' social connectivity during class could be 148 conducted in a large auditorium (open indoor space). A potential application for CPHS is to redesign 149 public open spaces based on patterns of measured usage and socialization to foster user sociability 150 and cohesion (open outdoor space). This variety not only allows downstream applications to lever-151 age DUET for investigating the direct and indirect impacts of ambient surroundings on algorithm 152 performance but also improves the resilience of deep learning models against background noise. We 153 intentionally collected data at various times and on different days to capture a range of environmental 154 conditions (e.g., lighting) and background noise, ensuring the dataset reflects real-world scenarios. 155 Another challenge is the limited number of views, which can affect system robustness (Perera et al., 156 2020). The lack of multiple views in existing literature (Table 1) undermines the generalizability 157 and view-invariance of video samples from different orientations. To address this, we propose a 158 novel data collection process that captures interactions from multiple angles using a single camera. 159 This low-cost approach captures activities from various orientations, something that even multiple cameras have struggled to achieve. We evaluate the performance of six open-source, state-of-the-art 160 algorithms using RGB, depth, and 3D skeleton joint data. This comparison not only highlights the 161 complexity of contextualizable dyadic interactions but also reveals the strengths of each modality.

Dataset	Modalities	#Videos	#Classes	#Loca- tions	#Views	Backgro- und noise	Year
UT Interaction (Ryoo et al., 2010)	RGB	160	6	2	1	No	2010
SBU Kinect (Yun et al., 2012)	RGB+D+J	300	8	1	1	No	2012
JPL Interaction (Ryoo & Matthies, 2013)	RGB	399	7	5	1	No	2013
G3Di (Bloom et al., 2016)	RGB+D+J	168	14	1	1	No	2015
M ² I (Liu et al., 2018)	RGB+D+J	1,760	9	1	2	No	2015
ShakeFive 2 (Van Gemeren et al., 2016)	RGB+J	153	8	1	1	No	2016
PKU-MMD (Liu et al., 2017)	RGB+D+J+IR	4225	10	1	3	No	2017
MMAct (Kong et al., 2019)	RGB+keypoints+ acceleration+ orientation+ Wi-FI+Pressure	2162	2	4	4+ego	No	2019
NTU RGB+D 120 (Liu et al., 2019)	RGB+D+J+IR	24,828	26	-	155	No	2019
Air Act2Act (Ko et al., 2021)	RGB+D+J	5,000	10	2	3	No	2020
DUET (our dataset)	RGB+D+J+IR	14,400	12	3	360	No	2024

Table 1: comparison of existing dyadic datasets shows that the proposed dataset has the *highest number of samples per class*, the most views, and a relatively high number of locations. Note: (1)
"Views" refer to different sensor orientations from which interactions are captured, and (2) "background noise" indicates the presence of random people's movement or cluttered environments.

189 190

191

192

193

The remainder of the paper is structured as follows. Section 2 details the taxonomy for classifying human interactions. Section 3 overviews the dataset, including modalities, format, acquisition configurations, biometrics, annotation, and data splits for cross-location and cross-subject evaluations. Section 4 benchmarks six open-source algorithms and their results. Finally, Section 5 presents conclusions, key takeaways, and future directions.

2 CONTEXTUALIZING HUMAN INTERACTIONS

A social interaction is an exchange of information between two or more individuals, and the deliv-199 ery can happen through various channels. Among all communication channels, bodily movement 200 represents a critical part of social interaction as instinctive actions convey unspoken cues of the con-201 versation or communication (Sharan et al., 2022). To study the context embedded in social interac-202 tions through bodily movement, generating a dataset that attempts to exhaust all existing interactions 203 would not be feasible. Similarly, selecting actions arbitrarily would detract from the dataset's ability 204 to preserve functional relationships and shared characteristics among actions, which is needed to 205 create semantic cohesion across classes. In this work, we propose a dataset, DUET, in which the 206 selection of interactions is not arbitrary but instead grounded in psychological principles.

A total of 12 kinesic interactions are drawn from a taxonomy developed by Ekman & Friesen (1969), which classifies human interactions into five groups based on their fundamental communication functions. This system provides a structured and effective approach for categorizing interactions and systematically extracting the information embedded in bodily movements. The categories include emblems, illustrators, affect displays, regulators, and adaptors:

212

Emblems: Emblems are gestures that have direct verbal translation and can be culturally specific.
The same gesture might be interpreted differently for different cultures (Hartman, 2024). For instance, a "thumbs up" indicates well done in most Western cultures, but is a derogatory sign in Middle Eastern countries. Interactions chosen are "waving in," "thumbs-up," and "hand waving."

- Illustrators: Bodily movements that illustrate the verbal message they accompany are called illustrators, which are used to clarify conversations and are context dependent (Chute et al., 2023). Interactions chosen are "*pointing*" and "*showing measurements*."
- Affect displays: Affect displays are gestures that reveal one's affective and emotional state. An example of an affect display is "arm crossing," which can signal defensiveness, insecurity, or anxiety. Interactions chosen are "*hugging*," "*laughing*," and "*arm crossing*."
- Regulators: During interactions, regulators determine the alternation of instigating and receiving. Interactions chosen are "nodding," "writing circles in the air," and "holding one's palms out." Nodding is a gesture of acceptance and acknowledgement used for the continuation of the conversation. "Drawing circles in the air" displays the need to expedite the conversation. "Holding one's palms out" is used to warn the other person to cease the conversation.
 - Adaptors: Adaptors are habitual movements that satisfy personal needs and can be used to increase or decrease emotional stability (Neff et al., 2011). The interaction chosen is *"twirling or scratching hair"* to moderate one's stress during contemplation.
- 230 231 232

235

228

229

3 DATA COLLECTION AND MANAGEMENT

3.1 DATA MODALITIES AND DATA FORMAT

For the data collection, we use the high-quality and multimodal Azure Kinect, equipped with an RGB camera, a depth sensor, and an IR sensor. These sensors all operate at 30 frames per second for three seconds for each video sample, yielding 91 frames per sample. The recorded data is saved in the Matroska ('.mkv') container format, allowing multiple tracks of data formats to be extracted through post-processing. Tracks of modalities used in this dataset are RGB, depth, IR, and 3D skeleton joint sequences.

242 The specification of each data format varies depending on the conventions commonly used in the 243 research community: each RGB frame is captured with a resolution of 1.920×1.080 and is stored 244 in a '.jpeg' format. We record depth and IR sequences with a resolution of 640×576 and store them 245 as 24-bit '.png' files. The skeleton joints of every sample video are stored in their corresponding '.csv' files. Each file contains a 91×193 array, where each row represents a frame, and each column 246 holds information related to that frame. The first column records the timestamp of the frame, and the 247 following 96 columns capture the x, y, and z coordinates of 32 joints of one subject (as illustrated 248 in Figure 2a), measured as the distance (in millimeters) from the joint to the camera. For instance, 249 the first three columns record the x, y, and z values of the first joint. The order of the joints follows 250 the joint index in (Microsoft, 2024b). The last 96 columns record the 32 joints of the other object. 251

Figure 1 presents sample frames from each action category across different modalities, each offering distinct strengths and weaknesses. RGB frames capture rich details such as interactions, locations, 253 and subject features, making them highly informative but lacking in user privacy protection. How-254 ever, since RGB frames compress the 3D world into a 2D plane, they often suffer from issues like occlusion and view variation. In contrast, 3D skeleton joints provide the spatial position of each 256 joint in a 3D space, offering a desirable view-invariant characteristic. Beyond joint positioning, 257 3D skeletons reveal little about the subject's identity, making this modality more privacy-friendly. 258 This privacy feature is particularly valuable in human-centered applications such as smart homes, 259 CPHS, and elder care management. Overall, the comparison of these modalities highlights an in-260 verse relationship between privacy and the amount of information conveyed-the more information 261 a modality provides, the less it typically protects user privacy. Our dataset includes four modalities 262 that span this entire spectrum, encouraging both the exploration of individual modalities and the fusion of multiple modalities to balance privacy preservation with information richness. 263

264 265

266

3.2 DATA ACQUISITION AND SETUP

After selecting the Azure Kinect as the sensing device, a setup for housing the sensor was needed
 to guarantee consistency throughout the experiments. We constructed a sensing module, shown in
 Figure 2b, which positions the Azure Kinect 215 cm above the ground and tilts it forward at a 37° angle. This setup allows for capturing interactions with a full field of view and minimal occlusions.

270

277 278 279

281

283

284 285



Figure 1: Sample data from 12 interactions. The modalities are, from top row to bottom row: RGB, IR, depth, and 3D skeleton joints. The 12 interactions are, from left to right: "waving in," "thumbs up," "waving," "pointing," "showing measurements," "nodding," "drawing circles in the air," "hold-ing one's palms out," "twirling or scratching hair," "laughing," "arm crossing," and "hugging."

287 An important aspect of the experiment is the selection of testbed locations. Rather than attempting to cover all possible environments, we chose three representative locations across a US university 288 campus: an open indoor area, a confined indoor space, and an outdoor area, as shown in Figure 289 3. These locations are selected to provide a variety of backgrounds and support the exploration of 290 the effects of the ambient environment on the sensors. A common limitation of HAR datasets is 291 the lack of diverse backgrounds, which can lead to deep learning models overfitting to background 292 noise. By conducting the experiment in three distinct locations, we aim to improve the generaliz-293 ability of background noise handling. We also acknowledge that a contextualizable dataset should be applicable across a range of environments, such as parks, schools, nursing facilities, and smart 295 homes. Collecting data in varied locations, especially outdoors, allows for the examination of how 296 the ambient environment directly and indirectly affects sensor performance and algorithm accuracy.

297 Since the experiment was conducted at three different locations, it was essential to ensure the data 298 collection process was consistent and repeatable. To achieve this, we designed a testbed setup, 299 shown in Figure 2c, which was used across all three environments. In this setup, volunteers were 300 asked to perform each interaction 40 times within a rectangular area marked on the ground. Af-301 ter each repetition, a beep signaled the participants to rotate either clockwise or counterclockwise 302 before proceeding to the next repetition. This structured process helped minimize labeling ambigu-303 ity by ensuring that subjects performed each action in a predefined sequence, one action at a time. 304 This approach allowed us to confidently associate specific images with their corresponding actions, effectively eliminating the potential for ambiguity or labeling errors. In less controlled settings, 305 where actions may overlap or occur simultaneously, we recommend incorporating contextual tags 306 to enhance label clarity and reduce ambiguity in the data. 307

308 The benefits of this innovative technique are two-fold. First, it enabled us to capture interactions from a wide range of orientations relative to the camera. As shown in Figures 6 and 7 in Section A, some frames capture the side profiles of the subjects, while in others, one subject faces the cam-310 era while the other has their back to it. This diversity in orientations enhances the view-invariance 311 of HAR algorithms. Second, our dataset includes samples with occlusions-a common challenge 312 in HAR tasks. Occlusion occurs when one subject fully or partially obstructs the other within the 313 camera's field of view. By incorporating occlusions, our dataset aims to help HAR algorithms ad-314 dress this issue more effectively. Furthermore, capturing multiple viewpoints using a single camera 315 reduces deployment costs, as achieving similar results would otherwise require multiple sensors. 316 Although the environments for this dataset were curated, similarly to other datasets in Table 1, we 317 intentionally collected data at different times of the day and on various days to capture a wide range 318 of environmental conditions. For example, in the outdoor setting, some participants performed dur-319 ing the early morning or late afternoon when the lighting was dim, while others were assigned to 320 midday sessions under bright sunlight. On several occasions, the sky was fully overcast, providing a 321 low-light environment. These variations in illumination are evident in Figures 6 and 7. In the indoor environments, we enriched the lighting conditions by opening curtains to allow natural light to filter 322 in and by configuring the overhead lights differently for each session. Additionally, the outdoor 323 setting introduced further variability, including breezes and higher winds that caused rustling in the



Figure 2: Using the Azure Kinect SDK (Microsoft, 2023), (a) 32 3D skeleton joints are extracted following this labeling scheme. (b) The sensing module configuration and (c) bird's-eye view of the testbed remain consistent across locations, with subjects confined to a rectangular area.



Figure 3: Testing locations: (a) open indoor area, (b) confined indoor space, and (c) outdoor area.

surrounding trees, creating varying levels of background noise. An active construction site located behind the testing area also contributed to the diversity of conditions, with noticeable changes in the site's layout and equipment placement from one test to another. By incorporating these variations in lighting, noise, and environmental dynamics, DUET more closely mirrors real-world scenarios, enhancing its relevance and robustness for human activity recognition tasks.

3.3 SUBJECTS

A total of 15 male and eight female subjects participated in the experiments. The subjects were randomly paired to perform actions across the three locations. The subjects' ages range from 23 to 42 years old with a mean age of 27 years and standard deviation of 4.01 years. Heights ranged from 165.1 cm to 185.4 cm with a mean height of 172.7 cm and a standard deviation of 8.46 cm. The subjects' weights ranged from 55 kg to 93 kg with a mean weight of 69 kg and a standard deviation of 10.1 kg. To further enhance the diversity and robustness of the dataset, users are encouraged to apply data augmentation techniques to create additional variations and improve the generalizability of machine learning models using this dataset.

368 369 370

339

340

341

351

352

353

354

355

356

357

358 359

360

3.4 DATA ANNOTATION

To simplify the file compilation, we organized the data into a folder structure, as shown in Figure 4. The folder structure comprises four hierarchical layers: (1) modality, (2) location combination, interaction label, and subject, (3) timestamps, and (4) image or '.csv' files. The first layer classifies files by modality, including RGB, depth, IR, and 3D skeleton joints. The next layer uses a six-digit code, *LLIISS*, to categorize the location, interaction label, and subject. In this code, *LL* represents the location: *CM* for the indoor open space, *CC* for the indoor confined space, and *CL* for the outdoor space. *II* refers to the numbered activities (1-12) listed in Table 2, and *SS* indicates the subject pair, ranging from 1–10. Note that the same subject pair number in different locations does not indicate

Labe ID	Dyadic interaction	Label ID	Dyadic interaction
1	Waving in	7	Drawing circles in the air
2	Thumbs up	8	Holding one's palms out
3	Waving	9	Twirling or scratching hair
4	Pointing	10	Laughing
5	Showing measuremen	ts 11	Arm crossing
6	Nodding	12	Hugging
	RGB CCC0101 CL0101 CM0101 Timestamps 0.jpg 90.jpg	epth IR The folder hierarchy is identical to that of RGB	3D joints 3D joints cco101 CL0101 CM0101 mestamps.csv

Table 2: Activity labels and their corresponding interactions.

Figure 4: The data folder structure for our dataset is designed to ensure easy access for users. The RGB, depth, and IR modalities follow the same hierarchical structure, while the 3D skeleton joint folders store all 3D coordinates for a sample video clip in a single '.csv' file.

the same pair; only the pairs CCII02 and CLII07, CCII01 and CMII10, and CCII03 and CMII05 represent the same individuals across locations. As mentioned earlier, each pair was asked to repeat an interaction 40 times, and all repetitions were recorded in a single video. To segment the video temporally, we organized each time window by start and end timestamps. For example, a folder named 40800222_43800211 contains a recording that begins at 40800222 and ends at 43800211 milliseconds after the Azure Kinect is connected. Inside each timestamp folder, the corresponding clip is stored frame by frame, with frames numbered sequentially from 0-90.

3.5 CROSS-LOCATION AND CROSS-SUBJECT EVALUATIONS

One of the key motivations for creating DUET is to encourage the research community to explore HAR in the context of dyadic, contextualizable interactions. To support this, we provide a base-line training and test data split for evaluating algorithm performance. In addition to the standard cross-subject evaluation, we also include a cross-location evaluation. We recognize that applica-tions involving dyadic, contextualizable interactions may take place in a variety of indoor and out-door settings, so the cross-location evaluation helps ensure HAR algorithms are resilient to location variation. For the cross-subject evaluation, we use CCII05, CCII07, CLII01, CLII05, CMII06, and *CMII09* for the test data, and the remainder for the training data. For cross-location evaluation, CCIISS is selected as the test data, while CLIISS and CMIISS are used as the training data.

BENCHMARKING STATE-OF-THE-ART HAR ALGORITHMS

In this section, we evaluate the performance of six open-source, state-of-the-art HAR models with publicly available code, as listed in Table 3. This work intentionally selects algorithms that are open-source to ensure that the implementation used in our benchmarking is consistent with the original benchmarking conducted by the algorithm's developers. This decision prioritizes reproducibility and transparency, both of which are essential for meaningful comparisons. Since DUET provides multiple modalities, the evaluation includes two RGB-based, two depth-based, and two skeleton-based algorithms. The results of the evaluation are presented in Table 3.

First, we analyze the effect of occlusion on the RGB modality, as its accuracy is relatively lower compared to other modalities. As previously mentioned, occlusion is a common challenge in vision-based HAR. To evaluate its impact, we train the two selected algorithms using only unoccluded joints. Note: the parenthesized values are accuracies for unoccluded samples.HAR algorithmModalityCross-location
accuracy (%)Cross-subject
accuracy (%)

m	Modality	accuracy (%)	accuracy (%)	
e et al., 2021)	RGB	9.65 (13.81)	17.85 (21.34)	
t al., 2020)	RGB	8.26 (18.58)	7.79 (34.68)	
(Xiaopeng et al., 2021)	Depth	13.15	18.77	
e et al., 2021)	Depth	14.94	23.18	
(Yang et al., 2020)	3D joints	30.73	36.65	
ı et al., 2021)	3D joints	38.17	41.57	
	m e et al., 2021) t al., 2020) (Xiaopeng et al., 2021) e et al., 2021) (Yang et al., 2020) n et al., 2021)	m Modality e et al., 2021) RGB t al., 2020) RGB (Xiaopeng et al., 2021) Depth e et al., 2021) Depth (Yang et al., 2020) 3D joints a et al., 2021) 3D joints	m Modality accuracy (%) e et al., 2021) RGB 9.65 (13.81) t al., 2020) RGB 8.26 (18.58) (Xiaopeng et al., 2021) Depth 13.15 e et al., 2021) Depth 14.94 (Yang et al., 2020) 3D joints 30.73 a et al., 2021) 3D joints 38.17	m Modality accuracy (%) accuracy (%) e et al., 2021) RGB 9.65 (13.81) 17.85 (21.34) t al., 2020) RGB 8.26 (18.58) 7.79 (34.68) (Xiaopeng et al., 2021) Depth 13.15 18.77 e et al., 2021) Depth 14.94 23.18 (Yang et al., 2020) 3D joints 30.73 36.65 n et al., 2021) 3D joints 38.17 41.57

Table 3: Cross-location and cross-subject accuracy comparison for RGB, depth, and 3D skeleton



Figure 5: Representative confusion matrices for cross-subject evaluation for (a) RGB (DB-LSTM (He et al., 2021)), (b) depth (DB-LSTM (He et al., 2021)), and (c) 3D skeleton joints (DR-GCN (Zhu et al., 2021)). Note: each label's interaction corresponds to the mapping in Table 2.

samples for both cross-location and cross-subject evaluations. The corresponding results (Table 3)
 are comparable to the benchmarking records of other datasets (Liu et al., 2017). The results indicate
 that both algorithms show improved performance when occluded samples are excluded. This experiment highlights not only the significant impact of occlusion on algorithm performance but also the
 critical importance of including occluded samples in datasets for comprehensive evaluation.

Overall, the cross-subject evaluation outperforms the cross-location evaluation across all modalities in the state-of-the-art algorithms, which can be explained by two key factors. First, RGB-based and depth-based algorithms are prone to learning view-dependent motion patterns, often correlating background with motion trajectories during training. In the cross-subject evaluation, the training set includes samples from three locations, whereas in the cross-location evaluation, only two locations are used for training. As a result, these models struggle to generalize to unseen backgrounds during testing, leading to lower accuracy in the cross-location evaluation. Second, the difference in the number of training samples also contributes to the performance gap. In the cross-subject evaluation, 80% of the dataset—approximately 11,520 samples—is used for training, while in the cross-location evaluation, only two-thirds of the dataset is available for training. Performance improves with a larger training sample size. These two phenomena are also present Liu et al. (2019)'s work.

Another observation is the gradual increase in accuracy of the state-of-the-art HAR algorithms tested in our study, progressing from RGB to depth, and then to 3D skeleton joints, which aligns with the expansion of dimensional information. RGB-based algorithms compress input into a 2D plane, lead-ing to lower accuracy since human interactions involve both 3D spatial and temporal coordination (Lee & Kim, 2022). This dimensional compression limits the system's ability to fully capture spatial dynamics. Adding depth information to each pixel in an image, as seen in depth-based algorithms, provides an additional layer of information. The improvement in performance with depth inputs is particularly clear when we compare the same model (i.e., DB-LSTM) using RGB and depth inputs separately. However, despite the increase in accuracy from RGB to depth modalities, both still leave room for improvement. This is due to the fact that both modalities operate in Euclidean space (i.e., images), making them more susceptible to view variations. DUET addresses this issue and improves accuracy by providing more robust data. Additionally, training in Euclidean space can be easily in-

fluenced by trivial features. As shown in Fig. 5a and Fig. 5b, RGB and depth models are confused
by common poses shared across activities—for example, standing is present in nearly all activities.
In contrast, skeleton-based algorithms perform HAR in non-Euclidean space (Peng et al., 2021),
representing human interactions in 3D space relative to the camera, leading to better accuracy.

490 Skeleton-based algorithms outperform other modalities because they capture activities in a 3D space 491 relative to the camera, which is well-suited for the spatial complexity of human interactions. These 492 algorithms can extract underlying motion patterns regardless of the viewpoint. Additionally, 3D 493 skeletons provide a sparse representation of the human body, which helps prevent the network from 494 learning irrelevant features. However, this sparsity can also hinder recognition in certain cases. In 495 our dataset, many dyadic interactions differ only in subtle ways. For example, both the "thumbs up" gesture and "holding one's palms out" (i.e., label ID 2 and 8, respectively) involve arm extension, 496 but the former requires raising the thumb, while the latter involves holding the hand vertically. 497 The simplified skeletal representation may not capture these fine distinctions using current HAR 498 algorithms. This is evident in Figure 5c, which shows these two actions are frequently confused by 499 the algorithm. While the nuances are more apparent in RGB and depth images, from which the 3D 500 skeleton joints are extracted, state-of-the-art skeleton-based algorithms still struggle to detect them. 501

502 503

504

5 DISCUSSION AND CONCLUSION

In this work, we introduce DUET, a contextualizable dataset consisting of 12 dyadic interactions, based on a psychological taxonomy that organizes human interactions into five groups according to their communication functions. This taxonomy advances the field of HAR beyond simple body movement tracking by extracting the embedded semantics in dyadic interactions. Moreover, contextualizing human activities enhances HAR models and paves the way for significant downstream applications, such as autonomous vehicles, urban infrastructure planning, and healthcare.

11,4,400 samples were collected across 12 interactions, resulting in 1,200 samples per activity—the
highest sample-to-class ratio published. The samples span four modalities: RGB, depth, IR, and
3D skeleton joints, each offering unique strengths. The multimodal dataset encourages both the
individual use of each modality to refine models for specific applications and the fusion of modalities
to combine complementary information, maximizing the value provided by each modality.

517 DUET also aims to improve the view and background invariance of HAR models. We introduce 518 a novel data collection procedure to capture human interactions from multiple angles using a sin-519 gle camera, something previously unachievable even with multiple sensors. This innovative setup 520 enhances resilience to variations in viewing angles, reflects real-life scenarios where observations 521 are not restricted to a specific angle, and reduces deployment costs. The choice of testbed locations 522 is carefully considered. Data was collected in three distinct environments: an open indoor area, a 523 confined indoor space, and an outdoor area. This variety not only improves generalizability but also 523 enables applications to assess how ambient environments affect system performance.

524 To establish baseline performance for DUET, we evaluate six HAR algorithms with open-source 525 code to ensure an accurate assessment of their capabilities-two RGB-based, two depth-based, and 526 two skeleton-based algorithms. While some previous work has attempted to recognize dyadic in-527 teractions using monadic algorithms, the performance reveals a persistent gap between recognizing 528 monadic and dyadic activities. In this study, we take a step further by benchmarking six dyadic algo-529 rithms with our dataset. The results highlight (1) the complexity of social interactions that remains 530 underexplored in existing literature, and (2) the vulnerability of HAR algorithms to changes in view and background, which presents new research opportunities for future investigation. 531

532 Future developments from this work can be broadly categorized into two areas: refining HAR al-533 gorithms and enhancing the taxonomy to better capture the embedded semantics of interactions. 534 As shown in Table 3, all modalities require improvement when it comes to contextualizable dyadic 535 interactions. In addition to developing more sophisticated HAR models capable of capturing the 536 nuances in these interactions, another way to enhance performance is through contextualization. 537 We have laid the groundwork for contextualizing human activities by integrating a psychological taxonomy with HAR. The next step is to further define this framework, mapping all interactions to 538 their corresponding levels of embedded meaning, which can benefit downstream applications such as CPHS, autonomous vehicles, smart homes, and healthcare.

540 REFERENCES 541

554

574

575

576 577

581

582

583

584

588

589

- Vida Adeli, Ehsan Adeli, Ian Reid, Juan Carlos Niebles, and Hamid Rezatofighi. Socially and 542 contextually aware human motion and pose forecasting. *IEEE Robotics and Automation Letters*, 543 5(4):6033-6040, 2020. 544
- Chaitanya Ahuja, Shugao Ma, Louis-Philippe Morency, and Yaser Sheikh. To react or not to react: 546 End-to-end visual pose forecasting for personalized avatar during dyadic conversations. In 2019 547 International Conference on Multimodal Interaction, pp. 74-84, 2019. 548
- Anuradha M Annaswamy, Karl H Johansson, and G Pappas. Control for societal-scale challenges: 549 Road map 2030. *IEEE Control Systems Magazine*, 44(3):30–32, 2023. 550
- 551 Authors. DUET Repository, Access on Oct 1 2024. URL https://huggingface.co/ 552 datasets/Anonymous-Uploader1/DUET. 553
- Victoria Bloom, Vasileios Argyriou, and Dimitrios Makris. Hierarchical transfer learning for online recognition of compound actions. Computer Vision and Image Understanding, 144:62–72, 2016. 555
- 556 Huili Chen, Sharifa Alghowinem, Soo Jung Jang, Cynthia Breazeal, and Hae Won Park. Dyadic affect in parent-child multimodal interaction: Introducing the dami-p2c dataset and its preliminary 558 analysis. IEEE Transactions on Affective Computing, 14(4):3345–3361, 2022. 559
- Andrea Chute, Sharon Johnston, and Brandi Pawliuk. 4.2 Types of nonverbal communication. 560 Professional Communication Skills for Health Studies, 2023. 561
- Claudio Coppola, Serhan Cosar, Diego R Faria, and Nicola Bellotto. Social activity recognition 563 on continuous rgb-d video sequences. International Journal of Social Robotics, 12(1):201-215, 2020. 565
- Maral Doctorarastoo, Katherine Flanigan, Mario Bergés, and Christopher McComb. Exploring the 566 potentials and challenges of cyber-physical-social infrastructure systems for achieving human-567 centered objectives. In Proceedings of the 10th ACM International Conference on Systems for 568 Energy-Efficient Buildings, Cities, and Transportation, pp. 385–389, 2023a. 569
- 570 Maral Doctorarastoo, Katherine Flanigan, Mario Bergés, and Christopher McComb. Modeling hu-571 man behavior in cyber-physical-social infrastructure systems. In Proceedings of the 10th ACM 572 International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation, 573 pp. 370-376, 2023b.
 - Paul Ekman and Wallace V Friesen. The repertoire of nonverbal behavior: Categories, origins, usage, and coding. *semiotica*, 1(1):49–98, 1969.
- Neal А Hartman. Nonverbal communication teaching Acnote. 578 cessed Jun 03 2024. URL https://ocw.mit.edu/courses/ on 15-279-management-communication-for-undergraduates-fall-2012/ 579 251fccce2dabe0f6ceafb86218d74c57_MIT15_279F12_nonVerbalComm.pdf. 580
 - Jun-Yan He, Xiao Wu, Zhi-Qi Cheng, Zhaoquan Yuan, and Yu-Gang Jiang. Db-lstm: Denselyconnected bi-directional lstm for human action recognition. Neurocomputing, 444:319–331, 2021.
- Woo-Ri Ko, Minsu Jang, Jaeyeon Lee, and Jaehong Kim. Air-act2act: Human-human interaction 585 dataset for teaching non-verbal social behaviors to robots. The International Journal of Robotics 586 Research, 40(4-5):691-697, 2021.
 - Quan Kong, Ziming Wu, Ziwei Deng, Martin Klinkigt, Bin Tong, and Tomokazu Murakami. Mmact: A large-scale dataset for cross modal human action understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 8658-8667, 2019.
- Gawon Lee and Jihie Kim. Improving human activity recognition for sparse radar point clouds: A 592 graph neural network model with pre-trained 3D human-joint coordinates. Applied Sciences, 12 (4):2168, 2022.

- 594 Cheyu Lin, Maral Doctorarastoo, and Katherine Flanigan. Your actions talk: Automated sociomet-595 ric analysis using kinesics in human activities. In Proceedings of the 11th ACM International 596 Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation, 2024a. 597 Cheyu Lin, John Martins, and Katherine A Flanigan. Read the room: Inferring social context through 598 dyadic interaction recognition in cyber-physical-social infrastructure systems. In Proceedings of the ASCE International Conference on Computing in Civil Engineering, I3CE 2024, 2024b. 600 601 Chunhui Liu, Yueyu Hu, Yanghao Li, Sijie Song, and Jiaying Liu. Pku-mmd: A large scale bench-602 mark for continuous multi-modal human action understanding. arXiv preprint arXiv:1703.07475, 603 2017. 604 Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. NTU 605 RGB+D 120: A large-scale benchmark for 3D human activity understanding. IEEE Transactions 606 on Pattern Analysis and Machine Intelligence, 42(10):2684-2701, 2019. 607 608 Tingting Liu, Zengzhao Chen, Hai Liu, Zhaoli Zhang, and Yingying Chen. Multi-modal hand 609 gesture designing in multi-screen touchable teaching system for human-computer interaction. In 610 Proceedings of the 2nd International Conference on Advances in Image Processing, pp. 198–202, 611 2018. 612 Microsoft. Azure Kinect Body Tracking SDK: Welcome, Accessed on Dec 03 2023. URL 613 https://microsoft.github.io/Azure-Kinect-Body-Tracking/release/ 614 1.1.x/index.html. 615 616 Microsoft. Buy the Azure Kinect Developer Kit, Accessed on Jul 02 2024a. URL https://www. 617 microsoft.com/en-us/d/azure-kinect-dk/8pp5vxmd9nhq. 618 Microsoft. Azure Kinect Body Tracking Tracking Joints, Accessed on Jun 03 2024b. URL https: 619 //learn.microsoft.com/en-us/azure/kinect-dk/body-joints. 620 621 Jacob Levy Moreno. Foundations of sociometry: An introduction. Sociometry, pp. 15–35, 1941. 622 623 Michael Neff, Nicholas Toothman, Robeson Bowmani, Jean E Fox Tree, and Marilyn A Walker. 624 Don't scratch! Self-adaptors reflect emotional stability. In Intelligent Virtual Agents: 10th International Conference, IVA 2011, Reykjavik, Iceland, September 15-17, 2011. Proceedings 11, pp. 625 398-411, 2011. 626 627 Friedrich Niemann, Stefan Lüdtke, Christian Bartelt, and Michael Ten Hompel. Context-aware 628 human activity recognition in industrial processes. Sensors, 22(1):134, 2021. 629 630 Wei Peng, Jingang Shi, Tuomas Varanka, and Guoying Zhao. Rethinking the ST-GCNs for 3d 631 skeleton-based human action recognition. *Neurocomputing*, 454:45–53, 2021. 632 Asanka G Perera, Yee Wei Law, Titilayo T Ogunwa, and Javaan Chahl. A multiviewpoint outdoor 633 dataset for human action recognition. IEEE Transactions on Human-Machine Systems, 50(5): 634 405-413, 2020. 635 636 Michael S Ryoo and Larry Matthies. First-person activity recognition: What are they doing to 637 me? In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 638 2730-2737, 2013. 639 Michael S Ryoo, Chia-Chih Chen, JK Aggarwal, and Amit Roy-Chowdhury. An overview of contest 640 on semantic description of human activities (sdha) 2010. Recognizing Patterns in Signals, Speech, 641 Images and Videos: ICPR 2010 Contests, Istanbul, Turkey, August 23-26, 2010, Contest Reports, 642 pp. 270–285, 2010. 643 644 Navya N Sharan, Alexander Toet, Tina Mioch, Omar Niamut, and Jan BF van Erp. The relative 645 importance of social cues in immersive mediated communication. In Human Interaction, Emerging Technologies and Future Systems V: Proceedings of the 5th International Virtual Conference 646 on Human Interaction and Emerging Technologies, IHIET 2021, August 27-29, 2021 and the 6th 647
 - IHIET: Future Systems (IHIET-FS 2021), October 28-30, 2021, France, pp. 491–498, 2022.

648	Coert Van Gemeren, Ronald Poppe, and Remco C Veltkamp. Spatio-temporal detection of fine-
649	grained dyadic human interactions. In Human Behavior Understanding: 7th International Work-
650	shop, HBU 2016, Amsterdam, The Netherlands, October 16, 2016, Proceedings 7, pp. 116–133,
651	2016.

- Ji Xiaopeng, Zhao Qingsong, Cheng Jun, and Ma Chenfei. Exploiting spatio-temporal representation for 3D human action recognition from depth map sequences. *KnowledgeBased Systems*, 227, 2021.
- Chao-Lung Yang, Aji Setyoko, Hendrik Tampubolon, and Kai-Lung Hua. Pairwise adjacency matrix
 on spatial temporal graph convolution network for skeleton-based two-person interaction recognition. In 2020 IEEE International Conference on Image Processing (ICIP), pp. 2166–2170,
 2020.
- Kiwon Yun, Jean Honorio, Debaleena Chattopadhyay, Tamara L Berg, and Dimitris Samaras. Twoperson interaction detection using body-pose features and multiple instance learning. In 2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, pp. 28–35, 2012.
- Shiwen Zhang, Sheng Guo, Weilin Huang, Matthew R. Scott, and Limin Wang. V4D: 4d convolutional neural networks for video-level representation learning. In *International Conference on Learning Representations*, 2020.
- Liping Zhu, Bohua Wan, Chengyang Li, Gangyi Tian, Yi Hou, and Kun Yuan. Dyadic relational graph convolutional networks for skeleton-based human interaction recognition. *Pattern Recognition*, 115, 2021.

А



Figure 6: Sample data from the first six interactions. The locations presented are, from left to right: the confined indoor space, the open indoor space, and the open outdoor space. The six interac-tions are, from the top to bottom rows: "waving in," "thumbs up," "waving," "pointing," "show-ing measurements," and "nodding." These images demonstrate the variation in lighting conditions, viewpoints, and occlusions.

757			
758			
759			
760			
761			
762			
763			
764			
765			
766			
767		all the second	
768			
769			
770			
771			
772			
773		all the second	/
774			
775			
776			And the second second
777			
778			
779	1 4		
780			
781			
782			
783		HILE PARTICIPATION	
784			A CONTRACTOR
785			
786		24 A AT TO ME	1
787			
788			
789			
790			
791		Christen -	
792		and when it	
793			
794			
795			and a sector
796			
797			
798		Contract Contraction	/
799		the the the second to get	

Figure 7: Sample data from the last six interactions. The locations presented are, from left to right: the confined indoor space, the open indoor space, and the open outdoor space. The six interactions are, from the top to bottom rows: "drawing circles in the air," "holding one's palms out," "twirling or scratching hair," "laughing," "arm crossing," and "hugging." These images demonstrate the variation in lighting conditions, viewpoints, and occlusions.