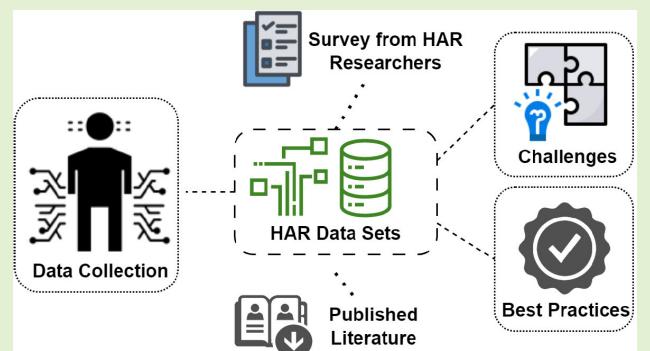


Open Datasets in Human Activity Recognition Research—Issues and Challenges: A Review

Gulzar Alam^{ID}, Ian McChesney, Peter Nicholl^{ID}, and Joseph Rafferty^{ID}

Abstract—Huge amounts of data are generated with the emergence of new sensor technologies. Human activity recognition (HAR) datasets are generated from cameras, such as video or still images, capturing human behavior through sensors such as gyroscopes, Bluetooth, sound sensors, and accelerometers. These generated data sources are collected by the researchers and formed into open datasets. However, these datasets often show issues during dataset construction, sharing, and searching, which could produce further challenges for the reuse of the data by others. The main objective of this research is to explore the current issues and challenges faced by researchers in the HAR domain. A detail literature review was conducted to extract information from the published literature. Similarly, a questionnaire survey was sent to selected researchers having expertise in the HAR domain, who work with open datasets. The main issues and challenges were identified and classified into a hierarchical structure. This research will help HAR researchers to be aware of the current issues and challenges in the field of HAR open datasets. It will help to promote important attributes applicable to many open datasets, such as privacy, anonymity, platform maintenance, datasets' descriptions, metadata, environmental conditions, resources, and training, while constructing and sharing new datasets.

Index Terms—Artificial intelligence (AI), dataset quality, datasets' issues and challenges, human activity recognition (HAR), open dataset lifecycle.



I. INTRODUCTION

WITH the emergence of new computing technologies that are capable of capturing huge amounts of data on human activities [1], there is a need to store the information in a structured, meaningful, and sharable way. Generated data are collected in the form of datasets. Across many areas of business and society, these datasets facilitate analysis, forecasting, and decision-making that can significantly impact quality of life, ranging from the optimization of supply chains to the transformation of healthcare systems. For example, within the realm of business, these datasets function as a dynamic tool for decision-makers, enabling them to develop targeted marketing plans, forecast future market patterns,

reduce time to market, and stimulate economic growth [2], [3], [4]. Within the domain of healthcare, these datasets serve as a catalyst for pioneering research, facilitating the development of tailored therapies, and accelerating advancements in the field of medicine [5], [6]. Other areas in which open datasets are having an impact include urban planning and public services [7], and the world of entertainment and virtual reality [8], [9]. Researchers and practitioners collect datasets for research objectives to understand the totality of an area of interest and develop a basis for making decisions. Similarly, the researchers aim to investigate and tackle the challenges associated with ensuring the findability, accessibility, interoperability, and reusability (FAIRness) of expanding biomedical datasets housed in diverse repositories. Their objective is to improve transparency, reproducibility, and the progress of research by promoting open science practices and the reuse of data [10]. The primary objective of constructing and sharing open datasets, metadata, related dataset publications, and results is to encourage open dataset benchmarking, replication, validation of research approaches, applied data analysis practices, detection of experimental errors, and exploration of novel hypotheses [11], [12], [13].

Manuscript received 20 August 2023; accepted 5 September 2023. Date of publication 4 October 2023; date of current version 14 November 2023. This work was supported in part by the Department for the Economy (DfE), Ulster University, and in part by the School of Computing, Ulster University, Belfast Campus, U.K. The associate editor coordinating the review of this article and approving it for publication was Prof. Chao Tan. (*Corresponding author: Gulzar Alam.*)

The authors are with the School of Computing, Ulster University, BT15 1AP Belfast, U.K. (e-mail: alam-g@ulster.ac.uk; ir.mcchesney@ulster.ac.uk; p.nicholl@ulster.ac.uk; j.rafferty@ulster.ac.uk).

This article has supplementary downloadable material available at <https://doi.org/10.1109/JSEN.2023.3317645>, provided by the authors.

Digital Object Identifier 10.1109/JSEN.2023.3317645

health, social, environmental, and economic improvement. The significance of open datasets is due to the advances in data-driven systems and research across the globe. For example, datasets from healthcare can help to enrich personal care, speed up diagnosis, disease prediction, and planning of treatment, as well as other benefits [14].

HAR systems can typically collect, store, process, and share data in any format over any network describing specific human activities and environments [15]. Researchers and data science practitioners are exploring new approaches to data visualization, building processes to link users with sensors/devices, understanding the contextual importance of sensors/devices, data annotation, and data management during open dataset creation and sharing [16]. Pervasive computing has made significant progress in providing insight into human activity in a range of natural settings through the data generated by intricate sensors and actuators, establishing a connected globe that produces massive amounts of data [17], [18], [19].

HAR has become an emerging research area due to the development of devices and sensors with minimum cost, low power consumption, and real-time streaming of data combined with sophisticated processing capabilities through technologies such as artificial intelligence (AI), machine learning (ML), and the Internet of Things (IoT) [20]. HAR is the detection of movements/activities (walking, running, standing, sitting, drinking, talking, sleeping, and so on) of an individual based on sensor/device data [21]. HAR data collection can involve data from a camera such as video or images comprising human sentiments and from sensors such as gyroscopes, Bluetooth, sound sensors, and accelerometers. The data received from devices/sensors are connected via a network or attached to the human body [22]. HAR can play a significant role in human daily life activities by understanding people's behaviors from collected sensor data. Researchers apply ML and other statistical techniques for extracting important features, activities, and patterns from video- and sensor-based HAR datasets [23], [24].

However, there are issues and challenges when using HAR open datasets such as recognizing similar and dissimilar actions, emotion recognition from video and images, detecting background noise in activity recording, missing values in a dataset, and dataset annotation. The number of issues and challenges grows due to the numerous and complex activities being recorded, sensor placement and orientation, camera movement, and the increasing number of activity types being studied such as person-only, person and object, and person-to-person activities. Further issues and challenges related to HAR open datasets are the annotation/labeling ground truth, activity detection among multiple participants, sensor orientation and heterogeneity, different data formats, rogue sensor values, missing values, imbalanced data, and background environment [25].

The objective of this article is to explore the current issues and challenges faced by researchers in the HAR domain when using open datasets. A dual approach has been taken involving a comprehensive literature review to extract themes from the published literature (January 2016–February 2023)

and a questionnaire survey of researchers having expertise in the use of open datasets in the HAR domain. The main contributions arising from this research work are given as follows:

- 1) a conceptual framework of the open dataset lifecycle;
- 2) greater exploration of the issues and challenges related to HAR open datasets through a questionnaire survey;
- 3) identification of HAR datasets' issues and challenges through a comprehensive literature review;
- 4) derivation of useful insight from the analysis and comparisons of both survey and literature review results;
- 5) evidence-based classification of issues and challenges in the use of open datasets in HAR.

In contrast to previously published literature studies, such as [23], [26], [27], [28], [29], [30], [31], [32], [33], [34], and [35], this work specifically examines the issues and challenges faced by researchers in the field of HAR with regards to open datasets. Similarly, a conceptual framework for the lifecycle of open datasets was developed, along with a classification of issues and challenges based on findings derived from surveys and published literature. In addition to referring to the existing literature, we conducted a focused questionnaire survey with experts in the HAR domain, thereby obtaining timely and firsthand information. Through the process of comparing literature findings with survey responses, we can gain a deeper understanding of the practical challenges that researchers endure. Our research presents an opportunity for significant progress in this domain. This has the potential to bring about groundbreaking developments in the utilization and interpretation of open datasets within HAR research, fundamentally modifying our perspectives and methodologies.

The remainder of this article is structured as follows. Section II describes the research background, and Section III presents the research methodology. Section IV discusses the results from both the conducted literature review and the survey, and then, Section V presents an analysis and comparison of the survey and literature review results. Section VI illustrates the classification of open issues and challenges related to HAR open datasets, and Section VII discusses the derived HAR datasets' future scope from the research study. Finally, Section VIII describes the conclusion and future work of the proposed research.

A. Motivation and Scope

HAR researchers are facing issues and challenges when interacting with open datasets. The main issues and challenges are discussed comprehensively in Section IV. By addressing these issues and challenges, the aim is to improve dataset structure and quality, and access and facilitate analysis, forecasting, and decision-making that can then be more readily and reliably shared within the research community. Long-term, good-quality datasets can benefit personal healthcare, rehabilitation, early disease detection, and globalization of datasets.

The scope of this article is related to open datasets in the HAR domain, covering both video- and wearable-based recognition systems. Similarly, the questionnaire survey was conducted by involving HAR researchers, and a

comprehensive literature review was performed considering HAR studies.

II. RESEARCH BACKGROUND

Open datasets consist of those data that are freely available and accessible for use and sharing. Various public/private sector organizations and individuals use these open datasets for health, social, environmental, and economic purposes. The significance of datasets is due to the increase in data-driven systems across the globe [36], [37]. Hence, open data can allow a deeper understanding of worldwide trends and common problems. It can play a role, in combination with ML and statistics, to solve problems in the domains of healthcare, engineering, and science. Similarly, it can encourage international collaboration and improve transparency [38]. They can modernize the development processes and systems built by governments, private organizations, and societies. A useful characterization of open datasets is given in the following as described by the Organization for Economic Co-operation and Development (OECD) [39], [40].

A. Redistribution and Reusability

Open datasets should be reusable and distributed for commercial and noncommercial use. To ensure reusability, they must be properly licensed, well-structured, and in a machine-readable format [25]. The important concern of redistribution and reusability intersects with the fundamental principle of data ownership in open datasets. One noteworthy fact is that data owners frequently restrict redistribution, resulting in a complex system in which accessibility and sharing coexist with ownership rights. The dynamic interaction of these factors influences the extent to which open datasets can be freely shared and used for a variety of purposes [41]. The complexities of data ownership delicately thread themselves into the fabric of data transmission, influencing the possibilities for expanded reuse and collaboration. Recognizing the traditional constraints surrounding redistribution, it becomes critical to strike a delicate balance between facilitating open access and respecting ownership privileges, fostering a discourse that champions equitable data utilization while recognizing the critical role of data curators [42].

B. Open Access and Availability

The datasets should have transparent open access arrangements and should be available for download, reuse, and sharing [43]. Various platforms such as the World Bank Open Data, the WHO Open Data Repository, the European Union Open Data Portal, the UCI ML Repository, the U.S. Census Bureau, Data.gov, DBpedia, the UNICEF dataset, and Kaggle [44] provide open access to datasets for research use. They must also be available with full documentation, user guidelines, and a modifiable format.

C. General Participation

Everybody should be able to utilize open datasets, reuse them, and further redistribute them regardless of the domain.

Every domain should use open data for its intended purpose with no discrimination and opposition to any person, team, or field of work [45]. A good example is how noncommercial use of datasets prevents commercial use of it and keeps restrictions on datasets used only for certain purposes that are not permitted for general use.

D. Interoperability

This is the ability to combine diverse datasets from different systems and organizations. This permits various components to work together and merge several datasets to develop larger and more complex systems for solving complicated problems [40].

HAR data collection is mainly conducted based on two methods, namely, video data collection and sensor data collection [46]. It has been successfully applied to individual behavior analysis [47], movement recognition [48], gait analysis [49], and video surveillance [50]. Researchers apply ML and other statistical techniques for extracting important features and activities from video- and sensor-based HAR datasets [24].

The utilization of HAR datasets specifically designed for individuals with disabilities has great importance. The utilization of these datasets holds the potential to greatly enhance the quality of life for those with disabilities, making them more independent and self-sufficient. The use of these datasets has facilitated the development of HAR systems by researchers, which, in turn, has the potential to assist those with disabilities in activities of daily living (ADLs). Developing these datasets has the capacity to enhance the quality of life for those with disabilities by fostering more independence and self-sufficiency. Kim et al. [51] collected the MyMove dataset spanning a duration of seven days, a group of 13 elderly individuals participated by gathering activity labels and wristwatch sensor data. Leving et al. [52] collected the Activ8 dataset by involving 16 non-disabled individuals who performed 16 various standardized 60s ADLs. Additional datasets have been created to cater to the needs of those with disabilities, encompassing KFall [53], CAUCAFall [54], and SisFall [55].

Sing et al. [11] surveyed the use of different datasets in HAR for the research community. They discussed dataset benchmarking, comparisons, and improvement of building datasets. Relevant features are explored, such as participant classes and data sources for dataset construction, focus area of the datasets, modality, annotations, and evaluation of HAR datasets. These are outlined in the following.

- 1) *Classes* such as person-only, person and object, and person-to-person [56].
- 2) *Focus* is the type of activity under observation such as sports, gaming, surveillance, and healthcare [57].
- 3) *Modality* is related to temporal- or spatial-based data [58], [59].
- 4) *Data source* concerns the origin of the data such as sensor placement, whether the data are recorded or scripted, data labeling procedure, and, finally, whether it is a generated dataset or crowdsourced [60], [61].



Fig. 1. Word cloud of HAR datasets explored in the literature review.

- 5) *Annotation* is concerned with the correct annotation of the dataset and confirmation of labeling actions/activities [62].
- 6) *Evaluation* is checking the accuracy of the dataset, such as identifying imbalanced data and missing values [63].

Established dataset repositories having HAR datasets are the UCI ML Repository [64], the Harvard Dataverse [65], the Dataset Search from Google [66], IEEE Dataport [67], Zonedo [68], and Figshare [69]. Some of the popular HAR datasets available include Hollywood [70], Action Similarity LAbelN (ASLAN) [71], YouCook [72], ActivityNet [56], CASAS [73], and Opportunity [74]. Datasets in HAR are illustrated in Fig. 1 in the form of a word cloud generated through an online tool¹ from the published SLR research studies. The largest fonts of words show the highest frequency of a dataset to the lowest fonts with low words frequency. Some datasets that are commonly used by the researchers are EmotiW [75], CASAS [73], WISDM [76], ARUBA [73], PAMAP [77], SHL [78], Opportunity [74], and WHARF [79].

However, HAR open datasets present researchers with issues and challenges, such as recognizing similar and dissimilar actions, emotion recognition from video and images, background noise in activity recording, missing values in a dataset, and dataset annotation, as shown in Section IV. These issues and challenges are growing due to the numerous and complex activities involved, sensor placement and orientation, camera movement, and the increasing number of classes and activities being studied.

There has been a rise in the use of network analysis as a method for exploring the interconnections between various concepts and phrases in the academic literature [80]. Unit analysis as a keyword and type analysis as co-occurrence allow researchers to find important phrases that are regularly linked with the unit of study, allowing them to begin to construct a more comprehensive knowledge of the issue at hand [81]. The constructed network analysis, as shown in Fig. 2, consists

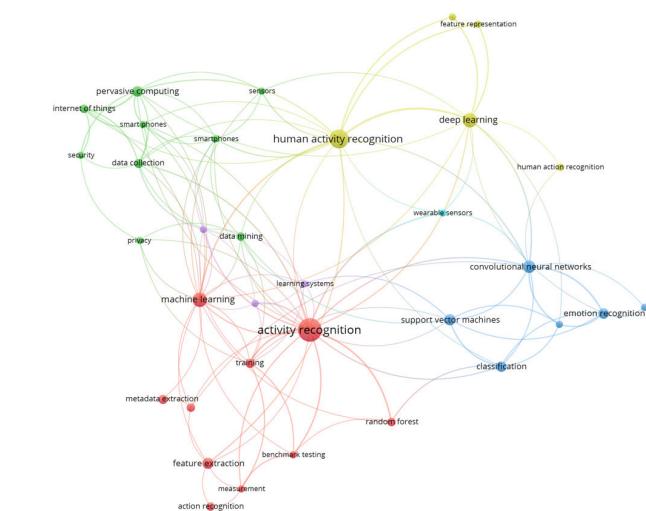


Fig. 2. Network analysis of the conducted literature review studies.

of the relationship between keyword and co-occurrence, developed via VOSviewer.² The colors represent the keyword categories.

As shown in Fig. 2, the terms ML, convolutional neural network (CNN), classification, support vector machine (SVM), feature representation, feature extraction, random forest (RF), learning systems, and training were received in the graph and are related to AI and ML. In the context of ML and AI, each of these concepts is fundamental to the study of unit analysis. ML is a subfield of computer science that uses statistical models and algorithms to enable computer systems to learn from data without being explicitly programmed [82]. The graph depicts the most prevalent ML methods, including CNNs, SVMs, and RFs.

CNNs are a form of artificial neural network (ANN) that is frequently used for image categorization and identification applications [83]. SVMs are a form of ML technique utilized often for classification problems [84]. They are particularly efficient when the data are high-dimensional and have the ability to effectively handle datasets with a large number of features being a notable benefit of SVMs, making them very capable in complex scenarios [85]. RFs, on the other hand, are an ensemble learning method for classification, regression, and other tasks that function by constructing a large number of decision trees (DTs) at training time and outputting the class that is the mode of the classes (classification) or the mean prediction (regression) of the individual trees [86]. Both feature representation and feature extraction are essential topics in ML. Feature representation is the manner in which data are represented, whereas feature extraction is the process of extracting features from data. In ML, learning systems and training are both fundamental ideas. Learning systems are the techniques and models used to train an ML system, whereas training is the act of adjusting the model's parameters to the input [87].

Similarly, human activity recognition (HAR), data collection, pervasive computing, activity action recognition, emotion recognition, benchmark testing, action recognition,

¹<https://worditout.com/word-cloud/create>

²<https://www.vosviewer.com/>

measurement, metadata extraction, data mining, privacy, sensors, smartphones, security, the IoT, and wearable sensors are all terms gathered from the network analysis of the graph in Fig. 2. In the context of recognizing human activities, unit analysis is important to each of these words.

HAR is a field of activity recognition that employs sensor data and other information sources to recognize and comprehend human actions. This is possible in several settings, including healthcare and fitness. Activity recognition is a similar term that applies to the identification and categorization of various activity types [88]. This may be accomplished using a variety of methods, such as ML algorithms and statistical models. Pervasive computing refers to the use of technology to build highly interconnected and interactive environments. This may be utilized in the context of human activity identification to produce precise and dependable systems [89]. Emotion recognition uses sensors and other data sources to recognize and comprehend human emotions. This is applicable in several settings, including healthcare and personal fitness [90].

Data collection is a vital stage in the process of recognizing human activities. This involves collecting data from a number of sources, such as sensors, cameras, and other devices, and then analyzing that data to derive actionable insights [91]. In the realm of HAR, benchmark testing and action recognition are both essential. Benchmark testing is the process of evaluating the accuracy and dependability of various algorithms and models, whereas action recognition is the identification of particular actions or gestures [92]. Important phases in the process of evaluating and interpreting sensor data include measurement, extraction of information, and data mining. The process of extracting important features from the data is followed by the utilization of ML algorithms to identify patterns and trends [28]. In the context of HAR, privacy and security are crucial issues, especially when data are collected from a growing number of sources and participants. Wearable sensors, smartphones, and the IoT are all key data sources that may be used in this context; nevertheless, it is crucial that these data are gathered and utilized in an ethical manner [93].

E. Open Dataset Lifecycle

An open dataset lifecycle is presented in the form of a conceptual framework consisting of four main phases, namely, construction, sharing, finding, and using a dataset, as shown in Fig. 2. We believe that this is the first open dataset lifecycle framework in computing and HAR domain [25], and this framework was influenced from the research work of Stampers et al. [94]. We use this lifecycle to organize the conducted questionnaire survey. As shown in Fig. 3, our proposed life cycle consists of the following four phases.

1) Construction: In this phase, researchers design a protocol for data collection according to the main goal of their research. This phase involves the detailed documentation and description of the dataset. The description includes specific elements such as dataset title, creation date, subject description, sensor and device information, actions and activities, and supported data format. Data are collected in various formats such as video, audio, images, and text. During the dataset construction phase,

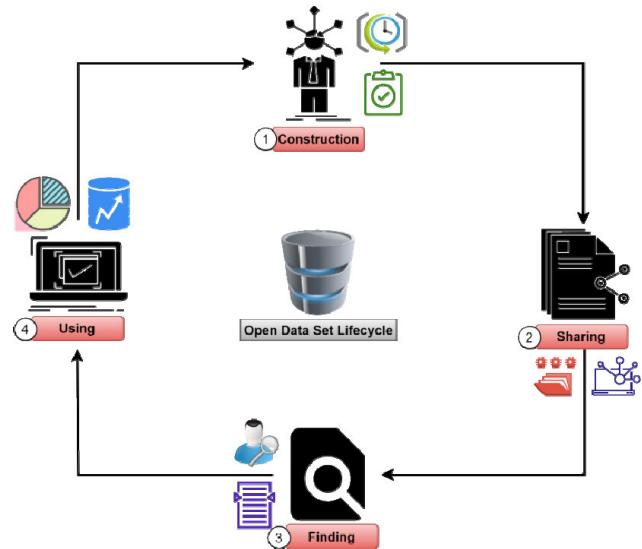


Fig. 3. Open datasets life cycle.

a significant emphasis is placed on carefully establishing the fundamental aspects of the datasets' integrity and usefulness. The annotation protocol assumes an important role in this context since it establishes explicit criteria and procedures for the systematic labeling and categorization of data points. Concurrently, it is crucial to prioritize the pursuit of validity, ensuring that every annotation accurately aligns with the facts of ground truth. The datasets' dependability is supported by the implementation of rigorous validation methods and quality assurance mechanisms, which effectively mitigates the risk of mistakes or inconsistencies. As the dataset undergoes transformation via the use of these techniques, its creation becomes a careful undertaking characterized by precision and diligence. This process lays the foundation for the following stages of data sharing, finding, and using.

2) Sharing: This phase is related to sharing the dataset and making it available to the research community. The dataset owner grants access to the dataset for research purposes. The datasets are stored in a centralized or local repository depending on the data storage and access arrangements as agreed in the research proposal and the distribution license determined for the dataset.

3) Finding: This phase is where dataset users (researchers) search and find a dataset for achieving their research objective. Researchers use different search terms and keywords to retrieve a dataset from a central repository or to retrieve it from a collaborative research group via a local repository.

4) Using: In this phase, users/researchers usually apply various ML techniques, statistical analysis tools, and frameworks to visualize the data and produce meaningful results for solving a problem. Users typically perform dataset preprocessing to make the dataset informative and improve dataset quality before use.

III. RESEARCH METHODOLOGY

A. Literature Review

The extracted information was collected from various well-known digital libraries, such as ACM, IEEE, Science

Direct, Springer, and Google Scholar. The following research questions were addressed in the conducted literature review.

RQ1: What are the unresolved challenges in the open dataset lifecycle in HAR research?

RQ2: What are the best practices in the open dataset lifecycle in HAR research?

RQ3: What ML techniques have been used for the analysis of open datasets?

RQ1 (What Are the Unresolved Challenges in the Open Datasets Lifecycle in HAR Research?): This concerns issues and challenges faced by the researcher during dataset construction, sharing, finding, and using, for example, challenges related to the dataset itself such as metadata, data representation, contents and structure, annotations, and documentation, and also challenges related to the datasets context such as reusability, privacy, societal concern, and usage policies [95].

RQ2 (What Are the Best Practices in the Open Datasets Lifecycle in HAR Research?): This might be practices such as what criteria and quality measures were used by the researchers while finding and using a dataset. Search best practices relating to Internet search engines, social media platforms, and other publicly available online websites and resources for finding and uploading datasets. The best practices are used for dataset documentation and annotations for real-world objects throughout the development process of data collection. Also, what criteria and practices are used by HAR researchers for dataset sharing in a dataset repository?

RQ3 (What ML Techniques Have Been Used for the Analysis of Open Datasets?): The importance of datasets for ML cannot be overlooked, and ML largely depends on the datasets and training algorithms/techniques to make decisions. This research question explores the applied algorithms and techniques on different datasets in the HAR domain. ML often relies on huge-size datasets at the center of model development and evaluation, and it depends heavily on data sources.

B. Questionnaire Survey

A questionnaire survey was chosen to explore the current issues and challenges faced by HAR researchers. This would enable a direct response from researchers on their current approach to working with open datasets. A comprehensive questionnaire survey was conducted consisting of three main phases of planning, performing, and reporting, as shown in Fig. 4.

1) Planning: The questionnaire consisted of 40 questions as shown in Appendix A; these included both open and closed questions.

Open Questions: An opportunity was given to the participants to express their expertise and opinions related to a particular question.

Closed Questions: The respondents were restricted to choosing an answer from the given options in the form of yes/no, multiple choice, ranking, and rating scale questions.

The overall aim of the survey was to elicit open issues and challenges, best practices and processes used by the researchers, and their views on the quality of existing datasets in the HAR domain.

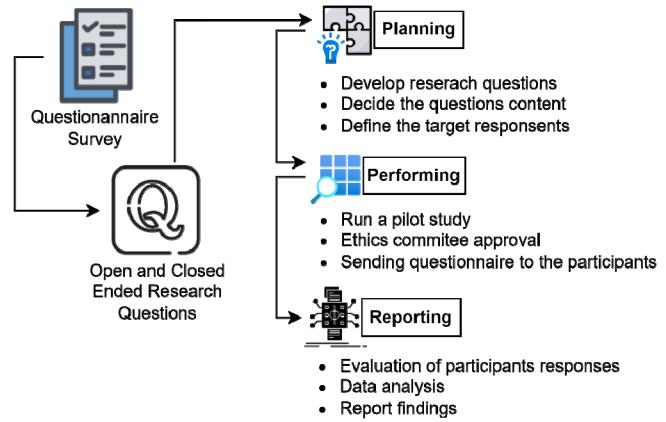


Fig. 4. Protocol for the conducted questionnaire survey.

After a successful pilot study involving four HAR researchers, the survey was sent to 98 participants who had some research experience in the HAR domain. The participants were selected from academia and research institutes, and included researchers, faculty, Ph.D. students, and practitioners from industry.

2) Performing: Ethical approval for the survey was granted by the Faculty Research Ethics Committee at Ulster University to ensure the integrity, trustworthiness, and authenticity of the research outcomes. The questionnaire survey was administered using the Jisc Online Survey Tool,³ and the link was sent to all participants through email.

3) Reporting: In this phase, the participants' responses were evaluated. The information was extracted from the open-ended questions through structural and thematic coding analysis [97], [98]. The findings of the survey are reported in Section IV-B.

IV. RESULT AND DISCUSSION

This research aims to provide a comprehensive analysis of available datasets in the field of HAR by integrating a literature review with a conducted questionnaire survey. This study conducts a comprehensive analysis of the existing literature to explore the complexities associated with the issues and challenges prevalent in the field of HAR.

A. From Literature Review

RQ1 (What Are the Unresolved Challenges in the Open Datasets Lifecycle in HAR Research?):

Table I shows the issues and challenges associated with HAR open datasets that were identified through the conducted literature review. The groups and subgroups were created from the authors' brainstorming and from the published work of Singh et al. [11]. The authors identified challenges related to red-green-blue (RGB) and RGB-depth (RGB-D) datasets, and then, they were divided into five groups based on application domain, environmental condition, occlusion, viewpoints variations, and similarity of actions.

However, these challenges are only related to human actions and external factors related to the datasets. The conducted literature review considered both external factors, such

³<https://www.onlinesurveys.ac.uk/>

TABLE I
OPEN ISSUES AND CHALLENGES IN HAR OPEN DATASETS

Group	Subgroup	In percent	Reference
<i>Activities/Actions Recognition</i>	<ul style="list-style-type: none"> • Similar actions • Dissimilar actions • Complex activities • Emotion recognition • Monitoring activities • Multiple occupants 	15.8%	[102], [103], [104], [105], [106], [107], [108], [109], [110], [111], [112], [113], [114], [115], [116], [117], [118], [119], [120], [121], [122], [123], [124], [125]
<i>Annotations</i>	<ul style="list-style-type: none"> • Data labelling problems 	14.5%	[126], [102], [104], [127], [128], [109], [111], [112], [113], [129], [114], [115], [130], [131], [116], [132], [133], [119], [134], [121], [135], [124]
<i>Noise/Imbalanced data set</i>	<ul style="list-style-type: none"> • Data interpolation • Temporal relationship (Video frames) • Irrelevant information 	11.8%	[126], [136], [137], [127], [105], [106], [138], [114], [114], [139], [131], [132], [118], [119], [140], [141], [142], [122]
<i>Background Condition</i>	<ul style="list-style-type: none"> • Other objects • Light condition • Image quality • Movement/posture pattern • Shadow • Dynamic background 	9.2%	[102], [103], [126], [136], [106], [108], [111], [138], [112], [139], [140], [143], [144], [125]
<i>Resource/Training</i>	<ul style="list-style-type: none"> • Subject/user training • Battery condition of device • Internet connection • Device variety • Sensor orientation • Camera variation 	7.9%	[136], [128], [108], [145], [146], [147], [138], [115], [148], [119], [124], [125]
<i>Missing values</i>	<ul style="list-style-type: none"> • Lack of efficient data representation 	7.2%	[149], [102], [103], [104], [105], [128], [145], [146], [147], [133], [135]
<i>Privacy</i>	<ul style="list-style-type: none"> • Data sensitivity • Ethics concern 	5.9%	[107], [150], [111], [113], [139], [116], [151], [118], [124]
<i>Feature Selection</i>	<ul style="list-style-type: none"> • Auto extraction of important features • Selection of active features • Removal of irrelevant features • Appropriate preprocessing techniques 	5.9%	[137], [146], [150], [111], [112], [115], [117], [119], [144]
<i>Heterogeneity (Device/Data/subject)</i>	<ul style="list-style-type: none"> • Various data format • Various devices installed • Subject heterogeneity 	5.9%	[137], [107], [146], [147], [112], [152], [116], [133], [141]
<i>Data set size</i>	<ul style="list-style-type: none"> • Less sample of data 	4.6%	[150], [112], [152], [117], [143], [121], [144]
<i>Other Issues and Challenges</i>	<ul style="list-style-type: none"> • Lack of sharing data • No data collection guidelines • Subject observation and monitoring • Subject annoying to wear device • Support new participants in data collection • Data set benchmarking 	11.2%	[131], [148], [119], [120], [143], [121], [125]

as activities and actions, background conditions, resources training, device and subject heterogeneity, and dataset sharing, and the internal factors of a dataset such as annotation,

noisy and imbalanced datasets, missing values, data privacy, features selection, and dataset size. The brief descriptions of the identified challenges are given in the following.

Activities/Action Recognition: During activity recognition, it can be challenging to differentiate similar and dissimilar activities, such as walking and running, and stairs up and down movement. Complex activities consist of more than two activities and actions, for example, exercise activities such as jumping, extending legs, and bending down. Complex activity recognition can be achieved by incorporating a model that addresses, for example, proper posture monitoring. 15.8% of papers in the conducted literature review reported activity recognition to be a challenge.

Annotations: Annotation is the labeling of data (to label an activity/action) in different formats, such as audio, video, text, and images. Annotated datasets are important for supervised ML for pattern interpretation, and accurate outcomes named entity recognition (NER) and sentiment analysis require annotated training data to detect emotions and opinions [100].

Noise/Imbalanced Datasets: Noisy and imbalanced datasets are created when participants tend to make mistakes, and sensors/devices record incorrect data or insert incorrect values to attributes while collecting data. Also, the collection of additional and irrelevant information can create noise in the data, which can impact the prediction and recognition of activities, and affect the overall quality of a dataset [101]; 11.8% of the papers reported such issues as challenging.

Background Condition: Background condition refers to the different objects of the observed environment, which can hinder activity recognition—objects such as trees, rain, water, and waves [99]. Similarly, moving objects in the background can also impact activity recognition. The datasets that are collected from social media and YouTube also contain challenges arising from background conditions and moving objects. Light condition refers to brightness and darkness. Image quality is also affected by camera movement. Object shadows are also an issue when detecting human activity from images. 9.2% of papers in the literature review reported such background conditions to be a challenge when working with an open HAR dataset.

Resource/Training: Training is where the recruited participants/users/subjects for the experiments are instructed on how to use devices correctly according to the defined protocol guidelines. Resource issues and challenges include the battery life of devices and reliable Internet connection for all experimental wearable devices. Also, the installation of different devices producing data in various formats can be a resource-intensive activity. Finally, datasets created from video recordings are large and require enough GPUs for model training. Human complexity is an issue when constructing an HAR dataset because humans are highly variable and unpredictable in their movements and actions. This can make it difficult to accurately capture and classify a wide range of different activities. In addition, factors such as lighting, camera angle, and background can also affect the quality of the data and make it difficult to generalize the dataset to different environments. In the literature review, 7.9% of papers reported such challenges.

Missing Values: Missing values are situations where data have not been recorded in relation to a significant or meaningful event, which should have been observed, for example, missing an activity, incomplete information, or feature due to a network problem, where participants forget to record an activity or the low-energy characteristics of a wearable device. This is challenging when recognizing elderly people's activities because it is difficult to predict the activities associated with missing data, such as falling down or any other serious activity. 7.2% of papers reported this to be a challenge in HAR dataset use.

Privacy: Secure storing and proper distribution of human-obtained data are crucial parts of data ethics [153]. Each dataset involves people who both collect data and provide data. It is important to consider that, in both cases, we have human beings [154]. International and cultural contexts should be respected where privacy concerns may have less governance in some areas, and the possibility of data manipulation is a real threat to the protection of data participants. This can be a challenge and a risk in very sensitive contexts such as personal information related to finance, medicine, and biometrics [155]. Using sensors and devices for elderly assistance and patient monitoring presents privacy issues. Installation of devices in the home for tracking could be considered a violation of intimacy and privacy [156], [157]; 5.9% of papers reported data handling and privacy to be a challenge.

Feature Selection: Feature selection is the process of choosing the most relevant features that are essential to perform machine and deep learning, and eliminate unnecessary features. This process is useful because it will reduce noise in a dataset. However, it is difficult to select the optimal features due to the high correlation among features and to select the most appropriate features that contribute to prediction outcomes.

The other issue related to feature selection is how to choose the appropriate statistical and ML approach and technique; 5.9% of papers reported this process as a challenge.

Heterogeneity (Devices/Data/Subjects): The use of sensors, RFID tags, and wearable devices for collecting activity data in the HAR research area is growing. The variations in sensors and the frequencies at which data have been collected create data in different formats. Subjects' heterogeneity due to different lifestyles and cultures, and different mechanisms for their recruitment is also a potential challenge in the HAR domain; 5.9% of papers reported such issues.

Dataset Size: When the size of activity data is small, the HAR model developed through training leads to anomalies and random noise. Therefore, it negatively affects the model's generalization capability. Furthermore, limited data mean that HAR models are unable to model new data and generalize to unseen (new) data leading to low model performance.

The activities'/actions' recognition and annotations are the highest mentioned groups of the identified issues and challenges. Similarly, noise/imbalanced datasets were mentioned by the researchers up to 11.8% and other issues and challenges to 11.2%. The remaining issues and challenges

are 9.2% and below. The dataset issues and challenges as described above are important to address and create datasets that are meaningful for the HAR research community.

Other Issues and Challenges: At a fundamental level, the lack of data sharing among HAR researchers is itself an issue. This may be due to constraints on dataset sharing from within the researcher's organization, the confidential nature of data, or a dataset having low quality such that sharing and reuse are not possible. This might arise when open datasets are collected without proper data collection guidelines and protocols. Poor quality data might arise because it is difficult to fully control the participants during data collection. Sometimes, it might be annoying for the subject to wear a device for recording activity, or they might have a privacy or other social concern. Similarly, introducing a new participant into an experiment during data collection or an existing participant leaving is a potential issue for consistent data collection. A benchmark dataset is a set of data that are widely used in a particular field to evaluate the performance of different models or algorithms. These datasets are often used as a standard for comparing the performance of different approaches and for measuring progress in the field. Benchmark datasets are taken from various sources, each with a distinct composition of copyright owners and agreements for their usage in training and evaluation in ML models [158]. Open dataset benchmarking is good for validating datasets and evaluation of a model or an experiment against internal or external standards related to HAR datasets.

RQ2 (What Best Practices Are Used by Researchers When Constructing, Sharing, Finding, and Using Datasets in HAR Research?):

Researchers developed their own practices according to the needs of their experiments. However, most of the researchers used dataset preprocessing techniques in the construction and using phase of a dataset. They extracted the prominent features and improved accuracy by using different ML classifiers. Some researchers replaced missing values from a dataset with a mean value.

Researchers also observed that not all extracted features are good for classification because some of them have a negative impact on classification performance. To reduce overfitting in a dataset and to make a dataset more generalizable, data augmentation was applied by some researchers [149], [102], [104], [137], [127].

Researchers are working to improve metrics for dataset benchmarking by including the full empirical evaluations, including negative results during evaluation and full sharing of additional experimental details [159]. For data labeling, researchers used data transformation and segmentation algorithms to distribute the data into different length windows, and then, human experts assign the labels by marking a time stamp for each activity [1], [160]. To overcome the problems of data scarcity and privacy concerns, researchers are developing synthetic datasets, which makes it easy to share data and increase model robustness. Synthetic datasets can be generated by learning the statistical properties of actual datasets [161], [162]. Detailed best practices are shown in Fig. 24.

RQ3 (What ML Techniques Have Been Used for the Analysis of Open Datasets?): In addition to investigating issues and

challenges, the conducted literature review identified the range of ML techniques that have been used in HAR analysis as it has been applied to open datasets. Currently, researchers are using ML techniques to identify and interpret activities in HAR in various domains, such as sports, healthcare, and falls of elderly people. The limitation of the conventional approaches such as basis transform coding, statistics of raw signals, and symbolic representation was the shallow learning that needed feature engineering from data, and it was largely dependent on human knowledge of the specific domain [163]. The conventional approaches of heuristic nature and human understanding make it difficult to capture complex actions having micro activities. Higher level activities include more semantic and contextual information, making it harder to discern their hierarchical structure. Existing approaches often overlook signal correlation, which limits their ability to provide satisfactory outcomes. Therefore, the widespread adoption of ML can be attributed to its ability to effectively learn complex patterns and behaviours through the utilization of CNNs [164].

For large-scale datasets, deeply learned feature methods are suitable for the HAR domain because deep learned features, also known as deep features, are representations of data that are learned by a deep neural network (DNN). Deep learning-based solutions consist of feature extraction and classification in video-based datasets. In addition to the development of high computational power and a growing volume of video-based datasets, a deep learning-based solution is useful for real-life applications [99]. Feature extraction is an important aspect of HAR to identify the most relevant and significant features from the data to reduce errors in classification and computational complexity. The efficient performance of HAR activity recognition depends on suitable feature representation. Achieving maximum performance depends on the selection of suitable techniques for feature extraction [119]. Researchers already work and apply ML to automatically extract features from raw data generated through sensors. ML techniques and classifiers that are used in HAR for automatic feature extraction, preprocessing, and classification of datasets are SVM, k-nearest neighbor (KNN), DT, CNN, long short-term memory (LSTM), and RF [24], [165].

Researchers are applying new trends of ML such as transfer learning, which is the ability to transfer knowledge from one model to another to train it with a minimum amount of data and to reduce computational complexity and effort [166]. Similarly, the other emerging ML method is active learning with the objective of reducing learning complexity and computational cost. It seeks to choose the relevant information from unlabeled data and ask the annotator for labeling information. The main advantages of active ML in HAR are to reduce annotation efforts and increase forecasting accuracy [167]. As shown in Fig. 5, ML techniques were used in the literature review papers for prediction, classification, feature extraction, and dataset labeling purposes.

B. Survey Results

The purpose of the survey of HAR researchers was to elicit their views on the use of open datasets. The open dataset

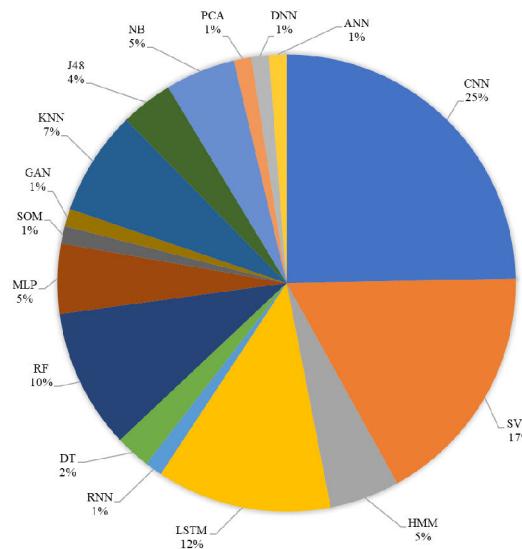


Fig. 5. Weightage ML techniques explored in SLR applied on open datasets CNN, SVM, hidden Markov model (HMM), LSTM, recurrent neural network (RNN), DT, RF, multilayer perceptron (MLP), self-organizing map (SOM), generative adversarial network (GAN), KNNs, naive Bayes (NB), principal component analysis (PCA), DNN, and ANN.

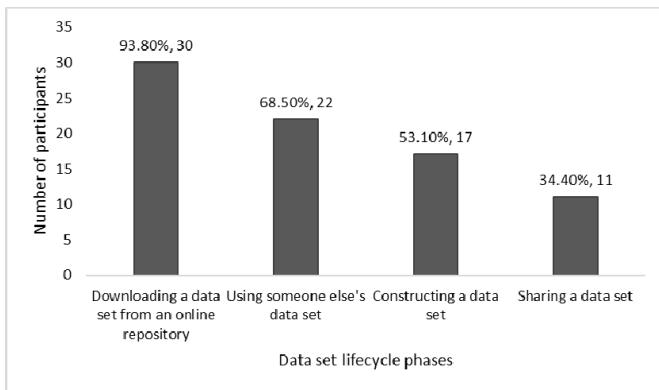


Fig. 6. Experience of the participant with respect to each phase of the dataset lifecycle.

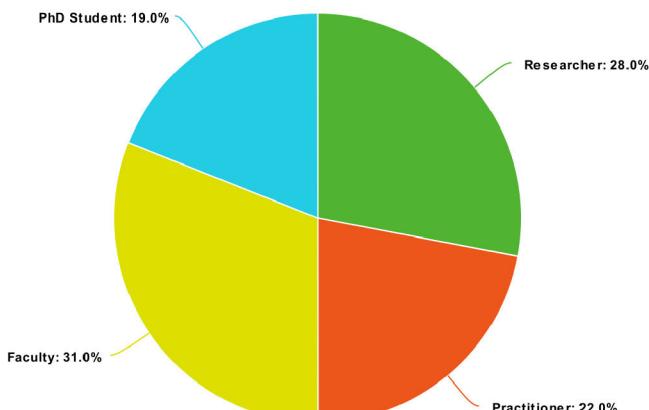


Fig. 7. Participant occupation.

lifecycle (see Fig. 3) was used as a conceptual framework for designing this survey. The survey was sent to 98 HAR researchers, and 32 (32%) responses were received. The experience of the participants with respect to each phase of the dataset lifecycle is shown in Fig. 6, and Fig. 7 illustrates the range of participant occupations.

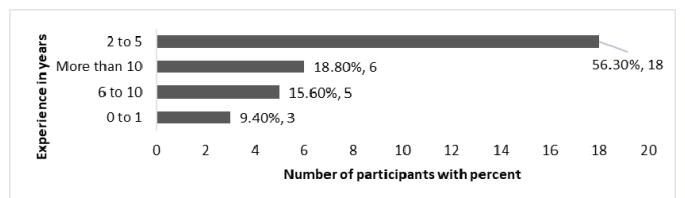


Fig. 8. Participants' experience in years.

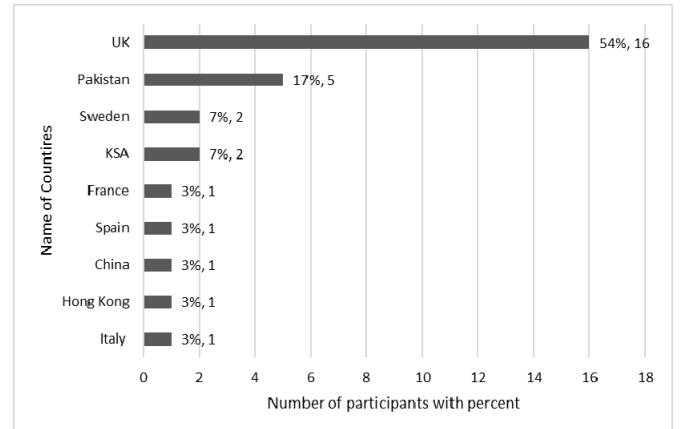


Fig. 9. Location of the participants.

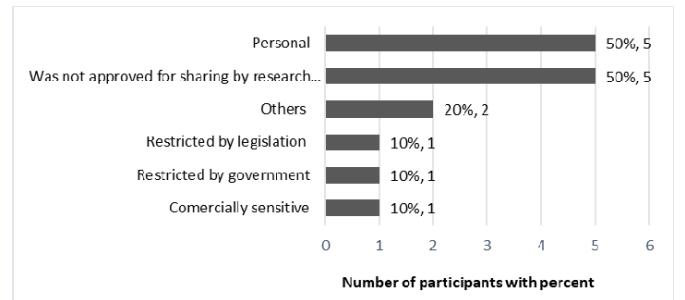


Fig. 10. Constructed datasets' sharing constraints.

Similarly, Fig. 8 shows the participants' experience by years, and Fig. 9 presents the location of the participants.

1) Construction: This is the most significant phase of the open dataset lifecycle as it involves the identification of what type of data is needed for analysis to achieve the research objectives. Furthermore, what type of data collection protocol will be used to collect data systematically and to build the overall guidelines for data collection? A total of 53.1% of participants who took part had experience in constructing a dataset.

From an initial analysis of the question “what is the main issue or challenge you have faced when generating a new dataset in HAR?” the major issues and challenges identified are shown in Table II.

In datasets' sharing after construction, the responses were taken from the participants regarding sharing restrictions of the constructed datasets from the HAR research community, as shown in Fig. 10.

Finally, participants were asked what is the main piece of advice they would give to another researcher when generating a new dataset in HAR. From the 32 survey participants, 18 responses were received, as summarized in Table III, with similar advice paraphrased and grouped as shown.

TABLE II
ISSUES AND CHALLENGES IN OPEN DATASETS' CONSTRUCTION PHASE

Main issues/challenges	Description	Frequency of participants
<i>Data set annotation</i>	Generating sufficiently precise annotation of the data. This is a very time-consuming task if labelling the data set after collection.	4
<i>Time synchronization with multiple sensors</i>	Synchronization of timings from multiple sensors and heterogeneous sources	2
<i>Ethical approval</i>	Obtaining Ethical approval from the host organization	1
<i>Data set size</i>	Data set size is important in determining the performance of a machine learning model. The consequence of using a test set with a limited sample size might lead to a large variation as well as a high error rate. Further, a small set of data set tends to overfitting and creates incorrect results.	1
<i>Privacy</i>	Collecting video data involves enough with sufficient contextual information can lead to privacy issues.	1
<i>Training of users and actors</i>	Providing proper training to the users/actors who help in generating new data by performing different activities.	1
<i>Missing values</i>	Missing data due to sensor malfunctioning or poor user engagement when collecting data with multiple users in free-living. The missing data could bias the final results because missing data adds ambiguity to the data.	1
<i>Finding volunteers</i>	Finding sufficient volunteers for collecting the data set while ensuring diversity of participants with respect to privacy and social concerns.	1
<i>Lack of funding</i>	Lack of funding specifically for data collection process.	1

Responses related to the normal practice for sharing constructed datasets in which the statement “always seeks permission to share” up to 50% and “shared only if required by the research sponsor” equal to 33%. This indicates that datasets are often not open-sourced but used only for their original research purpose. When asked about the normal practice for dataset sharing, the responses are shown in Fig. 11.

2) Sharing: This phase of the dataset life cycle is concerned with making the dataset available for researchers and users to use it for solving their HAR problems. Dataset owners must grant access to their shared datasets for research purposes. Researchers and dataset owners store their datasets in a range of storage repositories. Different categories of dataset

TABLE III
MAIN PIECE OF ADVICE PARTICIPANTS GIVES TO ANOTHER RESEARCHER WHEN GENERATING A NEW DATASET IN HAR

Participants responses	Frequency
Be thorough with annotation of the dataset.	3
Keeping participants anonymity	2
Give information about the process of data filtering and cleaning.	2
Plan well the formatting and synchronisation of data for timestamps. Using the same format facilitates cross-validation of approaches	2
Be prepared to spend quite a bit of time trying to make it clear what is in the data and documenting it, if you want other research groups to be able to use it.	2
Carefully clean the data prior to sharing	1
Provide context and ground truths if possible.	1
What kinds of activities were going on in the time between the collection and the sensor placement?	1
Ensure ethical approval before data collection	1
Include or reference supporting papers with result discussions	1
Data should remain open if shared online, links must not be deactivated.	1
It is important to design a protocol for generating new dataset in HAR	1

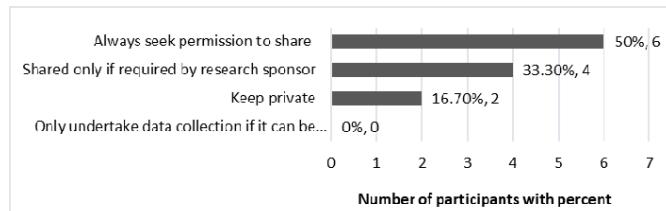


Fig. 11. Normal practice for constructing experimental datasets.

repositories are used based on the nature of the data and the researcher's licensing requirements. The main categories of open dataset repositories are discussed in the following.

a) Institutional dataset repositories: Institutional-based repositories are collections of datasets that are created and maintained by a specific institution, such as a university or research organization. These repositories may contain data from a wide range of research projects and disciplines, and may be used to share data across the institution or with external researchers. They may also have more robust data management and preservation capabilities, as they are typically maintained by dedicated staff. They distribute and manage datasets, research outcomes such as data analysis, source code, tools, and frameworks that are developed by the research community and postgraduate students. Examples include DataShare⁴ from The University of Edinburgh, the University of Cambridge Data repository⁵, and the UCL Research Data Repository.⁶

⁴<https://www.ed.ac.uk/information-services/research-support/research-data-service/after/data-repository>

⁵<https://www.data.cam.ac.uk/data-repository>

⁶<https://www.ucl.ac.uk/library/open-science-research-support/research-data-management/ucl-research-data-repository>

b) *Hosted repositories*: Hosted repositories refer to datasets that are stored and maintained on a remote server and can be accessed by the public over the Internet. These datasets can be easily discovered and downloaded by anyone with an Internet connection. UCI ML Repository, Figshare, DataPort, Zonedo, Kaggle, and so on also engage in the management and distribution of open datasets from educational organizations and industries.

c) *Government dataset repositories*: These repositories are managed by the government and store sensitive data for healthcare management, surveillance, and administrative purposes. These datasets are not shared commonly due to privacy concerns and intellectual property protection of the user data [168]. Examples of government dataset repositories are “data.europa.eu”⁷ and “data.gov.uk.”⁸

d) *Specific domain dataset repositories*: These repositories consist of datasets and metadata from a specific domain, such as healthcare, sports, engineering, and social sciences. The main advantages of these datasets are related to specific research domains. Examples are geriatric healthcare⁹ and physiotherapy.¹⁰

e) *Project-based repositories*: Project-based repositories are collections of datasets that are created and maintained by a specific research project or group. These repositories are typically focused on a specific topic or area of research and may contain data that are collected and analyzed by the project team, for example, the Centre for Data and Visualization Sciences, Duke University [169], managing research data by the University of Bristol [170], and ML and AI dataset managing by Carnegie Mellon University [171].

The percentage of participants who took part in dataset sharing was 46%. This clearly shows that most of the researchers are not taking part in dataset sharing.

Furthermore, the participants were also questioned whether the sharing of data was restricted because it contained information about an organization or participants. On a Likert scale that included strongly agree, agree, neutral, disagree, and strongly disagree, participants' opinions were measured. Overall, 13% of respondents were neutral, 46% agreed and strongly agreed with the statement, and the remaining respondents disagreed.

Furthermore, the participants were asked the question “sharing the dataset was not possible because of problems with the data?” In response, most of the participants agreed up to 93% and consent to the problem with data being the main barrier to datasets’ sharing.

Moreover, participants were asked to select reasons for not sharing their dataset. The top reasons were the missing values, errors during measurement, and the small sample size of the data. Also, 28% of participants agreed to the large size, the data are not in presentable form, and the lack of time and the lack of resources for datasets’ sharing are the main reasons. The detail responses are mentioned in Table IV with participants’ frequency and percentage.

⁷<https://data.europa.eu/en>

⁸<https://www.data.gov.uk/>

⁹<https://mira.mcmaster.ca/research/open-access-datasets-from-aging-studies>

¹⁰<https://pedro.org.au/>

TABLE IV
WHY THE SHARING OF DATA WAS LIMITED

Asked questions related to data set sharing	Participants responses in numbers	Participants responses in percentage
The sharing of data was limited because it contained identifying information about an organization or participants.	-Strongly disagree: 0 -Disagree: 6 -Neutral: 2 -Agree: 6 -Strongly agree: 1	-Strongly disagree: 0% -Disagree: 40% -Neutral: 13.3% -Agree: 40% -Strongly agree: 6.7%
Sharing the data set was not possible because of problems with the data?	-Yes: 1 -No: 14	-Yes: 6.7% -No: 93.3%
Please review the following list of reasons and select all that apply for your data set:	-Missing values: 1 -Outliers: 0 -Measurement errors: 1 -Overfitting is harder to avoid: 0 -Sample size too small: 1	-Missing values: 100% -Outliers: 0% -Measurement errors: 100% -Overfitting is harder to avoid: 0% -Sample size too small: 100%
Sharing the data set was not possible because it was too large?	-Yes: 4 -No: 10	-Yes: 28.6% -No: 71.4%
Please select all that apply	-Too large for selected repository: 1 -Raw data not in a presentable form: 1 -Lack of time and resources to share large data set: 4 -Lack of experience in data management: 0	-Too large for selected repository: 25% -Raw data not in a presentable form: 25% -Lack of time and resources to share large data set: 100% -Lack of experience in data management: 0%
I did not share the data set because I was unsure of the best approach?	Yes: 2 No: 13	Yes: 13.3% No: 86.7%

The other potential reasons that were asked by the participants for not sharing datasets are shown in Fig. 12.

When asked “what is the main issue or challenge you have faced during dataset sharing?” participants responded as shown in Table V.

Responses related to the license used for dataset sharing are shown in Fig. 13. Most of the participants did not use any license while sharing their datasets.

Participants were asked if they found the process of depositing and sharing a dataset time-consuming. As shown

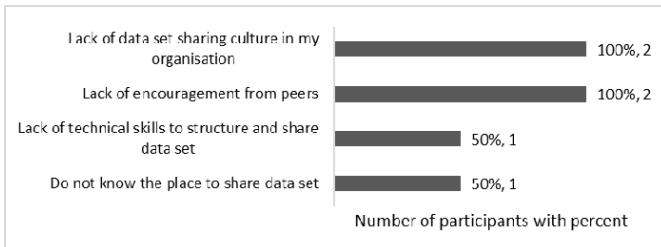


Fig. 12. Other reasons for not sharing datasets.

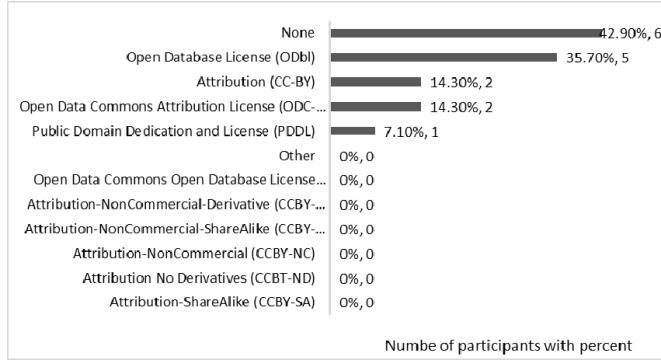


Fig. 13. License used by the researchers while sharing datasets.

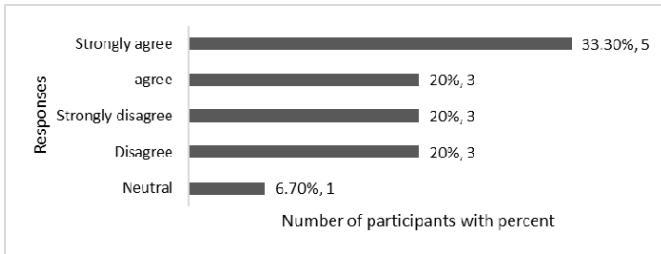


Fig. 14. Participants' responses regarding the process of depositing and sharing datasets are time-consuming.

in Fig. 14, most of the participants strongly agreed that this was the case.

Participants were also asked "what motivates you to share your dataset with the HAR research community?" Details of participants' responses with frequency are given in Table VI.

The final question in relation to dataset sharing is "what is the main piece of advice you would give to another researcher when sharing a new dataset in HAR?" Table VII shows the exact response of participants with frequency.

3) Finding: This phase of the open dataset life cycle is related to searching and finding a dataset for a specific problem domain. Responses received from the participants related to the questions such as "have you searched for and downloaded an open dataset from an open data online repository for experimental/research purposes?" A total of 96.9% of the participants obtained an open dataset for the purposes of study or experimentation by downloading it from an online open data repository. The next question is related to the criteria that are used by the researcher when searching a dataset—the results are shown in Fig. 15.

Participants were asked about the amount of preprocessing (none, a little, some, and a lot) required after downloading a dataset; 29% of participants indicated a lot, 38.7% some, 32% a little, and 0% with none. For the respondents, every dataset

TABLE V
ISSUES AND CHALLENGES IN OPEN DATASETS' SHARING PHASE

Main issues/challenges	Description	Frequency of participants
<i>Privacy</i>	Even when participant anonymity is a challenge, the collection and exchange of massive quantities of data may expose individuals to significant privacy violations.	4
<i>License</i>	License awareness and how to select the most appropriate license for data sharing	2
<i>Proper data collection</i>	Expensive to do a proper data collection by buying good quality resources, recruiting, training and monitoring of the participants	1
<i>Getting data visibility</i>	Difficult to interpret and to get the visibility of the data from the existing data sets.	1
<i>Ethics approval</i>	Obtaining ethics approval from the organization before data sharing	1
<i>Data format</i>	There is no standardized data format for data sharing	1
<i>Data presentation</i>	Presenting data in an understandable and easy way	1
<i>Resolving data set queries</i>	Resolving data set queries by checking access to the server, adding new data records and version control of the data set.	1
<i>Lack of time</i>	Need time to clean and share data	1
<i>Lack of resources</i>	Need resources for data sharing and to perform maintenance or version control of the data set.	1
<i>Data set size</i>	Low sample of data leads to bias in the results and effect's reliability due to higher variability in the data.	1
<i>Data quality</i>	Hard to maintain standard data format and quality before sharing by removing data bias and to deal with missing data.	1
<i>Platform for sharing</i>	No proper platform for researchers and practitioners to share data set, collaborate, benchmark, validate and improve data set.	1

required preprocessing after downloading to make it usable and informative.

As with previous stages of the open dataset lifecycle, participants were asked "what is the main issue or challenge

TABLE VI
MOTIVATING FACTORS FOR SHARING YOUR DATASET

Participants responses	Frequency
Further citations.	3
To allow for study replication, transparency of method and verification of results and contribution.	2
Serve the scientific community and industry.	2
Open research and funding grants.	2
To get more insights from the data and extract more useful information.	1
It motivates us to find insights from the human behaviour.	1
Having standard repository for storage and sharing	1
To see how different researchers manage and analyse the data.	1
Collecting a dataset requires is a huge effort. Making it available to other researchers maximize the return	1
The desire to benefit humankind by improving data analysis and algorithms.	

TABLE VII
MAIN PIECE OF ADVICE TO ANOTHER RESEARCHER WHEN SHARING A NEW DATASET IN HAR

Participants responses	Frequency
Ensure data set is shared on a repository that has longevity/permanence.	2
Be prepared to spend quite a bit of time making the dataset understandable to others.	2
To be aware of the size and privacy of data	2
Carefully clean the data	1
Publish your research to increase visibility of the dataset	1
Establishing clear and simple dataset guidance	1
Explore for datasets relevant to your research. This will facilitate future research.	1
Share in a standard repository	1
Ensure you capture the methods used to collect the dataset at the time of recording the data.	1
Research must be open.	1

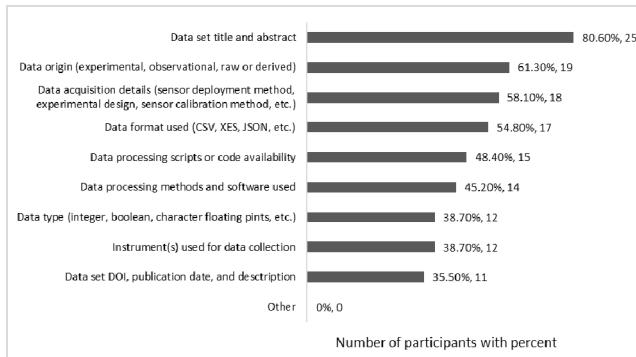


Fig. 15. Searching criteria for dataset searching.

you have faced when trying to find a suitable dataset?" The responses and their frequency are shown in Table VIII.

Similarly, participants were asked about their motivations for using an existing dataset instead of creating a new one. Responses are shown in Table IX.

TABLE VIII
ISSUES AND CHALLENGES IN OPEN DATASETS' FINDING PHASE

Main issue/challenge	Description	Frequency of participants
<i>Metadata</i>	Lack of data set description such as activities observed, devices, sensors, participants, etc.	5
<i>Data set completeness</i>	Missing data and missing annotations	5
<i>Specific domain data set</i>	Data sets are only limited to a specific domain	4
<i>Access restriction</i>	Data set access and use restriction from the data set owner	3
<i>Data set size</i>	Insufficient data for meaningful use	2
<i>Data set format and structure</i>	Lack of standard data format and structure	2
<i>Repository search optimization</i>	Hard to search relevant data set in the repository	2
<i>Data set authenticity</i>	Lack of trust in data collection methods and protocol used	1
<i>Time and effort</i>	Time and effort needed to search a data set and understand it	1

The participants' responses regarding the datasets' repository or directory used when searching are shown in Table X.

Finally, the main piece of advice from the participants for other researchers while searching for a dataset in HAR is shown in Table XI.

4) Using: This is the final stage of the open dataset lifecycle. After finding a dataset, researchers typically apply different ML and statistical analysis tools and frameworks to visualize the data and produce meaningful results. Researchers and practitioners use open datasets to solve their domain problems and make a decision for future solutions.

Dataset preprocessing is a key step to make the dataset informative and to improve dataset quality before using. It is a data-extracting procedure that includes converting raw data into a comprehensible and informative format. Datasets collected within the natural environment are often inadequate and inconsistent due to missing values and activities/actions labeling that leads to errors [172]. Data preprocessing is performed by the researchers to resolve challenges such as missing values, data annotations, handling errors and outliers, codes, and naming discrepancies [173]. The major steps involved while performing dataset preprocessing are: 1) data cleaning; 2) handling null values; 3) standardization; 4) handling categorical values; and 5) feature scaling and dependencies [174].

In the survey, the question was asked "have you used and evaluated someone else's open dataset for experimental/research purposes?" 78% of participants responded yes.

TABLE IX
MOTIVATION FACTORS FOR USING THE EXISTING DATASETS

Motivating factors	Description	Frequency
Save time/effort/resource	Cleansing data necessitates more time and resources before it can be utilised for experimental purposes. Already cleaned data will save more time/effort/resource.	14
Benchmarking	Data set able to be used as a benchmark	3
Meta data	Availability of a clear data description	2
Problem domain	Provide a description of the data together with the applicable problem domain.	1
Use already published data set	Data set already published in a research platform	1
Authentic data	Trusted data	1
Reproducibility	Ease of reproducibility of results and application of novel algorithms	1
Research topic	The topic and the challenges of the research are relevant?	1
Challenge	Creating new dataset is challenging! Existing public datasets are easy to use	1
Data set size	Data set provides enough amount of sufficient data for meaningful analysis.	1

TABLE X
DATASETS' REPOSITORY/DIRECTORY USED FOR FINDING A DATASET

Data set	Frequency
UCI	11
Google search	9
Kaggle	8
GitHub	5
Publisher website	5
Data search	3
IEEE Dataport	3
Direct contact	1
Sparkbankan	1
Data portal	1
Open ML	1
Planet lab	1
Zonedo	1

Two questions related to metadata description are given as follows: was the metadata description easy to understand? and was the metadata description accurate? Responses are shown in Figs. 16 and 17, respectively.

Respondents were asked if they encountered a dataset update issue after the initial sharing, with a response shown in Fig. 18.

The main issues and challenges identified by respondents while using someone else's dataset are shown in Table XII.

TABLE XI
MAIN PIECE OF ADVICE FOR OTHER RESEARCHERS WHILE SEARCHING A DATASET

Main piece of advice	Description	Frequency
Documentation	Look for supported documentation of a data set such as published papers, meta data, data collection protocol, sensors used, participants, format etc.	7
Data set check	Before using a data set, check for certain feature such as timestamps are being in order, missing values, no mismatch of columns, noise in data, data normalization etc.	4
Search terms and keywords	Identify the most relevant keywords and search terms for searching a data set according to the domain problem	3
Goal and objective	The main advice for searching HAR dataset is to first identify the goals you want achieve and also to arrange resources for processing the data in Advance.	3
Generic data set	Make datasets generic for a specific field of research so that several research questions can be answered using the datasets	3
Data quality	Look for data set completeness such as meta data and description	2
Libraries identification	Exploring available libraries for data set cleansing and pre-processing	1
Check data set license	Check licensing to ensure permission to use is within remit of experiment.	1
Available resources	Try every available resource (repository?) for searching and finding a data set	1
Benchmarked data set	Search for an already used and benchmarked data set	1
Contact	To contact related parties that may provide the data set	1

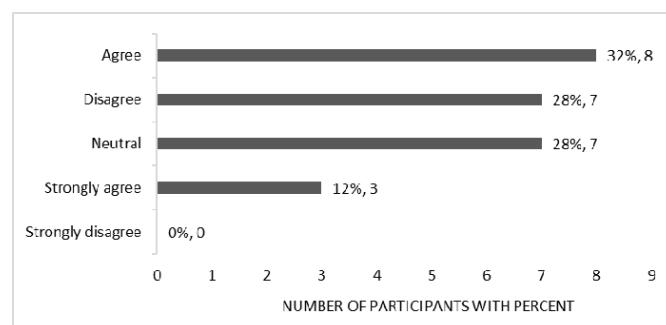


Fig. 16. Metadata description makes it easy to understand a dataset.

Finally, the main piece of advice from the participants to other researchers while using someone else's dataset in HAR is illustrated in Table XIII.

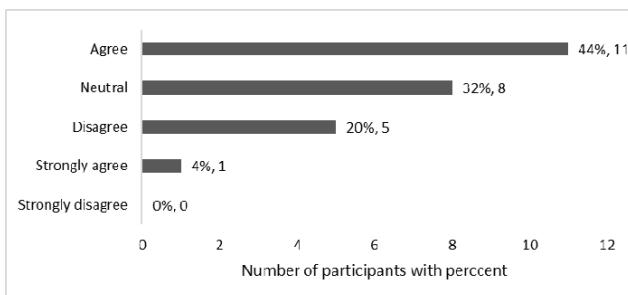


Fig. 17. Metadata description accuracy.

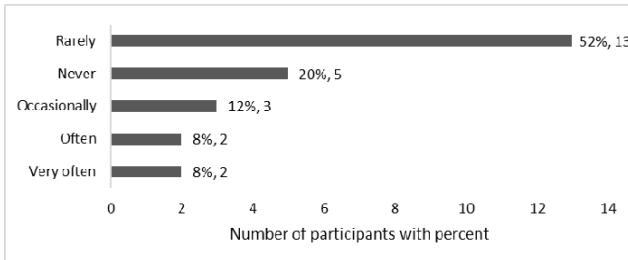


Fig. 18. Agreement of dataset update issue after the initial sharing.

TABLE XII

ISSUES AND CHALLENGES IN OPEN DATASETS USING PHASE

Main issue/challenge	Description	Frequency of participants
Documentation	A complete description of data set such as data collection protocol, devices used, guidelines for using the data set, etc.	6
Annotation	Understanding how data was captured, what the data shows (annotation/labelling of activity)	5
Format	Various formats of data sets	5
Missing values	Missing values and empty columns of a data set	4
Data noise and imbalanced data	Missing columns, irrelevant information, anomalies in a data set	3
Inaccurate assumption	The An inaccurate assumption from the researcher about the experiment relevance of the data set??.	2
Sensor	No information about sensors and data collection with limited sensors	1
Data restriction	Limited/restricted amount of data. For example, 7 days' work or a few months' work or lab settings only, etc.	1
Trust	Lack of trust and proper ground truth about a data set	1

Other relevant questions were asked “in your view, what factors improve dataset quality (one per line)?” The responses associated with the important factors of dataset quality are shown in Table XIV.

TABLE XIII
MAIN PIECE OF ADVICE WHILE USING SOMEONE ELSE'S DATASET

Main piece of advice	Description	Frequency
Data set analysis	Conduct a Complete analysis of data set and meta data for understanding data set	5
Time and effort	Provide more Set aside sufficient time and effort for understanding the data set and matching it to the requirements of problem domain	4
Documentation	Read supporting documentation provided with data set and also published papers on a data set	3
Data set pre-processing	Pre-process and cleaning of the data set to make it usable and informative	2
Issues reporting	During data set analysis, if found any issues or problems are found, report it them to the data set owner/research community	1
Trust	Be careful. If you don't trust the data/labels, then everything built on top of it cannot be trusted.	1
Open source	Priority should be to find open-source dataset first.	1
Validation	Validate your research experiment on multiple datasets and in different scenarios to increase achieve generalization.	1

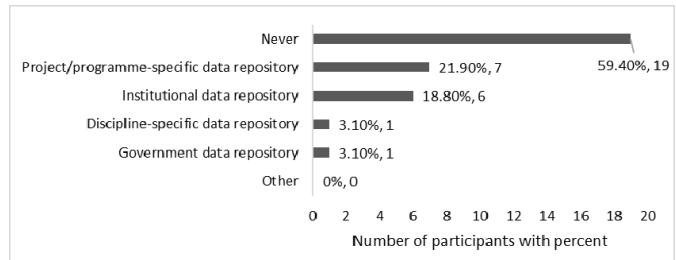


Fig. 19. Datasets' repository used for online registration.

Participants were asked “have you ever preregistered an experiment in any of the following online repository domains?” 59% indicated that they never used any online repository, 21% had preregistered on the project-/program-based repository, and 18.8% indicated an institutional dataset repository, as shown in Fig. 19.

The question was asked “how important to you is the replication of an HAR experiment using an open dataset?” Responses are shown in Fig. 20.

Finally, the participants were asked about the importance of dataset benchmarking with responses, as shown in Fig. 21.

V. ANALYSIS AND COMPARISON OF LITERATURE REVIEW AND QUESTIONNAIRE SURVEY

This section presents an analysis and comparison of the outcomes of both the literature review and questionnaire

TABLE XIV
IMPORTANT FACTORS FOR DATASET QUALITY

Improvement factors	Description	Frequency
<i>Annotation</i>	Annotating a data set involves adding metadata, such as descriptions or tags. This can be done in a variety of ways, for as by transcribing audio recordings, labelling certain elements in videos, or providing textual descriptions to images.	6
<i>Standard format and structure</i>	The term "standard format and structure data set" is used to describe a collection of information that has been organised in a way that makes it simple to read, analyse, and disseminate.	6
<i>Data cleaning and pre-processing</i>	Preparing a dataset for analysis necessitates cleaning and pre-processing it to remove or rectify errors, inconsistencies, and missing data. Data cleansing include eliminating duplicates, fixing typos, completing blanks, and standardising file formats.	4
<i>Data set noise</i>	The term "dataset noise" is used to describe the occurrence of meaningless, contradictory, or incorrect information inside a dataset. This may arise as a result of typos, bad data gathering, or the addition of extraneous information. The quality of any studies or conclusions made from a dataset might be diminished by the presence of noise in the dataset.	4
<i>Documentation</i>	The term "dataset documentation" is used to describe the information and documentation supplied about a dataset, such as the dataset's goal, origin, structure, format, relevant metadata, annotations, etc.	3
<i>Experimental setup</i>	An experimental setup is the arrangement of equipment, materials, and circumstances utilised in an experiment. It covers experiment design, data collection, and quality control. The experiment's validity and reliability depend on its setup.	3
<i>Typos errors and duplication</i>	A dataset is a collection of data that is structured and organised in a certain way. Typographical errors are mistakes produced while typing data, such as misspellings or grammatical errors. The occurrence of duplicate data inside a dataset,	2

TABLE XIV
(Continued.) IMPORTANT FACTORS FOR DATASET QUALITY

	where the same information is repeated many times, is referred to as duplication.	
<i>Activity description</i>	The dataset often contains sensor data obtained from wearable devices, such as accelerometer and gyroscope measurements, as well as annotations or labels indicating the activity being done at a given moment. The possible actions include walking, running, leaping, sitting, and standing, among others. The dataset may additionally contain extra information, such as participant demographics and environmental characteristics.	2
<i>Quality metrics</i>	A dataset's quality metrics are a set of standards by which its reliability and accuracy may be judged. Completeness, correctness, consistency, timeliness, and relevance are some examples of metrics that may be used. The authenticity and trustworthiness of the sources and the suitability of the data for its intended purpose are other crucial considerations when assessing the quality of a dataset. Quality metrics for datasets may also incorporate indicators of data governance, such as data lineage, data provenance, and data security.	1
<i>Accuracy</i>	Dataset accuracy refers to the degree to which the data inside a dataset are correct or reliable. It assesses how accurately the dataset's data matches the actual values or attributes of the items or individuals it represents. A dataset with high accuracy has few mistakes or inaccuracies, whereas a dataset with low accuracy contains a significant number of errors or inaccuracies.	1
<i>Meta data</i>	Dataset metadata includes the title, author, date generated, format, and other attributes about a dataset. The dataset or its metadata file normally contains this information. Metadata aids data discovery, reuse, and analysis by explaining the data's context and purpose.	1
<i>Feature engineering</i>	Making new features or modifying existing features in a dataset is called "feature engineering," and it may be	1

survey. Section V-A assesses the open issues and challenges identified. Section V-B outlines the best practices identified

from the literature review and mentioned by the survey participants.

TABLE XIV
(Continued.) IMPORTANT FACTORS FOR DATASET QUALITY

	used to enhance a machine learning model's ability to accurately represent the data. Dimensionality reduction, feature scaling, and feature extraction are all examples of such methods. The purpose of feature engineering is to enhance the performance of a machine learning model by identifying and utilising the most informative and pertinent features within the available data.	
<i>Missing data</i>	A dataset with missing data is a collection of information in which certain values are absent or have not been captured. This can arise for a variety of reasons, including mistakes in data collection, data input, or data source constraints. Missing data may have a substantial influence on the accuracy and use of a dataset, since it might result in findings that are skewed or insufficient.	1
<i>Data set size</i>	Dataset size is the number of observations or records. It can also mean a dataset's variables or properties. A dataset might comprise hundreds or millions of observations. Data processing, analysis, and tools depend on a dataset's size.	1

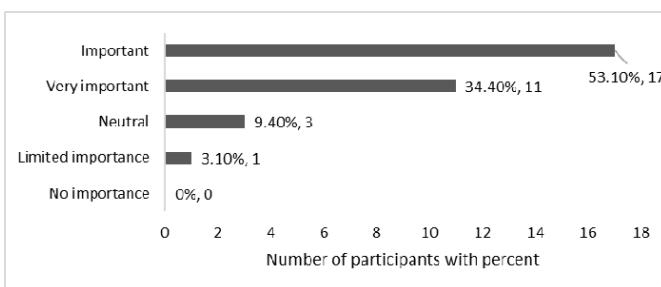


Fig. 20. Importance of dataset replication.

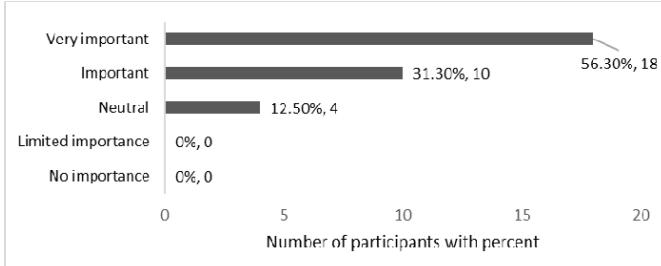


Fig. 21. Importance of benchmarking.

A. Issues and Challenges

Fig. 22 shows the main issues and challenges identified from the literature review and survey. Some issues and

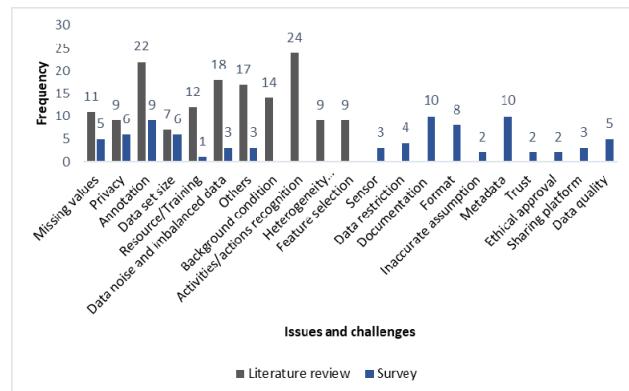


Fig. 22. Open datasets' issues and challenges from both literature review and questionnaire survey in the HAR domain.

challenges intersect the questionnaire survey and literature review, which are given as follows:

- 1) missing values;
 - 2) privacy;
 - 3) annotations;
 - 4) dataset size;
 - 5) resource and training;
 - 6) data noise and imbalanced datasets;
 - 7) others (dataset documentation, volunteer participants recruitment, lack of sharing datasets, subject annoying to wear device, and lack of protocol for data collection).
- Some issues and challenges that are identified only in the literature review are given as follows:
- 1) background condition;
 - 2) activities'/actions' recognition;
 - 3) device/data/subjects' heterogeneity;
 - 4) feature selection.

Those from the questionnaire survey only are given as follows:

- 1) sensor's information;
- 2) data restriction from organizations and research groups;
- 3) dataset documentation;
- 4) standard data format;
- 5) inaccurate assumption during dataset construction;
- 6) metadata;
- 7) trust in the dataset authenticity;
- 8) ethical approval from the hosted organizations;
- 9) sharing platform;
- 10) data quality.

As a result, all the issues and challenges that are obtained as an intersection form both literature review and questionnaire survey or separated from each need focusing to address by proposing approaches, frameworks, and tools to overcome.

Getting over the issues and challenges and coming up with solutions for the future might be beneficial in the following areas.

- 1) *To Improve Dataset Structure, Quality, and Access Improvement:* New knowledge and approaches can be used to improve the structure of the data and the quality of the data by removing irrelevant information and to increase data access due to open-source technology. Open-source technology improves the availability and transparency of the data [175], [176].
- 2) *Good Quality Open Datasets Can Improve Personal Fitness:* Good quality open datasets may contain

data on levels of physical activity, diet, sleep habits, and other health-related information. These data may be used to develop individualized workout routines, establish and track objectives, and measure progress over time. Leveraging technology and datasets with personal fitness can improve health related quality of life [177], [178].

- 3) *Good Quality Open Datasets Can Enhance Decision-Making by Providing a Broad and Diverse Range of Information That Can Be Used to Inform Decisions:* Decision-making can be improved and optimize by the access of accurate and up to date information and data. This will make easy to allow decision faster. HAHAR and medicine-related datasets assist healthcare to get more informed decision [175], [179], [180].
- 4) *Globalization of Datasets:* Good datasets play a critical role in the globalization of datasets by providing accurate, reliable, and relevant information that can be shared and used across different countries and cultures. The distribution of datasets plays an essential role in leveraging the practices and knowledge globally. Accurate and good quality datasets are widely accessible, and the contributors may retrieve new information from all regions globally [175], [181].
- 5) *Good Datasets Can Play a Crucial Role in Early Disease Detection by Providing the Information Needed to Identify Patterns and Trends That Can Indicate the Presence of a Disease:* HAR and healthcare-related datasets developed from wearable sensors allow to detect a disease early to assist medical experts in improving treatment and patient healthcare [179], [180], [181].
- 6) *Rehabilitation:* By providing valuable information that can inform the design, implementation, and evaluation of rehabilitation programs, for example, identifying patient needs, tracking progress, evaluating outcomes, and identifying risk factors and disparities, properly collected datasets can play an important role in rehabilitation to assess, treat, and manage an individual to leverage their social, physical, cognitive, and physiological functions and get back to normal condition [182], [183].

B. Best Practices

Fig. 23 presents the best practices used and recommended by HAR researchers from both the literature review published literature and survey. The best practices arising from the survey are related to all the dataset lifecycle phases. However, the best approaches from the literature review are mostly related to the using phase of the dataset lifecycle due to their use in experiments and subsequent preprocessing and improvement in order to fit the reported experiments.

The best practices identified from the literature review on the right-hand side of Fig. 23 are data transformation and segmentation to divide the data into various length windows and then assign labels manually based on marking timestamps for each activity to improve the metrics for datasets' benchmarking. Also, feature extraction and dataset

preprocessing to make datasets fit for their experiments. Researchers used mean values when dealing with missing values in a dataset. Deeply learned features are frequently more robust and accurate than hand-engineered features because the model can learn from the data itself, instead of relying on human-designed features. Furthermore, researchers used data augmentation to reduce overfitting and deeply learned features for dealing with large-scale video-based datasets. Finally, use synthetic data to overcome the problem of privacy and data scarcity to deal with sensitive and more personal information.

The recommended best practices from survey participants on the left-hand side of Fig. 23 are data cleaning and preprocessing of the dataset before sharing and providing data annotation while constructing a dataset. Also, they recommended using standardized data formats during data collection and performing dataset benchmarking to validate and improve datasets; datasets should be open source and should be easily accessible to all researchers. Researchers should design a protocol for data collection during dataset construction and provide complete metadata related to a dataset to improve data quality and reusability. Importantly, researchers should ensure ethical approval before data collection from the host organization and provide anonymity to participants to preserve privacy. Finally, providing complete documentation with shared datasets will assist other researchers prior to performing their experiments.

Regarding the data collection methods, such as controlled and uncontrolled in the HAR domain, the distinction between both methods plays a pivotal role in shaping the landscape of dataset quality and applicability [184]. Controlled data collection involves particularly structured environments and predefined activities, enabling precise data labeling and facilitating algorithm training. This method ensures a high degree of consistency and reproducibility that can be used for benchmarking and fine-tuning algorithms [26]. On the other hand, uncontrolled data collection unfolds in real-world contexts, capturing natural human behavior. While generating datasets that represent real scenarios, this approach introduces challenges such as data noise and annotation ambiguity [30]. Both methodologies offer unique insights and bear relevance in addressing diverse research problems. Exploring the interaction between these methods is crucial, as it drives advancements in HAR methodologies, enriches dataset diversity, and augments the field's overall robustness.

VI. CLASSIFICATION OF ISSUES AND CHALLENGES WITH OPEN DATASETS

The main objective of this section is to organize and classify all of the identified issues and challenges in a systematic manner to assist researchers in their understanding and interpretation of the results. From the conducted comprehensive literature review and survey, key issues and challenges relating to open datasets in HAR have been identified in terms of all phases of the open dataset lifecycle. In addition, we also make reference to external factors relevant to these issues and challenges, as shown in Fig. 24.

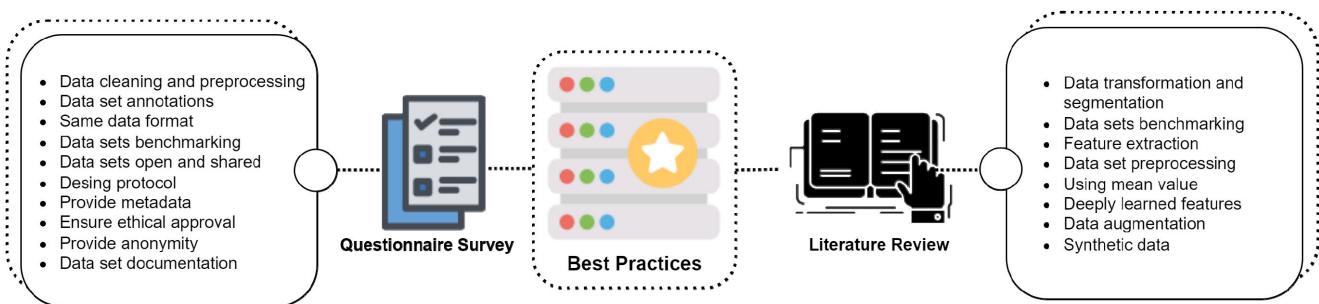


Fig. 23. Best practices explored from the literature review and recommended by survey participants.

These external factors were derived through the authors' brainstorming on the results of the literature review, survey, and related published research work [99].

A. Based on Dataset Lifecycle

This is divided into researcher and user perspectives. Researchers construct datasets and then share and deposit them into dataset repositories. Users find datasets by applying relevant keywords and search strings to download and use them according to their research purposes. The issues and challenges in this category have been discussed above and include dataset annotations/labeling, missing values, data format, dataset size, and data noise.

B. Based on Organization and Governance

This is further divided into eight major categories: security, sharing platform, governance, guidelines, activities, training, environment, and resources.

1) Security: This category is concerned with the security aspect of open datasets. Privacy concerns are the main issue while constructing and sharing datasets. Dataset owners must ensure anonymity during dataset construction and sharing to protect the personal data of the participants because some people in the dataset may not want their data to be made public. In addition, certain sorts of data, such as medical information or financial records, may be particularly sensitive and should be secured properly. Data must be appropriately anonymized and deidentified, and access restrictions and monitoring must be in place to avoid illegal access and exploitation of open datasets.

2) Sharing Platform: A sharing platform is an online community where people and groups may pool their resources and knowledge for the greater good. A sharing repository is a database or online storage space where users may deposit and retrieve files, documents, and other media. Making it simpler for users to access and share open datasets, a sharing repository on a platform might be invaluable. Open datasets are those that may be accessed, shared, and augmented by anybody who wants to do so. All sorts of studies, analyses, and ML projects might benefit from these datasets.

Several methods may be used by a sharing platform to make the most of available public datasets:

- 1) facilitating the discovery of useful datasets by providing search and filtering tools;
- 2) facilitating data analysis and interpretation by providing visual representations of data;

- 3) facilitating users' ability to work together and share their platform-based discoveries;
- 4) pushing users to add their own open datasets to the service as a means of contributing to it.

3) Governance: These are the emerging issues related to open datasets such as proper licensing of the data while sharing. Obtaining ethical approval from the institutions hosting the study. Proper licensing, ethical approval, and access limits are all components of open dataset governance. Data sharing and commercial usage are only two examples of the kinds of restrictions that might be stated in a license. Legal and ethical considerations, such as the need to preserve individuals' privacy, are considered throughout the ethical review process. Data can be made available to the public or kept private, depending on the access settings in place. Good governance of open datasets guarantees that data are utilized in a way that does not break the law or compromise ethical principles.

4) Guidelines: A common issue to arise during dataset construction is the need to design a protocol for data collection in a systematic way. Similarly, data consistency and mitigation of missing values, insufficient data samples, and metadata description of the dataset are common problems faced by the research community. Open dataset guidelines involve a protocol for data collection and sharing, which ensures that data are collected ethically and legally, and shared in a way that is accessible and understandable to others. Documentation is also important, as it provides information on the data's origin, context, and any limitations or biases present. Quality is crucial, as inaccurate or unreliable data can negatively impact research and decision-making. Metadata, or data about the data, is also important as it provides information on the data's structure, format, and any relevant information for understanding and using the data. Overall, open dataset guidelines aim to promote transparency, accessibility, and reliability in data collection and sharing.

5) Activities: Open HAR datasets are collections of data that are publicly available for researchers and developers to use in their work. These datasets are often used in activities related to recognizing similar, dissimilar, and complex activities. The main technical concern in HAR datasets is the recognition of dissimilar activities, such as eating, drinking, and making tea. For example, researchers may use open HAR datasets to train ML models to recognize similar activities, such as walking or running, based on sensor data. Similarly, the recognition of complex activities having subactivities and

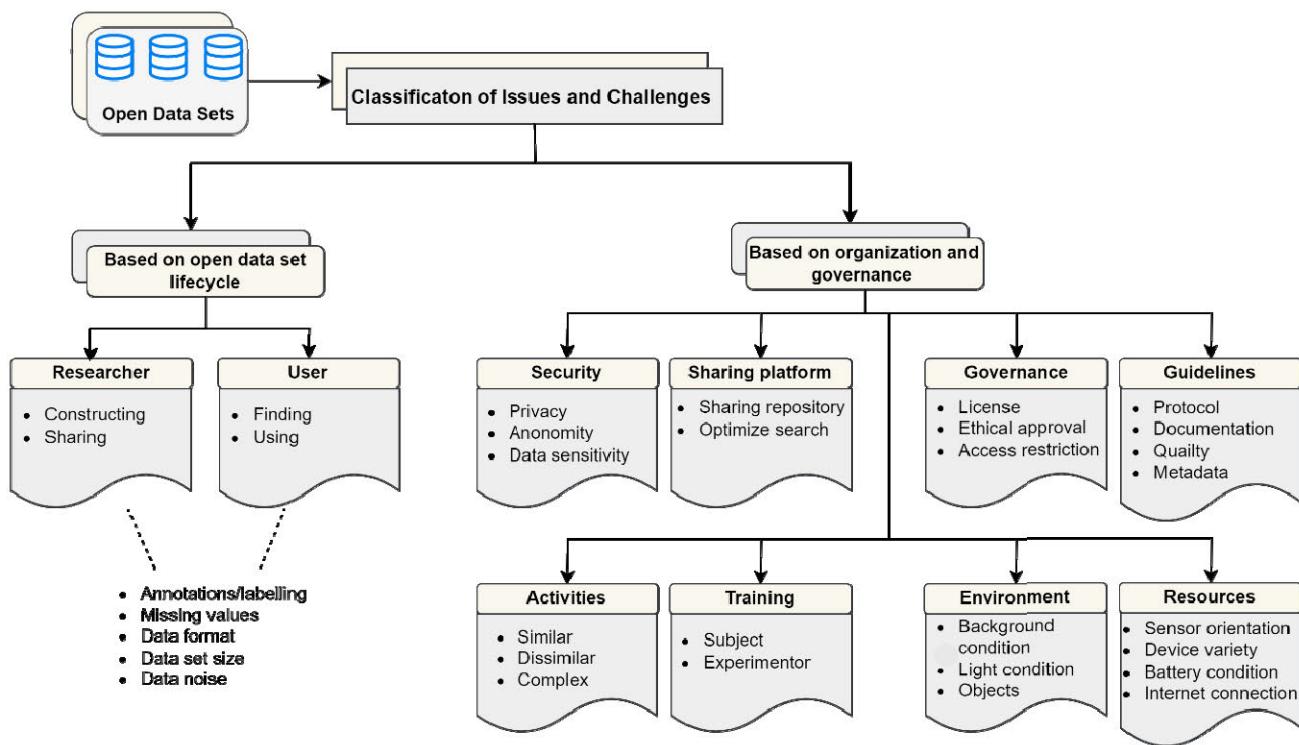


Fig. 24. Classification of open datasets' issues and challenges.

activities involves multiple participants. In addition, open HAR datasets can be used to analyze complex activities, such as multitasking or multitasking while performing different activities.

6) Training: HAR open datasets are collections of data collected from a variety of subjects and individuals, generally via wearable devices, such as smartwatches or activity trackers. These datasets are used to train ML models to identify and categorize various human activities, such as walking, running, and cycling. Dataset annotation is the main challenge for constructing a dataset, and it is often unnoticed due to insufficient or no training of the recruited participants/users for the experimental study. Awareness and training of the experimenter related to environmental factors involved are also needed while conducting an experiment.

7) Environment: Environmental and background conditions are the most challenging factors during dataset construction. The open HAR dataset collecting environment consists of a range of background conditions, lighting conditions, and other potential scene objects. This includes indoor and outdoor settings, varying levels of illumination, and diverse things such as furniture, people, and vehicles. The dataset is gathered in a naturalistic context, which means that the individuals are conducting activities of everyday life in their natural surroundings, unrestricted by any external factors. Researchers must be able to control and avoid irrelevant and noisy information and irrelevant objects during data collection.

8) Resources: This category includes the challenges related to sensor placement and orientation during data collection. Similarly, Internet connectivity with sensors, devices, and software applications. Also, a variety of different devices produce various data output formats and camera variations,

and the battery condition of wearable devices. The dataset resources contain information on sensor orientation, devices employed, battery health of devices and sensors, and Internet connectivity. This information is essential for comprehending the performance and limitations of various wearable devices and sensors, and it may be utilized to enhance the design and functioning of future devices. In addition, the dataset is continually updated to reflect new technologies and trends in the industry, making it an excellent resource for keeping up with the latest developments in wearable technology.

VII. HAR DATA EVOLUTION

Datasets are not always static artifacts, and the dataset lifecycle can be iterative in nature. Adopting this perspective, a number of promising avenues come to the forefront regarding dataset evolution.

A. Enhancing Real-World Diversity

Subsequent iterations of HAR datasets have the potential to incorporate a broader range of real-world settings, therefore reflecting the complex composition of human actions across varied surroundings. By broadening its scope to include uncontrolled environments, datasets have the potential to accurately represent the intricacies of everyday life, resulting in models that possess more resilience, adaptability, and alignment with real-world user encounters.

B. Cross-Domain Merging

The integration of HAR data with datasets originating from other domains, such as healthcare, environmental monitoring, or social interactions, presents a promising opportunity for cross-domain merging. The utilization of an interdisciplinary

approach has the potential to reveal new and unique perspectives and associations, facilitating a more profound comprehension of human behavior within various settings.

C. Multimodal Fusion

The integration of data from diverse sensors and modalities, including accelerometers, gyroscopes, audio, and video, has the potential to enhance the informative richness of the dataset. The integration of these two components has the potential to facilitate the development of comprehensive and precise activity recognition models, hence expanding the present limitations in this field.

D. Longitudinal and Contextual Data

The inclusion of longitudinal data, which tracks the progression of activities over time and contextual information such as environmental conditions and user emotions, has the potential to increase the datasets' temporal and situational values. This technological development has the potential to provide more sophisticated and contextually sensitive identification of human activities.

E. Privacy Preservation

Privacy preservation refers to methods that aim to protect individuals' privacy during data collection and analysis. Evolving HAR datasets have the potential to investigate novel approaches to preserving user privacy, all the while providing important insights. Methods such as differential privacy and secure multiparty computing have the potential to be included in data-collecting protocols.

F. Customization Focused on User Needs

The adaptation of HAR datasets to provide the distinctive actions and preferences of individual users has the potential to generate personalized models that more effectively correspond to users' distinct behavioral patterns, hence improving accuracy and usability.

G. Benchmarking Standards

The establishment of defined standards and assessment criteria for HAR datasets has the potential to facilitate a more uniform and comparative evaluation of various models and algorithms, hence expediting advancements in the area.

H. Collaborative Environment

The growth of open data sharing and collaboration among researchers has the potential to facilitate the generation of extensive and diverse datasets, hence facilitating the advancement of more robust and generalizable models for HAR.

VIII. CONCLUSION AND FUTURE WORK

The emergence of new computing technologies with the ability to collect more detailed and accurate data has led to the generation of huge amounts of data, and its

relevance to understanding and impact on decision-making is increasing. Generated data are typically collected in the form of datasets. Researchers and practitioners use datasets for research objectives to understand the totality of an area of interest and develop a basis for making decisions. The primary objective of constructing a dataset and making it available and open to others is to allow benchmarking, replication, and validation of research approaches, as well as exploration of novel hypotheses.

The main objective of this research study is to identify current issues and challenges faced by researchers in the HAR domain with respect to open datasets. A literature review and survey were conducted to identify these issues and challenges from the published literature and the research community. The identified issues and challenges were classified for ease of understanding and interpretation. This classification of issues and challenges will help HAR researchers to be aware of the open issues and challenges in HAR open datasets. This research has helped to identify and promote important attributes such as privacy, anonymity, platform maintenance, datasets' descriptions and metadata, environmental conditions, resources, and training while constructing and sharing new datasets. In future work, our own datasets will be shared using the recommendations and good practices identified. An evaluation workshop will be conducted involving other HAR researchers to explore the above issues and challenges along with other outcomes of our work in curating open datasets in HAR, including the analysis of datasets to automatically extract metadata and assess dataset quality.

APPENDIX A

The following is the question collection that was asked in the conducted questionnaire survey.

Demographic Information (to Help Us Understand the Type of Respondents to Our Survey):

- 1) Occupation: required.
- 2) Organization/research institute.

Your Experience Using Datasets:

- 3) Approximately how many years of experience do you have in working with open datasets in HAR?
- 4) What is your experience in using open datasets in HAR?

Constructing?

- 5) Have you taken part in constructing an open dataset in HAR?

Dataset Construction:

- 6) The dataset I constructed may not be shared because it contains information that is (select all that apply); if you selected Other, please specify.
- 7) For the research group/organization in which I work, normal practice in relation to experimental datasets is to (please indicate the statement that best describes your situation).
- 8) What are your preferred data formats for dataset construction and sharing?
- 9) Are there other data-sharing formats not on the above list that you use?
- 10) What is the main piece of advice you would give to another researcher when generating a new dataset in HAR?

- 11) What is the main issue or challenge you have faced when generating a new dataset in HAR?

Sharing?

- 12) Have you taken part in sharing open datasets in HAR?
Dataset Sharing:

- 13) The sharing of data was limited because it contained identifying information about an organization or participants.
- 14) Sharing the dataset was not possible because of problems with the data.
- 15) Sharing the dataset was not possible because it was too large.
- 16) I did not share the dataset because I was unsure of the best approach.
- 17) When I shared the dataset, I used the following license.
- 18) I found the process of depositing and sharing the dataset time-consuming.
- 19) As a researcher, what motivates you to share your dataset with the HAR research community?
- 20) What is the main issue or challenge you have faced during dataset sharing?
- 21) What is the main piece of advice you would give to another researcher when sharing a new dataset in HAR?

Finding?

22) Have you searched for and downloaded an open dataset from an open data online repository for experimental/research purposes?

Dataset Finding:

- 23) In searching for a dataset, I make a selection based on the following.
- 24) How much preprocessing did you need to perform on the dataset after downloading?
- 25) In your view, what factors improve dataset quality (one per line)?
- 26) As a researcher, what motivates you to gain access to an existing dataset instead of creating a new dataset?
- 27) What open dataset repositories/directories do you typically search when looking for a dataset?
- 28) What is the main issue or challenge you have faced when trying to find a suitable dataset?
- 29) What is the main piece of advice you would give to another researcher when searching for a dataset in HAR?

Using?

- 30) Have you used and evaluated someone else's open dataset for experimental/research purposes?

Dataset Using:

- 31) Was the metadata describing the dataset easy to understand?
- 32) Was the metadata describing the dataset accurate?
- 33) I encountered a dataset update issue after the initial sharing.
- 34) What is the main issue or challenge you have faced when trying to use someone else's dataset in your research?
- 35) What is the main piece of advice you would give to another researcher when using someone else's dataset in HAR?
- 36) Have you ever preregistered an experiment in any of the following online repository domains?

TABLE XV
PAPER METRICS FOR FIG. 5

Digital Libraries	Paper Title	ML technique
ACM	Classical Machine Learning Approach for Human Activity Recognition Using Location Data	RF
	Emotion Recognition in the Wild from Videos using Images	CNN, SVM
	Emotion Recognition in the Wild using Deep Neural Networks and Bayesian Classifiers	CNN, BN
	Ensemble Approach for Sensor-Based Human Activity Recognition	SVM, KNN, RF
	Face Recognition via Active Annotation and Learning	DNN
	Feature Based Random Forest Nurse Care Activity Recognition Using Accelerometer Data	RF
	From Individual to Group-Level Emotion Recognition: EmotiW 5.0	SVM
	UPIC: User and Position Independent Classical Approach for Locomotion and Transportation Modes Recognition	RF
	Modeling Multimodal Cues in a Deep Learning-Based Framework for Emotion Recognition in the Wild	CNN, LSTM
	Multi-modal Emotion Recognition using Semi-supervised Learning and Multiple Neural Networks in the Wild	CNN, LSTM
	Multi-view Common Space Learning for Emotion Recognition in the Wild	CNN
	Summary of the 2nd Nurse Care Activity Recognition Challenge Using Lab and Field Data	kNN
	HoloNet: Towards Robust Emotion Recognition in the Wild	CNN
	Deep Triplet Networks with Attention for Sensor-based Human Activity Recognition	LSTM
IEEE	A Semisupervised Recurrent Convolutional Attention Model for Human Activity Recognition	CNN, LSTM, SVM
	A Comparative Study on Missing Data Handling Using Machine Learning for Human Activity Recognition	SVN, RF
	Group Activity Description and Recognition based on Trajectory Analysis and Neural Networks	SOM
	Hidden Markov Model-Based Fall Detection With Motion Sensor Orientation Calibration: A Case for Real-Life Home Monitoring	HMM
	Recognition of Real-life Activities with Smartphone Sensors using Deep Learning Approaches	CNN, LSTM
	Transition-Aware Housekeeping Task Monitoring Using Single Wrist-Worn Sensor	SVM, NB, LSTM

- 37) As a researcher, how important to you is the replication of an HAR experiment using an open dataset?

TABLE XV
(Continued.) PAPER METRICS FOR FIG. 5

Springer	Unsupervised Recognition of Multi-Resident Activities in Smart-Homes	HMM
	The joint use of sequence features combination and modified weighted SVM for improving daily activity recognition	PCA, SVM, HMM
	Human Action Prediction with 3D-CNN	CNN, LSTM
	Wrapper Filter Approach for Accelerometer-Based Human Activity Recognition	RF, k-NN, GB
	A revised framework of machine learning application for optimal activity recognition	SVM, MLP
	Feature learning for Human Activity Recognition using Convolutional Neural Networks	RF, CNN
	Multi modal human action recognition for video content matching	SVM
	A Framework for Semi-Supervised Adaptive Learning for Activity Recognition in Healthcare Applications	BN
	Efcacy of Imbalanced Data Handling Methods on Deep Learning for Smart Homes Environments	LSTM, CNN
	Novel approaches to human activity recognition based on accelerometer data	CNN
Science Direct	A resource conscious human action recognition framework using 26-layered deep convolutional neural network	CNN, SVM, k-NN
	Human-Sensing: Low Resolution Thermal Array Sensor Data Classification of Location-Based Postures	J48
	A fall detection method based on a joint motion map using double convolutional neural networks	CNN
	Daily Human Activities Recognition Using Heterogeneous Sensors from Smartphones Cross-subject transfer learning in human activity recognition systems using generative adversarial networks	MLP GAN
	Attention induced multi-head convolutional neural network for human activity recognition	CNN
	A smartphone sensors-based personalized human activity recognition system for sustainable smart cities	DRNN
	Efficiency investigation from shallow to deep neural network techniques in human activity recognition	ANN, CNN
	GCHAR: An efficient Group-based Context-aware human activity recognition on smartphone	k-NN, RF, j48
	Human activity recognition based on smartphone and wearable sensors using multiscale DCNN ensemble	CNN
	Human activity recognition with smartphone sensors using deep learning neural networks	ANN, SVM, MLP, J48, NB

- 38) As a researcher, how important to you is the benchmarking of an HAR approach using an open dataset?

TABLE XV
(Continued.) PAPER METRICS FOR FIG. 5

Investigating the significance of adversarial attacks and their relation to interpretability for radar-based human activity recognition systems	CNN
Robust least squares twin support vector machine for human activity recognition	SVM
Robust Human Activity Recognition using smartwatches and smartphones	RF, CNN, HMM, MLP, LSTM
Online active learning for human activity recognition from sensory data streams	DT, SVM
Multi-label classification based ensemble learning for human activity recognition in smart home	NB, DT, k-NN
Modified deep residual network architecture deployed on serverless framework of IoT platform based on human activity recognition application	SVM, CNN, LSTM

- 39) In your view, what are the major issues and challenges in relation to using open datasets in HAR research?
 40) Have you any other comments regarding the survey?

APPENDIX B

See Table XV.

REFERENCES

- [1] S. Gupta, “Deep learning based human activity recognition (HAR) using wearable sensor data,” *Int. J. Inf. Manage. Data Insights*, vol. 1, no. 2, Nov. 2021, Art. no. 100046, doi: [10.1016/j.jjime.2021.100046](https://doi.org/10.1016/j.jjime.2021.100046).
- [2] A. Haleem, M. Javaid, M. A. Qadri, R. P. Singh, and R. Suman, “Artificial intelligence (AI) applications for marketing: A literature-based study,” *Int. J. Intell. Netw.*, vol. 3, pp. 119–132, Jan. 2022, doi: [10.1016/j.ijin.2022.08.005](https://doi.org/10.1016/j.ijin.2022.08.005).
- [3] D. Sheth and M. Shah, “Predicting stock market using machine learning: Best and accurate way to know future stock prices,” *Int. J. Syst. Assurance Eng. Manage.*, vol. 14, no. 1, pp. 1–18, Feb. 2023, doi: [10.1007/s13198-022-01811-1](https://doi.org/10.1007/s13198-022-01811-1).
- [4] Y.-T. Chang and N.-H. Fan, “A novel approach to market segmentation selection using artificial intelligence techniques,” *J. Supercomput.*, vol. 79, no. 2, pp. 1235–1262, Feb. 2023, doi: [10.1007/s11227-022-04666-2](https://doi.org/10.1007/s11227-022-04666-2).
- [5] S. P. Arnerić, V. D. Kern, and D. T. Stephenson, “Regulatory-accepted drug development tools are needed to accelerate innovative CNS disease treatments,” *Biochem. Pharmacol.*, vol. 151, pp. 291–306, May 2018, doi: [10.1016/j.bcp.2018.01.043](https://doi.org/10.1016/j.bcp.2018.01.043).
- [6] H. Liu, W. Zhang, B. Zou, J. Wang, Y. Deng, and L. Deng, “DrugCombDB: A comprehensive database of drug combinations toward the discovery of combinatorial therapy,” *Nucleic Acids Res.*, vol. 48, no. D1, pp. D871–D881, Jan. 2020, Oct. 2019, doi: [10.1093/nar/gkz1007](https://doi.org/10.1093/nar/gkz1007).
- [7] M. Sun and J. Zhang, “Research on the application of block chain big data platform in the construction of new smart city for low carbon emission and green environment,” *Comput. Commun.*, vol. 149, pp. 332–342, Jan. 2020, doi: [10.1016/j.comcom.2019.10.031](https://doi.org/10.1016/j.comcom.2019.10.031).
- [8] N. S. Suhami, J. Mountstephens, and J. Teo, “A dataset for emotion recognition using virtual reality and EEG (DER-VREEG): Emotional state classification using low-cost wearable VR-EEG headsets,” *Big Data Cognit. Comput.*, vol. 6, no. 1, p. 16, Jan. 2022, doi: [10.3390/bdcc6010016](https://doi.org/10.3390/bdcc6010016).

- [9] C. Y. Wang, Q. Zhou, G. Fitzmaurice, and F. Anderson, "VideoPoseVR: Authoring virtual reality character animations with online videos," *Proc. ACM Hum.-Comput. Interact.*, vol. 6, pp. 575:448–575:467, Nov. 2022, doi: [10.1145/3567728](https://doi.org/10.1145/3567728).
- [10] G. Tsueng et al., "Developing a standardized but extendable framework to increase the findability of infectious disease datasets," *Sci. Data*, vol. 10, no. 1, p. 99, Feb. 2023, doi: [10.1038/s41597-023-01968-9](https://doi.org/10.1038/s41597-023-01968-9).
- [11] R. Singh, A. Sonawane, and R. Srivastava, "Recent evolution of modern datasets for human activity recognition: A deep survey," *Multimedia Syst.*, vol. 26, no. 2, pp. 83–106, Apr. 2020, doi: [10.1007/s00530-019-00635-7](https://doi.org/10.1007/s00530-019-00635-7).
- [12] A. Das Antar, M. Ahmed, and M. A. R. Ahad, "Challenges in sensor-based human activity recognition and a comparative analysis of benchmark datasets: A review," in *Proc. Joint 8th Int. Conf. Informat. Electron. Vis. (ICIEV) 3rd Int. Conf. Imag., Vis. Pattern Recognit. (icIVPR)*, May 2019, pp. 134–139.
- [13] H. A. Piwowar and W. W. Chapman, "Public sharing of research datasets: A pilot study of associations," *J. Informetrics*, vol. 4, no. 2, pp. 148–156, Apr. 2010, doi: [10.1016/j.joi.2009.11.010](https://doi.org/10.1016/j.joi.2009.11.010).
- [14] A. Ambhaikar, "A survey on health care and expert system," *Math. Stat. Eng. Appl.*, vol. 72, no. 1, p. 1, Jan. 2023.
- [15] Y. Li, G. Yang, Z. Su, S. Li, and Y. Wang, "Human activity recognition based on multienvironment sensor data," *Inf. Fusion*, vol. 91, pp. 47–63, Mar. 2023, doi: [10.1016/j.inffus.2022.10.015](https://doi.org/10.1016/j.inffus.2022.10.015).
- [16] C. Nugent et al., "An initiative for the creation of open datasets within pervasive healthcare," in *Proc. EAI Int. Conf. Pervasive Comput. Technol. Healthcare*, May 2016, pp. 318–321. Accessed: Nov. 20, 2021. [Online]. Available: <http://urn.kb.se/resolve?urn=urn:nbn:se:ltu:diva-64771>
- [17] Z. Zhang, W. Wang, A. An, Y. Qin, and F. Yang, "A human activity recognition method using wearable sensors based on convtransformer model," *Evolving Syst.*, pp. 1–17, Jan. 2023.
- [18] N. Davies and S. Clinch, "Pervasive data science," *IEEE Pervasive Comput.*, vol. 16, no. 3, pp. 50–58, Jul. 2017, doi: [10.1109/MPRV.2017.2940956](https://doi.org/10.1109/MPRV.2017.2940956).
- [19] A. Bexheti, M. Langheinrich, and S. Clinch, "Secure personal memory-sharing with co-located people and places," in *Proc. 6th Int. Conf. Internet Things*, Nov. 2016, pp. 73–81, doi: [10.1145/2991561.2991577](https://doi.org/10.1145/2991561.2991577).
- [20] P. Kumar and S. Suresh, "Deep-HAR: An ensemble deep learning model for recognizing the simple, complex, and heterogeneous human activities," *Multimedia Tools Appl.*, vol. 82, no. 20, pp. 30435–30462, Feb. 2023, doi: [10.1007/s11042-023-14492-0](https://doi.org/10.1007/s11042-023-14492-0).
- [21] Y. Chen, Y. Gu, X. Jiang, and J. Wang, "OCEAN: A new opportunistic computing model for wearable activity recognition," ACM in *Proc. Int. Joint Conf. Pervasive Ubiquitous Comput., Adjunct (UbiComp)*. New York, NY, USA: Association for Computing Machinery, Sep. 2016, pp. 33–36, doi: [10.1145/2968219.2971453](https://doi.org/10.1145/2968219.2971453).
- [22] C.-Y. Huang et al., "Flexible pressure sensor with an excellent linear response in a broad detection range for human motion monitoring," *ACS Appl. Mater. Interface*, vol. 15, no. 2, pp. 3476–3485, Jan. 2023, doi: [10.1021/acsami.2c19465](https://doi.org/10.1021/acsami.2c19465).
- [23] M. G. Morshed, T. Sultana, A. Alam, and Y.-K. Lee, "Human action recognition: A taxonomy-based survey, updates, and opportunities," *Sensors*, vol. 23, no. 4, p. 2182, Feb. 2023, doi: [10.3390/s23042182](https://doi.org/10.3390/s23042182).
- [24] H. F. Nweke, Y. W. Teh, M. A. Al-garadi, and U. R. Alo, "Deep learning algorithms for human activity recognition using mobile and wearable sensor networks: State of the art and research challenges," *Expert Syst. Appl.*, vol. 105, pp. 233–261, Sep. 2018, doi: [10.1016/j.eswa.2018.03.056](https://doi.org/10.1016/j.eswa.2018.03.056).
- [25] G. Alam, I. McChesney, P. Nicholl, and J. Rafferty, "An approach to extract and compare metadata of human activity recognition (HAR) data sets," in *Proc. Int. Conf. Ubiquitous Comput. Ambient Intell. (UCAmI)*, J. Bravo, S. Ochoa, and J. Favela, Eds. Cham, Switzerland: Springer, 2023, pp. 717–728, doi: [10.1007/978-3-031-21333-5_71](https://doi.org/10.1007/978-3-031-21333-5_71).
- [26] V. Sharma, M. Gupta, A. K. Pandey, D. Mishra, and A. Kumar, "A review of deep learning-based human activity recognition on benchmark video datasets," *Appl. Artif. Intell.*, vol. 36, no. 1, Dec. 2022, Art. no. 2093705, doi: [10.1080/08839514.2022.2093705](https://doi.org/10.1080/08839514.2022.2093705).
- [27] J. M. Chaquet, E. J. Carmona, and A. Fernández-Caballero, "A survey of video datasets for human action and activity recognition," *Comput. Vis. Image Understand.*, vol. 117, no. 6, pp. 633–659, Jun. 2013, doi: [10.1016/j.cviu.2013.01.013](https://doi.org/10.1016/j.cviu.2013.01.013).
- [28] D. R. Beddar, B. Nini, M. Sabokrou, and A. Hadid, "Vision-based human activity recognition: A survey," *Multimedia Tools Appl.*, vol. 79, nos. 41–42, pp. 30509–30555, Nov. 2020, doi: [10.1007/s11042-020-09004-3](https://doi.org/10.1007/s11042-020-09004-3).
- [29] S. Vishwakarma and A. Agrawal, "A survey on activity recognition and behavior understanding in video surveillance," *Vts. Comput.*, vol. 29, no. 10, pp. 983–1009, Oct. 2013, doi: [10.1007/s00371-012-0752-6](https://doi.org/10.1007/s00371-012-0752-6).
- [30] L. Onofri, P. Soda, M. Pechenizkiy, and G. Iannello, "A survey on using domain and contextual knowledge for human activity recognition in video streams," *Expert Syst. Appl.*, vol. 63, pp. 97–111, Nov. 2016, doi: [10.1016/j.eswa.2016.06.011](https://doi.org/10.1016/j.eswa.2016.06.011).
- [31] M. H. Arshad, M. Bilal, and A. Gani, "Human activity recognition: Review, taxonomy and open challenges," *Sensors*, vol. 22, no. 17, p. 6463, Aug. 2022, doi: [10.3390/s22176463](https://doi.org/10.3390/s22176463).
- [32] A. Lentzas and D. Vrakas, "Non-intrusive human activity recognition and abnormal behavior detection on elderly people: A review," *Artif. Intell. Rev.*, vol. 53, no. 3, pp. 1975–2021, Mar. 2020, doi: [10.1007/s10462-019-09724-5](https://doi.org/10.1007/s10462-019-09724-5).
- [33] L. M. Dang, K. Min, H. Wang, M. J. Piran, C. H. Lee, and H. Moon, "Sensor-based and vision-based human activity recognition: A comprehensive survey," *Pattern Recognit.*, vol. 108, Dec. 2020, Art. no. 107561, doi: [10.1016/j.patcog.2020.107561](https://doi.org/10.1016/j.patcog.2020.107561).
- [34] A. Ferrari, D. Micucci, M. Mobilio, and P. Napoletano, "Trends in human activity recognition using smartphones," *J. Reliable Intell. Environ.*, vol. 7, no. 3, pp. 189–213, Sep. 2021, doi: [10.1007/s40860-021-00147-0](https://doi.org/10.1007/s40860-021-00147-0).
- [35] S. Zolfaghari, M. R. Keyvanpour, and R. Zall, "Analytical review on ontological human activity recognition approaches," *Int. J. E-Bus. Res.*, vol. 13, no. 2, pp. 58–78, Apr. 2017, doi: [10.4018/IJEBR.2017040104](https://doi.org/10.4018/IJEBR.2017040104).
- [36] I. H. Sarker, "Data science and analytics: An overview from data-driven smart computing, decision-making and applications perspective," *Social Netw. Comput. Sci.*, vol. 2, no. 5, p. 377, Jul. 2021, doi: [10.1007/s42979-021-00765-8](https://doi.org/10.1007/s42979-021-00765-8).
- [37] Y. Wang et al., "A novel deep multifeature extraction framework based on attention mechanism using wearable sensor data for human activity recognition," *IEEE Sensors J.*, vol. 23, no. 7, pp. 7188–7198, Apr. 2023, doi: [10.1109/JSEN.2023.3242603](https://doi.org/10.1109/JSEN.2023.3242603).
- [38] D. Helbing et al. (Nov. 26, 2021). *Will Democracy Survive Big Data and Artificial Intelligence*. Scientific American. Accessed: Nov. 26, 2021. [Online]. Available: <https://www.scientificamerican.com/article/will-democracy-survive-big-data-and-artificial-intelligence/>
- [39] Open Government Data OECD. Accessed: Nov. 26, 2021. [Online]. Available: <https://www.oecd.org/gov/digital-government/open-government-data.htm>
- [40] What is Open Data. Accessed: Nov. 26, 2021. [Online]. Available: <https://opendatahandbook.org/guide/en/what-is-open-data/>
- [41] P. Andanda, "Towards a paradigm shift in governing data access and related intellectual property rights in big data and health-related research," *IIC Int. Rev. Intellectual Property Competition Law*, vol. 50, no. 9, pp. 1052–1081, Nov. 2019, doi: [10.1007/s40319-019-00873-2](https://doi.org/10.1007/s40319-019-00873-2).
- [42] X. Zhu, C. Thomas, J. C. Moore, and S. Allen, "Open government data licensing: An analysis of the U.S. state open government data portals," in *Diversity, Divergence, Dialogue (Lecture Notes in Computer Science)*, K. Toeppe, H. Yan, and S. K. W. Chu, Eds. Cham, Switzerland: Springer, 2021, pp. 260–273, doi: [10.1007/978-3-030-71305-8_21](https://doi.org/10.1007/978-3-030-71305-8_21).
- [43] C. Pernet, C. Svarer, R. Blair, J. D. Van Horn, and R. A. Poldrack, "On the long-term archiving of research data," *Neuroinformatics*, vol. 21, no. 2, pp. 243–246, Feb. 2023, doi: [10.1007/s12021-023-09621-x](https://doi.org/10.1007/s12021-023-09621-x).
- [44] (Jan. 10, 2019). *These are the Best Free Open Data Sources Anyone Can Use*. freeCodeCamp.org. Accessed: Nov. 19, 2021. [Online]. Available: <https://www.freecodecamp.org/news/https-medium-freecodecamp-org-best-free-open-data-sources-anyone-can-use-a65b514b0f2d/>
- [45] T. Davidson, E. Wall, and J. Mace, "A qualitative interview study of distributed tracing visualisation: A characterisation of challenges and opportunities," *IEEE Trans. Vis. Comput. Graphics*, early access, Feb. 1, 2023, doi: [10.1109/TVCG.2023.3241596](https://doi.org/10.1109/TVCG.2023.3241596).
- [46] D. Cook, K. D. Feuz, and N. C. Krishnan, "Transfer learning for activity recognition: A survey," *Knowl. Inf. Syst.*, vol. 36, no. 3, pp. 537–556, Sep. 2013, doi: [10.1007/s10115-013-0665-3](https://doi.org/10.1007/s10115-013-0665-3).
- [47] P. Vepakomma, D. De, S. K. Das, and S. Bhansali, "A-wristocracy: Deep learning on Wrist-Worn sensing for recognition of user complex activities," in *Proc. IEEE 12th Int. Conf. Wearable Implant. Body Sens. Netw. (BSN)*, Jun. 2015, pp. 1–6, doi: [10.1109/BSN.2015.7299406](https://doi.org/10.1109/BSN.2015.7299406).
- [48] Y. Kim and Y. Li, "Human activity classification with transmission and reflection coefficients of on-body antennas through deep convolutional neural networks," *IEEE Trans. Antennas Propag.*, vol. 65, no. 5, pp. 2764–2768, May 2017, doi: [10.1109/TAP.2017.2677918](https://doi.org/10.1109/TAP.2017.2677918).

- [49] N. Y. Hammerla, S. Halloran, and T. Ploetz, "Deep, convolutional, and recurrent models for human activity recognition using wearables," 2016, *arXiv:1604.08880*.
- [50] J. Qin, L. Liu, Z. Zhang, Y. Wang, and L. Shao, "Compressive sequential learning for action similarity labeling," *IEEE Trans. Image Process.*, vol. 25, no. 2, pp. 756–769, Feb. 2016, doi: [10.1109/TIP.2015.2508600](https://doi.org/10.1109/TIP.2015.2508600).
- [51] Y.-H. Kim et al., "MyMove: Facilitating older adults to collect in-situ activity labels on a smartwatch with speech," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, Apr. 2022, pp. 1–21, doi: [10.1145/3491102.3517457](https://doi.org/10.1145/3491102.3517457).
- [52] M. T. Leving, H. L. D. Horemans, R. J. K. Vegter, S. de Groot, J. B. J. Bussmann, and L. H. V. van der Woude, "Validity of consumer-grade activity monitor to identify manual wheelchair propulsion in standardized activities of daily living," *PLoS ONE*, vol. 13, no. 4, Apr. 2018, Art. no. e0194864, doi: [10.1371/journal.pone.0194864](https://doi.org/10.1371/journal.pone.0194864).
- [53] X. Yu, J. Jang, and S. Xiong, "A large-scale open motion dataset (KFall) and benchmark algorithms for detecting pre-impact fall of the elderly using wearable inertial sensors," *Frontiers Aging Neurosci.*, vol. 13, Jul. 2021, Art. no. 692865, doi: [10.3389/fnagi.2021.692865](https://doi.org/10.3389/fnagi.2021.692865).
- [54] J. C. E. Guerrero, E. M. Espa  a, M. M. Af  asco, and J. E. P. Lopera, "Dataset for human fall recognition in an uncontrolled environment," *Data Brief*, vol. 45, Dec. 2022, Art. no. 108610, doi: [10.1016/j.dib.2022.108610](https://doi.org/10.1016/j.dib.2022.108610).
- [55] A. Sucerquia, J. L  pez, and J. Vargas-Bonilla, "SisFall: A fall and movement dataset," *Sensors*, vol. 17, no. 12, p. 198, Jan. 2017, doi: [10.3390/s17010198](https://doi.org/10.3390/s17010198).
- [56] F. C. Heilbron, V. Escorcia, B. Ghanem, and J. C. Niebles, *ActivityNet: A Large-Scale Video Benchmark for Human Activity Understanding*. Accessed: Nov. 20, 2021. [Online]. Available: https://openaccess.thecvf.com/content_cvpr_2015/html/Heilbron_ActivityNet_A_Large-Scale_2015_CVPR_paper.html
- [57] H. Zheng, D. Liu, and Y. Liu, "RETRACTED: Design and research on automatic recognition system of sports dance movement based on computer vision and parallel computing," *Microprocessors Microsyst.*, vol. 80, Feb. 2021, Art. no. 103648, doi: [10.1016/j.micpro.2020.103648](https://doi.org/10.1016/j.micpro.2020.103648).
- [58] Y. Liu, L. Nie, L. Liu, and D. S. Rosenblum, "From action to activity: Sensor-based activity recognition," *Neurocomputing*, vol. 181, pp. 108–115, Mar. 2016, doi: [10.1016/j.neucom.2015.08.096](https://doi.org/10.1016/j.neucom.2015.08.096).
- [59] X. Dai, B. Singh, G. Zhang, L. S. Davis, and Y. Q. Chen, "Temporal context network for activity localization in videos," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 5727–5736.
- [60] D. Wu, N. Sharma, and M. Blumenstein, "Recent advances in video-based human action recognition using deep learning: A review," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, May 2017, pp. 2865–2872, doi: [10.1109/IJCNN.2017.7966210](https://doi.org/10.1109/IJCNN.2017.7966210).
- [61] T. Hassner, "A critical review of action recognition benchmarks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2013, pp. 245–250.
- [62] N. Ikizler-Cinbis, R. G. Cinbis, and S. Sclaroff, "Learning actions from the Web," in *Proc. IEEE 12th Int. Conf. Comput. Vis.*, Sep. 2009, pp. 995–1002, doi: [10.1109/ICCV.2009.5459368](https://doi.org/10.1109/ICCV.2009.5459368).
- [63] A. Torralba and A. A. Efros, "Unbiased look at dataset bias," in *Proc. CVPR*, Jun. 2011, pp. 1521–1528, doi: [10.1109/CVPR.2011.5995347](https://doi.org/10.1109/CVPR.2011.5995347).
- [64] UCI Machine Learning Repository: *Human Activity Recognition Using Smartphones Data Set*. Accessed: Nov. 20, 2021. [Online]. Available: <http://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>
- [65] Harvard Dataverse. Accessed: Nov. 20, 2021. [Online]. Available: <https://dataverse.harvard.edu/>
- [66] Dataset Search. Accessed: Nov. 20, 2021. [Online]. Available: <https://datasetsearch.research.google.com/>
- [67] IEEE DataPort. IEEE DataPort. Accessed: Nov. 20, 2021. [Online]. Available: <https://ieee-dataport.org/> (accessed Nov. 20, 2021).
- [68] Zenodo Research Shared. Accessed: Nov. 20, 2021. [Online]. Available: <https://zenodo.org/>
- [69] Figshare Credit for All Your Research. Accessed: Nov. 20, 2021. [Online]. Available: <https://figshare.com/>
- [70] Hollywood Data on Data.World | 10 Datasets. data.world. Accessed: Nov. 26, 2021. [Online]. Available: <https://data.world/datasets/hollywood>
- [71] Action Similarity Labeling Challenge. Accessed: Nov. 27, 2021. [Online]. Available: <https://talhassner.github.io/home/projects/ASLAN/ASLAN-main.html>
- [72] Papers With Code YouCook Dataset. Accessed: Nov. 27, 2021. [Online]. Available: <https://paperswithcode.com/dataset/youcook>
- [73] Welcome to CASAS. Accessed: Nov. 27, 2021. [Online]. Available: <http://casas.wsu.edu/datasets/>
- [74] UCI Machine Learning Repository: OPPORTUNITY Activity Recognition Data Set. Accessed: Nov. 27, 2021. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/opportunity+activity+recognition>
- [75] A. Dhall, R. Goecke, S. Ghosh, J. Joshi, J. Hoey, and T. Gedeon, "From individual to group-level emotion recognition: EmotiW 5.0," in *Proc. 19th ACM Int. Conf. Multimodal Interact.* New York, NY, USA: Association for Computing Machinery, Nov. 2017, pp. 524–528, doi: [10.1145/3136755.3143004](https://doi.org/10.1145/3136755.3143004).
- [76] UCI Machine Learning Repository: WISDM Smartphone and Smartwatch Activity and Biometrics Dataset Data Set. Accessed: Aug. 5, 2022. [Online]. Available: <http://archive.ics.uci.edu/ml/datasets/WISDM+Smartphone+and+Smartwatch+Activity+and+Biometrics+Dataset+>
- [77] UCI Machine Learning Repository: PAMAP2 Physical Activity Monitoring Data Set. Accessed: Aug. 5, 2022. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/pamap2+physical+activity+monitoring>
- [78] Sussex-Huawei Locomotion Dataset. Accessed: Aug. 5, 2022. [Online]. Available: <http://www.shl-dataset.org/>
- [79] F. Mastrogiovanni. (Dec. 11, 2014). Fulviomas/WHARF. Accessed: Aug. 5, 2022. [Online]. Available: <https://github.com/fulviomas/WHARF>
- [80] M. Soliman, T. Fatnassi, I. Elgammal, and R. Figueiredo, "Exploring the major trends and emerging themes of artificial intelligence in the scientific leading journals amidst the COVID-19 era," *Big Data Cognit. Comput.*, vol. 7, no. 1, p. 12, Jan. 2023, doi: [10.3390/bdcc7010012](https://doi.org/10.3390/bdcc7010012).
- [81] J. Cui, Z. Wang, S.-B. Ho, and E. Cambria, "Survey on sentiment analysis: Evolution of research methods and topics," *Artif. Intell. Rev.*, vol. 56, no. 8, pp. 8469–8510, Jan. 2023, doi: [10.1007/s10462-022-10386-z](https://doi.org/10.1007/s10462-022-10386-z).
- [82] J. Osterrieder, *A Primer on Artificial Intelligence and Machine Learning for the Financial Services Industry*. Rochester, NY, USA: SSRN, Feb. 2023, doi: [10.2139/ssrn.4349078](https://doi.org/10.2139/ssrn.4349078).
- [83] M. Chhabra, K. K. Ravulakollu, M. Kumar, A. Sharma, and A. Nayyar, "Improving automated latent fingerprint detection and segmentation using deep convolutional neural network," *Neural Comput. Appl.*, vol. 35, no. 9, pp. 6471–6497, Mar. 2023, doi: [10.1007/s00521-022-07894-y](https://doi.org/10.1007/s00521-022-07894-y).
- [84] D. Garcia-Gonzalez, D. Rivero, E. Fernandez-Blanco, and M. R. Luaces, "New machine learning approaches for real-life human activity recognition using smartphone sensor-based data," *Knowl.-Based Syst.*, vol. 262, Feb. 2023, Art. no. 110260, doi: [10.1016/j.knosys.2023.110260](https://doi.org/10.1016/j.knosys.2023.110260).
- [85] L. Bai, H. Li, W. Gao, J. Xie, and H. Wang, "A joint multiobjective optimization of feature selection and classifier design for high-dimensional data classification," *Inf. Sci.*, vol. 626, pp. 457–473, May 2023, doi: [10.1016/j.ins.2023.01.069](https://doi.org/10.1016/j.ins.2023.01.069).
- [86] S. Buyruko  lu and S. Sava  , "Stacked-based ensemble machine learning model for positioning footballer," *Arabian J. Sci. Eng.*, vol. 48, no. 2, pp. 1371–1383, Feb. 2023, doi: [10.1007/s13369-022-06857-8](https://doi.org/10.1007/s13369-022-06857-8).
- [87] H. H. Ali, H. M. Moftah, and A. A. A. Youssif, "Depth-based human activity recognition: A comparative perspective study on feature extraction," *Future Comput. Informat. J.*, vol. 3, no. 1, pp. 51–67, Jun. 2018, doi: [10.1016/j.fcij.2017.11.002](https://doi.org/10.1016/j.fcij.2017.11.002).
- [88] A. Ray, M. H. Kolekar, R. Balasubramanian, and A. Hafiane, "Transfer learning enhanced vision-based human activity recognition: A decade-long analysis," *Int. J. Inf. Manage. Data Insights*, vol. 3, no. 1, Apr. 2023, Art. no. 100142, doi: [10.1016/j.ijimi.2022.100142](https://doi.org/10.1016/j.ijimi.2022.100142).
- [89] K. Henricksen, J. Indulska, and A. Rakotonirainy, "Modeling context information in pervasive computing systems," in *Pervasive Computing* (Lecture Notes in Computer Science), F. Mattern and M. Naghshineh, Eds. Berlin, Germany: Springer, 2002, pp. 167–180, doi: [10.1007/3-540-45866-2_14](https://doi.org/10.1007/3-540-45866-2_14).
- [90] B.   umak, S. Brdnik, and M. Pu  nik, "Sensors and artificial intelligence methods and algorithms for human-computer intelligent interaction: A systematic mapping study," *Sensors*, vol. 22, no. 1, p. 20, Dec. 2021, doi: [10.3390/s22010020](https://doi.org/10.3390/s22010020).
- [91] Y. Wang, S. Cang, and H. Yu, "A survey on wearable sensor modality centred human activity recognition in health care," *Expert Syst. Appl.*, vol. 137, pp. 167–190, Dec. 2019, doi: [10.1016/j.eswa.2019.04.057](https://doi.org/10.1016/j.eswa.2019.04.057).

- [92] S. K. Challa, A. Kumar, and V. B. Semwal, "A multibranch CNN-BiLSTM model for human activity recognition using wearable sensor data," *Vis. Comput.*, vol. 38, no. 12, pp. 4095–4109, Dec. 2022, doi: [10.1007/s00371-021-02283-3](https://doi.org/10.1007/s00371-021-02283-3).
- [93] Q. Shen, H. Feng, R. Song, D. Song, and H. Xu, "Federated meta-learning with attention for diversity-aware human activity recognition," *Sensors*, vol. 23, no. 3, p. 1083, Jan. 2023, doi: [10.3390/s23031083](https://doi.org/10.3390/s23031083).
- [94] J. C. Stamper et al., "Managing the educational dataset lifecycle with DataShop," in *Artificial Intelligence in Education* (Lecture Notes in Computer Science), G. Biswas, S. Bull, J. Kay, and A. Mitrovic, Eds. Berlin, Germany: Springer, 2011, pp. 557–559, doi: [10.1007/978-3-642-21869-9_100](https://doi.org/10.1007/978-3-642-21869-9_100).
- [95] A. Paullada, I. D. Raji, E. M. Bender, E. Denton, and A. Hanna, "Data and its (dis)contents: A survey of dataset development and use in machine learning research," *Patterns*, vol. 2, no. 11, Nov. 2021, Art. no. 100336, doi: [10.1016/j.patter.2021.100336](https://doi.org/10.1016/j.patter.2021.100336).
- [96] K. Kelley, "Good practice in the conduct and reporting of survey research," *Int. J. Quality Health Care*, vol. 15, no. 3, pp. 261–266, May 2003, doi: [10.1093/intqhc/mzg031](https://doi.org/10.1093/intqhc/mzg031).
- [97] (Jun. 8, 2018). *Open-Ended Questions: How to Code & Analyze for Insights [2018]*. Thematic. Accessed: Aug. 9, 2022. [Online]. Available: <https://getthematic.com/insights/code-open-ended-questions-in-surveys-to-get-deep-insights>
- [98] A. Castleberry and A. Nolen, "Thematic analysis of qualitative research data: Is it as easy as it sounds?" *Currents Pharmacy Teaching Learn.*, vol. 10, no. 6, pp. 807–815, Jun. 2018, doi: [10.1016/j.cptl.2018.03.019](https://doi.org/10.1016/j.cptl.2018.03.019).
- [99] T. Singh and D. K. Vishwakarma, "Video benchmarks of human action datasets: A review," *Artif. Intell. Rev.*, vol. 52, no. 2, pp. 1107–1154, Aug. 2019, doi: [10.1007/s10462-018-9651-1](https://doi.org/10.1007/s10462-018-9651-1).
- [100] F. Cruciani, I. Cleland, C. Nugent, P. McCullagh, K. Synnes, and J. Hallberg, "Automatic annotation for human activity recognition in free living using a smartphone," *Sensors*, vol. 18, no. 7, p. 2203, Jul. 2018, doi: [10.3390/s18072203](https://doi.org/10.3390/s18072203).
- [101] F. Alharbi, L. Ouarbya, and J. A. Ward, "Comparing sampling strategies for tackling imbalanced data in human activity recognition," *Sensors*, vol. 22, no. 4, p. 1373, Feb. 2022.
- [102] S. A. Bargal, E. Barsoum, C. C. Ferrer, and C. Zhang, "Emotion recognition in the wild from videos using images," in *Proc. 18th ACM Int. Conf. Multimodal Interact.*, Oct. 2016, pp. 433–436, doi: [10.1145/2993148.2997627](https://doi.org/10.1145/2993148.2997627).
- [103] L. Surace, M. Patacchiola, E. B. Sönmez, W. Spataro, and A. Cangelosi, "Emotion recognition in the wild using deep neural networks and Bayesian classifiers," in *Proc. 19th ACM Int. Conf. Multimodal Interact.*, Nov. 2017, pp. 593–597, doi: [10.1145/3136755.3143015](https://doi.org/10.1145/3136755.3143015).
- [104] S. Brajesh and I. Ray, "Ensemble approach for sensor-based human activity recognition," in *Proc. Adjunct ACM Int. Joint Conf. Pervasive Ubiquitous Comput. ACM Int. Symp. Wearable Comput.*, Sep. 2020, pp. 296–300, doi: [10.1145/3410530.3414352](https://doi.org/10.1145/3410530.3414352).
- [105] D. H. Kim, M. K. Lee, D. Y. Choi, and B. C. Song, "Multimodal emotion recognition using semi-supervised learning and multiple neural networks in the wild," in *Proc. 19th ACM Int. Conf. Multimodal Interact.*, Nov. 2017, pp. 529–535, doi: [10.1145/3136755.3143005](https://doi.org/10.1145/3136755.3143005).
- [106] A. Yao, D. Cai, P. Hu, S. Wang, L. Sha, and Y. Chen, "HoloNet: Towards robust emotion recognition in the wild," in *Proc. 18th ACM Int. Conf. Multimodal Interact.*, Oct. 2016, pp. 472–478, doi: [10.1145/2993148.2997639](https://doi.org/10.1145/2993148.2997639).
- [107] B. Khaertdinov, E. Ghaleb, and S. Asteriadis, "Deep triplet networks with attention for sensor-based human activity recognition," in *Proc. IEEE Int. Conf. Pervasive Comput. Commun. (PerCom)*, Mar. 2021, pp. 1–10, doi: [10.1109/PERCOM50583.2021.9439116](https://doi.org/10.1109/PERCOM50583.2021.9439116).
- [108] R. Mohamed, T. Perumal, M. N. Sulaiman, N. Mustapha, and M. N. Razali, "Conflict resolution using enhanced label combination method for complex activity recognition in smart home environment," in *Proc. IEEE 6th Global Conf. Consum. Electron. (GCCE)*, Oct. 2017, pp. 1–3, doi: [10.1109/GCCE.2017.8229477](https://doi.org/10.1109/GCCE.2017.8229477).
- [109] K. Chen, L. Yao, D. Zhang, X. Wang, X. Chang, and F. Nie, "A semisupervised recurrent convolutional attention model for human activity recognition," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 5, pp. 1747–1756, May 2020.
- [110] J. Azorin-Lopez, M. Saval-Calvo, A. Fuster-Guiló, J. García-Rodríguez, M. Cañizola, and M. T. Signes-Pont, "Group activity description and recognition based on trajectory analysis and neural networks," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jul. 2016, pp. 1585–1592.
- [111] H. Bi, M. Perello-Nieto, R. Santos-Rodriguez, and P. Flach, "Human activity recognition based on dynamic active learning," *IEEE J. Biomed. Health Informat.*, vol. 25, no. 4, pp. 922–934, Apr. 2021.
- [112] K.-C. Liu, C.-Y. Hsieh, and C.-T. Chan, "Transition-aware housekeeping task monitoring using single Wrist-Worn sensor," *IEEE Sensors J.*, vol. 18, no. 21, pp. 8950–8962, Nov. 2018, doi: [10.1109/JSEN.2018.2868278](https://doi.org/10.1109/JSEN.2018.2868278).
- [113] D. Riboni and F. Murru, "Unsupervised recognition of multi-resident activities in smart-homes," *IEEE Access*, vol. 8, pp. 201985–201994, 2020, doi: [10.1109/ACCESS.2020.3036226](https://doi.org/10.1109/ACCESS.2020.3036226).
- [114] S. Jha, M. Schiemer, F. Zambonelli, and J. Ye, "Continual learning in sensor-based human activity recognition: An empirical benchmark analysis," *Inf. Sci.*, vol. 575, pp. 1–21, Oct. 2021, doi: [10.1016/j.ins.2021.04.062](https://doi.org/10.1016/j.ins.2021.04.062).
- [115] Z. N. Khan and J. Ahmad, "Attention induced multi-head convolutional neural network for human activity recognition," *Appl. Soft Comput.*, vol. 110, Oct. 2021, Art. no. 107671, doi: [10.1016/j.asoc.2021.107671](https://doi.org/10.1016/j.asoc.2021.107671).
- [116] J. Sena, J. Barreto, C. Caetano, G. Cramer, and W. R. Schwartz, "Human activity recognition based on smartphone and wearable sensors using multiscale DCNN ensemble," *Neurocomputing*, vol. 444, pp. 226–243, Jul. 2021, doi: [10.1016/j.neucom.2020.04.151](https://doi.org/10.1016/j.neucom.2020.04.151).
- [117] C. A. Ronao and S.-B. Cho, "Human activity recognition with smartphone sensors using deep learning neural networks," *Expert Syst. Appl.*, vol. 59, pp. 235–244, Oct. 2016, doi: [10.1016/j.eswa.2016.04.032](https://doi.org/10.1016/j.eswa.2016.04.032).
- [118] M. Jethanandani, A. Sharma, T. Perumal, and J.-R. Chang, "Multi-label classification based ensemble learning for human activity recognition in smart home," *Internet Things*, vol. 12, Dec. 2020, Art. no. 100324, doi: [10.1016/j.iot.2020.100324](https://doi.org/10.1016/j.iot.2020.100324).
- [119] B. M. Abidine, L. Fergani, B. Fergani, and M. Oussalah, "The joint use of sequence features combination and modified weighted SVM for improving daily activity recognition," *Pattern Anal. Appl.*, vol. 21, no. 1, pp. 119–138, Feb. 2018, doi: [10.1007/s10044-016-0570-y](https://doi.org/10.1007/s10044-016-0570-y).
- [120] L. Al-Frayd and A. Al-Taei, "Wrapper filter approach for accelerometer-based human activity recognition," *Pattern Recognit. Image Anal.*, vol. 30, no. 4, pp. 757–764, Oct. 2020, doi: [10.1134/S1054661820040033](https://doi.org/10.1134/S1054661820040033).
- [121] P. Gupta, R. McClatchey, and P. Caleb-Solly, "Tracking changes in user activity from unlabelled smart home sensor data using unsupervised learning methods," *Neural Comput. Appl.*, vol. 32, no. 16, pp. 12351–12362, Aug. 2020, doi: [10.1007/s00521-020-04737-6](https://doi.org/10.1007/s00521-020-04737-6).
- [122] R. A. Hamad, M. Kimura, and J. Lundström, "Efficacy of imbalanced data handling methods on deep learning for smart homes environments," *Social Netw. Comput. Sci.*, vol. 1, no. 4, p. 204, Jun. 2020, doi: [10.1007/s42979-020-00211-1](https://doi.org/10.1007/s42979-020-00211-1).
- [123] A. Jordao, L. A. B. Torres, and W. R. Schwartz, "Novel approaches to human activity recognition based on accelerometer data," *Signal, Image Video Process.*, vol. 12, no. 7, pp. 1387–1394, Oct. 2018, doi: [10.1007/s11760-018-1293-x](https://doi.org/10.1007/s11760-018-1293-x).
- [124] B. Pontes, M. Cunha, R. Pinho, and H. Fuks, "Human-sensing: Low resolution thermal array sensor data classification of location-based postures," in *Distributed, Ambient and Pervasive Interactions* (Lecture Notes in Computer Science), N. Streitz and P. Markopoulos, Eds. Cham, Switzerland: Springer, 2017, pp. 444–457, doi: [10.1007/978-3-319-58697-7_33](https://doi.org/10.1007/978-3-319-58697-7_33).
- [125] L. Yao, W. Yang, and W. Huang, "A fall detection method based on a joint motion map using double convolutional neural networks," *Multimedia Tools Appl.*, vol. 81, no. 4, pp. 4551–4568, Feb. 2022, doi: [10.1007/s11042-020-09181-1](https://doi.org/10.1007/s11042-020-09181-1).
- [126] H. Ye et al., "Face recognition via active annotation and learning," in *Proc. 24th ACM Int. Conf. Multimedia*, Oct. 2016, pp. 1058–1062, doi: [10.1145/2964284.2984059](https://doi.org/10.1145/2964284.2984059).
- [127] S. Pini, O. B. Ahmed, M. Cornia, L. Baraldi, R. Cucchiara, and B. Huet, "Modeling multimodal cues in a deep learning-based framework for emotion recognition in the wild," in *Proc. 19th ACM Int. Conf. Multimodal Interact.*, Nov. 2017, pp. 536–543, doi: [10.1145/3136755.3143006](https://doi.org/10.1145/3136755.3143006).
- [128] S. S. Alia et al., "Summary of the 2nd nurse care activity recognition challenge using lab and field data," in *Proc. Adjunct ACM Int. Joint Conf. Pervasive Ubiquitous Comput. ACM Int. Symp. Wearable Comput.*, Sep. 2020, pp. 378–383, doi: [10.1145/3410530.3414611](https://doi.org/10.1145/3410530.3414611).
- [129] E. Soleimani and E. Nazerfard, "Cross-subject transfer learning in human activity recognition systems using generative adversarial networks," *Neurocomputing*, vol. 426, pp. 26–34, Feb. 2021, doi: [10.1016/j.neucom.2020.10.056](https://doi.org/10.1016/j.neucom.2020.10.056).

- [130] A. R. Javed, R. Faheem, M. Asim, T. Baker, and M. O. Beg, "A smartphone sensors-based personalized human activity recognition system for sustainable smart cities," *Sustain. Cities Soc.*, vol. 71, Aug. 2021, Art. no. 102970, doi: [10.1016/j.scs.2021.102970](https://doi.org/10.1016/j.scs.2021.102970).
- [131] L. Cao, Y. Wang, B. Zhang, Q. Jin, and A. V. Vasilakos, "GCHAR: An efficient group-based context-aware human activity recognition on smartphone," *J. Parallel Distrib. Comput.*, vol. 118, pp. 67–80, Aug. 2018, doi: [10.1016/j.jpdc.2017.05.007](https://doi.org/10.1016/j.jpdc.2017.05.007).
- [132] R. Khemchandani and S. Sharma, "Robust least squares twin support vector machine for human activity recognition," *Appl. Soft Comput.*, vol. 47, pp. 33–46, Oct. 2016, doi: [10.1016/j.asoc.2016.05.025](https://doi.org/10.1016/j.asoc.2016.05.025).
- [133] R. San-Segundo, H. Blunck, J. Moreno-Pimentel, A. Stisen, and M. Gil-Martín, "Robust human activity recognition using smartwatches and smartphones," *Eng. Appl. Artif. Intell.*, vol. 72, pp. 190–202, Jun. 2018, doi: [10.1016/j.engappai.2018.04.002](https://doi.org/10.1016/j.engappai.2018.04.002).
- [134] S. Mohamad, M. Sayed-Mouchaweh, and A. Bouchachia, "Online active learning for human activity recognition from sensory data streams," *Neurocomputing*, vol. 390, pp. 341–358, May 2020, doi: [10.1016/j.neucom.2019.08.092](https://doi.org/10.1016/j.neucom.2019.08.092).
- [135] P. Gupta and P. Caleb-Solly, "A framework for semi-supervised adaptive learning for activity recognition in healthcare applications," in *Engineering Applications of Neural Networks (Communications in Computer and Information Science)*, E. Pimenidis and C. Jayne, Eds. Cham, Switzerland: Springer, 2018, pp. 3–15, doi: [10.1007/978-3-319-98204-5_1](https://doi.org/10.1007/978-3-319-98204-5_1).
- [136] C. Lübbe, B. Friedrich, S. Fudickar, S. Hellmers, and A. Hein, "Feature based random forest nurse care activity recognition using accelerometer data," in *Proc. Adjunct ACM Int. Joint Conf. Pervasive Ubiquitous Comput. ACM Int. Symp. Wearable Comput.*, Sep. 2020, pp. 408–413, doi: [10.1145/3410530.3414340](https://doi.org/10.1145/3410530.3414340).
- [137] M. S. Siraj et al., "UPIC: User and position independent classical approach for locomotion and transportation modes recognition," in *Proc. Adjunct ACM Int. Joint Conf. Pervasive Ubiquitous Comput. ACM Int. Symp. Wearable Comput.*, Sep. 2020, pp. 340–345, doi: [10.1145/3410530.3414343](https://doi.org/10.1145/3410530.3414343).
- [138] S. Mekruksavanich and A. Jitpattanakul, "Recognition of real-life activities with smartphone sensors using deep learning approaches," in *Proc. IEEE 12th Int. Conf. Softw. Eng. Service Sci. (ICSESS)*, Aug. 2021, pp. 243–246, doi: [10.1109/ICSESS52187.2021.952231](https://doi.org/10.1109/ICSESS52187.2021.952231).
- [139] J. Suto and S. Oniga, "Efficiency investigation from shallow to deep neural network techniques in human activity recognition," *Cognit. Syst. Res.*, vol. 54, pp. 37–49, May 2019, doi: [10.1016/j.cogsys.2018.11.009](https://doi.org/10.1016/j.cogsys.2018.11.009).
- [140] R. Alfaifi and A. M. Artoli, "Human action prediction with 3D-CNN," *Social Netw. Comput. Sci.*, vol. 1, no. 5, p. 286, Aug. 2020, doi: [10.1007/s42979-020-00293-x](https://doi.org/10.1007/s42979-020-00293-x).
- [141] M. Bilal, F. K. Shaikh, M. Arif, and M. F. Wyne, "A revised framework of machine learning application for optimal activity recognition," *Cluster Comput.*, vol. 22, no. S3, pp. 7257–7273, May 2019, doi: [10.1007/s10586-017-1212-x](https://doi.org/10.1007/s10586-017-1212-x).
- [142] F. Cruciani et al., "Feature learning for human activity recognition using convolutional neural networks," *CCF Trans. Pervasive Comput. Interact.*, vol. 2, no. 1, pp. 18–32, Mar. 2020, doi: [10.1007/s42486-020-00026-2](https://doi.org/10.1007/s42486-020-00026-2).
- [143] J. Guo, H. Bai, Z. Tang, P. Xu, D. Gan, and B. Liu, "Multi modal human action recognition for video content matching," *Multimedia Tools Appl.*, vol. 79, nos. 45–46, pp. 34665–34683, Dec. 2020, doi: [10.1007/s11042-020-08998-0](https://doi.org/10.1007/s11042-020-08998-0).
- [144] M. A. Khan, Y.-D. Zhang, S. A. Khan, M. Attique, A. Rehman, and S. Seo, "A resource conscious human action recognition framework using 26-layered deep convolutional neural network," *Multimedia Tools Appl.*, vol. 80, nos. 28–29, pp. 35827–35849, Nov. 2021, doi: [10.1007/s11042-020-09408-1](https://doi.org/10.1007/s11042-020-09408-1).
- [145] T. Wu et al., "CARMUS: Towards a general framework for continuous activity recognition with missing values on smartphones," in *Proc. IEEE 42nd Annu. Comput. Softw. Appl. Conf. (COMPSAC)*, Jul. 2018, pp. 850–859, doi: [10.1109/COMPSAC.2018.00148](https://doi.org/10.1109/COMPSAC.2018.00148).
- [146] T. Hossain and S. Inoue, "A comparative study on missing data handling using machine learning for human activity recognition," in *Proc. Joint 8th Int. Conf. Informat., Electron. Vis. (ICIEV) 3rd Int. Conf. Imag., Vis. Pattern Recognit. (icIVPR)*, May 2019, pp. 124–129.
- [147] S. Yu, H. Chen, and R. A. Brown, "Hidden Markov model-based fall detection with motion sensor orientation calibration: A case for real-life home monitoring," *IEEE J. Biomed. Health Informat.*, vol. 22, no. 6, pp. 1847–1853, Nov. 2018.
- [148] A. Keshavarzian, S. Sharifian, and S. Seyedin, "Modified deep residual network architecture deployed on serverless framework of IoT platform based on human activity recognition application," *Future Gener. Comput. Syst.*, vol. 101, pp. 14–28, Dec. 2019, doi: [10.1016/j.future.2019.06.009](https://doi.org/10.1016/j.future.2019.06.009).
- [149] S. H. Arib, R. Akter, O. Shahid, and M. A. R. Ahad, "Classical machine learning approach for human activity recognition using location data," in *Proc. Adjunct ACM Int. Joint Conf. Pervasive Ubiquitous Comput. ACM Int. Symp. Wearable Comput.*, Sep. 2021, pp. 340–345, doi: [10.1145/3460418.3479376](https://doi.org/10.1145/3460418.3479376).
- [150] E. Kwon, H. Park, S. Byon, E.-S. Jung, and Y.-T. Lee, "HaaS(human activity analytics as a service) using sensor data of smart devices," in *Proc. Int. Conf. Inf. Commun. Technol. Converg. (ICTC)*, Oct. 2018, pp. 1500–1502.
- [151] U. Ozbulak, B. Vandersmissen, A. Jalalvand, I. Couckuyt, A. Van Messem, and W. De Neve, "Investigating the significance of adversarial attacks and their relation to interpretability for radar-based human activity recognition systems," *Comput. Vis. Image Understand.*, vol. 202, Jan. 2021, Art. no. 103111, doi: [10.1016/j.cviu.2020.103111](https://doi.org/10.1016/j.cviu.2020.103111).
- [152] M.-S. Dao, T.-A. Nguyen-Gia, and V.-C. Mai, "Daily human activities recognition using heterogeneous sensors from smartphones," *Proc. Comput. Sci.*, vol. 111, pp. 323–328, Jan. 2017, doi: [10.1016/j.procs.2017.06.030](https://doi.org/10.1016/j.procs.2017.06.030).
- [153] N. M. Richards and J. H. King, "Big data ethics," *Wake Forest L. Rev.*, vol. 49, p. 393, Jan. 2014.
- [154] T. Simonite, *Google's AI Guru Wants Computers to Think More Like Brains*. San Francisco, CA, USA: Wired, 2018.
- [155] S. Mohamed, M.-T. Png, and W. Isaac, "Decolonial AI: Decolonial theory as sociotechnical foresight in artificial intelligence," *Philosophy Technol.*, vol. 33, no. 4, pp. 659–684, Dec. 2020.
- [156] P. Chahuara, A. Fleury, M. Vacher, and F. Portet, "Méthodes SVM et MLN pour la reconnaissance automatique d'activités humaines dans les habitats perceptifs: Tests et perspectives," in *Proc. Reconnaissance des Formes et Intell. Artificielle (RFIA)*, Lyon, France, Jan. 2012, p. 978-2-9539515-2-3.
- [157] N. Mollet and R. Chellali, "Détection et interprétation des gestes de la main," in *Proc. 3rd Int. Conf. (SETIT)*, 2005, pp. 1–7.
- [158] R. Merkley, *Use and Fair Use: Statement on Shared Images in Facial Recognition AI*. Mountain View, CA, USA: Creative Commons, 2019.
- [159] D. Sculley, J. Snoek, A. Wiltschko, and A. Rahimi, "Winner's curse? On pace, progress, and empirical rigor," in *Proc. ICLR Workshop*, Vancouver, BC, Canada: Vancouver Convention Center, 2018.
- [160] G. Bhat, N. Tran, H. Shill, and U. Y. Ogras, "W-HAR: An activity recognition dataset and framework using low-power wearable devices," *Sensors*, vol. 20, no. 18, p. 5356, Sep. 2020.
- [161] S.-J. van Els, D. Graus, and E. Beauxis-Aussalet, "Improving fairness assessments with synthetic data: A practical use case with a recommender system for human resources," in *Proc. 1st Int. Workshop Comput. Jobs Marketplace (CompJobs)*, Feb. 2022.
- [162] B. Xin et al., "Federated synthetic data generation with differential privacy," *Neurocomputing*, vol. 468, pp. 1–10, Jan. 2022, doi: [10.1016/j.neucom.2021.10.027](https://doi.org/10.1016/j.neucom.2021.10.027).
- [163] J. Yang, M. N. Nguyen, P. P. San, X. L. Li, and S. Krishnaswamy, "Deep convolutional neural networks on multichannel time series for human activity recognition," in *Proc. 24th Int. Joint Conf. Artif. Intell.*, 2015, pp. 3995–4001.
- [164] J. Wang, Y. Chen, S. Hao, X. Peng, and L. Hu, "Deep learning for sensor-based activity recognition: A survey," *Pattern Recognit. Lett.*, vol. 119, pp. 3–11, Mar. 2019.
- [165] Y. LeCun, F. J. Huang, and L. Bottou, "Learning methods for generic object recognition with invariance to pose and lighting," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2004, p. 104, doi: [10.1109/CVPR.2004.1315150](https://doi.org/10.1109/CVPR.2004.1315150).
- [166] J. J.-C. Ying, B.-H. Lin, V. S. Tseng, and S.-Y. Hsieh, "Transfer learning on high variety domains for activity recognition," in *Proc. ASE BigData & SocialInformatics*. New York, NY, USA: Association for Computing Machinery, Oct. 2015, pp. 1–6, doi: [10.1145/2818869.2818890](https://doi.org/10.1145/2818869.2818890).
- [167] H. M. S. Hossain, M. A. A. H. Khan, and N. Roy, "Active learning enabled activity recognition," *Pervas. Mobile Comput.*, vol. 38, pp. 312–330, Jul. 2017, doi: [10.1016/j.pmcj.2016.08.017](https://doi.org/10.1016/j.pmcj.2016.08.017).
- [168] V. Xafis and M. K. Labude, "Openness in big data and data repositories," *Asian Bioethics Rev.*, vol. 11, no. 3, pp. 255–273, Sep. 2019, doi: [10.1007/s41649-019-00097-z](https://doi.org/10.1007/s41649-019-00097-z).

- [169] *Center for Data and Visualization Sciences | Duke University Libraries*. Accessed: Aug. 9, 2022. [Online]. Available: <https://library.duke.edu/data>
- [170] U. of Bristol. *Managing Research Data*. Accessed: Aug. 9, 2022. [Online]. Available: <http://www.bristol.ac.uk/staff/researchers/data/>
- [171] H. Wang. (Aug. 9, 2022). *LibGuides: Machine Learning and AI: Home*. [Online]. Available: <https://guides.library.cmu.edu/machine-learning/home>
- [172] X. Chu, I. F. Ilyas, S. Krishnan, and J. Wang, "Data cleaning: Overview and emerging challenges," in *Proc. Int. Conf. Manag. Data (SIGMOD)*. New York, NY, USA: Association for Computing Machinery, Jun. 2016, pp. 2201–2206, doi: [10.1145/2882903.2912574](https://doi.org/10.1145/2882903.2912574).
- [173] S. Singhal and M. Jena, "A study on WEKA tool for data preprocessing, classification and clustering," *Int. J. Innov. Technol. Exploring Eng.*, vol. 2, no. 6, pp. 250–253, May 2013.
- [174] S. A. Alasadi and W. S. Bhaya, "Review of data preprocessing techniques in data mining," *J. Eng. Appl. Sci.*, vol. 12, no. 16, pp. 4102–4107, 2017.
- [175] K. Jee and G.-H. Kim, "Potentiality of big data in the medical sector: Focus on how to reshape the healthcare system," *Healthcare Inform. Res.*, vol. 19, no. 2, pp. 79–85, 2013.
- [176] M. Mancini, "Exploiting big data for improving healthcare services," *J. E-Learn. Knowl. Soc.*, vol. 10, no. 2, pp. 1–11, 2014.
- [177] E. Baro, S. Degouf, R. Beuscart, and E. Chazard, "Toward a literature-driven definition of big data in healthcare," *BioMed Res. Int.*, vol. 2015, pp. 1–9, Jun. 2015.
- [178] M. B. Howren, M. W. V. Weg, and F. D. Wolinsky, "Computerized cognitive training interventions to improve neuropsychological outcomes: Evidence and future directions," *J. Comparative Effectiveness Res.*, vol. 3, no. 2, pp. 145–154, Mar. 2014.
- [179] W. Raghupathi and V. Raghupathi, "Big data analytics in healthcare: Promise and potential," *Health Inf. Sci. Syst.*, vol. 2, no. 1, pp. 1–10, Dec. 2014.
- [180] L. M. Fernandes, M. O'Connor, and V. Weaver, "Big data, bigger outcomes," *J. AHIMA*, vol. 83, no. 10, pp. 38–43, 2012.
- [181] J.-C. Hsieh, A.-H. Li, and C.-C. Yang, "Mobile, cloud, and big data computing: Contributions, challenges, and new directions in telecardiology," *Int. J. Environ. Res. Public Health*, vol. 10, no. 11, pp. 6131–6153, Nov. 2013.
- [182] A. B. Khanghah, G. Fernie, and A. R. Fekr, "Design and validation of vision-based exercise biofeedback for tele-rehabilitation," *Sensors*, vol. 23, no. 3, p. 1206, Jan. 2023, doi: [10.3390/s23031206](https://doi.org/10.3390/s23031206).
- [183] J. Wu, S. G. Faux, I. Harris, C. J. Poulos, and T. Alexander, "Record linkage is feasible with non-identifiable trauma and rehabilitation datasets," *Austral. New Zealand J. Public Health*, vol. 40, no. 3, pp. 245–249, Jun. 2016, doi: [10.1111/1753-6405.12510](https://doi.org/10.1111/1753-6405.12510).
- [184] F. Cruciani et al., "A public domain dataset for human activity recognition in free-living conditions," in *Proc. IEEE SmartWorld, Ubiquitous Intell. Comput., Adv. Trusted Comput., Scalable Comput. Commun., Cloud Big Data Comput., Internet People Smart City Innov. (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI)*, Aug. 2019, pp. 166–171, doi: [10.1109/SmartWorld-UIC-ATC-SCALCOM-IOP-SCI2019.00071](https://doi.org/10.1109/SmartWorld-UIC-ATC-SCALCOM-IOP-SCI2019.00071).



Gulzar Alam received the M.Sc. degree in software engineering from the King Fahd University of Petroleum and Minerals (KFUPM), Dhahran, Saudi Arabia, in 2020. He is currently pursuing the Ph.D. degree in computing science with the School of Computing, Ulster University, Belfast, U.K.

Following the completion of his master's degree, he worked on several cybersecurity and software engineering-related research projects. He also served as an Undergraduate Teaching Assistant at KFUPM. He has devoted his career to enhancing the field of computer science through his contributions to data curation and the quality of datasets for human activity recognition (HAR). He obtained valuable experience conducting research, collaborating with teams, and publishing academic papers and a patent during this time. He has also collaborated with academic institutions on various funded projects to advance the current state of knowledge in computer science. His current research focuses on enhancing the quality of datasets in the realm of HAR.



Ian McChesney received the degree in computer science in 1987, and the Ph.D. degree in software engineering from Ulster University, Coleraine, Northern Ireland, in 1998.

He is a senior lecturer in computing science. He has higher education experience in research, teaching, and knowledge transfer with industry. In pervasive computing, he has worked on experimental human activity recognition and support tools for managing open datasets. His knowledge transfer experience covers areas such as software project management and requirements engineering. His research interests in software engineering have focused on sociotechnical themes, such as software estimation, program comprehension, and software team coordination, using methods such as field surveys, empirical software engineering, and eye tracking.

Dr. McChesney is a Fellow of The British Computer Society, a Chartered Engineer, and a Senior Fellow of the Higher Education Academy.



Peter Nicholl received the B.Eng. degree in electronic systems and the Ph.D. degree in feature encoding from Ulster University, Belfast, U.K., in 1991 and 1994, respectively.

He is a Senior Lecturer in computing science with the School of Computing, Ulster University. He has higher education experience in research, teaching, and knowledge transfer with industry. His research interests include intelligent transportation, computer vision, and deep learning.

Dr. Nicholl is a Senior Fellow of the Higher Education Academy.



Joseph Rafferty received the B.Eng. degree in computer science from Queen's University Belfast, Belfast, U.K., in 2010, and the M.Sc. degree in computing and the Ph.D. degree in computer science from Ulster University, Belfast in 2011 and 2016, respectively.

He is currently a Lecturer with the School of Computing, Ulster University. His research interests include intention recognition, smart environments, agent-based systems, connected health, sensor technology, and planning and intelligent systems.