

Let's Play for Action: Recognizing Activities of Daily Living by Learning from Life Simulation Video Games

Alina Roitberg*

David Schneider*

Aulia Djamal

Constantin Seibold

Simon Reiß

Rainer Stiefelhagen

Institute for Anthropomatics and Robotics
 Karlsruhe Institute of Technology, Germany
 {firstname.lastname}@kit.edu

Abstract—Recognizing Activities of Daily Living (ADL) is a vital process for intelligent assistive robots, but collecting large annotated datasets requires time-consuming temporal labeling and raises privacy concerns, e.g., if the data is collected in a real household. In this work, we explore the concept of constructing training examples for ADL recognition by playing life simulation video games and introduce the SIMS4ACTION dataset created with the popular commercial game THE SIMS 4. We build SIMS4ACTION by specifically executing actions-of-interest in a “top-down” manner, while the gaming circumstances allow us to freely switch between environments, camera angles and subject appearances. While ADL recognition on gaming data is interesting from the theoretical perspective, the key challenge arises from transferring it to the real-world applications, such as smart-homes or assistive robotics. To meet this requirement, SIMS4ACTION is accompanied with a GAMING→REAL benchmark, where the models are evaluated on real videos derived from an existing ADL dataset. We integrate two modern algorithms for video-based activity recognition in our framework, revealing the value of life simulation video games as an inexpensive and far less intrusive source of training data. However, our results also indicate that tasks involving a mixture of gaming and real data are challenging, opening a new research direction. We will make our dataset publicly available at <https://github.com/aroitberg/sims4action>.

I. INTRODUCTION

With growing elderly population¹, assistive household robots increasingly gain attention, but such technologies also require raising the levels of robot’s awareness. Understanding human behaviour enables better assistance, improves human-robot-communication, situation-dependent planning and safety [1], [2], [3]. This perception task is often referred to as recognizing Activities of Daily Living (ADL) and is addressed by training deep neural networks on carefully curated and manually labelled datasets [4], [5], [6], [7], [8].

The quality of recognition is directly linked to the amount of labelled examples available during training [9]. However, collecting large video datasets tailored for a specific use-case is expensive and requires application-dependent sen-

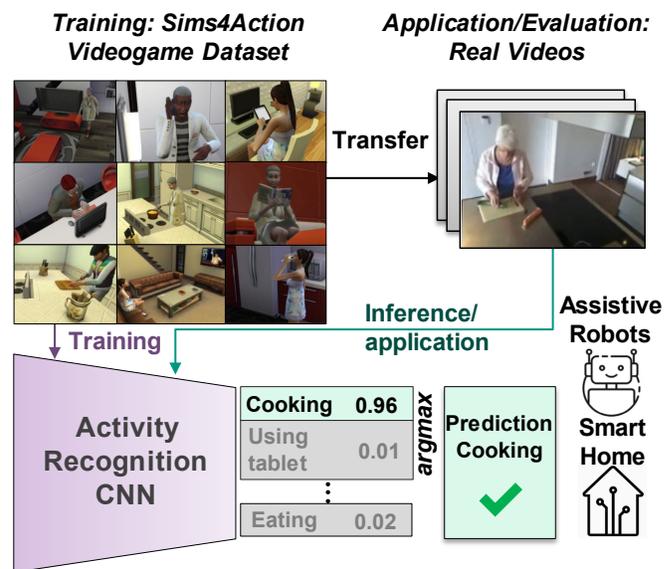


Fig. 1: We leverage the life simulation video game THE SIMS 4 for inexpensive ADL training data which we then utilize in the GAMING→REAL setting in order to suit real-world applications, such as smart homes or assistive robots.

sory setups, recruitment of diverse participants and time-consuming accurate temporal annotation of activities. Even after this effort is put in, the collected datasets often become insufficient, since the environment or the desired behaviour types to be recognized may change and the costly process has to be revisited. On top of this, data protection and privacy concerns make recruitment of participants for recording training videos in real domestic environments even more difficult, as seen, e.g., by the blurred faces in the Toyota Smarthome dataset [8] for ADL recognition.

Video games constitute a multi-billion dollar industry², where developers put great effort into build highly realistic worlds. Intuitively we connect gameplay exclusively

*David Schneider and Alina Roitberg contributed equally to this work.

¹Demography - Elderly population - OECD Data: data.oecd.org/pop/elderly-population.htm

²Statistics on the videogame market obtained from www.statista.com/topics/868/video-games/

with entertainment, but the impressive amount of details and diversity of objects, environments, appearances, human body motion and character traits coupled with the power to control the virtual world render it an excellent source of training data. While adventure-driving games, such as Grand Theft Auto, have been recently leveraged as a data source for semantic segmentation in autonomous vehicles applications [10], [11], the potential of life simulation games for robot learning in assisted living scenarios has been left untouched and is the main motivation of our work.

We view life simulation video games as an opportunity for inexpensive collection of visual training examples for smart homes and domestic robots and introduce the SIMS4ACTION dataset built with the popular commercial game THE SIMS 4³. Our dataset features ten hours of virtual subjects engaged in different household activities, such as cooking or working on a computer. As THE SIMS 4 is specifically designed to give the player full control of everyday situations, we are able to freely design the environments (six different locations), camera angles (four fixed cameras per location and a moving camera mode) and human appearances (four elderly people, two adults and two young adults and equal distribution of males and females). We capture ten different behaviour categories, which we have triggered by “playing” the desired ADL, and accompany our dataset with the GAMING→REAL benchmark, where the models are evaluated on real videos derived from Toyota Smarthome, an existing ADL dataset [8]. As our visual model, we adopt two off-the-shelf video classification architectures, which we train with the SIMS4ACTION data only. Our results constitute a challenging benchmark for future research and have two-fold conclusions: (1) video game supervision has potential for ADL recognition in real settings, as all models surpass the random baseline by a large margin, although (2) The GAMING→REAL evaluation regime is very challenging as modern recognition algorithms are highly susceptible to domain shifts. We hope that our dataset detaches ADL recognition from heavy physical data collection cycles and creates a new avenue for future research focusing on capturing the essence of human behaviour in both, synthetic and real domains.

II. RELATED WORK

A. Recognizing Activities of Daily Living

The performance of models recognizing domestic activities is directly linked to the progress in general video classification, which has undergone a sudden shift from machine learning approaches operating on hand-crafted features, such as Hidden Markov Models or Improved Dense Trajectories [12], [13] to end-to-end Convolutional Neural Networks (CNNs) [14]. Today, spatiotemporal 3D CNNs [15], [14], [16] are considered front-runners in activity classification and are used in a wide range of applications, such as robotics, autonomous driving or healthcare [17], [18], [19], [20], [21].

³The Sims 4 is a commercial game by Electronic Arts (EA): ea.com/en-gb/games/the-sims/the-sims-4

Research of ADL recognition for domestic robotics or smart homes often follows similar classification frameworks and focuses on constructing annotated datasets with the specific application requirements in mind [22], [8], [4], [23], [5], [6]. The creators of the Toyota Smarthome dataset [8], [22], for example, equipped apartments of elderly people with multiple Kinect cameras and leverage I3D enhanced with spatio-temporal attention to distinguish 31 activity classes. While Toyota Smarthome is an excellent benchmark, recording and annotating such data requires heavy effort and extending it with new behaviours, environments or appearance types is harder than reproducing desired situations inside rich simulation games such as THE SIMS 4. We aim to overcome this labour-intensive constraint and explore ADL recognition by learning from life simulation video games, which, to the best of our knowledge, has not been considered yet. While our main contribution is the SIMS4ACTION dataset with domestic activities recorded during gameplay, we establish correspondences between our classes and multiple Toyota Smarthome categories in order to validate if such videogame-based training examples can be used to classify real data.

B. Leveraging Video Games and Simulations for Data

Recent works have considered video games as a source for inexpensive training data especially focusing on autonomous driving applications. Richter et al. [10], [11] presented a method to collect data and ground truth for visual perception tasks in driving scenarios (semantic segmentation, visual odometry, optical flow, object detection & tracking) from the video game Grand Theft Auto V (GTA V). Krähenbühl [24], presented an improved data collection method: fully automated ground truth extraction from video games, without human annotation effort. Through a wrapper around the DirectX 11 API the rendered data is obtained in real-time, i.e. during actual gameplay. The authors collected 220 thousand training- and 60 thousand test images across three different games (Far Cry Primal, The Witcher 3, GTA V) and were able to gather rich image meta information such as albedo, depth, instance segmentation, semantic labeling, optical flow and occlusion boundaries. Our dataset is created with a different type of game - the extensive life simulator THE SIMS 4 and targets human behaviour, while previous work has strong focus on the underlying scene and video features, such as segmentation maps or optical flow.

While learning behaviours from synthetic data is an under-researched area, few recent works have introduced simulation tools explicitly targeting human activity recognition [25], [26], [27]. The research most similar to ours is presumably the concurrent work of Varol et al. [25], which leverage Blender-based simulations to improve activity recognition from viewpoints not present during training. A different line of work considers activity simulations from the perspective of body poses, synthesizing skeleton trajectories with neural networks [28], [29]. In contrast to the previous research, we focus on learning new categories of daily activities with *videogame supervision*, which, to the best of our knowledge, is explored for the first time.

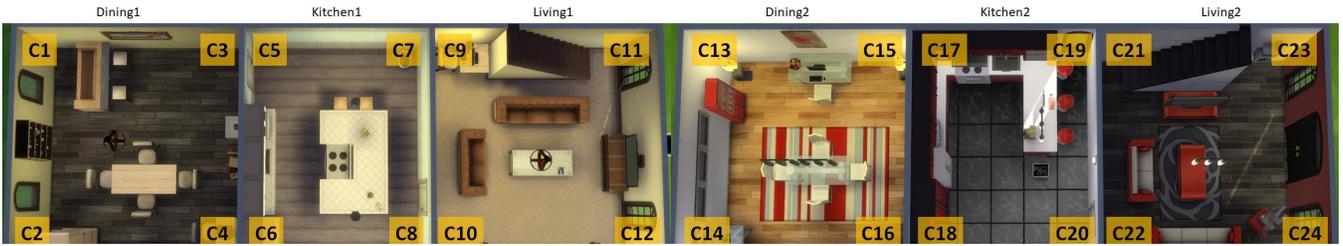


Fig. 2: Snapshots of the virtual apartments we designed to simulate household situations. The two homes cover a dining room, a living room and a kitchen. *C1 - C24* marks the 24 camera placements used in the static camera mode recordings.



Fig. 3: Eight subjects (four elderly people, two adults and two young adults) are covered in the SIMS4ACTION dataset.

Property	Train	Validation	Test	Total
Duration (minutes)	388.5	79.6	157.45	625.6
Number of samples*	8217	1684	3331	13232
Number of subjects	5	1	2	8
Male/Female	2/3	1/0	1/1	4/4
Static camera	✓	✓	✓	✓
Moving camera	✓	✓	✓	✓
Number of static views	24	24	24	24
Number of activity classes	10	10	10	10
Min. nr. samples per class	769	150	324	1265
Avg. nr. samples per class	821.7	168.4	333.1	1323.2

TABLE I: Characteristics of the SIMS4ACTION dataset.

*A sample is a three second snippet with its assigned label.

III. PLAYING SIMS 4 FOR ADL TRAINING DATA

We introduce SIMS4ACTION – a dataset designed for visually recognizing human behaviour at home by learning from life simulation video games. Our dataset covers ten different behaviours carried out by eight designed subjects, while the gaming situation allows us to control their appearance and environments according to our needs. The main characteristics of our dataset are provided in Table I.

A. Data Collection

SIMS4ACTION dataset is collected using THE SIMS 4 life-simulation video game, where players guide the virtual humans through everyday activities of all aspects leading to complex character traits, varying emotional states, diverse subject appearances, and ages. We leverage this environment to induce the situations we want to recognize and use the built-in video capture tool for the recordings (the tool allows us to select the time frame, capture size and image resolution, which is 30 FPS and 640×368 pixels in our case). The main stages of our data collection are now described in detail and can be summarized in 1) design of the virtual human 2) design of the home environment, and 3) actively “playing” behaviour categories we want to recognize while recording the gameplay.

Virtual Subjects. Recruitment of diverse subjects is essential to overcome dataset biases, but is often difficult in practice. Complementing real datasets with samples from life simulation video games might mitigate this issue, as THE SIMS 4 enables detailed subject personalization even capturing details such as walking style. Keeping broad range of appearances in mind, we have designed eight virtual

humans diversifying their look, genders (four male and four female) and age groups (four elderly people, two young adults, two adults). An overview of the created subjects is provided in Figure 6.

Inducing Virtual Behaviours. In THE SIMS 4, the virtual subject (Sim) follows given instructions. Thus, we explicitly enforce the behaviors-of-interest in a top-down manner, by repeatedly guiding the Sim to engage in the desired activities, which range from cooking or taking pills to rare events such as having a wedding party or morning sickness during pregnancy. In this work, we focus on ten everyday household activities popular in previous works on recognition in real smart homes. Examples for each of the classes are provided in Figure 6. While balanced datasets are rare in real-life [8], synthetic data collection enables us to easily design equal distribution of classes (ca. one hour footage per category), although the number of videos for each activity varies as the activities have different durations (e.g., around 15 to 30 seconds for *drink*, and around 30 seconds to one minute for *cook* and *eat*). The duration of some activities can be controlled by the player (*get up/sit down, use computer, use tablet, walk, watch TV*), whereas other activities will be completed automatically after a certain amount of time (*cook, drink, eat, read book, use phone*). Figure 4 illustrates sample statistics for the individual categories.

Virtual Homes. As users can create customized interiors, we can generate a large variety of testing grounds. We chose to lean close to the Toyota Smarhome [8], [22] framework and recreate digital versions of their households with the game’s build mode. Our dataset covers two simulated apartments

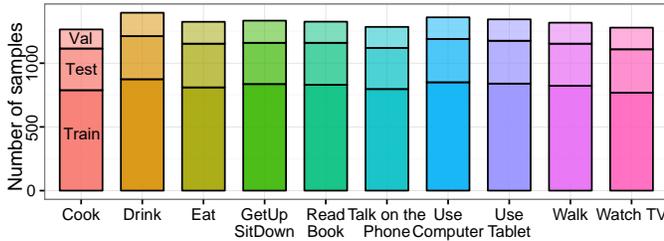


Fig. 4: Sample frequency of the activities in *train*, *test* and *val*. In contrast to real-life data collection [8], life simulation video games give us control of the data distribution resulting in a well-balanced dataset. (A sample corresponds to a 3s snippet with the assigned label.)

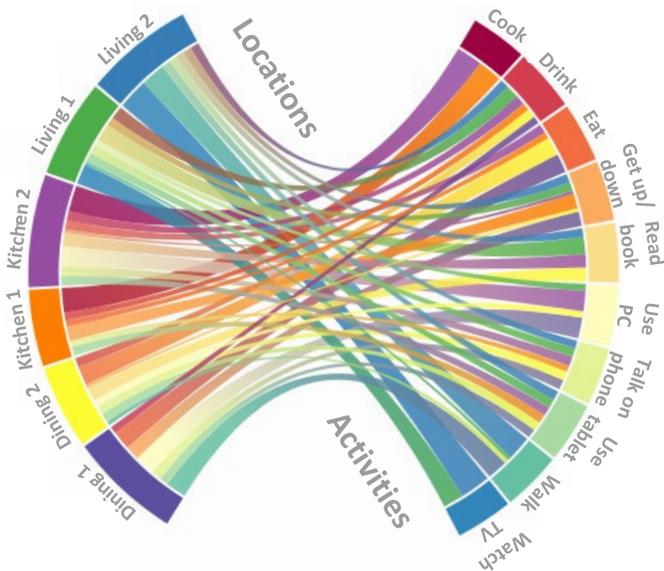


Fig. 5: *Where did the activities take place?* The cord diagram visualizes the correspondences of the captured behaviours and the rooms it took place in. Best viewed in color.

with three rooms: a kitchen, a dining-, and a living room. While some activities (such as *talking on the phone*) happen at various locations, other events, such as *cooking* have natural location biases (the location distribution of the captured activities is provided in Figure 5).

The simulation environment allows us to freely choose the camera angle. We record the videos with two different cameras: (1) the *fixed camera* mode, where the camera stays in one determined position for the entire video, and (2) the *moving camera* mode, where the camera is set in various random angles and moves for every few seconds in the video. For the *fixed camera* mode, we define four different camera angles in each room, resulting in a total of 24 distinct camera views. An overview of the modeled homes and the camera views is shown in Figure 2.

B. Dataset and Splits Statistics

We collect 942 videos with a total length of 10 hours, 25 minutes, and 32 seconds and a roughly equal representation

of activity classes (ca. 1 hour each) and subjects (ca. 1 hour, 17 minutes each).

SIMS4ACTION Splits. The generalization to new humans is essential in ADL recognition. We therefore divide our dataset *subject-wise* into *train* (5 subjects), *val* (1 subjects), and *test* (2 subjects). We learn model parameters on *train*, select checkpoints and hyperparameters on *val*, and present final evaluation on *test* (overview of the split statistics provided in Figure 4). Since behaviour complexity and duration vary strongly, we divide the videos into chunks of three seconds and use them as samples in our benchmark. Hence, the classification task on our dataset amounts to assigning the correct activity label to a three second video.

GAMING→REAL Benchmark. Since neural networks heavily rely on processing the data distribution they were trained on, the task to transfer recognition capabilities from the GAMING→REAL domain is challenging. To enable transfer from video games in actual robotics or smart home applications long-term, SIMS4ACTION is accompanied by a real data benchmark. To achieve this, we build correspondences between our categories and activities captured in the real Toyota Smarthome dataset [8] and use the Toyota Smarthome cross-subject test set with the corresponding classes for evaluation. Note, that we do not directly take Toyota Smarthome classes as-is, and certain modifications and consolidations are required to match the SIMS4ACTION classes. For example, we consolidate different variants of cooking activities into a single cooking class. Figure 6 compares the videogame examples and their real-life counterparts for every category. Overall, the GAMING→REAL test set comprises 4412 videos with the corresponding ten activities inside a real smart home, where each video is used as a test sample, as in [8]. This real-life test data is not evenly distributed, with the number of samples ranging between 15 for *use tablet* and 1225 for the *walk* category.

IV. NEURAL ARCHITECTURES

To evaluate our idea of ADL recognition by learning from life simulation video games and provide a competitive benchmark based on the SIMS4ACTION dataset, we adapt two off-the-shelf CNNs for general activity recognition.

Inflated 3D CNN. The Inflated 3D ConvNet (I3D) by Carreira and Zisserman [14] is a wildly successful video classification architecture utilized for both, general- and ADL recognition [14], [8]. The model benefits from the object recognition architecture Inception-v1 [30] and modifies it to deal with the temporal domain via expansion of 2D- to 3D kernels and ablations on how to integrate pooling operations. Inherited by the ancestor network, I3D is 27 layers deep and stacks three convolution layers at the beginning, nine inception modules which themselves are two layers deep, one fully-connected layer at the end as well as five max- and average-pooling layers.

Separable 3D CNN. While I3D offers excellent classification accuracy, it is a heavy architecture with > 12 million parameters and might be impractical in robotic applications.

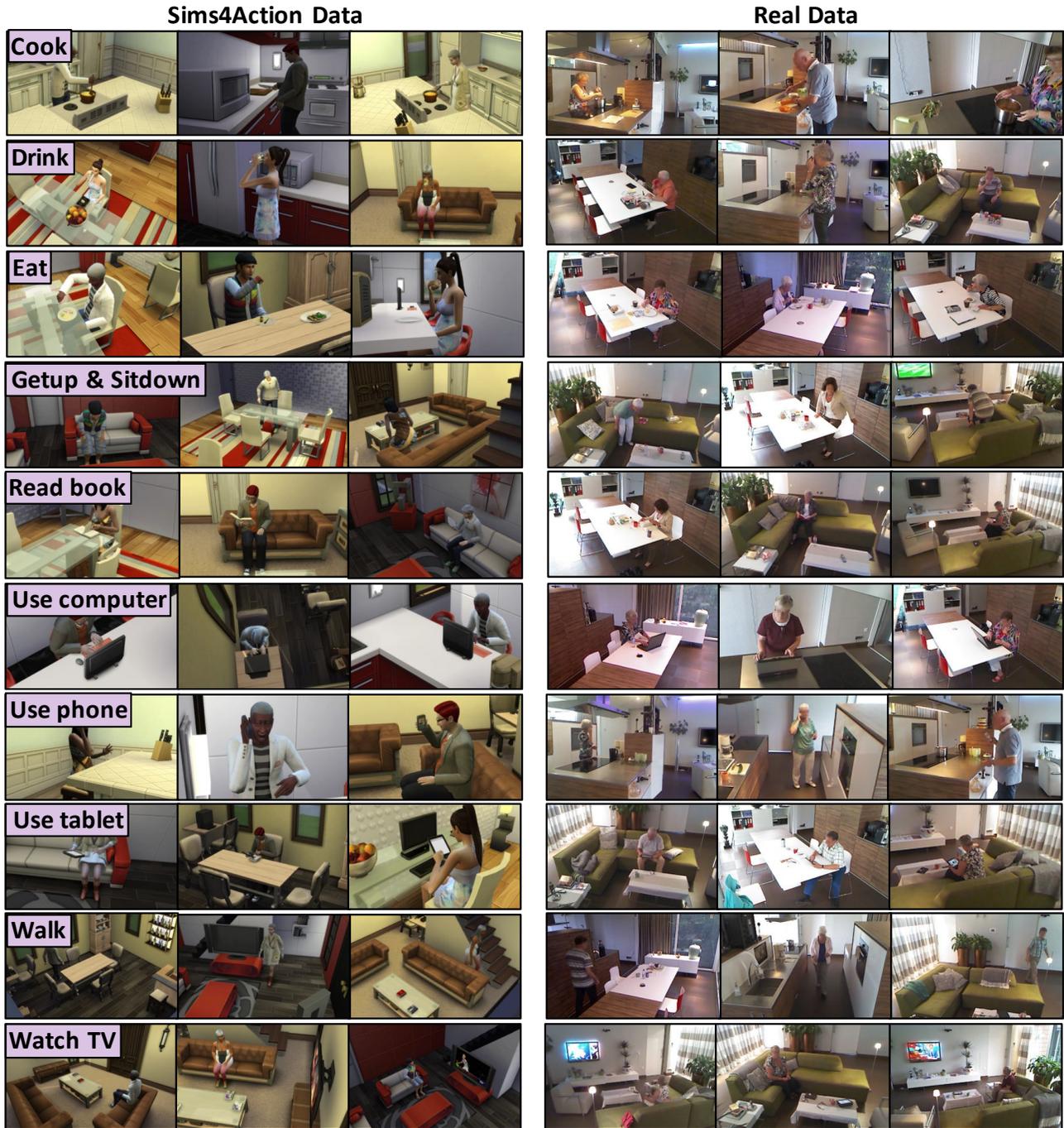


Fig. 6: Examples of subjects being engaged into different domestic activities in the SIMS4ACTION dataset and the real-life counterparts used for the GAMING→REAL Benchmark, which are derived from Toyota Smarhome [8].

With this in mind, Xie et al. [16] redesigned the I3D architecture by considering (1) switching to a top-heavy model, (2) temporally separable convolutions, and (3) spatio-temporal feature gating. These alterations led to a major reduction in floating point operations per second (FLOPs) and culminated in the new Separable 3D architecture (S3D). The so-called top-heavy model property reduces FLOPs by switching early 3D convolutional blocks in the I3D architecture with standard 2D convolutions and only in later layers involve temporally

separable 3D convolutions. The high accuracy of S3D can be attributed to a gating mechanism after temporal convolutions in the separable convolution blocks, which resembles the mechanics of *self-attention*.

Implementation and Training Details. In all experiments except for the baselines, we train our networks end-to-end on SIMS4ACTION for 100 epochs. We use the I3D architecture with an input frame size of 224×224 or the S3D architecture with an input frame size of 128×128 . The net-

Model	Weight initialization	Normal Acc [%]			Balanced Acc [%]		
		Top 1	Top 3	Top 5	Top 1	Top 3	Top 5
Random Baseline	None	10.0	33.33	50.00	10.0	33.33	50.00
Separable 3D Net	Random	56.74	79.36	91.18	56.52	79.33	91.25
Separable 3D Net	Kinetics400	84.67	96.65	98.67	84.61	96.67	98.68
I3D Net	Random	67.0	87.81	94.29	66.91	87.78	94.33
I3D Net	Kinetics400	89.06	98.35	99.55	89.12	98.34	99.54

TABLE II: Evaluation of the ADL recognition in the GAMING→GAMING setting of our benchmark. The performance metrics are standard- and balanced accuracies.

True Activity Class	Prec. [%]	Recall [%]	F1 [%]
Cook	98.99	93.61	96.22
Drink	84.74	94.44	89.33
Eat	93.90	95.65	94.77
Getup or Sitdown	78.44	85.37	81.76
Read book	97.15	88.35	92.54
Use computer	99.34	94.06	96.63
Talk on phone	78.17	96.67	86.44
Use tablet	96.43	85.71	90.76
Walk	100.00	57.37	72.91
Watch TV	76.74	100.00	86.84

TABLE III: Per-class performance analysis of I3D on the SIMS4ACTION dataset: Precision, Recall and F1 (their harmonic mean) score are calculated for each individual class.

work architecture-related parameters (e.g. dropout, input size) were set according to the original architectures [14], [16]. We apply multiple augmentations during training including random rotations ($\pm 20^\circ$), hue changes ($\pm 180^\circ$), saturation changes ($\pm 100\%$), brightness changes ($\pm 80\%$) and random cropping. For testing, we only scale the input clips to the desired input size and take a center crop. Conceptually, we divide all videos into 90 frame chunks. During training, we randomly select single 30 frame clips from each chunk and classify them separately. For the final classification on videos, we only classify the centre chunk of each video. We use the Adam optimizer with a learning rate of $1e-4$ and weight decay $1e-4$. For training on the Sims dataset, we used a mini-batch size of 48 for S3D and 20 for I3D. For our transfer learning experiments on ADL, we exchanged the last layer to match the number of 31 classes and then trained on that layer for one epoch in order to fine-tune the parameters. We used a batch size of 64 chunks for S3D and a batch size of 32 chunks for I3D, all other settings are equal to the Sims training.

V. EXPERIMENTS

We now investigate how well off-the-shelf activity recognition models described in Section IV can leverage our synthetic SIM4ACTION data to recognize activities of daily living. To gather this insight, we present results for two scenarios: (1) recognition of ADL in the video game domain and (2) recognition transfer from video game-trained models to real-life ADL classification. We use the standard accuracy and the balanced multi-class accuracy (average of the individual class recalls) as our main evaluation metrics.

Model	Weight initialization	Normal Acc [%]			Balanced Acc [%]		
		Top 1	Top 3	Top 5	Top 1	Top 3	Top 5
Random Baseline	None	10.0	33.33	50.00	10.0	33.33	50.00
Separable 3D Net	Random	19.95	45.85	64.62	12.4	34.39	58.7
Separable 3D Net	Kinetics400	34.08	58.64	75.94	23.23	43.85	62.29
I3D Net	Random	18.02	44.39	66.06	10.91	31.49	53.58
I3D Net	Kinetics400	22.75	47.76	64.99	23.25	47.02	64.41

TABLE IV: Recognition results in the GAMING→REAL setting. Such transfer is a hard task for current recognition models, although all methods surpass the random baseline.

A. ADL Recognition on Video Game Data

First, we explore how well the described CNN architectures can handle the video game data in the GAMING→GAMING setting, where we train and evaluate in the synthetic domain. In Table II we present the results of the S3D and I3D networks both trained on the SIMS4ACTION dataset. Similarly to the Toyota Smarthome research [8], we initialize the networks using the Kinetics-400 classification models [31], but also report the results when training from scratch. As expected, the accuracy of the pretrained models is considerably higher, since 3D CNNs are known to strongly benefit from pretraining. When granting the models a small margin of error in the *Top 3* evaluation scenarios, they already come close to perfectly recognizing the synthetic ADLs. In addition to the overall accuracy, we examine model performance for the individual classes, with the exact precision, recall and F1-score (harmonic mean of the precision and recall) provided in Table III. Interestingly, the best recognition in terms of *F1*-score is achieved for the activities *cook* and *use computer*. This is likely due to the distinctiveness of the cooking utensils and the computer making them easily distinguishable patterns in those activities. *Get up/sit down* is quite ambiguous and hard to disentangle from other activities like *watch TV* or *eating* where the virtual human is also in a seated position resulting in lower precision.

B. ADL Recognition in the REAL-WORLD from Video Games

Our next area of investigation is the *transfer* scenario, where models trained on SIMS4ACTION are applied on real videos not present during training. The complementary real-world benchmark was derived from the Toyota Smarthome dataset [8] (see Figure 6 for images from both datasets), as described in Section III-B. The results of this experiment are reported in Table IV. Direct cross-domain recognition is an important but very hard task for modern data-driven algorithms and the performance also drops in our case. Still, all models clearly outperform the random baseline and 34% of cases were correctly identified by the I3D Net in the case of Kinetics-initialization. This experiment results provide encouraging evidence, that a cheap, synthetic data collection from sources such as life simulation video games is a promising direction for training ADL models. We believe, that such frameworks provide a pathway towards economical data collection in the future. As our results also demonstrate the sensitivity of modern ADL recognition

models to changes in data distribution, we hope that our dataset facilitates development of domain-agnostic models, which is an active topic, e.g., in autonomous driving [32], [33] but has been rather overlooked in ADL recognition.

VI. CONCLUSION

In this paper, we present SIMS4ACTION, a dataset for recognizing activities of daily living by learning from synthetic data recorded from life simulation video games. Our dataset covers over 10 hours of annotated video material with individual video clips set in different simulated environments, subjects of different age groups and ethnicities as well as multiple camera view-points. We adopt two prominent architectures for video classification to our task and train them to distinguish ten domestic behaviours with the videogame data only present during training. To provide an extensive benchmark for comparison, we formalize the GAMING→GAMING and GAMING→REAL evaluation settings, which we use to validate the models. Our dataset will be publicly released upon publication and we believe that SIMS4ACTION has strong potential to foster research on not only ADL itself but also on domain adaptation between synthetic and real-world activities and motivate the needed development of generalizable video understanding models.

Acknowledgements. This work was supported by the Competence Center Karlsruhe for AI Systems Engineering (CC-KING, www.ai-engineering.eu) sponsored by the Ministry of Economic Affairs, Labour and Housing Baden-Württemberg. This research was also supported by JuBot project sponsored by the Carl Zeiss Stiftung.

REFERENCES

- [1] E. G. Christoforou, A. S. Panayides, S. Avgousti, P. Masouras, and C. S. Pattichis, "An overview of assistive robotics and technologies for elderly care," in *Mediterranean Conference on Medical and Biological Engineering and Computing*. Springer, 2019, pp. 971–976.
- [2] M. Leo, A. Furnari, G. G. Medioni, M. Trivedi, and G. M. Farinella, "Deep learning for assistive computer vision," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018.
- [3] L. Marco and G. M. Farinella, *Computer vision for assistive healthcare*. Academic Press, 2018.
- [4] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+d: A large scale dataset for 3d human activity analysis," in *CVPR*, 2016.
- [5] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta, "Hollywood in homes: Crowdsourcing data collection for activity understanding," *ArXiv*, vol. abs/1604.01753, 2016.
- [6] J. Jang, D. Kim, C. Park, M. Jang, J. Lee, and J. Kim, "Etri-activity3d: A large-scale rgb-d dataset for robots to recognize daily activities of the elderly," *arXiv preprint arXiv:2003.01920*, 2020.
- [7] S. Das, S. Sharma, R. Dai, F. Bremond, and M. Thonnat, "Vpn: Learning video-pose embedding for activities of daily living," in *European Conference on Computer Vision*. Springer, 2020, pp. 72–90.
- [8] S. Das, R. Dai, M. Koperski, L. Minciullo, L. Garattoni, F. Bremond, and G. Francesca, "Toyota smarhome: Real-world activities of daily living," in *ICCV*, 2019.
- [9] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in *ICCV*, 2017.
- [10] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *European conference on computer vision*. Springer, 2016, pp. 102–118.
- [11] S. R. Richter, Z. Hayder, and V. Koltun, "Playing for benchmarks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2213–2222.
- [12] H. Wang and C. Schmid, "Action recognition with improved trajectories," in *International Conference on Computer Vision (ICCV)*, 2013, pp. 3551–3558.
- [13] A. Roitberg, N. Somani, A. Perzylo, M. Rickert, and A. Knoll, "Multimodal human activity recognition for industrial manufacturing processes in robotic workcells," in *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, 2015.
- [14] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [15] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *International Conference on Computer Vision*. IEEE, 2015.
- [16] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, "Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 305–321.
- [17] X. Wen, H. Chen, and Q. Hong, "Human assembly task recognition in human-robot collaboration based on 3d cnn," in *2019 IEEE 9th Annual International Conference on CYBER Technology in Automation, Control, and Intelligent Systems (CYBER)*. IEEE, 2019.
- [18] A. Roitberg, M. Haurilet, S. Reiß, and R. Stiefelhofen, "Cnn-based driver activity understanding: Shedding light on deep spatiotemporal representations," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2020, pp. 1–6.
- [19] P. Gebert, A. Roitberg, M. Haurilet, and R. Stiefelhofen, "End-to-end prediction of driver intention using 3d convolutional neural networks," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019.
- [20] A. Roitberg, C. Ma, M. Haurilet, and R. Stiefelhofen, "Open set driver activity recognition," in *2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2020, pp. 1048–1053.
- [21] A. Sharghi, H. Haugerud, D. Oh, and O. Mohareri, "Automatic operating room surgical activity recognition for robot-assisted surgery," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 385–395.
- [22] R. Dai, S. Das, S. Sharma, L. Minciullo, L. Garattoni, F. Bremond, and G. Francesca, "Toyota smarhome untrimmed: Real-world untrimmed videos for activity detection," *arXiv preprint arXiv:2010.14982*, 2020.
- [23] M.-Y. Wang, A. A. Kogkas, A. Darzi, and G. P. Mylonas, "Free-view, 3d gaze-guided, assistive robotic system for activities of daily living," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018.
- [24] P. Krähenbühl, "Free supervision from video games," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2018, pp. 2955–2964.
- [25] G. Varol, I. Laptev, C. Schmid, and A. Zisserman, "Synthetic humans for action recognition from unseen viewpoints," *International Journal of Computer Vision*, vol. 129, no. 7, pp. 2264–2287, 2021.
- [26] X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba, "Virtualhome: Simulating household activities via programs," in *CVPR*, 2018, pp. 8494–8502.
- [27] H. Hwang, C. Jang, G. Park, J. Cho, and I.-J. Kim, "Eldersim: A synthetic data generation platform for human action recognition in eldercare applications," *arXiv preprint arXiv:2010.14742*, 2020.
- [28] C. Guo, X. Zuo, S. Wang, S. Zou, Q. Sun, A. Deng, M. Gong, and L. Cheng, "Action2motion: Conditioned generation of 3d human motions," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2021–2029.
- [29] A. Blattmann, T. Milbich, M. Dorkenwald, and B. Ommer, "Behavior-driven synthesis of human dynamics," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*. IEEE, 2021, pp. 12 236–12 246.
- [30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [31] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.
- [32] E. Alberti, A. Tavera, C. Masone, and B. Caputo, "Idda: a large-scale multi-domain dataset for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5526–5533, 2020.
- [33] C. Ma, J. Zhang, K. Yang, A. Roitberg, and R. Stiefelhofen, "Densepass: Dense panoramic semantic segmentation via unsupervised domain adaptation with attention-augmented context exchange," in *International Conference on Intelligent Transportation Systems (ITSC)*, 2021.