

Inferring Human Intentions from Predicted Action Probabilities

LEI SHI, Institute for Visualisation and Interactive Systems, University of Stuttgart, Germany

PAUL-CHRISTIAN BÜRKNER, Cluster of Excellence SimTech, University of Stuttgart, Germany

ANDREAS BULLING, Institute for Visualisation and Interactive Systems, University of Stuttgart, Germany

Predicting the next action that a human is most likely to perform is key to human-AI collaboration and has consequently attracted increasing research interests in recent years. An important factor for next action prediction are human intentions: If the AI agent knows the intention it can predict future actions and plan collaboration more effectively. Existing Bayesian methods for this task struggle with complex visual input while deep neural network (DNN) based methods do not provide uncertainty quantifications. In this work we combine both approaches for the first time and show that the predicted next action probabilities contain information that can be used to infer the underlying intention. We propose a two-step approach to human intention prediction: While a DNN predicts the probabilities of the next action, MCMC-based Bayesian inference is used to infer the underlying intention from these predictions. This approach not only allows for independent design of the DNN architecture but also the subsequently fast, design-independent inference of human intentions. We evaluate our method using a series of experiments on the Watch-And-Help (WAH) and a keyboard and mouse interaction dataset. Our results show that our approach can accurately predict human intentions from observed actions and the implicit information contained in next action probabilities. Furthermore, we show that our approach can predict the correct intention even if only few actions have been observed.

Additional Key Words and Phrases: Bayesian inference, deep neural network, intention

ACM Reference Format:

Lei Shi, Paul-Christian Bürkner, and Andreas Bulling. 2023. Inferring Human Intentions from Predicted Action Probabilities. 1, 1 (August 2023), 22 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Everyday human activities are composed of individual actions that, if performed in the right sequence, allow humans to successfully perform even complex tasks. There is an increasing interest in artificial intelligence (AI) research to develop computational agents that are able to support and collaborate with humans on such tasks [7, 25, 38]. A key requirement for effective collaboration is that artificial agents understand human actions and can plan their own actions accordingly, such as to maximise utility [12]. A number of recent works have therefore focused on predicting future actions based on observed previous action sequences [1, 10, 11, 13, 17]. Predicting future actions is challenging because even if observed past actions are similar, the next action can differ due to a

Authors' addresses: Lei Shi, Institute for Visualisation and Interactive Systems, University of Stuttgart, Germany, lei.shi@vis.uni-stuttgart.de; Paul-Christian Bürkner, Cluster of Excellence SimTech, University of Stuttgart, Germany, paul.buerkner@gmail.com; Andreas Bulling, Institute for Visualisation and Interactive Systems, University of Stuttgart, Germany, andreas.bulling@vis.uni-stuttgart.de.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

XXXX-XXXX/2023/8-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

different underlying intention [9, 30]. Consequently, understanding human intentions is crucially important for next action prediction, and thus effective human-AI collaboration [23, 26, 45].

One line of previous work on intention prediction and next action prediction has focused on Bayesian approaches [2, 3, 34]. A key limitation is that these probabilistic models were geared for the game domain and do not easily generalise to other settings and tasks. If the method should be applied to a different task domain, the models had to be changed substantially to adapt to the new task. Additionally, probabilistic models would not be effective if the input is more complex form such as images and videos. On the other hand, Deep Neural Networks (DNNs) have also been used to for intention prediction and next action prediction [8, 15]. Deploying DNNs can in principle handle the intention prediction in very complex scenarios. However, one shortcomings of DNNs is that they cannot easily quantify the epistemic uncertainty in the prediction, which is crucial in safety-critical applications where wrong predictions could lead to severe consequences. Moreover, the uncertainty of the predicted intentions can also provide additional information to help planning future actions and collaboration in Human-AI collaboration scenarios.

Despite the importance of understanding intentions in collaborative settings, previous work has mainly focused on either using Bayes-only or DNN-only approaches. Both types of approach integrate next action prediction and intention prediction together. A model that combining DNN-based and Bayesian-based method together could have the advantages of both. This can benefit practical applications in two aspects. First, collaborative AI system needs to operate in the real world and it needs to deal with the data with high (visual) complexity. Deploying DNNs can easily adapt to the complex input. And second, the uncertainty quantification about the prediction of intention can help better providing decisions on the future actions in Human-AI collaboration and reduce the risk of a potential wrong prediction.

The DNN-based next action prediction is typically treated as a classification task, i.e., the next action predicted is the one with the highest probability across the full action space [29]. However, the probabilities of all actions in the full action space may contain implicit information that can be used for uncertainty-aware intention prediction. This has not been explored in DNN-based methods so far. In this work, we propose a novel two-step procedure to infer human intentions from sequences of actions. Our approach combines DNNs to obtain the probabilities of the next action with MCMC-based Bayesian inference for subsequently inferring intentions. The DNNs are trained to classify the next action from all possible actions in the action space. Specifically, given action data from N different intention, we train N DNN models for next action prediction. At test time, one action sequence data is fed to all N models to obtain the action probabilities, the action sequence represents the true intention. The output from N DNNs represents the probabilities assuming N intentions are applied. Next, we use Markov Chain Monte Carlo (MCMC) sampling to train a Bayesian model with all action probabilities from all DNNs to infer the intentions. Our two-step method decouples the next action prediction (DNNs) and actual intention prediction (Bayesian model). We do not have any requirement on the DNN input format and network architecture. They can be modified and optimised according to different tasks. The Bayesian model is independent on the DNNs, it only takes the action probabilities from DNNs and predict the intentions with uncertainties.

We demonstrate the effectiveness of our method through experiments on two datasets: Watch-And-Help (WAH) [25] and keyboard and mouse interaction dataset [44]. In the WAH dataset, a computational agent is performing household tasks, such as putting food in a fridge, in a virtual environment. These tasks are composed of actions, such as walk to the table, take an apple etc. Then intention is the task the agent is performing. In the keyboard and mouse interaction dataset, participants are interacting with computers by formatting text generated randomly. The participants follow the given rules to format the text. A rule consists of performing actions (e.g. make font

bold, italic) on title, sub-title and paragraph. The intention in this dataset is the rule given to the participant. We choose the two datasets to showcase that our method works in two different settings of scenario, i.e. computational agents in virtual environments and humans in real life interaction. Our experimental results show that our method can correctly predict the intentions of users. We further evaluate the performance of our method with 10% to 100% of observed actions in one action sequence on both datasets. The results show that using 20% of actions in a sequence, is often sufficient for the true intention to have clearly higher posterior probability than all other intentions, although with substantially uncertainty. This demonstrates that our method can infer the users' intention already in an early stage where only few actions have been observed.

The main contribution of this work is the two-step method to infer human intention. We use action probabilities of the next action prediction from DNNs with a Bayesian model to infer the intention. To our best knowledge, we are the first ones to propose the joint use of DNNs and Bayesian models to decouple the next action prediction and intention prediction. Our method has three advantages. First, the DNNs and the Bayesian inference are decoupled. The inference of intention does not depend on the DNN architecture. One can optimise the DNN architecture for classifying the next action separately. Second, training the Bayesian model requires less time and Bayesian inference provides a fast prediction on intention. Third, our method can predict the intention correctly and efficiently when using few observed actions in the series of actions.

2 RELATED WORK

Human-AI collaboration has attracted increasing interest recently. Several works focused on developing computational agents for collaboration in virtual environments [20, 26, 39]. The virtual environments possess near-realistic scenes and objects and supports different types of actions. Puig et al. [25] pointed out that understanding human intention and predicting next actions are important in Human-AI collaboration. They developed a two-stage collaboration scheme where an helping agent watches another agent demonstrating an activity in the first stage and collaborating with it in the second stage. The importance of intention prediction in collaboration has also been shown in real-life scenarios [42]. Correctly predicted human intentions lead to more effective collaborations in robotic shared autonomy [19], Human-Robot handover [41] and cooperative assembly [24].

Several prior works have focused on action anticipation based on videos, i.e., the task of predicting future actions based on observed behaviour in the past [13, 14, 27]. Different types of models have been used, e.g. two-stream CNN [14], LSTM [13], video transformer [16], or graph neural networks [43]. In [5], the authors further used label smoothing technique to improve the work [13]. Other works have also used goal in anticipating the future actions. In [30], the authors used latent goal to improve the action anticipation. The latent goal was defined as the visual representation after the final action. The latent goal computation was integrated into the model architecture to help anticipate the next action. In [9], procedure planning was introduced in video-based action anticipation. Given the current action and the visual goal (visual representation of final action), the authors encoded the visual observation and the action into two separate latent space and planned the future actions based on Markov decision process (MDP). These two works have brought the intention components in predicting the future actions. However they do not predict the intention, rather they use the intention to improve the prediction of future actions.

In [18], the intention was defined as the intended ingredients for a sandwich. SVMs were used to predict the intention from human gaze data. SVM was also used in intention prediction during human interactions [4]. Other approaches include MDP [21], probabilistic graphical model [36], k-nearest neighbour (kNN) [28]. In [44], the authors investigated the task of predicting user intents from mouse and keyboard input as well as gaze behaviour. In another line of work, gaze behaviour

was also identified as a rich source of information for predicting users' search intents [31–33] and even visually reconstructing it [37]. Perhaps the works in [22, 35] are the most similar to ours. In [35], a Bayesian model to infer intentions from ontic actions and gaze action. The ontic actions are the actions which change the state of the world and the gaze actions are the regions where an agent is looking at with regard to the ontic action. The work in [22] further introduced a deceptive component into the Bayesian model for the scenario where the human might perform ambiguous actions on purpose. Although Bayesian models were used for intention prediction, the actions were not predicted by neural networks. Rather they were pre-processed and then used in the Bayesian models.

3 METHOD

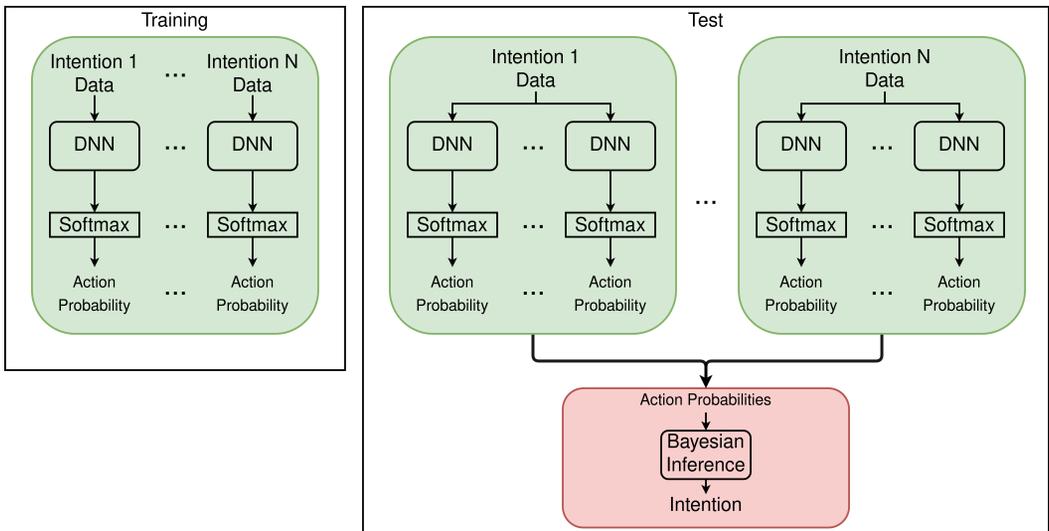


Fig. 1. Overview of the proposed method. We train one deep neural network (DNN) for each intention. Softmax layers are used to obtain the probability of next actions. The Bayesian inference model takes the action probabilities from all activities and infers the intention.

Figure 1 provides an overview of our method. For each of the N possible intentions, we train a separate DNN on data where the ground-truth intention is known. The task of the DNNs is to predict the next action from all previous actions. Since our method works with arbitrary DNN architectures that perform this prediction task, we are not focusing on the specific architecture of the DNN here. It is important, however, that each DNN has a final Softmax (or equivalent) layer to obtain the predicted next-action probabilities. All actions probabilities are then used to train a Bayesian model to predict the intention from the set of predicted next-action probabilities (see below for details). At test time, the data of each intention is forwarded to all DNNs to obtain N next-action probabilities representing N intentions. We refer the action sequence forwarded to the DNNs with known intention label as the true intention. The N DNNs are trained with N intentions and we interpret each DNN as an assumed intention, i.e. given one action sequence, the DNNs do not know the true intention, the i th DNN assume it is from the i th intention. The Bayesian model then uses all action probabilities jointly to infer the posterior distribution over the N assumed intention. Specifically, the Bayesian only use the action probabilities from one action sequence to

infer user intention. All DNNs can be trained separately as they do not share weights for our procedure to work (but they could if this was beneficial).

Formally, an intention \mathcal{I} consists of a series of actions, denoted as

$$\mathcal{I}_{ij} = [a_0, \dots, a_L], 0 < i < N, 0 < j < M, 0 < k < L, \quad (1)$$

where a is action, M is the number of action series belonging to i^{th} intention and L is the total number of actions in \mathcal{I}_{ij} . \mathcal{I}_{ij} represents an action series instance in the i th intention. In the training of i th DNN, we use all instances in \mathcal{I}_i as training data. The input of for the DNN are the \mathcal{I}_{ij} , whereas the ground-truth next action $y_i = a_{k+1}$ constitutes the target variable. The loss function is then

$$\mathcal{L} = f(y, \hat{y}), \quad (2)$$

where \hat{y} is the DNN prediction and $f(\cdot)$ is a cross entropy loss. After all networks are trained, each intention data is passed to all N networks and obtain N Softmax outputs. The i^{th} Softmax outputs produced by the i^{th} DNN represents the action probability assuming i^{th} intention is applied.

To infer the intention of human based on the series of actions and their DNN predictions serving as a surrogate likelihood, we set-up the following Bayesian model:

$$\begin{aligned} a_k &\sim \text{categorical}(\theta_k), \forall k = 1, \dots, L, \\ \theta_{km} &= \sum_{i=1}^N P(a_{km} | \mathcal{I}_i) P(\mathcal{I} = \mathcal{I}_i), \forall m = 1, \dots, M, \\ P(\mathcal{I}) &\sim \text{Dirichlet}(\alpha), \end{aligned}$$

where $P(\mathcal{I} = \mathcal{I}_i)$ is the i^{th} element of the intention probability $P(\mathcal{I})$ to be inferred by the model, and $\alpha \in \mathbb{R}_+^N$ is the concentration vector of the Dirichlet prior on $P(\mathcal{I})$, which we set to $\alpha = 1$ to obtain an uninformative prior. The action probabilities $P(a_{km} | \mathcal{I}_i)$ of the m th possible action to occur at the k th position in the sequence are obtained from the output of DNNs. To predict the intention $P(\mathcal{I})$, we use the probabilistic programming language Stan [6], which employs a state-of-the-art Markov-chain Monte-Carlo (MCMC) sampler to $P(\mathcal{I})$.

4 EXPERIMENT

4.1 Datasets

4.1.1 Watch-And-Help Dataset. WAH is a dataset for social intelligence and human-AI collaboration [25]. In the dataset, an AI agent Bob helps another human-like agent Alice perform household activities. The world is a 3D virtual environment. There are two stage of collaboration, i.e., the *Watch* stage and the *Help* stage. In the *Watch* stage, Bob observes Alice demonstrating an activity and Bob helps Alice with the same activity in the *Help* stage. In this work we only consider the *Watch* state given that we are interested in inferring the intentions of Alice. We understand the activities are defined by a set of sub-goals represented by predicates. For instance, the sub-goal *ON*(cupcake, coffeetable) means that the cupcake should be on the coffee table and it belongs the activity read book. Both agents can perform different actions to accomplish their goals. An action can either be navigation (e.g., walk to fridge) or interaction with an object (e.g., grab an apple). An activity is accomplished once the states of all sub-goal predicates are reached. In total, there are five types of activities with each activity having two to eight sub-goals. The total number of potential actions is 79. The dataset has one training and and two test sets, i.e., test set 1 and test set 2. The training set contains data of 1,011 action sequences with different goals and both test sets contain data of 100 action sequences each.

	Put fridge	Put dishwasher	Read book
Training set	52,752	24,306	40,786
Test set 1	3,443	8,591	11,557
Test set 2	275	2,149	2,670

Table 1. Final size of the train set and test set in three intentions in the Watch-And-Help dataset.

To evaluate our method we only need information on the activity and actions and thus leave the sub-goals aside. Furthermore, to keep the activity category consistent, we focus only on those types of activities that are present in training set, test set 1, and test set 2: (*put fridge*, *put dishwasher*, and *read book*). We treat the activity as the intention and predict the intentions from sequence of actions. Since we use the DNN to predict the next action in an action sequence, we modify the original action sequences for the use of next action prediction. For an action sequence $[a_0, \dots, a_L]$, when a new action is observed, we create a new action sequence. The final sizes of the training set, test set 1 and test set 2 are shown in Table 1.

4.1.2 Keyboard and Mouse Interaction Dataset. To complement the household activities performed in the virtual environment in the WAH dataset, we also evaluated our method on keyboard and mouse interaction dataset introduced in [44]. In the dataset, 16 participants were asked to format text according to several formatting rules (the interaction intentions). The evaluation task on this dataset was to predict these interaction intentions from mouse and keyboard input. The text consisted of titles, subtitles and paragraphs and a rule contained instructions on how to format it using the mouse and keyboard (e.g. "make the title bold"). Participants could perform seven different actions for formatting the text, i.e. bold, italic, underline, text size, font family, indentation and alignment. The dataset contains data from two types of formatting tasks: First, participants were asked to perform formatting according to seven predefined formatting rules. Each rule was repeated five times. Second, each participant was asked to create a custom rule themselves and to format the text according to this rule. We only used data from the first part of the dataset for our experiment since there is only one intention for each participant in the second part. We used the data from participants one to 11 for training and the data from participants 12 to 16 for testing.

4.2 Experimental Settings

We first evaluated our method on inferring users' intentions on the WAH dataset. We used 10% to 90% of the actions in an action sequence with a 10% step to infer the intentions. Since the WAH dataset is created in a virtual environment and the action sequences do not belong to any user, we created virtual users by randomly grouping the data in test set 1 and test set 2. As a result, test set 1 had 92 artificial users, each user had one action sequence in *put fridge*, two actions sequences in *put dishwasher*, and three action sequences in *read book*. Test set 2 had nine users, each user had one action sequence in *put fridge*, five actions sequences in *put dishwasher*, and five action sequences in *read book*.

The architecture of our DNN model is based on the one in [26]. The architecture of DNN is shown in Figure 2. We padded the action sequence to the maximum length of action sequences in the dataset. The number of embeddings was set to the length of the action sequence. The features of the last cell of LSTM layer in the mini-batch were stacked. Finally two FC layers were used as prediction head. We read-out both the predicted next action and the Softmax probability. To train the DNNs on the WAH dataset, we used 2,000 epochs, a batch size of 32 and a learning rate of $3e^{-4}$.

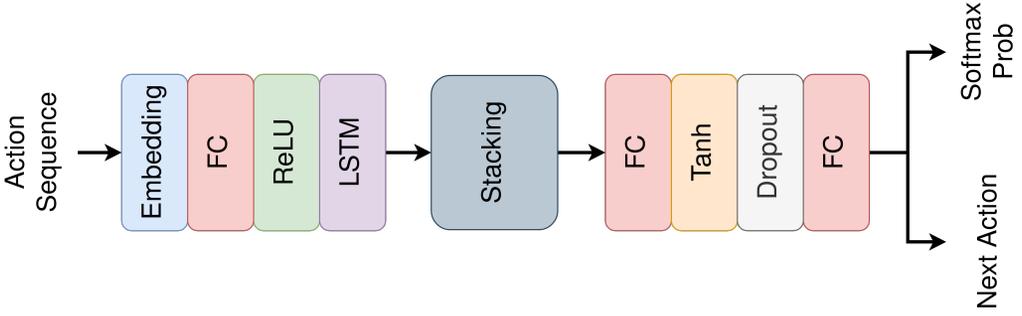


Fig. 2. The architecture of the DNN used in experiments.

For the keyboard and mouse interaction dataset, we trained for 100 epochs, the batch size was eight, the learning rate was $1e^{-4}$.

To train the Bayesian model we used the same training strategy for both datasets. For each action sequence, we performed Bayesian inference via four MCMC chains, each with 2,000 iterations of which the first 1,000 were discarded as warmup. All Bayesian models converged well according to standard convergence criteria [40]. In section 5, we exemplarily show one randomly selected action sequence of each user from each true intention to perform experiments. Results for all individual action sequences are shown in section 8 in the appendix. Note that in both datasets the baseline methods [25, 44] are based on machine learning models that do not provide uncertainty-aware predictions. For fairness, we therefore did not compare with these baselines in this work. Additionally, there are not variant components in the Bayesian model that are sensitive to the results, hence we did not perform ablation study.

5 RESULTS

5.1 User Intention Prediction

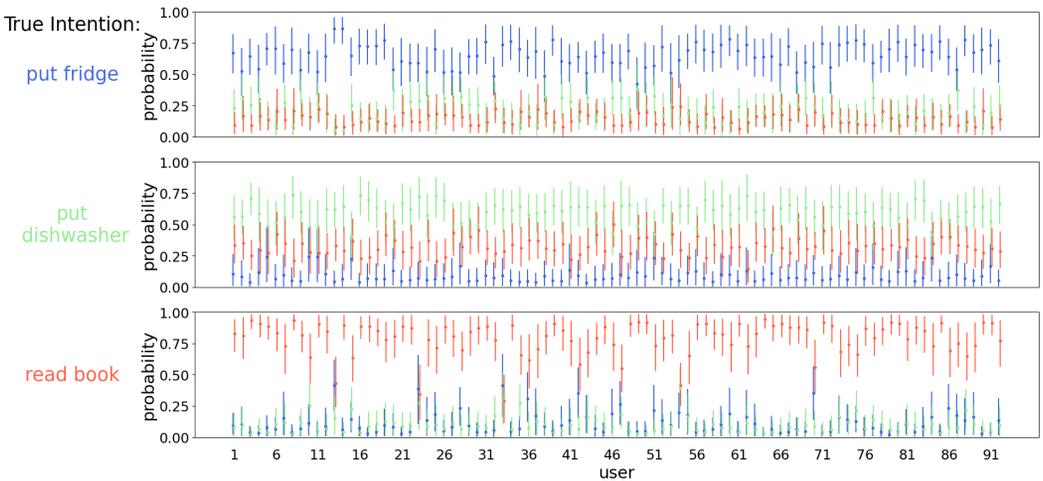


Fig. 3. Result of intention prediction of users on test set 1 in WAH dataset.

Figure 3 shows the result of user intention prediction on the test set 1 of the WAH dataset. We report the posterior mean and 90% credible intervals (CIs) of the probabilities of all assumed intentions. The top, middle and bottom plot shows the results when the true intention is *put fridge*, *put dishwasher* and *read book*. For the true intention *put fridge*, for most users our method can predict the correct intention, meaning that the assumed intention with the highest posterior mean probability is the same as the true intention. In a few cases, the posterior mean probability of put fridge is close to put dishwasher or read book. From Table 2 we can see that the posterior mean of put fridge and put dishwasher for user 49 is 0.43 and 0.39. When the true intention is *read book*, our method can also predict the correct intentions of most users with a few exceptions, i.e. user 13, 23, and 33 (Table 2). For user 13, the posterior mean of read book 0.43, only 0.02 higher than put fridge. For user 23 and 33, the posterior mean of put fridge is slightly higher than read book. For the true intention *put dishwasher*, although the posterior mean of put dishwasher are the highest in most users predictions, the difference between put dishwasher and the other two assumed intentions are smaller compared to the cases in true intention *put fridge* and *read book*. For user 6, 19, 20, 44, 49, 51, 53, the difference between the posterior mean of put dishwasher and the posterior mean of read book are around 0.1. For user 15, 27, 30, 46, 50, 56, 65, 71, and 84, the differences are below 0.07. Overall, our model can predict the correct true intention *put fridge*, however the prediction for user 49 is rather uncertain. For true intention *read book*, predictions for most users are correct but more predictions are more uncertain. The model can predict users true intention *put fridge* and *read book* better than *put dishwasher*.

User	<i>Put Fridge</i>			<i>Put Dishwasher</i>			<i>Read Book</i>		
	put fridge	put dishwasher	read book	put fridge	put dishwasher	read book	put fridge	put dishwasher	read book
6	0.71	0.09	0.20	0.06	0.51	0.42	0.06	0.11	0.83
13	0.87	0.05	0.08	0.05	0.62	0.33	0.41	0.15	0.43
15	0.65	0.25	0.10	0.23	0.41	0.37	0.15	0.22	0.63
19	0.77	0.12	0.10	0.07	0.52	0.41	0.10	0.13	0.78
20	0.54	0.37	0.09	0.13	0.49	0.38	0.08	0.11	0.81
23	0.59	0.29	0.12	0.07	0.72	0.21	0.39	0.27	0.34
27	0.52	0.30	0.17	0.08	0.48	0.43	0.08	0.11	0.81
30	0.65	0.26	0.09	0.05	0.49	0.45	0.04	0.09	0.87
33	0.74	0.15	0.11	0.09	0.62	0.29	0.42	0.30	0.29
44	0.64	0.15	0.20	0.06	0.52	0.42	0.07	0.14	0.79
46	0.60	0.31	0.09	0.07	0.43	0.50	0.19	0.07	0.74
49	0.43	0.38	0.19	0.12	0.49	0.39	0.05	0.03	0.92
50	0.56	0.24	0.20	0.09	0.52	0.40	0.05	0.03	0.92
51	0.57	0.33	0.10	0.23	0.42	0.34	0.22	0.05	0.73
52	0.72	0.17	0.10	0.11	0.65	0.23	0.14	0.07	0.79
53	0.51	0.25	0.24	0.07	0.51	0.42	0.10	0.09	0.82
54	0.61	0.14	0.24	0.05	0.65	0.30	0.20	0.39	0.41
56	0.76	0.15	0.09	0.13	0.47	0.40	0.04	0.06	0.90
65	0.65	0.17	0.18	0.11	0.42	0.47	0.06	0.04	0.89
71	0.74	0.17	0.09	0.06	0.49	0.45	0.04	0.04	0.92
84	0.68	0.16	0.16	0.23	0.42	0.35	0.16	0.11	0.73

Table 2. Intention prediction results on test set 1 in WAH dataset. We show the posterior mean probabilities of the users which are uncertain or the highest posterior mean is not the correct intention. *Italic* indicates the true intention. Under each true intention is the assumed intention.

Figure 4 shows the posterior mean and 90% CIs of the probabilities of user intention prediction on the test set 2 of WAH dataset. For the true intention *put fridge*, the posterior mean of put fridge are higher than put dishwasher and read book. Except user 2 and 9, the posterior mean probabilities

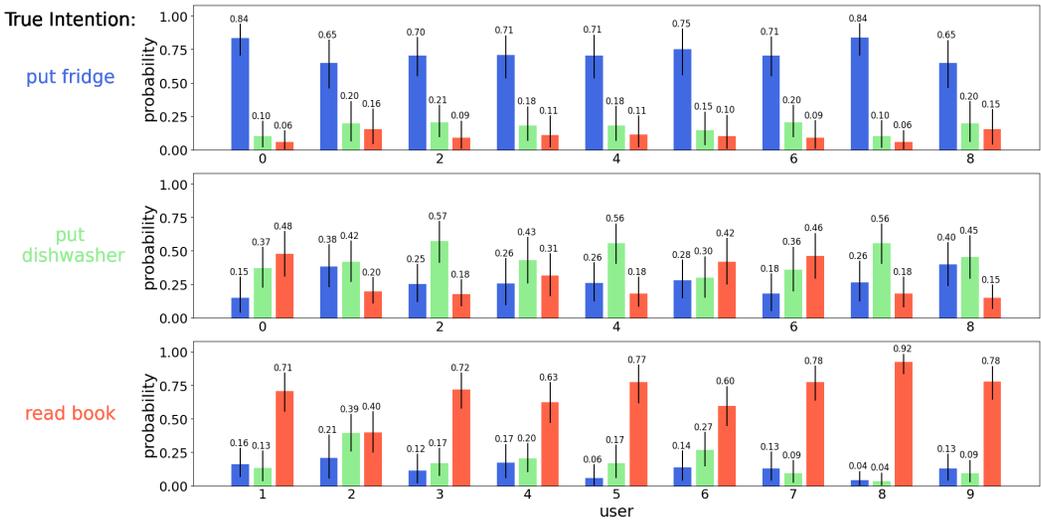


Fig. 4. Result of intention prediction of users on test set 2 in WAH dataset.

of the rest users are all over 0.7. The highest posterior mean of put fridge and read book is 0.21 and 0.15 respectively. The Bayesian model can accurately distinguish the correct intention from the assumed intentions. For the case of true intention *read book*, the posterior mean probabilities of user 1, 3, 5, 7, 8, and 9 are all above 0.7, while the numbers of user 4 and 6 are slightly lower. The prediction of user 2 is not as good as the rest, the posterior mean is 0.4 for read book and 0.39 for put dishwasher. For the true intention *put dishwasher*, the predictions are much more uncertain. Three (user 1, 6, and 7) out of nine users' posterior mean of read book are higher than put dishwasher. For user 2 and 9, although the posterior mean of put dishwasher is still the highest, it is very close to the posterior mean of put fridge (0.42 vs. 0.38 for user 2 and 0.45 vs. 0.4 for user 9). Only three (3, 5, and 8) users' posterior mean probabilities of put dishwasher are higher the 0.5. Overall, our method can predict the correct intentions of most users except user 2, when the true intentions are *put fridge* and *read book*. When identifying the intention in the case of *read book*, the prediction is more uncertain. For the true intention *put dishwasher*, the model only predicts the correct intention relatively confident for user 3, 5, and 8. The rest are either uncertain about the prediction or the wrong intention is predicted.

Figure 5 shows the posterior mean and 90% CIs of the predicted intentions in the keyboard and mouse interaction dataset. The prediction on all user data on all rules are correct in term of the highest posterior mean of the assumed intention being the true intention. For user 1, 2,3 and 5, the differences of the posterior mean between the correctly predicted intention and the rest intentions in all true intentions are quite large. For user 4, the posterior of the correct intention for true intention *rule 4* and *rule 7* are more uncertain than the other true intentions. For true intention *rule 4*, the posterior mean of rule 4 is 0.41 while the posterior mean of rule 1 is 0.27. For true intention *rule 7*, the posterior mean probabilities of rule 7 and rule 6 are 0.45 and 0.22 respectively.

5.2 Different Lengths of Observed Actions in Action Sequences

Figure 6 shows the results when different percentages of observed actions in one action sequence are used for inferring intention on test set 1 in WAH dataset. The percentage of the action in a sequence is varied from 10% to 100% in steps of 10%. For instance, 50%, we extract the first 50% of

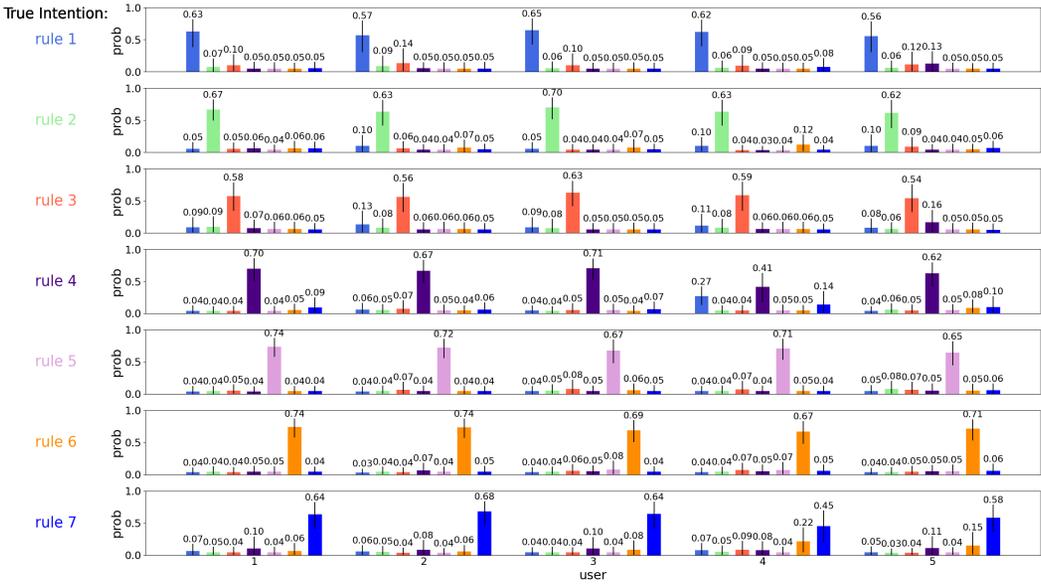


Fig. 5. Result of intention prediction of users in keyboard and mouse interaction dataset.

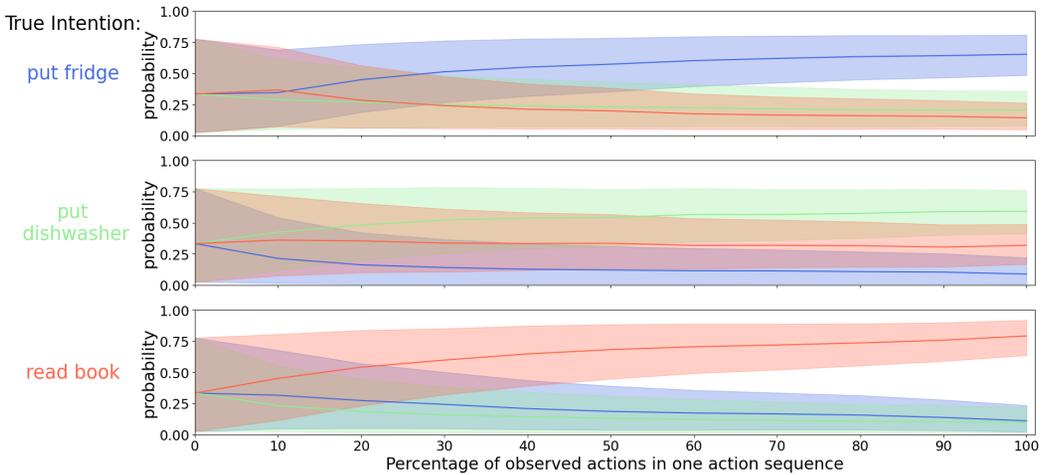


Fig. 6. Posterior mean probabilities and CI bounds when different percentages of observed actions in an action sequence are used for inference. The results on test set 1 in WAH dataset are shown.

the actions in the sequence and use it to infer the intention. We plot the average posterior mean probabilities and 90% CIs, i.e. at each percentage of observed actions in one sequence, the values of posterior mean and CI bounds are the average value from all users. For all true intentions, the posterior mean probabilities of the correct intentions are already relatively higher than the other assumed intentions when 20% of action in a sequence has been observed. The posterior mean probabilities increase with the increase of observed actions in action sequences. The CI bounds decrease as the percentage of observed actions increases. This shows that the more actions have

been observed in one action sequence, the more certain the Bayesian model is about its intention predictions.

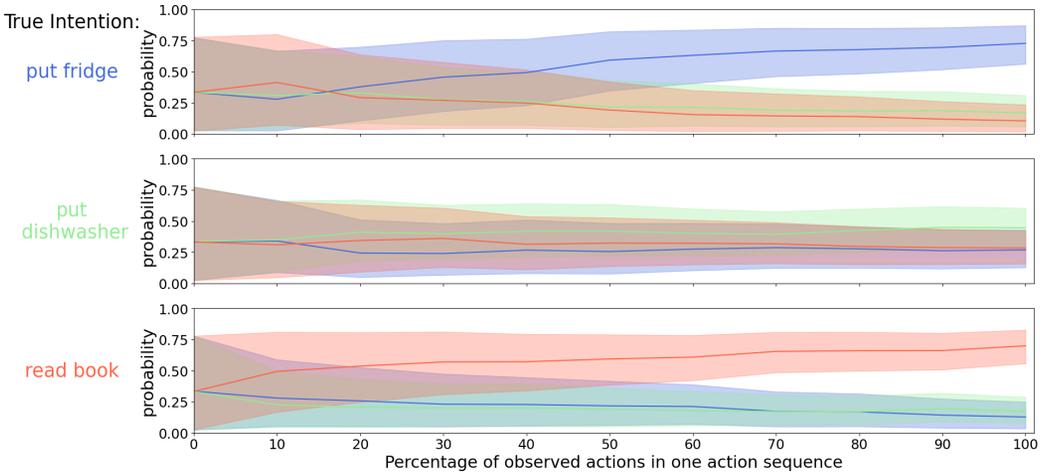


Fig. 7. Posterior mean probabilities and CI bounds when different percentages of observed actions in an action sequence are used for inference. The results on test set 2 in WAH dataset are shown.

Figure 7 shows the results when different percentage of observed actions in an action sequence are used on test set 2 in WAH dataset. The average posterior mean and 90% CIs of all users are plotted for each increase of the percentage of actions in an action sequence. When 10% of actions have been observed in one sequence, only in the true intention *read book* the Bayesian model is able to infer the true intention correctly. With more actions observed in the action sequence, the posterior mean probabilities in true intention *put fridge* and *read book* also increase, and the CI comes smaller. For the true intention *put dishwasher*, the posterior mean probabilities from 20% to 100% are comparable. The CIs however become smaller even though the posterior mean keeps on the same level.

In both test set 1 and test set 2, when the true intention is *put dishwasher*, the predictions of Bayesian models are more uncertain than the other two true intentions, especially in the test set 2. This can be observed when predicting user intention using full action sequence (Figure 4) and partially observed action sequence (Figure 7). We interpret that it is due to the noisier distribution of the actions in action sequences and the predictions of the DNNs. Figure 8 shows the distribution of representative action labels in train set, test set 1 and test set 2 of WAH dataset. By representative actions we mean the actions from which the intention can be easily interpreted. For instance, *open fridge* is a representative action for the intention *put fridge*. For the true intention *put dishwasher*, the number of action *walk dishwasher* in test set 2 is much less than the number in train set. It means *walk dishwasher* happens less in the action sequences in test set 2. This could be the reason that the predicted intention is more uncertain. Taking a look at the true intention *put dishwasher* in test set 1 (Figure 6) and test set 2 (Figure 7), the performance on test set 1 is better than on test set 2. And the number of action *walk dishwasher* in test set 1 is comparable to the number in train set.

Figure 9 shows the result on the keyboard and mouse interaction dataset. We show the average posterior mean probabilities and the 90% CIs of all seven intentions when different percentage of actions in one action sequence are used for inference. For all true intentions, the posterior mean probabilities of the correct intentions increase with more actions having been observed in the

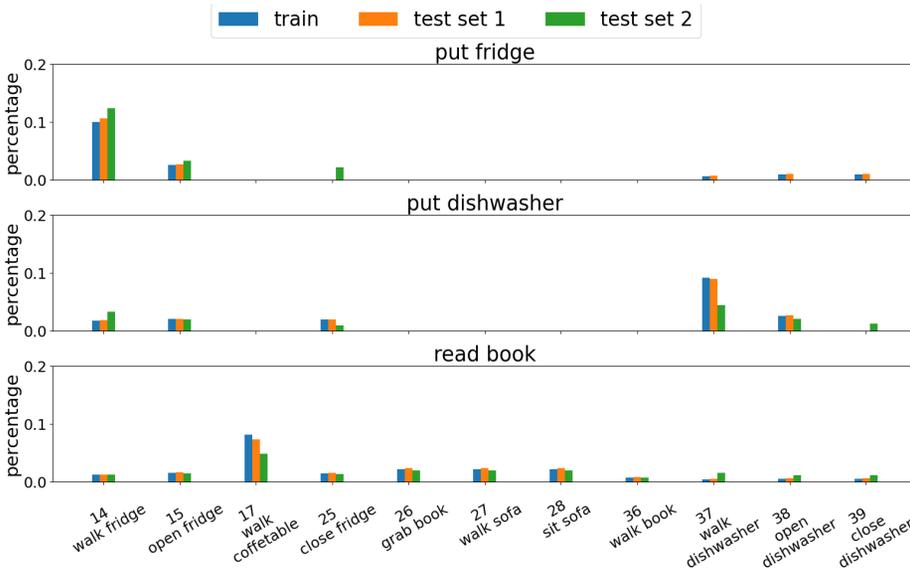


Fig. 8. Distribution of the representative action labels on train set, test set 1 and test set 2 in WAH dataset.

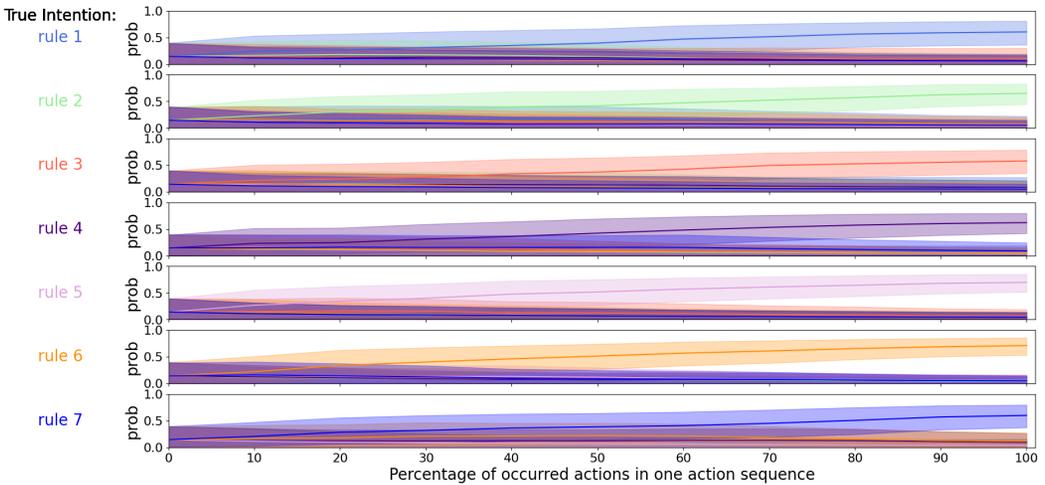


Fig. 9. Posterior mean probabilities and CI bounds when different percentages of observed actions in an action sequence are used for inference. The results in keyboard and mouse interaction dataset are shown.

action sequences. At 10%, the assumed intention with highest posterior mean probability is the same as true intention for all rules but the differences are small. The predictions at this are still quite uncertain. At 20%, the differences in true intention *rule 2*, *rule 5*, and *rule 6* become larger, but the CI bounds still remain at wide ranges. For true intentions *rule 5* and *rule 6*, the probabilities of the correct intentions are close to 0.5, however, the CIs has not decreased a lot. At 50%, the CIs of the correct intention already no longer overlaps with the CIs of the other intentions.

6 DISCUSSION

We proposed a two-step procedure to predict intentions from a sequence of user actions. We first trained DNNs to predict the next action from a sequence of observed actions. Based on the DNN output action probabilities, a Bayesian model then inferred the intention. The training of the DNNs and the Bayesian model is independent, so any of the two components can be exchanged without affecting the other. Moreover, the time for training and inference of the Bayesian model is negligible (on the level a few seconds) compared to the time to train DNNs. We evaluated our proposed method on two datasets i.e. WAH dataset and keyboard and mouse interaction dataset. The WAH dataset was generated in a virtual environment and the interactions were between computational agents. In keyboard and mouse interaction dataset, the interaction data between humans and computers were collected. The interactions were between real humans and computers.

In the evaluation of WAH dataset, we manually created artificial users by assigning action sequences to users. On test set 1, most intentions of users were inferred correctly. On test set 2, our method performed well on true intention *put fridge* and *read book*, the predictions were more uncertain for true intention *put dishwasher*. In the keyboard and mouse interaction dataset, all predictions of all users for all true intentions were correct. We used one action sequence from one user to perform Bayesian inference. This shows the Bayesian model is efficient for inferring intention in terms of the number of observations of action sequence. It does not have to see multiple action sequences to infer the correct intention, only seeing one action sequence is adequate.

We were also interested in how the Bayesian model performs on the data with fewer actions being observed in the action sequences. An intuition is that the model is more confident about the inferred intention when more actions have been observed. This is confirmed by experiments in two aspects. First, the posterior mean probabilities increase when more actions are observed. Second, the ranges of CI bounds become smaller meaning the Bayesian model is more certain about its predictions. Additionally, the Bayesian model can predict the true intentions correctly even at an early stage in the action sequence. Being able to predict human/agent intention in an early stage can benefit agent-agent and human-agent interaction. For instance, in the WAH scenario, an agent can help with the other agent finishing a task by complete other actions in the same task. In the scenario of keyboard and mouse interaction, the computer/agent can optimise the user interface or give suggestions while the human is formatting the text. It is necessary that the agent has enough time to plan and deploy collaboration and interaction. To be able to predict intentions when only partial actions in an action sequence allows the agent to have sufficient time for planning. It is worth noting that the uncertainties of the predicted intention in early stages are relatively high and this should be taken into consideration when designing the interaction with a human.

7 CONCLUSION

In this work we proposed a two-step procedure to infer human intentions from a series of actions based on DNNs and Bayesian inference. In a first step we trained DNNs to obtain the predicted probabilities of the next action in a sequence. In a second step we used MCMC-based Bayesian inference to infer the human intention from the predicted next-action probabilities. We performed experiments on the WAH and keyboard and mouse interaction datasets to validate our approach on both virtual environment and real life scenario. The results demonstrate that we can accurately infer the intentions correctly on both datasets even when only one action sequences from one user is available at inference time. This suggests that the implicit information contained in the next action probabilities generated by DNNs can be used to infer the intention using a Bayesian model. In addition, we demonstrated that our approach still provides correct predictions even if only few

actions have been observed. Prediction uncertainty then decreases further as more actions in the sequence become available.

REFERENCES

- [1] ABU FARHA, Y., RICHARD, A., AND GALL, J. When will you do what?-anticipating temporal occurrences of activities. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 5343–5352.
- [2] ALBRECHT, D. W., ZUKERMAN, I., AND NICHOLSON, A. E. Bayesian models for keyhole plan recognition in an adventure game. *User modeling and user-adapted interaction* 8 (1998), 5–47.
- [3] BAKER, C. L., JARA-ÉTTINGER, J., SAXE, R., AND TENENBAUM, J. B. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour* 1, 4 (2017), 0064.
- [4] BEDNARIK, R., VRZAKOVA, H., AND HRADIS, M. What do you want to do next: a novel approach for intent prediction in gaze-based interaction. In *Proceedings of the symposium on eye tracking research and applications* (2012), pp. 83–90.
- [5] CAMPORESE, G., COSCIA, P., FURNARI, A., FARINELLA, G. M., AND BALLAN, L. Knowledge distillation for action anticipation via label smoothing. In *2020 25th International Conference on Pattern Recognition (ICPR)* (2021), IEEE, pp. 3312–3319.
- [6] CARPENTER, B., GELMAN, A., HOFFMAN, M. D., LEE, D., GOODRICH, B., BETANCOURT, M., BRUBAKER, M., GUO, J., LI, P., AND RIDDELL, A. Stan: A probabilistic programming language. *Journal of Statistical Software* 76, 1 (2017), 1–32.
- [7] CARROLL, M., SHAH, R., HO, M. K., GRIFFITHS, T., SESHIA, S., ABBEEL, P., AND DRAGAN, A. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems* 32 (2019).
- [8] CASAS, S., LUO, W., AND URTASUN, R. Intentnet: Learning to predict intention from raw sensor data. In *Conference on Robot Learning* (2018), PMLR, pp. 947–956.
- [9] CHANG, C.-Y., HUANG, D.-A., XU, D., ADELI, E., FEI-FEI, L., AND NIEBLES, J. C. Procedure planning in instructional videos. In *European Conference on Computer Vision* (2020), Springer, pp. 334–350.
- [10] CRAMER, M., KELLEN, K., AND DEMEESTER, E. Probabilistic decision model for adaptive task planning in human-robot collaborative assembly based on designer and operator intents. *IEEE Robotics and Automation Letters* 6, 4 (2021), 7325–7332.
- [11] DARVISH, K., SIMETTI, E., MASTROGIOVANNI, F., AND CASALINO, G. A hierarchical architecture for human-robot cooperation processes. *IEEE Transactions on Robotics* 37, 2 (2020), 567–586.
- [12] FREEDMAN, R., AND ZILBERSTEIN, S. Integration of planning with recognition for responsive interaction using classical planners. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2017), vol. 31.
- [13] FURNARI, A., AND FARINELLA, G. M. Rolling-unrolling lstms for action anticipation from first-person video. *IEEE transactions on pattern analysis and machine intelligence* 43, 11 (2020), 4021–4036.
- [14] GAMMULLE, H., DENMAN, S., SRIDHARAN, S., AND FOOKES, C. Predicting the future: A jointly learnt model for action anticipation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 5562–5571.
- [15] GEBERT, P., ROITBERG, A., HAURILET, M., AND STIEFELHAGEN, R. End-to-end prediction of driver intention using 3d convolutional neural networks. In *2019 IEEE Intelligent vehicles symposium (IV)* (2019), IEEE, pp. 969–974.
- [16] GIRDHAR, R., AND GRAUMAN, K. Anticipative video transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 13505–13515.
- [17] HAWKINS, K. P., BANSAL, S., VO, N. N., AND BOBICK, A. F. Anticipating human actions for collaboration in the presence of task and sensor uncertainty. In *2014 IEEE international conference on Robotics and automation (ICRA)* (2014), IEEE, pp. 2215–2222.
- [18] HUANG, C.-M., ANDRIST, S., SAUPPÉ, A., AND MUTLU, B. Using gaze patterns to predict task intent in collaboration. *Frontiers in psychology* 6 (2015), 1049.
- [19] JAIN, S., AND ARGALL, B. Probabilistic human intent recognition for shared autonomy in assistive robotics. *ACM Transactions on Human-Robot Interaction (THRI)* 9, 1 (2019), 1–23.
- [20] KOLVE, E., MOTTAGHI, R., HAN, W., VANDERBILT, E., WEIHS, L., HERRASTI, A., GORDON, D., ZHU, Y., GUPTA, A., AND FARHADI, A. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474* (2017).
- [21] KOPPULA, H. S., JAIN, A., AND SAXENA, A. Anticipatory planning for human-robot teams. In *Experimental robotics* (2016), Springer, pp. 453–470.
- [22] LE, T., SINGH, R., AND MILLER, T. Goal recognition for deceptive human agents through planning and gaze. *Journal of Artificial Intelligence Research* 71 (2021), 697–732.
- [23] LEVINE, S. J., AND WILLIAMS, B. C. Watching and acting together: Concurrent plan recognition and adaptation for human-robot teams. *Journal of Artificial Intelligence Research* 63 (2018), 281–359.
- [24] LIU, T., LYU, E., WANG, J., AND MENG, M. Q.-H. Unified intention inference and learning for human-robot cooperative assembly. *IEEE Transactions on Automation Science and Engineering* 19, 3 (2021), 2256–2266.
- [25] PUIG, X., SHU, T., LI, S., WANG, Z., LIAO, Y.-H., TENENBAUM, J. B., FIDLER, S., AND TORRALBA, A. Watch-and-help: A challenge for social perception and human-ai collaboration. *arXiv preprint arXiv:2010.09890* (2020).

- [26] PUIG, X., SHU, T., TENENBAUM, J. B., AND TORRALBA, A. Nopa: Neurally-guided online probabilistic assistance for building socially intelligent home assistants. *arXiv preprint arXiv:2301.05223* (2023).
- [27] QI, Z., WANG, S., SU, C., SU, L., HUANG, Q., AND TIAN, Q. Self-regulated learning for egocentric video activity anticipation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [28] QU, C., YANG, L., CROFT, W. B., ZHANG, Y., TRIPPAS, J. R., AND QIU, M. User intent prediction in information-seeking conversations. In *Proceedings of the 2019 Conference on Human Information Interaction and Retrieval* (2019), pp. 25–33.
- [29] RODRIGUEZ, C., FERNANDO, B., AND LI, H. Action anticipation by predicting future dynamic images. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops* (2018), pp. 0–0.
- [30] ROY, D., AND FERNANDO, B. Action anticipation using latent goal learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2022), pp. 2745–2753.
- [31] SATTAR, H., BULLING, A., AND FRITZ, M. Predicting the category and attributes of visual search targets using deep gaze pooling. In *Proc. IEEE International Conference on Computer Vision Workshops (ICCVW)* (2017), pp. 2740–2748.
- [32] SATTAR, H., FRITZ, M., AND BULLING, A. Deep gaze pooling: Inferring and visually decoding search intents from human gaze fixations. *Neurocomputing* 387 (2020), 369–382.
- [33] SATTAR, H., MÜLLER, S., FRITZ, M., AND BULLING, A. Prediction of search targets from fixations in open-world settings. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 981–990.
- [34] SHERGADWALA, M. N., TENG, Z., AND EL-NASR, M. S. Can we infer player behavior tendencies from a player’s decision-making data? integrating theory of mind to player modeling. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment* (2021), vol. 17, pp. 195–202.
- [35] SINGH, R., MILLER, T., NEWN, J., SONENBERG, L., VELLOSO, E., AND VETERE, F. Combining planning with gaze for online human intention recognition. In *Proceedings of the 17th international conference on autonomous agents and multiagent systems* (2018), pp. 488–496.
- [36] SONG, D., KYRIAZIS, N., OIKONOMIDIS, I., PAPAIOZOV, C., ARGYROS, A., BURSCHKA, D., AND KRAGIC, D. Predicting human intention in visual observations of hand/object interactions. In *2013 IEEE International Conference on Robotics and Automation* (2013), IEEE, pp. 1608–1615.
- [37] STROHM, F., SOOD, E., MAYER, S., MÜLLER, P., BÂCE, M., AND BULLING, A. Neural photofit: Gaze-based mental image reconstruction. In *Proc. IEEE International Conference on Computer Vision (ICCV)* (2021), pp. 245–254.
- [38] STROUSE, D., MCKEE, K., BOTVINICK, M., HUGHES, E., AND EVERETT, R. Collaborating with humans without human data. *Advances in Neural Information Processing Systems* 34 (2021), 14502–14515.
- [39] SZOT, A., CLEGG, A., UNDERSANDER, E., WIJMANS, E., ZHAO, Y., TURNER, J., MAESTRE, N., MUKADAM, M., CHAPLOT, D. S., MAKSYMETS, O., ET AL. Habitat 2.0: Training home assistants to rearrange their habitat. *Advances in Neural Information Processing Systems* 34 (2021), 251–266.
- [40] VEHTARI, A., GELMAN, A., SIMPSON, D., CARPENTER, B., AND BÜRKNER, P.-C. Rank-normalization, folding, and localization: An improved \hat{R} for assessing convergence of MCMC (with discussion). *Bayesian Analysis* 16, 2 (2021), 667–718.
- [41] WANG, W., LI, R., CHEN, Y., SUN, Y., AND JIA, Y. Predicting human intentions in human–robot hand-over tasks through multimodal learning. *IEEE Transactions on Automation Science and Engineering* 19, 3 (2021), 2339–2353.
- [42] WANG, Z., MÜLLING, K., DEISENROTH, M. P., BEN AMOR, H., VOGT, D., SCHÖLKOPF, B., AND PETERS, J. Probabilistic movement modeling for intention inference in human–robot interaction. *The International Journal of Robotics Research* 32, 7 (2013), 841–858.
- [43] WU, X., ZHAO, J., AND WANG, R. Anticipating future relations via graph growing for action prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2021), vol. 35, pp. 2952–2960.
- [44] ZHANG, G., HINDENNACH, S., LEUSMANN, J., BÜHLER, F., STEUERLEIN, B., MAYER, S., BÂCE, M., AND BULLING, A. Predicting next actions and latent intents during text formatting. In *Proceedings of the CHI Workshop Computational Approaches for Understanding, Generating, and Adapting User Interfaces* (2022), pp. 1–6.
- [45] ZHANG, Y., AND WILLIAMS, B. Adaptation and communication in human-robot teaming to handle discrepancies in agents’ beliefs about plans. In *Proceedings of the International Conference on Automated Planning and Scheduling* (2023), vol. 33, pp. 462–471.

8 APPENDIX

8.1 Additional Results of User Intention Prediction

Figure 10 and Figure 11 show the additional results of user intention prediction on WAH test set 1 and test set 2. Since each user has only one action sequence in true intention *put fridge*, we show the results of *put dishwasher* and *read book*. Each plot shows the posterior mean and CI bounds when a different action sequence is used.

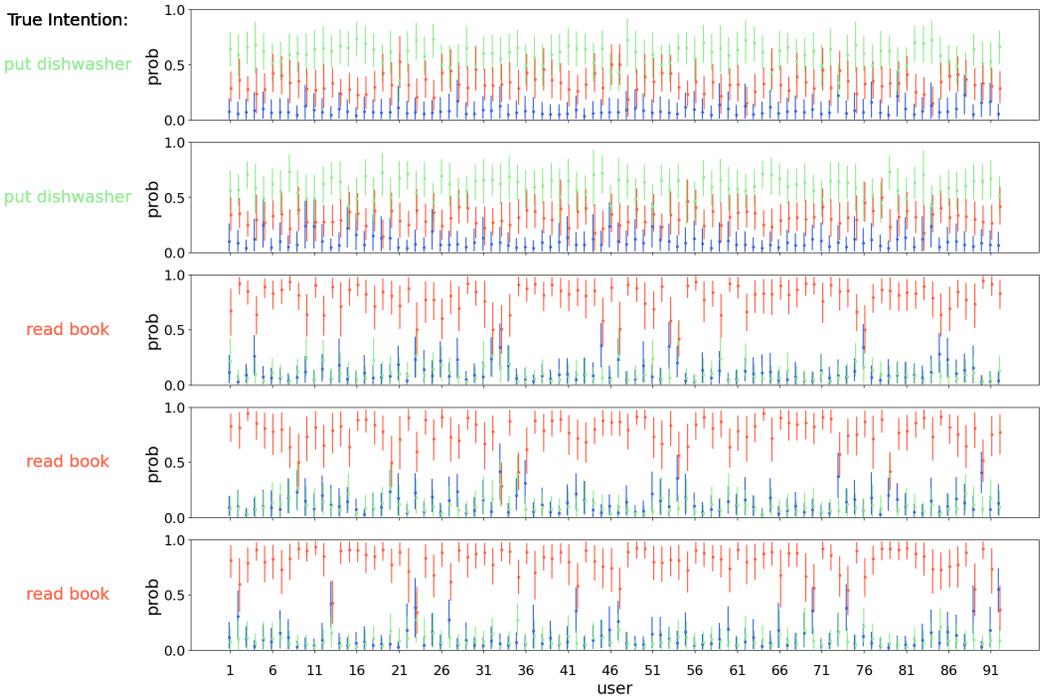


Fig. 10. Additional result of user intention prediction on WAH test set 1.

Figure 12 to Figure 15 show the additional results of user intention prediction on the keyboard and mouse interaction dataset. Each figure shows the results of using one different action sequence from one user.

8.2 Additional Results in Different Lengths Observed Actions in Action Sequences

Figure 16 and Figure 17 show the additional results of on WAH test set 1 and test set 2. we show the results of *put dishwasher* and *read book*. Each plot shows the posterior mean and CI bounds when a different action sequence is used.

Figure 18 to Figure 21 show the additional results of user intention prediction using different percentages of observed actions on the keyboard and mouse interaction dataset. Each figure shows the results of using one different action sequence from one user.

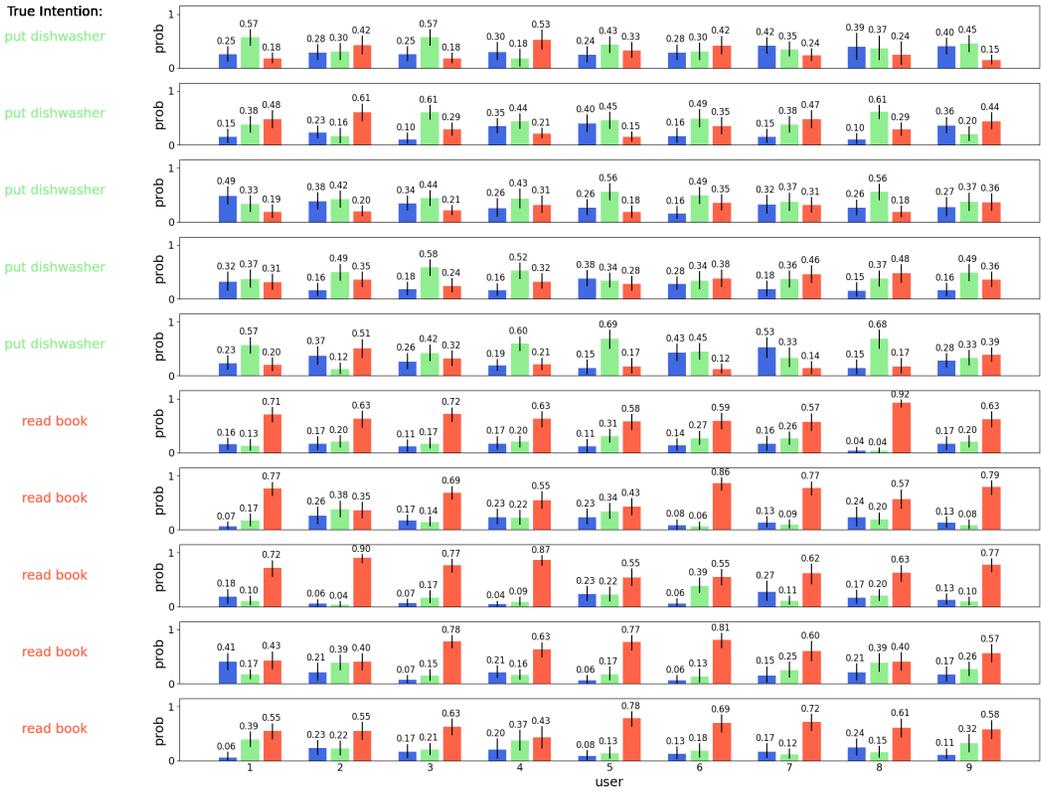


Fig. 11. Additional result of user intention prediction on WAH test set 2.

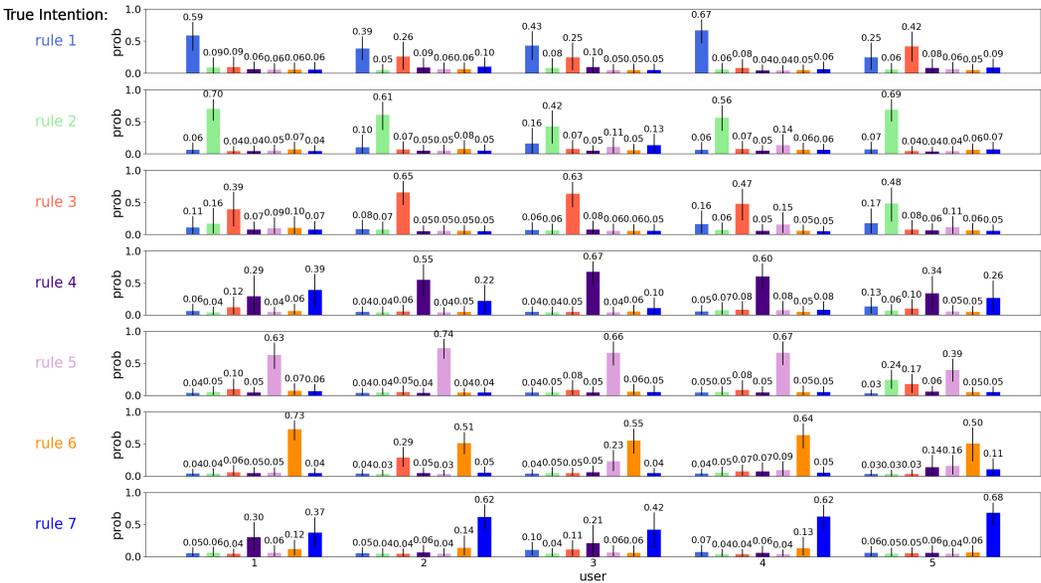


Fig. 12. Additional result 1 of user intention prediction on keyboard and mouse interaction dataset.

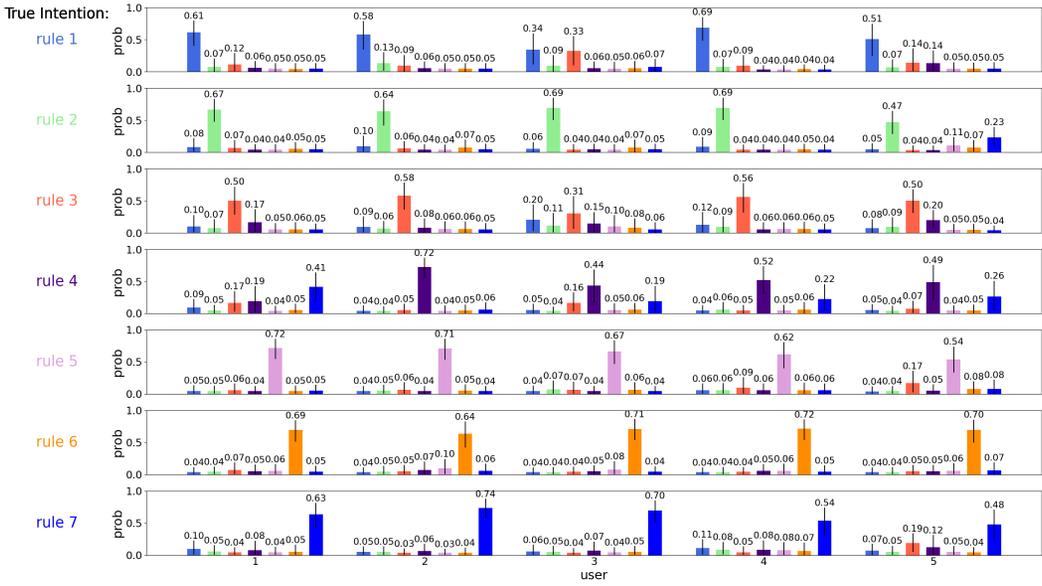


Fig. 13. Additional result 2 of user intention prediction on keyboard and mouse interaction dataset.

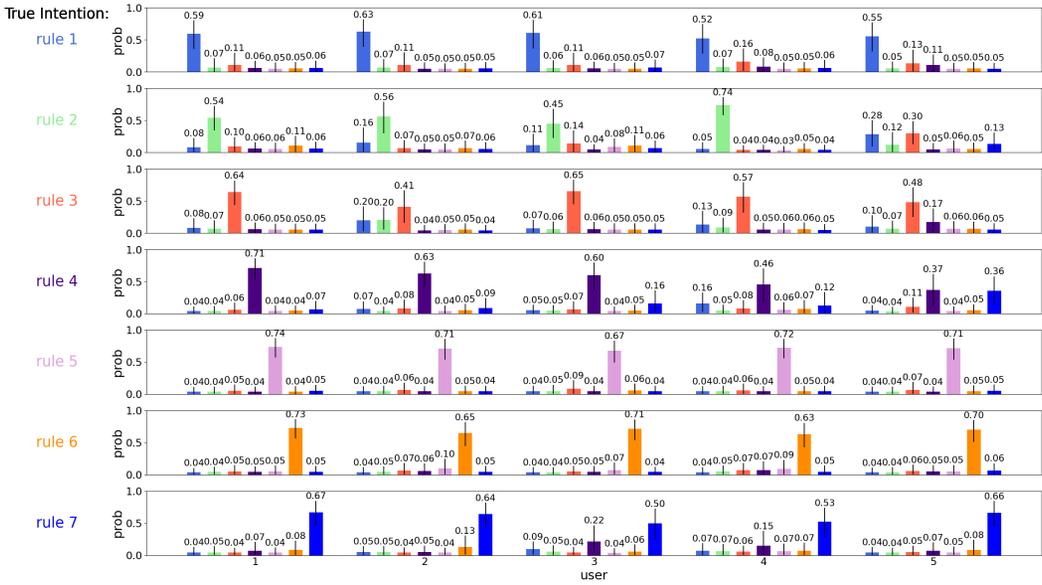


Fig. 14. Additional result 3 of user intention prediction on keyboard and mouse interaction dataset.

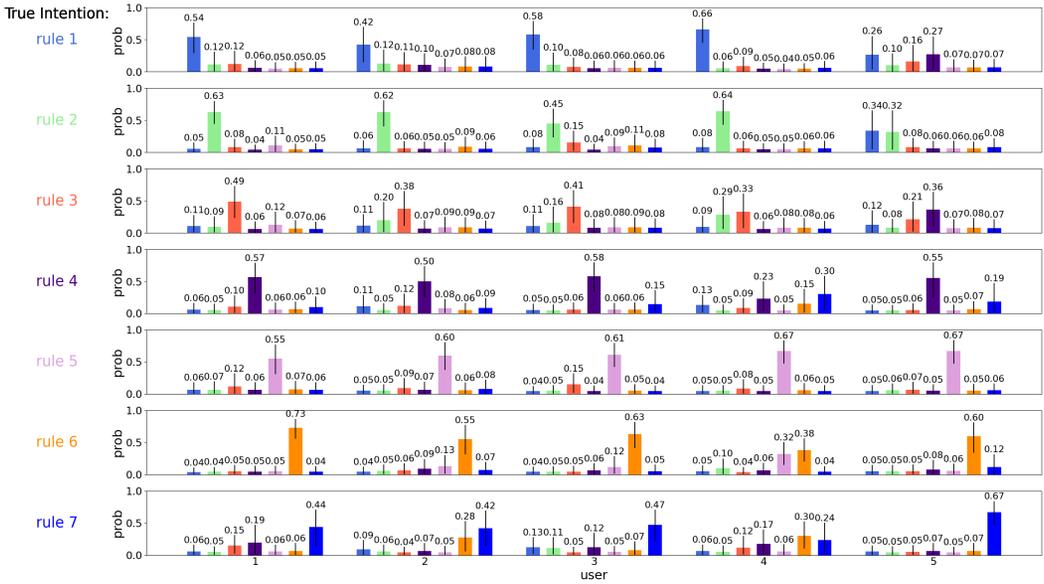


Fig. 15. Additional result 4 of user intention prediction on keyboard and mouse interaction dataset.

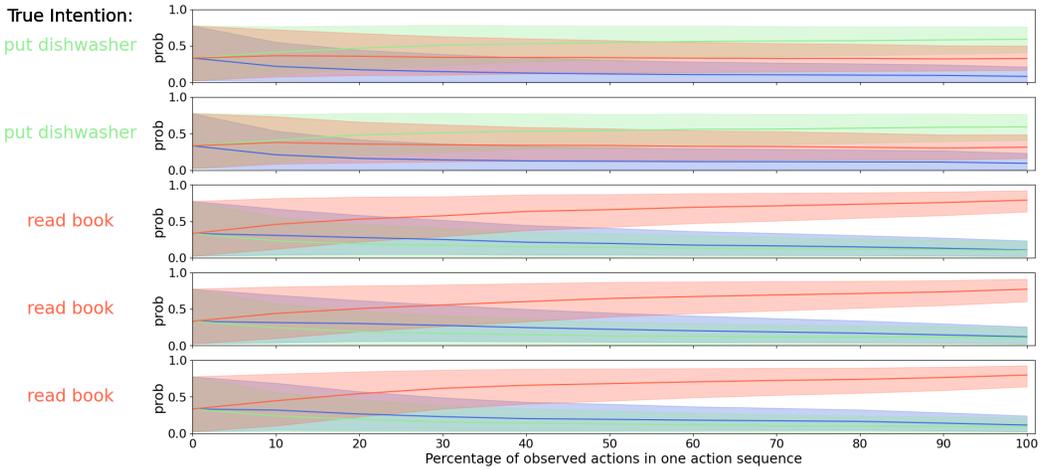


Fig. 16. Additional result of user intention prediction with different percentages of occurred actions in actions sequences on WAH test set 1.

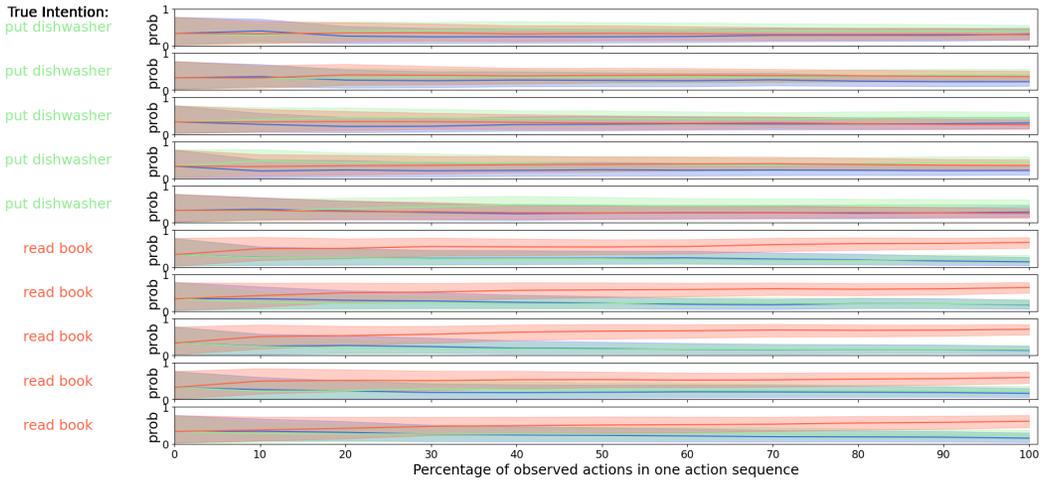


Fig. 17. Additional result of user intention prediction with different percentages of observed actions in actions sequences on WAH test set 2.

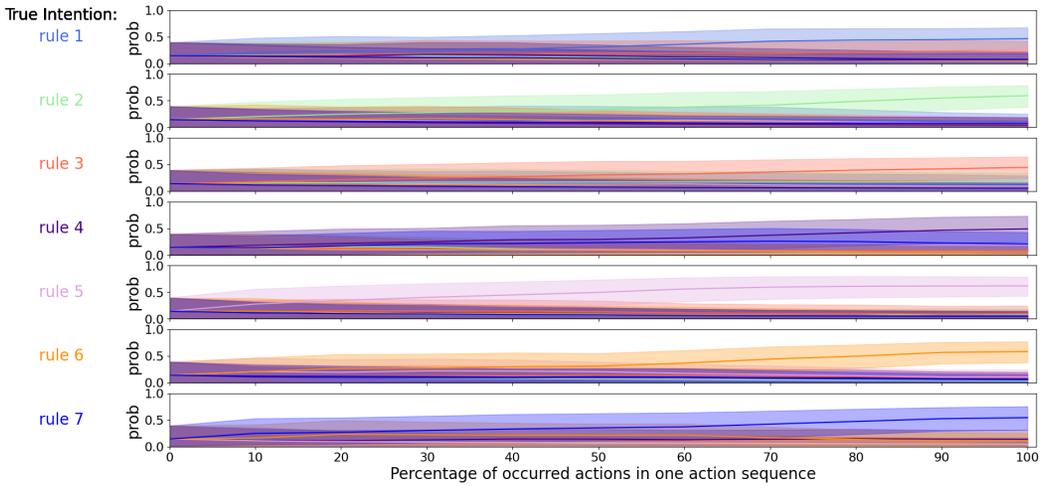


Fig. 18. Additional result 4 of user intention prediction using different percentage of observed action on keyboard and mouse interaction dataset.

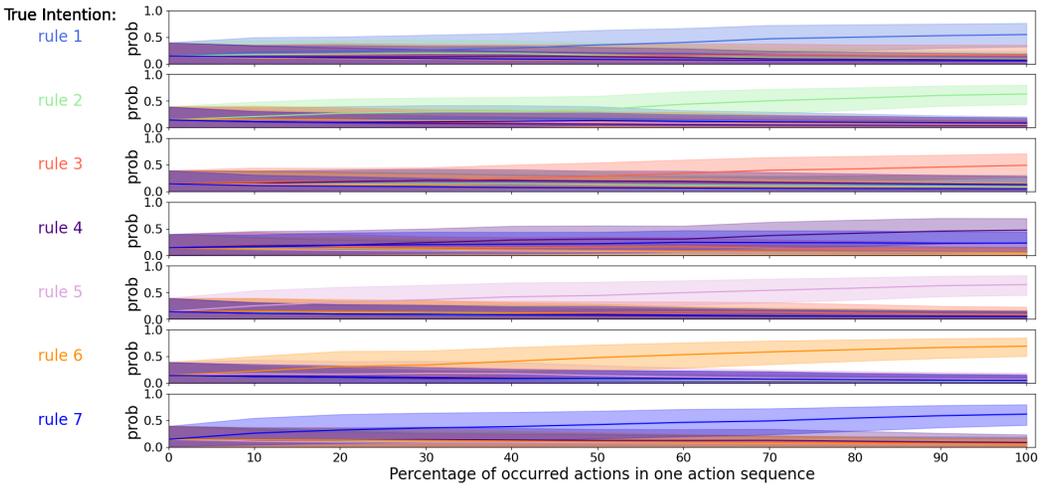


Fig. 19. Additional result 2 of user intention prediction using different percentage of observed action on keyboard and mouse interaction dataset.

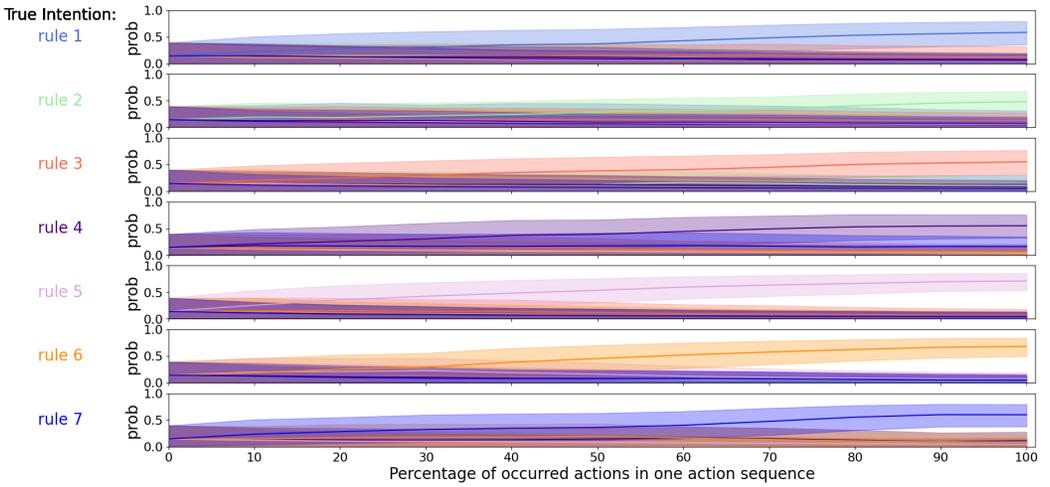


Fig. 20. Additional result 3 of user intention prediction using different percentage of observed action on keyboard and mouse interaction dataset.

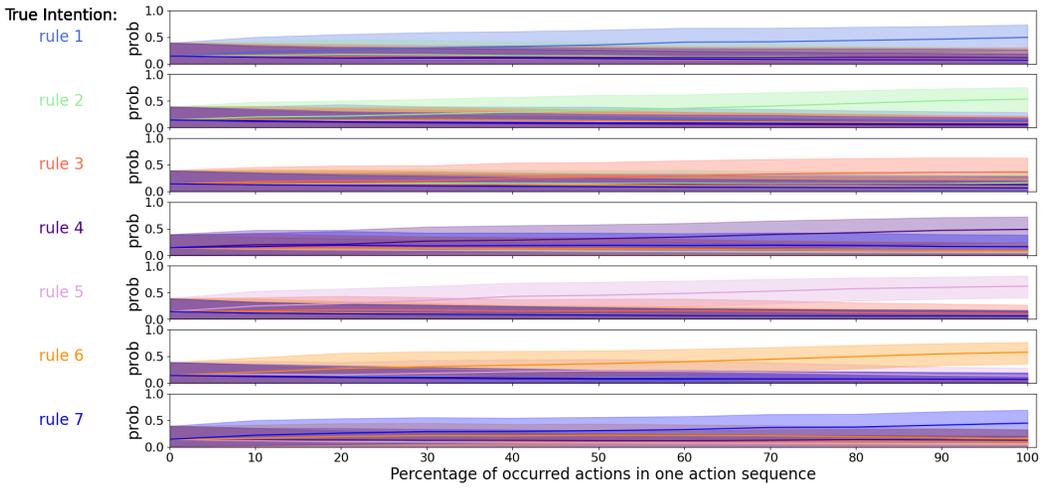


Fig. 21. Additional result 4 of user intention prediction using different percentage of observed action on keyboard and mouse interaction dataset.