

Action Anticipation for Collaborative Environments: The Impact of Contextual Information and Uncertainty-Based Prediction

Clebeson Canuto^a, Plinio Moreno^b, Jorge Samatelo^a, Raquel Vassallo^a, José Santos-Victor^b

^aRoom 20, CT-II, Department of Electrical Engineering, Federal University Espírito Santo, Av. Fernando Ferrari, 514, Goiabeiras, 29075-910, Vitória - ES - Brazil

^bFloor 7, North Tower, Institute for Systems and Robotics, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1, 1049-001, Lisbon, Portugal

Abstract

To interact with humans in collaborative environments, machines need to be able to predict (i.e., anticipate) future events, and execute actions in a timely manner. However, the observation of the human limb movements may not be sufficient to anticipate their actions unambiguously. In this work, we consider two additional sources of information (i.e., context) over time, gaze, movement and object information, and study how these additional contextual cues improve the action anticipation performance. We address action anticipation as a classification task, where the model takes the available information as the input and predicts the most likely action. We propose to use the uncertainty about each prediction as an online decision-making criterion for action anticipation. Uncertainty is modeled as a stochastic process applied to a time-based neural network architecture, which improves the conventional class-likelihood (i.e., deterministic) criterion. The main contributions of this paper are four-fold: (i) We propose a novel and effective decision-making criterion that can be used to anticipate actions even in situations of high ambiguity; (ii) we propose a deep architecture that outperforms previous results in the action anticipation task when using the Acticipate collaborative dataset; (iii) we show that contextual information is important to disambiguate the interpretation of similar actions; and (iv) we also provide a formal description of three existing performance metrics that can be easily used to evaluate action anticipation models. Our results on the Acticipate dataset showed the importance of contextual information and the uncertainty criterion for action anticipation. We achieve an average accuracy of 98.75% in the anticipation task using only an average of 25% of observations. Also, considering that a good anticipation model should perform well in the action recognition task, we achieve an average accuracy of 100% in action recognition on the Acticipate dataset, when the entire observation set is used.

Keywords: Action Anticipation, Early Action Prediction, Context Information, Bayesian Deep Learning, Uncertainty.

1. Introduction

Humans have the natural ability to interact with each other and perform joint tasks. Part of this ability is due to their capacity of perceiving the environment and recognizing patterns that help them anticipate the actions of others, and thus make better

Email addresses: clebeson.canuto@gmail.com (Clebeson Canuto), plinio@isr.tecnico.ulisboa.pt (Plinio Moreno), jorge.samatelo@ufes.br (Jorge Samatelo), raquel@ele.ufes.br (Raquel Vassallo), jasv@isr.tecnico.ulisboa.pt (José Santos-Victor)

decisions. Similarly, artificial machines need this capacity of anticipating actions, to act accordingly and achieve an effective interaction with humans [1].

Action anticipation and action recognition are two different tasks. The “action recognition” task is based on a model that uses an entire sequence of information, which represents one performed action, to associate the observed action to one possible action class [2]. If the decision-making depends on the entire action, it can only be performed after the action is completely executed. However, this approach is not suitable for systems that manage risks or perform joint tasks with humans. For instance, in a situation where a self-driving car approaches a pedestrian, it must perceive whether the pedestrian will cross the road in time, in order to safely stop or deviate the car if necessary. In this scenario, the model must not only recognize actions but, more importantly, must anticipate them [3].

Action anticipation consists of classifying an action even before it occurs, by using the partial information provided up to a certain moment in time. Usually, an anticipation model is more complex than a recognition one. This comes from its capacity to classify actions based on an incomplete sequence of data, which makes the choice of the correct class more uncertain. Ideally, every anticipation model should be capable of recognizing actions; on the other hand, not every recognition model would be able to anticipate them.

In the last few years, deep learning has achieved the state-of-the-art results in many tasks, such as image recognition [4, 5, 6], natural language processing [7, 8] and action/activity recognition [9, 10, 11]. Some works, like [12, 13, 14], represent an action by estimating the movement of the involved actors (i.e., users). In the case of simple and unambiguous actions, the movement can be sufficient for a successful recognition/anticipation task. However, in the case of more complex and ambiguous actions, it would not be enough to recognize/anticipate successfully, mainly when the information about objects, persons, environment configuration, movements performed previously, are important for recognizing or anticipating actions. Furthermore, some details during action anticipation, such as objects’ position, the relation between hands and object/person and the type of object

manipulated, can offer as much or even more information than only movement. As such, using only movement, the model rules out the context, a critical information that can help characterize the actions.

Regarding the action recognition task, the two-stream approaches [10, 11, 15] are the most successful, because they use movement as the main source of information to describe each action, and they use the context as additional information that can help characterize each class individually. In these solutions, the movement is the optical flow calculated between sequential images, and the contextual information [16] is extracted implicitly by CNNs (Convolutional Neural Networks)[17]. However, to obtain the implicit contextual information from images in a self-supervised manner, the training procedure of the CNN models requires large datasets to achieve good results. As a consequence, the two-stream approaches are not effective when solving problems provided by small datasets, such as those commonly used for human-human or human-robot collaboration.

Analyzing from another perspective, even achieving satisfactory results in their experiments, the aforementioned works are not crystal clear about how one could use their solutions in a real-time situation, once they measure the model performance using accuracy or observation ratio. They do not discuss how to handle action anticipation or what kind of function must be used as the decision-making criterion. Due to the absence of such discussion, it is unclear how to use this approach in a real application, where the data is continuously generated, as in a video streaming.

Another problem of most deep learning solutions is their overconfidence in their predictions. A deterministic model will always provide a prediction, even when there is a high uncertainty about the correct class, and the final decision becomes unclear. A trustworthy model should assess its uncertainty about each prediction and provide the system with the possibility of making more reliable decisions.

In this work, we focus on context-based action anticipation, but with small datasets. Thus, instead of implicitly learning the visual context, we define the contextual information in an action anticipation problem. With this in mind, we used the Acticipate

dataset [1], where one person hands an object over to another one, and receives it back. For this dataset, the dyadic interaction task requires the future prediction (i.e., anticipation) of the arm and head motion, gaze and object position. A previous work [18] has shown that, using the eye gaze and the 3D pose of the main character in the Acticipate dataset, a time-based deep learning architecture is able to anticipate his actions. As defined in [19], context is any information that can be used to characterize an entity. Therefore, when considering the 3D pose/movement as the entity that represents an action, the eye gaze in [18] can be seen as context information.

Now, to further investigate the importance of context in the task of anticipating actions, we have increased the complexity in the Acticipate dataset, extending its number of actions. To do that, we divided previous actions to create the new actions *receive* and *pick*, which add ambiguity into the actions *give* and *place*, correspondingly. We also consider an additional element of context information, the position of the handled object. Finally, instead of using 3D pose and gaze as in [18], we use only information taken from RGB images. Such restriction makes our proposal more general and less dependent on intrusive and/or expensive sensors.

Also, to investigate the possibility of using the uncertainty to provide a more reliable decision, we propose a context-aware model based on a recurrent neural network with an adaptive threshold. This threshold is calculated via an uncertainty metric and represents a decision-making criterion for action anticipation. The use of uncertainty significantly contributes to attenuate the overconfidence problem often faced by models trained with small datasets.

In summary, the main contributions of this paper are the following:

- We propose a novel and effective decision-making criterion that can be used to anticipate actions even in situations of high ambiguity. The proposed approach aims to minimize the model’s uncertainty instead of maximizing its class probabilities. Therefore, by applying a proper threshold over the uncertainty, the decision about whether an action should be anti-

pated or not can be done.

- We show the importance of context information to disambiguate similar actions.
- We propose a deep architecture that uses less information than [18], and outperforms the results in action anticipation task using Acticipate dataset. This result holds even in the case of its extended number of actions, which are more ambiguous than the original ones [1].
- We also provide a formal description of three existing performance metrics that can be easily used to evaluate action anticipation models.

To build a better understanding about our proposal, the next sections will cover, respectively: the related works (Sec. 2); action anticipation background and related problems (Sec. 3); the methodology of this work, including the hypotheses raised and its main contributions (Sec. 4); Bayesian neural networks and uncertainty (Sec. 5); our proposed approach (Sec. 6); experiments (Sec. 7 and 8, for results and discussions); and finally, conclusions and future works (Sec. 9).

2. Related Works

In the last few years, action anticipation has been addressed by many researchers [20, 21, 22, 23, 24, 25] due to its importance to perform an effective interaction.

In [23], the authors proposed to decrease the dimensionality on RNNs by allowing the sharing of weights, and improve the temporal representation of an action by using an RBF kernel (Radial Base Function) over the hidden-state of an LSTM network. They proposed to feed an LSTM with features extracted by a CNN. Next, they applied an RBF over the LSTM hidden states, and lastly, the RBF outcome is given as input to a Multilayer Perceptron (MLP). The authors use between 20% and 50% of a video to predict the next features and then perform the anticipation.

In [13], the authors use a convolutional auto-encoder network to predict the next movement of a

video. Such movement is generated by a ranking loss function, applied over the difference between consecutive images in a sequence, and is stored in a still RGB image called Dynamic Image [26]. With a Markov assumption, after generating a sequence of dynamic images using S frames for each one, the model generates the next k dynamic images, where $k \geq 1$. Further, those images feed a model that outputs the probability distribution over action classes. A drawback of the two previous works is to use movement as the only source of information to represent an action, which can harm the prediction of actions that are related not only to movement but also to context information.

In [27] is proposed a model to anticipate actions based only on RGB images. The authors use as feature extractor the pre-trained CNN VGG16 and, as the classifier, two LSTMs that predict the classes corresponding to each input frame of a video. A similar approach is also presented in [28].

LSTM is also used to anticipate actions of car drivers by using only RGB images [29] or in combination with GPS information [30]. Other approaches, as [31], use Generative Adversarial Networks (GAN) to predict future images and then anticipate the action, or more sophisticated architectures, as in [3], that uses Convolutional Graphical Models (CGM) to predict when a pedestrian will cross the road.

Despite these works present good results in terms of accuracy at each observation time, none of them explains how action anticipation should be performed in a real scenario, when it is not possible to know the size of the input sequence. They did not discuss what kind of decision-making criteria could be used in such a situation.

Even in works as [32], which aim to anticipate action in online videos, the authors only reported the accuracy at each observation, but nothing about how to make decisions. Only a couple of works address this question. For instance, in [33, 34], the authors use a threshold over the probability distribution provided by an HMM (Hidden Markov Model) to anticipate maneuvers of drivers. However, as discussed in [33], this approach faces problems in ambiguous situations, where it is not possible to be sure about the action to be anticipated, even when the probability

exceeds the specified threshold.

Many of the approaches mentioned above are not suitable for small datasets, since the high capacity of their models can lead to overfitting. Therefore, [18] proposes a different method to anticipate action in the Acticipate dataset - a small collaborative dataset used to understand the role of gaze on action anticipation [1], as discussed in Sec. 4. Their approach consists of feeding an LSTM cell with a 3D pose (Motion Capture-MoCap information) and gaze (fixation points), and then pass the LSTM output through a softmax classifier. They trained two models with different observations: one with only 3D pose and another with 3D pose plus eye gaze. When the model uses the pose and gaze information, the authors concluded that the actions in the dataset could be anticipated 92ms before. This result showed the importance of using not only movement information (here, the evolution of the pose in time) to anticipate actions. However, the authors did not notice that their model did not recognize all the actions (100% of action recognition accuracy) even after seeing the whole sequence. Their results for action anticipation were shown based only on one action sample. More conclusive results should present statistics for all classes in the entire dataset. In complement, they also did not provide an answer to when a model must anticipate an action. From their comments, we presume that it may be done using a threshold on the probability value, as mentioned in [33, 34].

After these explanations, our main objectives in this work are:

- propose a model that improves results in [18] even when using only RGB images;
- present how context can be used in a neural network architecture to improve action anticipation;
- present in detail how to anticipate an action using a threshold value as a decision-making criterion; and
- propose the use of uncertainty as an effective threshold value that improves action anticipation.

3. Action Anticipation Background

In this section, we describe the definition adopted here for action anticipation, its main properties, and how we address the problem. We can divide the works that try to solve the anticipation task into two main groups: (i) early action prediction, where an action must be predicted before it is fully executed [35, 24, 36, 37, 38]; and (ii) event anticipation, where an event must be predicted before it starts [39, 40, 40]. In this work, “action anticipation” is understood as in the first set of works: early action prediction by using sequential features.

3.1. Problem Definition

First of all, it is essential to formally define the action anticipation task. Let $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \mid \mathbf{x}_t \in \mathbb{R}^{d \times 1}\}$ be a sequence with N observations that represents the execution of a specific action $y \in \mathcal{Y}$, where \mathcal{Y} is a set with d action classes. Here, \mathbf{x}_t represents an observation taken at time t . Now, considering that $\mathcal{X}_{t_1:t_2}$ represents an indexed sequence composed by the observations taken between time t_1 and t_2 , we define a model M for action classification problem as a mapping function parametrized by θ that receives as input $\mathcal{X}_{1:t}$ (t observations from \mathcal{X}) and return as output the vector of probability scores $\mathbf{s} \in [0, 1]^{d \times 1}$, representing the probability that sequence \mathcal{X} belongs to each action class.

$$\mathbf{s} = M(\mathcal{X}_{1:t}, \theta). \quad (1)$$

In action recognition tasks, the model M has all the observations of the sequence \mathcal{X} ($t = N$) available to generate the probability score $\hat{\mathbf{s}}$. On the other hand, for an action anticipation task, the action is not completely executed, thus only an initial part of \mathcal{X} is available ($t < N$) so that M can infer $\hat{\mathbf{s}}$.

In Eq. (1) the parameter θ can be found by solving the following optimization problem:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \{ \mathcal{L}(\theta, \mathcal{D}) \} \quad (2)$$

where $\mathcal{D} = \{(\mathcal{X}^{(1)}, y^{(1)}), (\mathcal{X}^{(2)}, y^{(2)}), \dots, (\mathcal{X}^{(k)}, y^{(k)})\}$ is the training set, with each pair $(\mathcal{X}^{(i)}, y^{(i)})$ representing an action sequence and its respective label,

K is the number of sequences in the training set, and \mathcal{L} is a loss function.

During the prediction time, we do not know the value of N , and thus we do not know when the action will end. Therefore, at each time t , M only uses the observed current sequence, $\mathcal{X}_{1:t}$, and a function g is in charge of predicting the action class at instant t .

$$\begin{aligned} \hat{\mathbf{s}} &= M(\mathcal{X}_{1:t}, \hat{\theta}) \\ \hat{y} &= g(\hat{\mathbf{s}}). \end{aligned} \quad (3)$$

For action recognition tasks, the discriminant function g can be defined as:

$$g(\hat{\mathbf{s}}) = \operatorname{argmax}(\hat{\mathbf{s}}), \quad (4)$$

because the model M is more confident about the probability score assigned to $\hat{\mathbf{s}}$. On the other hand, for action anticipation tasks, since M uses only part of the observations, when the distribution $\hat{\mathbf{s}}$ is close to a uniform distribution, one can not be certain about the correct class. Hence, Eq. (4) is not an adequate discriminant function to anticipate actions.

In this way, a better option is to use a discriminant function with a threshold parameter p , as presented in Eq. (5).

$$g(\hat{\mathbf{s}}, p) = \begin{cases} \operatorname{argmax}(\hat{\mathbf{s}}), & h(\hat{\mathbf{s}}) > p \\ -1, & \text{otherwise} \end{cases} \quad (5)$$

Once p is specified as a probability value, h can be defined as:

$$h(\hat{\mathbf{s}}) = \max(\hat{\mathbf{s}}) \quad (6)$$

In Eq. (5), a value of $p \geq 0.9$ means that the model is highly certain about its prediction, and the action can be anticipated, which favors the use of such a model in real-time. On the other hand, when it returns -1 means that it is not certain about the correct class and needs more observation to improve its certainty.

3.2. Evaluation metrics

After determining how to anticipate an action, it is essential to decide how to ascertain the quality of the model M . Therefore, we formally describe three existing metrics that can be used in anticipation benchmark experiments: (i) accuracy at each observation

ratio, (ii) anticipation accuracy and (iii) expected observation ratio.

Accuracy at each observation ratio. Considering that each sequence \mathcal{X} can have a different length N , this metric helps evaluate all sequences in a normalized time scale. Thereby, the success ratio when anticipating an action after a observation ratio r , with an anticipation threshold p , can be calculated as follows:

$$ACC(r) = \frac{1}{K} \sum_{i=1}^K \text{pred}(\mathcal{X}_{1:\lceil r \times N \rceil}^{(i)}, y^{(i)}, p), \quad (7)$$

where,

$$\text{pred}(\mathcal{X}_{1:t}, y, p) = \begin{cases} 1, & g(M(\mathcal{X}_{1:t}), p) = y \\ 0, & \text{otherwise} \end{cases}. \quad (8)$$

In Eq. (7), in terms of r , $t = \lceil r \times N \rceil \forall r \in (0, 1]$, where N is the number of observations in a sequence $\mathcal{X}^{(i)}$. However, in terms of t , $r = t/N \forall t \in \{1, 2, \dots, N\}$.

Anticipation accuracy. In a real-time situation, the model can not access the label of each observation. So, the evaluation of the anticipation model during training must be performed when the model makes its first prediction for each sequence. In this sense, this classification metric measures the success ratio of the model M when anticipating actions by the first time. It is calculated as the average accuracy of each classification. Therefore, when using this metric, we do not regard in which observation the action was predicted but whether it was predicted correctly. Eq. 9 presents how it is calculated,

$$ACC_{act} = \frac{1}{K} \sum_{i=1}^K \sum_{t=1}^{N-1} I(t) \text{pred}(\mathcal{X}_{1:t}^{(i)}, y^{(i)}, p), \quad (9)$$

where,

$$I(t) = \begin{cases} 0, & ((t=1) \wedge (g(M(\mathcal{X}_{1:t}), p) = -1)) \\ & \vee (I(t-1) = 1) \\ 1, & \text{otherwise} \end{cases} \quad (10)$$

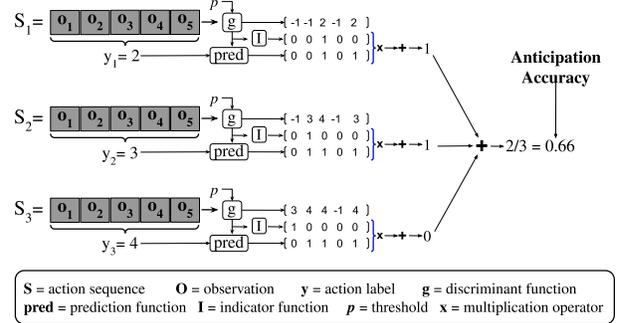


Figure 1: Graphical example of how to calculate the anticipation accuracy using Eq. (9).

where I is an indicator function that disable predictions based on whether the anticipation has already occurred or not.

A detailed example of how to apply this metric is given in Fig. 1.

Expected observation ratio. This measurement focuses on the expected amount of observations necessary to anticipate an action correctly. It can be implemented according to Eq. (11). Note that when the model is correct, it receives the value t , which corresponds to the observation where the prediction is executed. However, when it misses the anticipation, it is penalized by receiving the sequence size N .

$$E_{obs} = \frac{1}{K} \sum_{i=1}^K \text{obs}(\mathcal{X}^{(i)}, y^{(i)}, p), \quad (11)$$

where,

$$\text{obs}(\mathcal{X}, y, v) = \frac{1}{N} \min(\{f_{\text{pred}}(\mathcal{X}_{1:t}, y, p, t, N)\}_{t=1}^N)$$

$$f_{\text{pred}}(\mathcal{X}_{1:t}, y, p, t, N) = \begin{cases} t, & \text{pred}(\mathcal{X}_{1:t}, y, p) = 1 \\ N, & \text{otherwise} \end{cases}.$$

4. Methodology

This work aims to show the influence of context in the anticipation task and to use uncertainty as a decision-making criterion in a collaborative environment. To do this, we use a controlled dataset that contains, by each frame, the action performed and the

corresponding context information. Therefore, to understand how our intuitions have arisen and resulted in our proposal, it is necessary to analyze the used dataset and thus realize how the questions came up.

4.1. Acticipate dataset

The chosen dataset is the Acticipate¹, which was acquired to study the influence of gaze in action and/or intention anticipation [1, 18] in a collaborative environment. It comprises 120 trials, distributed into 6 classes. During the acquisition, the actor was wearing an eye gaze tracker binocular glasses (Pupil Labs eye-tracker [41]) and a suit with 25 markers. He should perform 6 different actions: *give* an object (left, middle, or right) and *place* an object (right, middle, or left). In the *give* actions, he should give an object (in this case, a small red ball) to one of the three volunteers located on: his right side, left side, or in front of him (middle). In the *place* actions, he should place the same object in one of the three points on the table located at his right, middle (in front of him), or his left. Each action starts with the object placed in a point near to the actor and finishes when the object returns to the same point. As showed by [1] the gaze is an essential source of information when one wants to anticipate actions. Besides, as we will see in the next section, the object plays a fundamental role when the action becomes more ambiguous. So it is possible to know what kind of context information must be taken into account for each action class during the anticipation process. Fig. 2 presents a sample of each action and the object starting point.

Each trial consists of 3-dimensional data corresponding to the positions of the markers on the actor’s suit, captured by an OptiTrack² MoCap system, at 120Hz; 2D gaze fixation point captured by the eye tracker glasses at 60Hz; and an RGB video captured by a camera facing the actor, at 30Hz. The dataset is unbalanced, because every class has a different number of samples: 17 (place right), 23 (place middle),

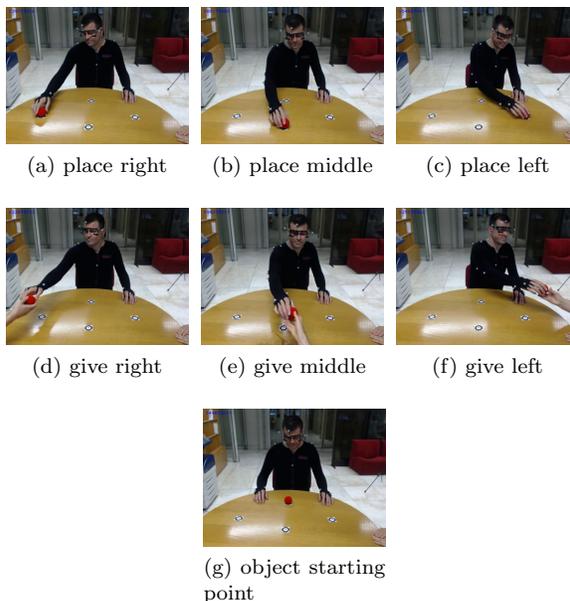


Figure 2: Sample of each action in Acticipate Dataset and the object starting point

20 (place left), 24 (give right), 19 (give middle) and 17 (give left).

In this work, when referring to the Acticipate dataset, we call movement the change of position of both arms.

4.2. Dataset analysis

By analyzing the dataset, it is possible to notice that, in many cases, the movement does not have enough information about action to provide good anticipation. For instance, each action of *place* and *give* has similar movements depending on its direction (left, middle or right). However, after taking into account gaze information, one can notice that the action can be anticipated long before. Gaze indicates whether the user will place the object on the table or give it to a volunteer. As discussed before, if we take the movement as the principal entity of each action, we can consider gaze as a context information (additional information that helps characterize the entity). Now, gaze and movement provide enough data to anticipate actions in this dataset. However,

¹Download: <http://vislab.isr.ist.utl.pt/datasets/>

²<https://optitrack.com/>

if the dataset was divided into more actions, would the gaze be a sufficient source of information to anticipate them?

An interesting but not considered characteristic of this dataset is that, once the interaction involves only one object, the actor must place it at its starting point after performing each action. Thus, when he places the object somewhere on the table, he must *pick* it up, and when he gives the object to someone, he must *receive* it back. A simple example of this behavior can be seen in Fig. 3. In that way, we can extend the dataset from 6 actions to 12 actions: *give*, *place*, *pick* and *receive* (each one with the directions left, right and middle).



(a) Place middle



(b) Give right

Figure 3: Sample of two actions (place middle and give right) from Anticipate Dataset

Considering now the extended dataset, if we constrain the analyzes to the movement and gaze (Fig. 4 (a-d)), we notice that, even with gaze, it is not possible to perform right anticipation between actions *give/receive* or *place/pick* when they are toward the same direction. In this case, it is necessary to wait for more observations.

On the other hand, when applying no constraint on what we can analyze in each image (Fig. 4 (e-h)), we can anticipate actions of the extended dataset as fast as in its original configuration. In some cases, as in Fig. 4 (f), the action can be anticipated after observing the first frame. This is possible because we take

object information into account as another essential context information. For instance, the starting position of the object makes it possible to anticipate a *pick* action after observing only one image.

Something similar occurs with *receive* actions, where the object is usually out of the scene, being held by a volunteer. For such actions, after seeing the first frame, it is not possible to assure which is the action, once it depends on the direction. However, we can tell that it will be a *receive* action. Therefore, the correct anticipation comes after perceiving the gaze or the movement direction. This helps us to eliminate less likely actions and allows us to focus on information that helps to find the right action. Fig. 4 illustrates four situations where there are significant ambiguities between actions, and the object information is critical to reduce it.

4.3. Anticipation

Although we are able to anticipate actions in the extended dataset, in some cases, there are issues about the anticipation that must be taken into account. Even people can have their prediction capacity compromised by overconfidence. In a particular case, as presented in Fig. 5, the volunteer wrongly anticipated an action after observing a movement similar to another one. Her confidence in her prediction deceived her. Thus, even people, in some situations, need to be more sure about the action before making a decision. If a person can be fooled by his/her overconfidence, this problem is possibly more significant in a computational model.

The overconfidence about a prediction could lead the model to make a wrong decision in a real-time situation. This problem can be mitigated by providing the model with the ability to estimate the uncertainty about its prediction. A deterministic model, even with a high value of probability threshold (e.g., $p > 0.9$), could wrongly anticipate an action when it is overconfident about its prediction. This overconfidence in prediction can be provoked by a lack of data to prevent the model from ambiguous classes.

Thus, as mentioned in [33], a possible solution to increase the model certainty is to lead it to make more z predictions before deciding on the correct action class. In this way, if the predicted class remains



(a) any action is possible (b) any action is possible (c) place or pick left (d) place or pick left



(e) give or place (any direction) (f) pick middle (g) place left (h) pick left

Figure 4: Situations in the extended dataset with great ambiguities when it is analyzed only gaze and movement. Any action is possible in (a) and (b). It is necessary to wait for the movement to infer the direction but, even after knowing the direction, it is necessary to observe almost the complete action to distinguish between the actions *give/receive* and *place/pick*. In (c) and (d), the movement starts toward the left side, simultaneously, the gaze is directed to the table. Therefore, the possible action is a *place* or *pick* toward the left direction. For the actions shown in (e)-(h), because the object position is taken into account, the ambiguities can be reduced or even eliminated. In (e) and (f), the number of possible actions is reduced after knowing the object position. In (e), *pick* and *receive* actions are not possible. On the other hand, in (f), only the action *pick middle* is possible. The same occurs in (g) and (h). In (g), the most likely action is *place left*. Finally, in (h), the only possibility is *pick left*. Notice that in (f) and (h), the action is anticipated after observing only one image.

for the next z observations, the model can be more confident about the correct class and can anticipate the action. However, even though it looks like a good solution, what is the best size for z ? An inaccurate choice of this new parameter can postpone the anticipation of actions that have no ambiguity problem in z observations. Additionally, z may not be enough for actions with more ambiguities.

A better solution is to use as threshold an uncertainty value rather than a probability value. Thus, the model can anticipate an action when it is more certain about its prediction. Therefore, ambiguous actions, which likely provide more uncertainty to the



(a) Place left



(b) Give left

Figure 5: Two action samples from the Acticipate dataset. In (a), the volunteer wrongly anticipated the action, thinking it would be a give left action (shown in (b)) instead of a place left.

model, would need more observations to be anticipated properly. On the other hand, those with less ambiguities could be anticipated previously. This solution can be taken as a tailored z value for each action chosen by the model during training.

4.4. Hypotheses and contributions

These observations led us to raise three main hypotheses regarding the Acticipate dataset:

1. more actions are likely to cause more ambiguities;
2. context information can help to distinguish different actions represented by similar movements;
3. uncertainty is a more reliable and effective threshold to anticipate actions than probability values.

In this work, the gaze and the object's position represent the context of each action. So, we propose a model based on Artificial Neural Networks (ANNs) that anticipates actions represented by sequences of data with varying lengths. The proposed model has two versions: a deterministic and a stochastic one.

5. Bayesian Neural Networks and Uncertainty

Deep neural networks are usually trained by optimization algorithms based on Stochastic Gradient Descent (SGD). As SGD uses the gradient of the weights, it needs the loss function to be differentiable for all weights, which implies the weights must be

deterministic variables. In consequence, most deep neural network models are deterministic, so they are unable to provide their uncertainty about their predictions. Thus, to measure uncertainty in this type of model, we can create a Bayesian Neural Network (BNN).

In a Bayesian model a posterior distribution must be inferred by applying the Bayes rule:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}, \quad (12)$$

where $p(\boldsymbol{\theta}|\mathcal{D})$ is the posterior distribution over $\boldsymbol{\theta}$ after observing data \mathcal{D} ; $p(\mathcal{D}|\boldsymbol{\theta})$ is the likelihood of \mathcal{D} ; $p(\boldsymbol{\theta})$ is the prior belief about the distribution of $\boldsymbol{\theta}$; and $\int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$ is the the normalization term (a.k.a evidence or marginal likelihood).

In many cases, the evidence term in Eq. (12) turns the posterior inference intractable. However, some works attempted to solve this problem via Variational Inference (VI) [42]. In 2011, [43] proposed in detail how to use VI in Bayesian Neural Networks so that a Gaussian distribution with known parameters could approximate its posterior distribution. Although effective, VI was not yet an easy task to accomplish. Therefore, in 2013, [44] proposed a way to train a BNN with VI thought a technique called reparametrization trick, which consists of drawing the activation Z of a layer l from a standard factorized Gaussian distribution.

In this sense, the layer l outputs two values, μ and σ , which represents the mean and variance of a factorized Gaussian distribution $N(\mu, \sigma)$, respectively. Next, the activation of l is drawn from $Z \sim N(\mu, \sigma)$. Aiming to approximate Z by a standard factorized Gaussian distribution ($N(0, 1)$), the authors use variational inference. However, as Z is now stochastic, SGD algorithms can not be used to train the parameters of l . To solve this problem, they proposed to parametrize Z so that μ and σ being deterministic with respect to Z , and, by consequence, differentiable with respect to a cost function. In that way, an SGD algorithm can be used to train the parameters of layer l . Eq. (13) presents this approach, so-called reparametrization trick.

$$\begin{aligned} Z &\sim N(\mu, \sigma) \\ Z &= \mu + \epsilon\sigma. \end{aligned} \quad (13)$$

Here, the noise $\epsilon \sim N(0, 1)$ is responsible for the stochasticity in Z .

Even with a significant contribution, the authors in [44] used Z as the last layer of an encoder, not in all network activations or weights. Hence, in 2015, [45] proposed to use this approach to create a BNN considering each weight as a distribution instead of a deterministic variable. The reparametrization trick allowed them to use the SGD algorithm to train the model, and to use VI to approximate the factorized weight posterior distribution to a distribution with known parameters. This approach is called Bayes By BackProp (BBB).

Other approaches, as MC dropout[46] and Variational dropout[47], use dropout to obtain an approximation of a Bayesian model.

In MC dropout, the model must have a dropout function before each weight layer. Thus, the Bayesian approximation is achieved by randomly deactivating weights based on a Bernoulli distribution with the probability of $1-p$, where p is a hyperparameter. The name MC dropout is given once the model prediction is calculated by the average of S Monte Carlo (MC) samples on the model with the dropout enabled.

Variational Dropout uses the local reparametrization and VI to train and to approximate the neural network model of a Bayesian model. With the reparametrization trick in BBB (Eq. (13)), after a layer i receives \mathbf{x}_i as input, it first samples the weights $\boldsymbol{\theta}$ from a Gaussian distribution $N(\mu, \sigma)$ and then computes the activation $\hat{y} = \boldsymbol{\theta}^T \mathbf{x}$ as the inner product between \mathbf{x} and $\boldsymbol{\theta}$. On the other hand, in local reparametrization, the activations are sampled directly from a factorized Gaussian distribution, as shown in Eq. (14):

$$\begin{aligned} \mu &= \boldsymbol{\theta}^T \mathbf{x} \\ \sigma &= (\boldsymbol{\theta}^2)^T \mathbf{x}^2 \\ \hat{y} &\sim N(\mu, \sigma). \end{aligned} \quad (14)$$

where, $\mathbf{b}^2 = \mathbf{b} \circ \mathbf{b}$, where \circ represents the pointwise multiplication operator.

This local reparametrization technique can be used in conjunction with a noise $\xi \sim N(1, \alpha)$ in order to get the posterior $p(\omega|D) = N(\theta, \alpha\theta^2)$, where ω is the variational parameter, θ is the model weight and $\alpha = p/(1-p)$. Eq. (15) presents the variational dropout approach.

$$\hat{y} = \theta^T (\mathbf{x} \circ \xi). \quad (15)$$

As ξ is drawn from a Gaussian distribution, the marginal distribution $\hat{y} = p(\hat{y}|\mathbf{x})$ is also a Gaussian distribution. Thus, one can sample \hat{y} directly from its marginal distribution $p(\hat{y}|\mathbf{x})$, as presented in Eq. (16).

$$\begin{aligned} \mu &= \theta^T \mathbf{x} \\ \sigma &= \alpha(\theta^2)^T \mathbf{x}^2 \\ \hat{y} &\sim N(\mu, \sigma). \end{aligned} \quad (16)$$

Even though p in MC dropout is a hyperparameter, α in variational dropout can be taken as a trained parameter, giving different importance for each element in $(\theta^2)^T \mathbf{x}^2$.

In a Bayesian model, regardless of the particular approach employed to infer the posterior distribution, the prediction of an observation \mathbf{x}^* is calculated by integrating the likelihood of \mathbf{x}^* over the entire posterior distribution (Eq. (17)). As this process involves an intractable integration, an unbiased approximation can be obtained by a Monte Carlo simulation, as presented in Eq. (18).

$$p(y^*|\mathbf{x}^*) = \int p(D|\theta)p(\theta|D)d\theta \quad (17)$$

$$\approx \frac{1}{S} \sum_{s=1}^S p(y^*|\mathbf{x}^*; \theta_s). \quad (18)$$

Here, S is the number of samples, y^* is the probability distribution of classes given \mathbf{x}^* , and $\theta_s \sim p(\theta|D)$ is the s^{th} parameter θ drawn from the posterior $p(\theta|D)$. For the variational dropout model, this posterior is $p(\omega|D)$. However, for MC dropout, this posterior distribution is represented by the dropout function inside each network layer.

5.0.1. Uncertainty

There are two main types of uncertainty in Bayesian modeling: aleatoric and epistemic. Aleatoric is the uncertainty of an event (a.k.a irreducible uncertainty). In a classification problem, this uncertainty is related to the event that generates the data. Therefore, even though some works propose ways to assess the aleatoric uncertainty of a model [48, 49], it is not an easy task to perform, as in most cases, one can not know how the data was sampled or which event generated them.

Epistemic uncertainty assesses the model uncertainty about the data and can be easily calculated when the model is stochastic. This type of uncertainty can be decreased by observing more data. Thus, it is important when one wants to know which class needs more data to improve model prediction. A detailed explanation about uncertainties for Bayesian Deep Neural Networks can be find in [50].

In this work, we are interested in determining the uncertainty of the model's prediction, which corresponds to its epistemic uncertainty. In this sense, the more data it receives during training, the more confident it would be about its predictions. Thereby, actions with fewer samples data would lead the model to uncertain predictions. In this case, it is possible to use epistemic uncertainty to realize when the model should wait for more observations to increase its certainty about prediction, and then anticipate the action correctly.

The epistemic uncertainty of a Bayesian Neural Network model can be estimated by the entropy or the mutual information metrics [50]. In the case of an MC simulation with S samples, a model with C actions can calculate the entropy of theses predictions (samples) by using Eq. (20) and the mutual information by Eq. (21).

$$\mathbb{E}_{pred}(x, c) = \frac{1}{S} \sum_{s=1}^S p(y = c|x; \theta_s), \quad (19)$$

$$\mathbb{H}(x) = - \sum_{c=1}^C \mathbb{E}_{pred}(x, c) \log(\mathbb{E}_{pred}(x, c)), \quad (20)$$

$$\mathbb{I}(x) = \mathbb{H}(x) + \frac{1}{S} \sum_{c=1}^C \sum_{t=1}^S p(y = c|x; \theta_s) \log p(y = c|x; \theta_s). \quad (21)$$

6. Proposal

In this section, we present our proposed architecture, which is divided into four main steps: (i) feature extraction and selection, (ii) feature embedding, (iii) classification model, and (iv) a novel decision-making criterion. The next topics will cover each step in detail.

6.1. Feature extraction and selection

This work proposes to use gaze and object position as the context information, and the evolution of the 2D body joints features as the movement information to perform action anticipation. Our approach aims to use only RGB images, where gaze and skeleton joints information are not straight available.

Therefore, to obtain the gaze and skeleton joints of the people present in the images, we consider to use the Openpose model[51] over each RGB image to extract such information. We used the Openpose version trained for COCO dataset that provides 19 2D joints for the body and 25 2D points for each hand.

It is important to mention that in [18], the authors used gaze and 3D body joints since they had glasses and a MoCap system, while, in our approach, we have only 2D joints to use as data, because we are considering just RGB images. Also, because the actor wore glasses during data acquisition, algorithms for 2D gaze estimation did not work. For this reason, we decided to use the head joints as information that likely may represent head direction or even gaze. However, this representation is a task to be assumed by the model. Aiming to reduce dimensionality, we calculated the central point of each hand instead of using their 25 2D points directly.

For the object information, we extracted the central point of the red ball for each frame using a segmentation method. This pre-processing procedure is summarized in Fig. 6 and described as follows:

1. Openpose model receives an RGB image representing an observation. This operation results in 19 joints and 25 hand points from each user present in the image.
2. A filter to remove false-positive users is applied.
3. Select the most important joints (arms, shoulders, and head)
4. Use hand points to calculate the central point of each hand
5. Give the same RGB image as input to segmentation function to extract the central point of the object.

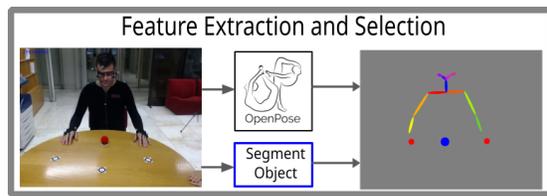


Figure 6: Summary of the feature extraction and selection step.

6.2. Feature Embedding

After the pre-processing step, head information is represented by: (i) five 2D points ($\mathbf{v}_h \in \mathbb{R}^{10 \times 1}$); (ii) object information by one 2D point ($\mathbf{v}_o \in \mathbb{R}^{2 \times 1}$); and (iii) user pose (movement) by nine 2D points ($\mathbf{v}_m \in \mathbb{R}^{18 \times 1}$), where the first seven points represent arms and shoulders, and the last two points represent the hands. Notice that, movement, head, and object have different quantity of points, which generate an unbalanced feature vector. Because of that, the model may consider the movement more important than the other features. Therefore, we propose to balance the input source by using an embedding structure, in such a way that movement and context features have the same dimension. Besides that, to represent the context, head and object features were also defined with the same dimension, so they had the same importance. The embedding process is explained below.

1. Embed head information: $\mathbf{e}_h = f(\mathbf{W}_h^T \mathbf{v}_h + \mathbf{b}_h)$, where $\mathbf{e}_h \in \mathbb{R}^{16 \times 1}$, $\mathbf{W}_h \in \mathbb{R}^{10 \times 16}$ and $\mathbf{b}_h \in \mathbb{R}^{16 \times 1}$.

2. Embed object information: $\mathbf{e}_o = f(\mathbf{W}_o^T \mathbf{v}_o + \mathbf{b}_o)$, where $\mathbf{e}_o \in \mathbb{R}^{16 \times 1}$, $\mathbf{W}_o \in \mathbb{R}^{2 \times 16}$ and $\mathbf{b}_o \in \mathbb{R}^{16 \times 1}$.
3. Embed movement information: $\mathbf{e}_m = f(\mathbf{W}_m^T \mathbf{v}_m + \mathbf{b}_m)$, where $\mathbf{e}_m \in \mathbb{R}^{16 \times 1}$, $\mathbf{W}_m \in \mathbb{R}^{18 \times 16}$ and $\mathbf{b}_m \in \mathbb{R}^{16 \times 1}$.
4. Embed context information: $\mathbf{e}_c = f(\mathbf{W}_c^T [\mathbf{e}_o | \mathbf{e}_h] + \mathbf{b}_c)$, where $\mathbf{e}_c \in \mathbb{R}^{16 \times 1}$, $\mathbf{W}_c \in \mathbb{R}^{32 \times 16}$, and $\mathbf{b}_c \in \mathbb{R}^{16 \times 1}$ and $|$ is a vector concatenation operator.
5. Create the embedded input vector: $\mathbf{e}_{cm} = [\mathbf{e}_c | \mathbf{e}_m]$, where $\mathbf{e}_{cm} \in \mathbb{R}^{32 \times 1}$.

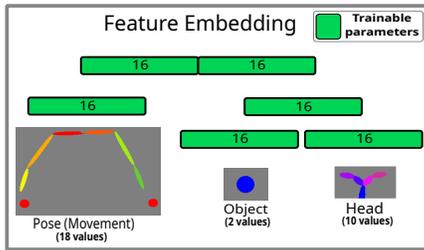


Figure 7: Feature embedding process for each observation. The connections between joints as well as the object shape are showed just for the sake of visualization. However, only their points are used.

Each \mathbf{W}_* and \mathbf{b}_* represents weights that are trained by the model, whereas $f(\bullet)$ represents the ReLU activation function. Thus, during the training phase, the model simultaneously learns to incorporate observations and classify actions.

6.3. Classification model

Since the problem addressed here has a sequential nature, and we assume that there exist dependencies between observations of different timesteps, we can treat the problem in two ways: considering that all sequences are limited to a fixed size of M observations or assuming the original size of sequences.

The problem with the first approach is to disregard the variance in the size of all sequences, besides introducing a new hyperparameter M . Thus, a sequence with T observations ($T > M$) must be truncated at observation $T - M$, meanwhile a sequence with L observations ($L < M$) must be padded with $M - L$ values (an illustration of this process can be

seen in Fig. 8). Therefore, aiming to acquire a score that corresponds to the chance of a sequence up to time t ($t \leq M$) belongs to one action class, we can feed the model with a sequence of size M , where the first t observations came from the real sequence, and the last $M - t$ are padding values. An advantage of this approach is to enable the use of non-sequential models, like Naive Bayes or Multilayer Perceptron, to classify the sequence, since the dependence between observations can be treated as dependence between features. The main problem with this approach is the waste of processing power when the sequence is starting (since the majority of the input data is padded with default values), besides the possible poor performance when the dependence between long sequences must be considered.

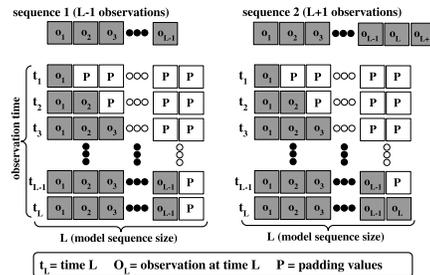


Figure 8: Representation of two sequences with different sizes for a model that receives a fixed input sequence of size M . At each time t , the t^{th} observation from the sequence is added to the fixed model sequence. Therefore, the model can predict the action represented by the t observations.

For the second case, only a sequential model can be used, as the size of the sequence is not available. In this case, models as HMM (Hidden Markov Model) and CRF (Conditional Random Fields) are possible candidates. However, these models assume the Markovian condition: a given observation depends only on the previous one. This assumption might not capture long dependencies on a sequence, which occurs during action execution. Therefore, we decided to use LSTM (Long-Short Term Memory)[52], a variant of RNN that can capture long dependencies in a sequence of observations.

An LSTM contains four trainable gates. These gates are responsible for capturing long and short dependencies in a sequence. An LSTM cell receives

as input an observation vector, a hidden state, and an echo cell. The input vector represents the actual observation; the hidden state represents the short-term memory and chooses what information should be paid attention in the next observation. The echo cell represents the long-term memory. At each new observation, the echo cell stores important pieces of information about the actual observation and forget part of its past when it considers less significant.

LSTM has been used mainly for NLP [8, 7] but in the last few years recognition tasks in videos are commonly using it as well. The Eq. (22)-(27) represent all LSTM gates and activations, respectively:

- Forget gate, \mathbf{f}_t , forget part of the memory stored in the echo cell.

$$\mathbf{f}_t = \sigma_g(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f). \quad (22)$$

- Input gate, \mathbf{i}_t , select part of the observation to be stored into the next echo cell.

$$\mathbf{i}_t = \sigma_g(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i). \quad (23)$$

- Output cell, \mathbf{o}_t , select what part of the input will be propagated to the next observation by the hidden state.

$$\mathbf{o}_t = \sigma_g(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o). \quad (24)$$

- Update gate, \mathbf{g}_t , normalize the observation in order to store it into the next echo cell. Part of this information will be forgot by using the input gate.

$$\mathbf{g}_t = \tanh(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o). \quad (25)$$

- Next echo cell, \mathbf{c}_t , forget part of the past observations and store part of the new one.

$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \mathbf{g}_t. \quad (26)$$

- next hidden state, \mathbf{h}_t , select a part of the normalized echo cell by using the output gate.

$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t). \quad (27)$$

where $\sigma_g(\bullet)$ and $\tanh(\bullet)$ are, respectively, *sigmoid* and *hyperbolic-tangent* activation functions.

Fig. 9 presents the proposed model. It comprises two LSTM layers followed by a softmax classifier. The first LSTM cell receives as input at timestep t the embedded input vector $\mathbf{e}_{cm}^{(t)}$, the hidden state $\mathbf{h}_1^{(t)} \in \mathbb{R}^{64 \times 1}$ and the echo cell $\mathbf{c}_1^{(t)} \in \mathbb{R}^{64 \times 1}$. The second LSTM cell receives as input the hidden state $\mathbf{h}_1^{(t)}$ resulted from the first layer, the hidden state $\mathbf{h}_2^{(t)} \in \mathbb{R}^{64 \times 1}$, and the echo cell $\mathbf{c}_2^{(t)} \in \mathbb{R}^{64 \times 1}$. Next, a fully connected layer receives as input $\mathbf{h}_2^{(t)}$, applies a transformation using a matrix $\mathbf{W}_{fc} \in \mathbb{R}^{d \times 64}$ (d is the number of actions) and normalize it using a *softmax* function. So,

$$\hat{\mathbf{s}}^{(t)} = \text{softmax}(\mathbf{W}_{fc} \mathbf{h}_2^{(t)}),$$

where $\hat{\mathbf{s}}^{(t)} \in \mathbb{R}^{d \times 1}$, and each element of $\hat{\mathbf{s}}^{(t)}$ can be interpreted as the probability of an action given the embedded input vector in the timestep t . Now, With this result and choosing a probability threshold value p , by using Eq. (5), the anticipation can be accomplished.

6.4. Decision-making criterion

As mentioned before, this work proposes a novel decision-making criterion based on the model uncertainty. In this sense, because the proposed model is deterministic, we propose to use three new stochastic versions of it: A Bayesian LSTM using Bayes By BackProp (BLSTM_{BBB}), an MC dropout Bayesian LSTM based on [53] (BLSTM_{MC}), and a Variational Dropout Bayesian LSTM (BLSTM_{VD}).

With these stochastic models, the uncertainty is obtained by running the architecture of Fig. 9 S times (MC sampling) using the same input sequence. As the models are stochastics, they must give a different value for each prediction. Thus, the Mutual Information (MI) over the S predictions give us the epistemic model uncertainty about the class prediction for the respective observation. MI is calculated using Eq. (21). Therefore, we propose to use a threshold over the mutual information (our new decision-making criterion) to anticipate actions. For this, Eq. (5) must be redefined as:

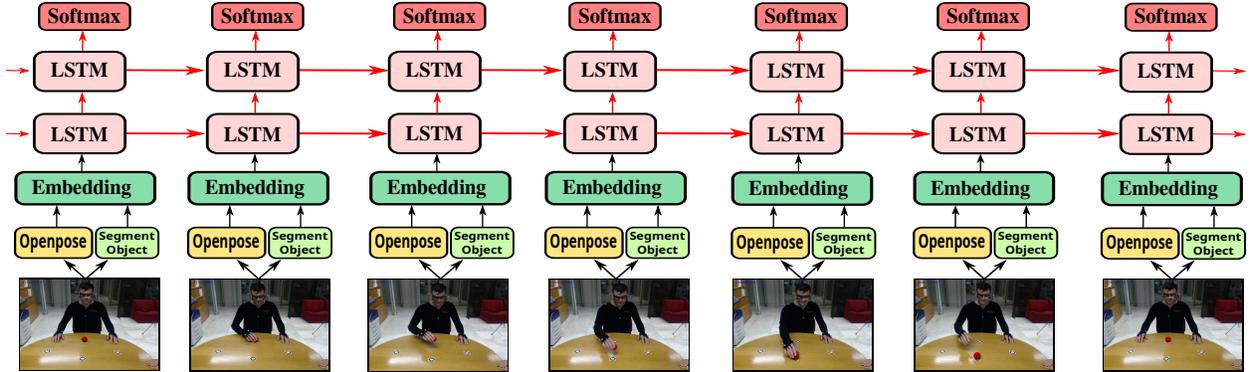


Figure 9: Proposed model architecture

$$g(\hat{\mathbf{s}}^{(t)}, u) = \begin{cases} \operatorname{argmax}(m(\hat{\mathbf{s}}^{(t)})), & h(\hat{\mathbf{s}}^{(t)}) < u \\ -1, & \text{otherwise} \end{cases}, \quad (28)$$

where u is an uncertainty value, h is the mutual information function (Eq. (21)), and m is the average of the S predictions.

Even though one can use the entropy (Eq. (20)) to measure the uncertainty, we chose to use MI because it takes into account not only entropy between classes (averaged over the S predictions) but also the mean entropy between them all.

7. Experiments

As discussed in Sec. 4, all experiments in this work used the Acticipate dataset. From the dataset, we extracted four different kinds of data by using the procedure described in Sec. 6: head points, object position, arm joints, and hand position. The head points and object position forms the context; arm joints and hand positions form a pose, which evolution in time represents the movement (also called here main entity). To better compare results and reach more reliable conclusions about how each source of information influences the action anticipation, we decided to carry out experiments using different combinations of the context (head and object) and movement, for the original version of the dataset (6 actions) and its extended version (12 actions).

The most common approaches, presented in Sec. 2, need a large dataset to be trained. Hence, they are not suitable to be used here, with the Acticipate dataset. However, to better compare and discuss our results, we used as baselines the proposal in [18], that uses Acticipate dataset, in addition to five classical models: Naive Bayes (NB), Multilayer Perceptron (MLP), 1D Convolutional Neural Network (CONV-1D), Support Vector Machine (SVM) and HMM. For the first four models, as mentioned before (Sec. 6), we used a fixed sequence size (as illustrated in Fig. 8), while for HMM we used the sequences with their original sizes.

Considering the architecture of the baseline models:

- NB uses a Gaussian Model to predict its conditional probability;
- MLP is composed of only one hidden layer;
- CONV-1D has three stacked 1D Convolutional layers followed by an MLP as classifier;
- SVM implements its non-linear version by using an RBF kernel; and
- HMM has its emissivity probability drawn from a Gaussian Distribution, which enables continuous observations.

The five models were trained with the original dataset (6 actions) and its extended version (12 actions). For the experiments with these five models,

we did not apply the proposed embedding technique. So each observation is represented by a vector with 18 values corresponding to the arms, two values for the object, and 10 for head points (see Sec. 6.2). By combining these features (movement, object, and head), we carried out a total of 45 different baseline experiments. However, for the sake of explanations, we will present only the most conclusive ones. We also carried out more 11 experiments with different versions of the proposed model, which will be explained next. Table 1 summarizes the main 16 experiments.

The experiments 6 and 7 (original dataset) in Table 1 provide results that can be compared with [18] and the baseline experiments (5 first experiments), which will validate our proposal against the other models. The experiments from 8 to 13 provide results to show the ambiguities between actions and the importance of context to anticipate them. The last three experiments show the importance of uncertainty in an anticipation model, and the contribution of the proposed decision-making criterion based on the uncertainty. For this reason, three Bayesian models were implemented: MC dropout, Variational Dropout, and Bayes by Backprop. For variational dropout, as mentioned before, we opt to use α (Eq. (16)) as a trainable parameter. With these three models, we can identify which model is best for this kind of application.

7.1. Experiment setups

For each RGB image in the video, the pre-processing procedure described in Sec. 6 was applied. Every missing data related to joints, hands and the object position were set to -1. Additionally the padding value used in the fixed sequence size for the four first experiments was also considered -1. To evaluate the model’s quality, the dataset was divided into 80% for training and 20% for testing. In each experiment, a 10-fold cross-validation over the training set was performed, where nine folds were used to train and one fold to validate. For each fold (round), the training process was finished when the recognition accuracy (the accuracy achieved at the last observation of the sequence) over the nine training folders in the previous five epochs was higher than 98% (early

stop) or when the iterations exceeded the maximum number of epochs.

The 16 experiments in Table 1 can be divided into four main categories: baselines with 6 actions (from 1 to 5), deterministic with 6 actions (6 and 7), deterministic with 12 actions (from 8 to 13) and stochastic with 12 actions (14,15 and 16). In both deterministic experiments, different configurations of the input data (movement, head and object) were achieved by assigning -1 to \mathbf{v}_m , \mathbf{v}_o , and/or \mathbf{v}_h in the entire dataset. For instance, by assigning -1 to \mathbf{v}_o , the model considers only movement (\mathbf{v}_m) and head (\mathbf{v}_h) information. For each one of the 16 experiments, proper hyperparameters were chosen by using a Bayesian Optimization process. The hyperparameters used in the tests are presented in Table 2 and 3.

7.2. Software and hardware environments

The models, including Openpose (a deep neural network), were implemented in Pytorch v1.0 and scikit-learn v0.22. The Bayesian Optimization was implemented using the library Hyperopt³. Additional parts, as object segmentation, filters, and chart plot scripts, were implemented using OpenCV v4.1, Python v3.7, Numpy v1.16.4, and Matplotlib v2.2.3. The computer used in the experiments had the following configuration:

- Linux Operating System, distribution Ubuntu Server 16.04;
- Intel Core i7-7700 processor, 3.60 GHz with four physical cores;
- 32 GB of RAM;
- 1 TB of storage unit (hard drive);
- Nvidia Titan V graphic card.

8. Results and Discussions

For each one of the 16 experiments, after running the Bayesian Optimization over the 10-fold cross-validation, the best hyperparameter configuration

³<https://github.com/hyperopt/hyperopt>

Table 1: List of all experiments

Id	Model	Type	Movement	Context		Dataset Version	
				Head	Object	6 actions	12 actions
1	<i>NB</i>	baseline	✓	✓		✓	
2	<i>CONV-1D</i>	baseline	✓	✓		✓	
3	<i>MLP</i>	baseline	✓	✓		✓	
4	<i>SVM</i>	baseline	✓	✓		✓	
5	<i>HMM</i>	baseline	✓	✓		✓	
6	<i>DLSTM_{6m}</i>	Deterministic	✓			✓	
7	<i>DLSTM_{6mh}</i>	Deterministic	✓	✓		✓	
8	<i>DLSTM_{12m}</i>	Deterministic	✓				✓
9	<i>DLSTM_{12h}</i>	Deterministic		✓			✓
10	<i>DLSTM_{12o}</i>	Deterministic			✓		✓
11	<i>DLSTM_{12mh}</i>	Deterministic	✓	✓			✓
12	<i>DLSTM_{12mo}</i>	Deterministic	✓		✓		✓
13	<i>DLSTM_{12mho}</i>	Deterministic	✓	✓	✓		✓
14	<i>BLSTM_{MC}</i>	MC Dropout	✓	✓	✓		✓
15	<i>BLSTM_{VD}</i>	Variational Dropout	✓	✓	✓		✓
16	<i>BLSTM_{BBB}</i>	Bayes By BackProp	✓	✓	✓		✓

Table 2: List of hyperparameters used in each experiment with the proposed models

Hyperparameters	Model Configurations			
	<i>DLSTM*</i>	<i>BLSTM_{MC}</i>	<i>BLSTM_{VD}</i>	<i>BLSTM_{BBB}</i>
Batch size	216	216	216	32
Truncate sequence	100	100	100	128
Sequence size	100	100	100	64
Max epochs	100	100	100	200
LR	1e-2	1e-2	1e-2	1e-2
LR decay (per epoch)	1%	1%	1%	1%
Weight decay	1e-5	1e-5	-	-
Dropout (keep prob)	0.7	0.2	-	-
Optimizer	Adam	Adam	Adam	Adam

Table 3: List of hyperparameters used in each baseline experiment

Hyperparameters	NB	MLP	SVM	CONV-1D	HMM
Architecture	Gaussian	hidden layer (48)	Kernel (RBF)	Conv Layers (1x7x32, 1x5x32, 1x5x32) Hiddens layer (64)	States (5)
Sequence Size	72	64	64	96	64
Max Epoch	-	50	-	50	-
Batch Size	-	32	-	32	-
Learning Rate	-	1e-3	-	1e-3	-
Optimizer	-	Adam	-	Adam	-

found was used to train the model with the complete training set. Then, each model was tested with the test set, which was never seen by the model during training. The results over the test set were plotted in many charts and carefully analyzed. Some of these charts are presented and discussed in this section.

8.1. Baseline models in the original dataset

As can be seen in Table 4, the classical models, used as baselines, offer a satisfactory result when one analyzes only the accuracy obtained at each observation ratio. However, as mentioned before, this type of analysis brings poor conclusions about the model’s capacity to anticipate the actions. A good anticipation model should increase the differentiation between classes, while the number of observations also increases. Therefore, let’s look at the graphs in Fig. 10, where we have the distribution of probabilities between classes for each observation ratio of an action *give middle*. Note that the baseline models fail to accumulate knowledge about the performing actions (NB, HMM, MLP, CONV-1D, and SVM) or even respond over-confidently (NB) about prediction. In HMM the distribution is almost uniform, making it difficult to determine a proper probability to be used as a decision-making threshold. Naive Bayes responds with high confidence even at the beginning of the action. So it is not a reliable model to be used in anticipation tasks. The other three models, MLP, SVM, and CONV-1D are quite noisy and do not give confidence about the correct action. In contrast, the proposed DLSTM was capable of representing the evolution of the action, in such a way that clearly differentiates one action from the others while more observations are provided.

For a more complete analysis, Table 5 brings the anticipation accuracy of each model calculated by Eq. (9). As we can see, the baseline models, even been quite effective in the recognition task, fail when applied to anticipation task. The max accuracy achieved by the HMM was 76.19% when using a threshold of 0.38. With such a small threshold, it is difficult to trust in this model, once at least one of the remaining classes can reach close values (e.g., 0.37). The best among the other baseline models (MLP) achieved only 61.90% with a threshold of 0.94. While

Naive Bayes could not even anticipate an action due to its over-confidence. On the other hand, DLSTM was capable of anticipating 95.25% of actions when applying a threshold of 0.90. All these pieces of evidence show us the superiority of the proposed model compared with the baselines.

Another important result is that our proposal, even using 2D skeleton joints extracted from images, outperforms [18], that used eye gaze and 3D pose (Tab. 4). With *movement + head* information we achieved 90% of accuracy with less than 40% of observations. On the other hand, the authors’ model, in [18], achieves 90% of accuracy after more than 50% of observations. As discussed before, an effective anticipation model must also be an effective recognizer. Our model $DLSTM_{6mh}$ recognized all actions at the last observation (100% of average accuracy); meanwhile, their model achieved a maximum of 97.5% with an observation ratio 0.60 and decreased to 87% with an observation ratio 1.0. Therefore, their model did not recognize all actions in the dataset.

In addition to the results above, $DLSTM_{6mh}$ can anticipate an action three frames before than $DLSTM_{6m}$, on average. As the video has a sample rate of 30Hz, this anticipation corresponds to 100ms, which is greater than the 92ms presented in [18] when comparing *pose* with *pose+gaze*. As such, besides outperforming [18], our proposal was able to solve the action recognition problem in the Anticipate dataset and improve the action anticipation results.

8.2. Deterministic models in the extended dataset

Fig. 11 shows the results for the 6 experiments (from 8 to 13) using the extended version of the dataset. We can observe that $DLSTM_{12m}$, $DLSTM_{12h}$, and $DLSTM_{12mh}$ did not achieve 100% accuracy at the last observation. This result shows that they were unable to separate actions properly, even when using head information. The 6 new actions are the only difference between $DLSTM_{6mh}$ and these three models. As such, when the dataset was divided into more actions, more ambiguities were generated among them. Therefore, these results support our first hypothesis: more actions are likely to cause more ambiguities.

Table 4: Comparing the results obtained by the models (baselines and DLSTM) at different observation ratios. The results for [18] were provided by the authors. The subscription *6mh* indicates that the correspondent model was trained with the original dataset (6 actions) using movement (m) and head (h) as source of information.

Models	Percentage of observation									
	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
<i>HMM</i> _{6mh}	19.05%	19.05%	38.10%	57.14%	90.48%	95.24%	100.00%	100.00%	100.00%	100.00%
<i>NB</i> _{6mh}	19.05%	9.52%	9.52%	9.52%	19.05%	33.33%	71.43%	95.24%	95.24%	95.24%
<i>CONV</i> _{6mh}	9.52%	14.29%	19.05%	19.05%	33.33%	47.62%	42.86%	61.90%	61.90%	80.95%
<i>MLP</i> _{6mh}	19.05%	19.05%	28.57%	28.57%	33.33%	47.62%	71.43%	90.48%	85.71%	95.24%
<i>SVM</i> _{6mh}	19.05%	19.05%	19.05%	19.05%	28.57%	71.43%	95.24%	95.24%	95.24%	95.24%
<i>DLSTM</i> _{6m}	14.29%	28.57%	38.10%	71.43%	85.71%	95.24%	95.24%	95.24%	95.24%	95.24%
<i>DLSTM</i> _{6mh}	14.29%	38.10%	47.62%	90.48%	95.24%	100.00%	100.00%	100.00%	100.00%	100.00%
[18] (Pose 3D)	16.25%	21.25%	28.75%	51.25%	76.25%	85.00%	86.25%	86.25%	86.25%	85.00%
[18] (Pose 3D + Gaze)	15.00%	16.25%	40.00%	73.75%	86.25%	97.50%	96.25%	92.50%	91.25%	87.50%

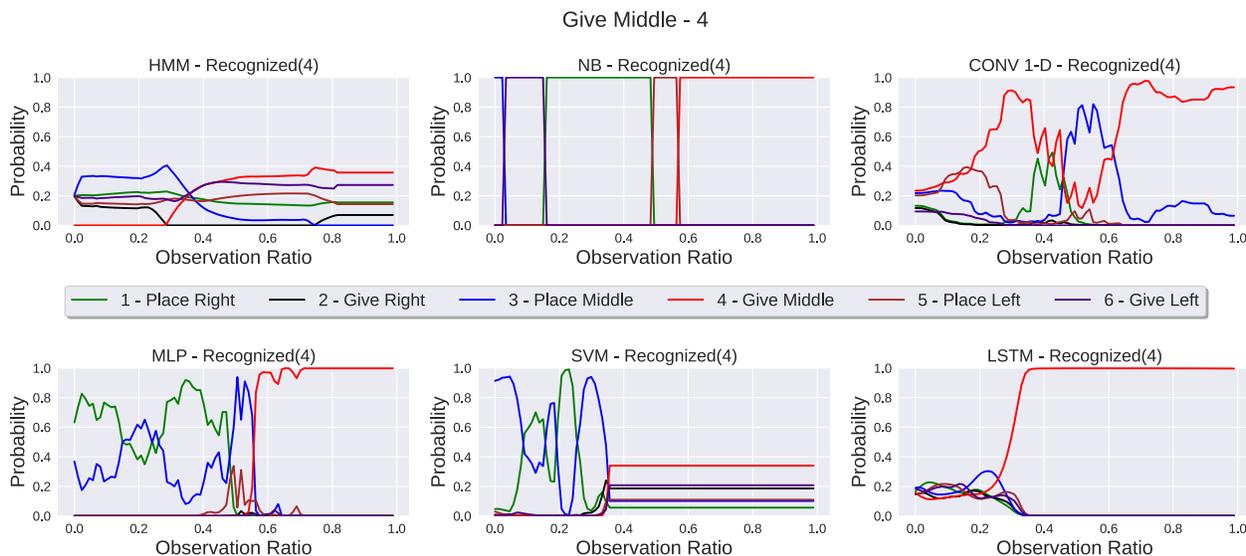


Figure 10: Evolution of an action *give middle* with its respective prediction for each baseline model and the proposed DLSTM. All models are trained on the original dataset (6 actions) with movement and head information.

Table 5: Maximum anticipation accuracy (second column) obtained by each model (first column) when applying the probability threshold (last column). To achieve the corresponding accuracy the model needed on average an observation ratio such that specified in third column. As the NB model is extremely over-confident in its predictions, it is not suitable to anticipate actions.

Model	Anticipation Accuracy	Mean Observation Ratio Needed	Probability Threshold
HMM	76.19%	0.45	0.38
NB	It is not possible to anticipate		
CONV-1D	47.62%	0.69	0.91
SVM	57.14%	0.55	0.68
MLP	61.90%	0.60	0.94
DLSTM	95.25%	0.40	0.90

The models that use object information (*DLSTM*_{12o}, *DLSTM*_{12mo}, *DLSTM*_{12mho}) were able to recognize all actions (accuracy at 100% of observations). Further, they achieved better results in the anticipation task. They start with more than 40% of accuracy at the first observation, and the model with complete context (*DLSTM*_{12mho}) achieves 98% of accuracy at the observation ratio of 0.42. Therefore, actions with similar movements can be distinguished better when using context information. So, this result supports our second hypothesis: context information can help to distinguish different actions represented by similar movements. We can also see how this last model was able to extract

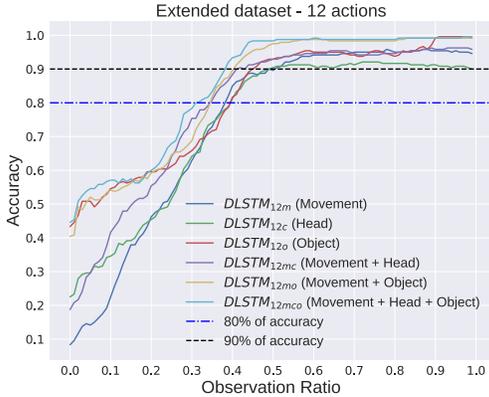


Figure 11: Results for deterministic models in the extended dataset (12 actions). Each experiment is the same model trained with different input data.

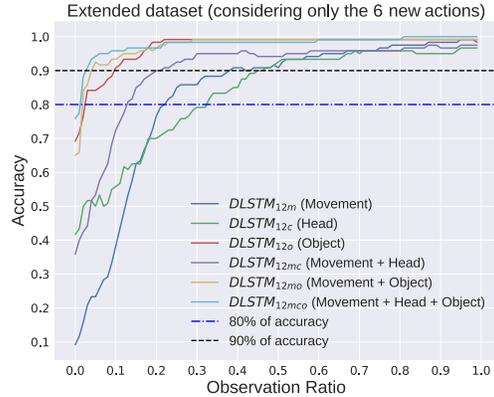


Figure 12: Results for deterministic models in the extended dataset but considering only the 6 new actions.

relevant information from the object position and head points. Even though the object is represented by only two values (a 2-dimensional point), as we suppose, it provided an important information about the actions to the model. This result shows the efficacy of our feature embedding process.

As mentioned above, after using the object information, the model might be able to anticipate some actions after a few observations. For better visualization, Fig. 12 illustrates the accuracy of the same 6 models for the 6 new added actions (receive and pick (left, middle, right)). The models that use object context start with a classification accuracy greater than 65% and achieve 90% of accuracy after less than 10% of observations. The best model reaches 95% of accuracy with less than 5% of observations, on average. In terms of frames, for the Acticipate dataset, that corresponds to an average of 4 frames. These results support our statement about the importance of object information for these 6 new actions.

To measure the anticipation accuracy, we use the Eq. 9 with a threshold $p = 0.9$ for the 6 models. Fig. 13 presents the evolution of an action *pick right* after passing throughout the 6 models. The charts illustrate how the model that uses only movement ($DLSTM_{12m}$) made a mistake in its anticipation. This mistake can be caused by the overconfidence of the model when anticipating ambiguous actions.

Other models, which use part/complete context information, anticipated the action correctly. Notice that the model with full context (*head + object*), anticipated the action after observing only 2% of the data sequence (two frames in its corresponding video). Another interesting result is that the models confused those classes we supposed they would. After analyzing the videos, one can notice that action *pick right* has similar movement to *place right*, *give right* and *receive right*, and similar gaze to *place right*. Thus, $DLSTM_{12m}$ *pick right* mistook for *receive right* and $DLSTM_{12h}$ was not certain about *pick right* and *place right*. These characteristics appear in almost all predictions.

To highlight the trade-off between the threshold and the anticipation accuracy, the chart in Fig. 14 presents the variation of anticipation accuracy and the percentage of observations w.r.t threshold (p). In the chart, we see that when $p = 0.9$, $DLSTM_{12mho}$ can anticipate correctly 95.42% of actions by using, on average, 19% of the video sequence. Each action in Acticipate dataset has an average of 79 images. Thus, this 19% in observation ratio corresponds to an average of 15 frames of a video. On the other hand, by aiming the minimum number of observations, with $p = 0.8$, the model would anticipate correctly 92.50% of actions by using, on average, 18% of observation (14 frames).

With $p = 0.9$, function g in Eq. (5) is not able

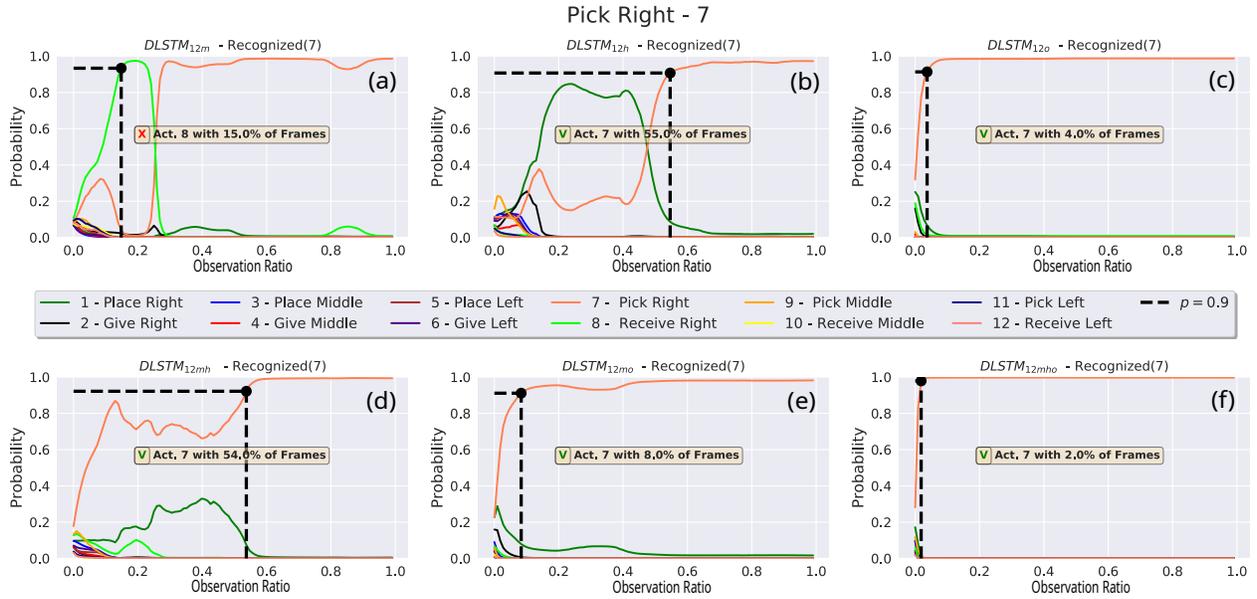


Figure 13: Evolution of sample of a *pick right* action for the 6 deterministic models.

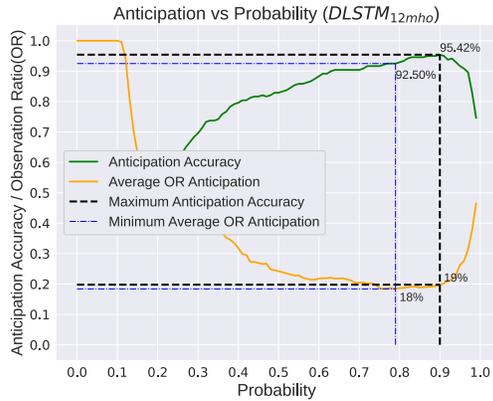


Figure 14: Variation of anticipation accuracy and average observation ratio w.r.t probability threshold.

to consider overconfidence in the model prediction, which may generate many false-positives. As discussed in Sec. 4, a possible solution to reduce the number of false-positives is to force the model to wait for more z observations to reaffirm its prediction. The problem with this approach is that z is a new parameter that can harm the anticipation, and must

be chosen carefully. Fig. 15 illustrates how anticipation accuracy ($p = 0.9$) and average observation ratio vary w.r.t z , where z is the additional observation ratio after anticipation. The best anticipation accuracy (97.08%) is achieved when $z = 0.18$. In other words, the model needs to wait on average for more 18% of observations to achieve an anticipation accuracy of 97.08%. Comparing with previous results, the gain of less than 2% in the accuracy cost an increase of more than the double of observations to the anticipation time (passing from 19% to 43%). Furthermore, the minimum observation ratio necessary to anticipate any action is now 18%, even for less ambiguous actions, such as those presented previously in Fig. 4 and Fig. 13. As conclusion, besides the fact that the choice of z inserts a new trade-off in the project (accuracy *vs* observation ratio), it does not provide an effective way to improve the action anticipation task.

8.3. Stochastic models

The results of the Bayesian models $LSTM_{MC}$, $BLSTM_{VD}$ and $BLSTM_{BBB}$ will be compared to $DLSTM_{12mho}$, our best deterministic model for the extended dataset. During prediction time, we fed

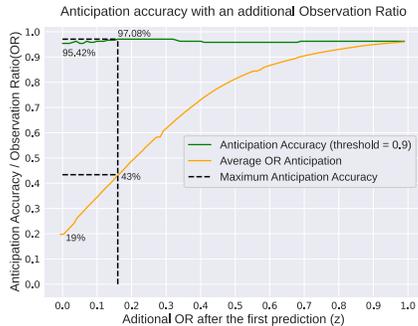


Figure 15: Variation of anticipation accuracy and average observation ratio w.r.t additional observation ratio after anticipation. If in time t the max probability exceeds the 0.9, the model must wait for more z observations in order to conform its prediction.

each Bayesian model 20 times with the same observation \mathbf{x}_t , which corresponds to an MC sampling with $S = 20$. Next, by applying Eq. (21) over the S predictions, we measured the epistemic uncertainty of each model prediction, concerning the observation \mathbf{x}_t . Then, we could use Eq. (28) to anticipate the action or not.

The Bayesian models also recognized all actions in the extended dataset. Furthermore, they achieved better results in anticipation accuracy than $DLSTM_{12mho}$, even if it waits for an observation ratio of $z = 0.18$. By applying the same procedure in Fig. 14, we could choose a threshold value to be used in each model. Therefore, for each model, the anticipation threshold was chosen by analyzing the variation of anticipation accuracy and the average observation time w.r.t the uncertainty value. Fig. 16 shows this comparison for $BLSTM_{MC}$.

Table 6 compares the results of the Bayesian models with our best deterministic model ($DLSTM_{12mho}$). Note that $BLSTM_{MC}$ achieves the best anticipation accuracy (98.75%) using the uncertainty threshold $u = 0.5$. However, $BLSTM_{VD}$ and $BLSTM_{BBB}$ also achieves satisfactory results: with $u = 0.5$, $BLSTM_{VD}$ achieved 98.33% of anticipation accuracy, and with $u = 0.3$ $BLSTM_{BBB}$ achieves 97.08%.

Considering the minimum number of observa-

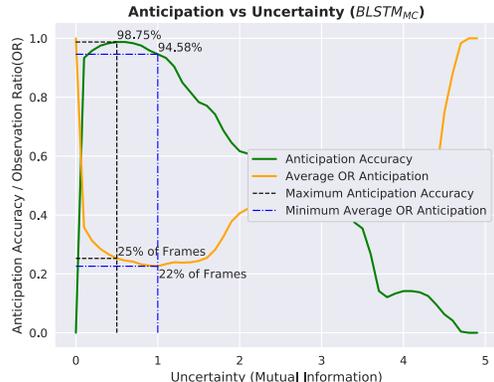


Figure 16: Variation of anticipation accuracy and average observation ratio w.r.t uncertainty threshold for MC dropout model

tions necessary for a good anticipation accuracy, $DLSTM_{12mho}$ gives the best result. On average, with $p = 0.79$, it needs to receive 18% of observations to achieve an anticipation accuracy of 92.50%. However, even needing less observations (18%) with $p = 0.79$, it presents less accuracy (92.50%) than considering more observations at a threshold of $p = 0.9$, which achieves an accuracy of 95.42%. Therefore, we can see that there is a tradeoff of accuracy for a less number of observations. In summary, the best model achieved 98.75% after observing, on average, 25% of the action ($BLSTM_{MC}$). An increase of 6.67% in the accuracy with a cost of only 7% in extra observations. Much better than using z in the deterministic model, where an increase of less than 2% costs 24% on extra observations. Besides, we do not need to choose more than one hyperparameter, only the uncertainty threshold u .

8.4. Discussions

As we mentioned in the previous sections, the model must have a short anticipation time for human-machine interaction and be accurate in its prediction. For $BLSTM_{MC}$, it needs 25% of observations to achieve its best prediction value, which indeed is not a high value. For instance, in a system based on images sampled at 30Hz (ordinary cameras), an action that lasts 2 seconds would be anticipated by such a model after elapsed on average 0.5s from its

Table 6: Results of stochastic models and the best deterministic model.

Model	Parameter	Anticipation Accuracy	Average Observation Ratio
<i>DLSTM</i> _{12mho}	$p = 0.9 / z = 0.0$	95.42%	19%
<i>DLSTM</i> _{12mho}	$p = 0.79 / z = 0.0$	92.50%	18%
<i>DLSTM</i> _{12mho}	$p = 0.9 / z = 0.18$	97.08%	43%
<i>BLSTM</i> _{MC}	$u = 0.5$	98.75%	25%
<i>BLSTM</i> _{MC}	$u = 1.5$	94.58%	22%
<i>BLSTM</i> _{VD}	$u = 0.5$	98.33%	26%
<i>BLSTM</i> _{VD}	$u = 1.5$	85.42%	20%
<i>BLSTM</i> _{BBB}	$u = 0.3$	97.08%	25%
<i>BLSTM</i> _{BBB}	$u = 1.3$	93.33%	20%

first frame. In other words, it might anticipate an action after the system observes, on average, 15 frames. Therefore, once the model can be considered accurate in its prediction, the system has about 1.5s to make a right decision.

As expected, our Bayesian models provided better results than deterministic ones with a small cost in additional observations. The overconfidence in model prediction decreases when waiting for more observations. However, as we could see, for deterministic models, this is a new parameter to be chosen (z) and did not provide satisfactory results. On the other hand, by using uncertainty as a threshold, we have only one parameter to be chosen, and the model can achieve better results of accuracy with a small cost in the observation ratio. These results support our last hypothesis that: uncertainty is a more reliable and effective threshold to anticipate actions than probability values.

In our opinion, the MC dropout [46, 53] and variational dropout [47] were the best models implemented in this work. Once dropout and local reparametrization can provide a different sample for each observation, a mini-batch with S observations correspond to an MC sampling of size S , which helps the model posterior distribution inference. Besides, for prediction, we only need to create a mini-batch of size S , repeating the same observation, that favors parallel prediction in GPUs. On the other hand, the reparametrization trick does not take advantage of the mini-batch to make samples. Every observation in the mini-batch uses the same sampled weight.

As a consequence, in our experience, BBB models train slower than Bayesian dropout approaches, and, during prediction, it needs to run the model S times with the same observation, which makes impractical to parallelize the prediction in GPUs. However, it seems that a significant advantage of BBB is the possibility of pruning the model by analyzing each parameter. As they are Gaussian distributions, the relation mean-variance can indicate if a parameter is required or maybe discarded [43, 45].

Finally, we could see that the proposed model outperformed the baselines, including [18], even using less accurate information (2D vs. 3D pose and head joints vs. eye gaze). The results supported the raised hypotheses and showed how the uncertainty provided by Bayesian models is vital for action anticipation. Our proposal can be used in other datasets even though the presented results were acquired in a small collaborative dataset. In this sense, it is necessary to analyze the possible sources of context for each class and adapt our embedding layer to represent all the context data.

9. Conclusions and Future Works

Machines need the capacity of anticipating actions to achieve effective interaction with humans. As such, the problem of action anticipation is drawing substantial research attention in recent years. Although many works have explored the issue, they do not provide a concise explanation about the importance of context in anticipating actions. They do not discuss

how to handle the problem of the uncertainty inherent in this kind of task and how to make decisions in a real-time situation.

We propose a decision-making criterion based on the uncertainty provided by a stochastic (Bayesian) LSTM model that can practically be used for action anticipation tasks. By selecting the action that minimizes the uncertainty, our model improves the action anticipation performance compared with the conventional class-likelihood maximization (i.e., deterministic model).

Considering arm motion as the primary source of information for action anticipation, we evaluate the influence of two additional (contextual) sources of information in the Acticipate dataset: gaze and object attributes. When considering all information sources in our stochastic LSTM, we achieved 100% of average accuracy in the action recognition task and 98.78% of average accuracy in the action anticipation task, outperforming previous results. Thus, our model serves both action recognition and anticipation purposes, while needing only 25% of the observations, on average, to anticipate each action. The results also show the evident importance of context for the anticipation task, since the actions that depend on the eye gaze information or the object position had impressive improvement in their anticipation time contrasted with [18]. For instance, actions that depend exclusively on the object information are anticipated precociously, some of them with only two observations.

Our work extends the current state-of-the-art and results in action anticipation, for small collaborative datasets. Also, our proposal uses context information to improve the classification probability, and the uncertainty as the decision-making criterion that can be used with any other probabilistic model.

As future work, we aim to increase the collaborative setup complexity by adding more objects to each action and designing a collaborative scenario where the performed actions depend on more than one object. Another important issue to be addressed is how to extract proper context from general datasets. In such a way, we could use this proposal to solve more complex anticipation tasks.

Acknowledgments

This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES), PDSE/Process n 88881.188840/2018-01, Finance Code 001. The authors also would like to acknowledge the support from NVIDIA Corporation through the donation of the Titan V GPU used in this research.

Bibliography

References

- [1] N. F. Duarte, M. Raković, J. Tasevski, M. I. Coco, A. Billard, J. Santos-Victor, Action anticipation: reading the intentions of humans and robots, *IEEE Robotics and Automation Letters* 3 (4) (2018) 4132–4139. doi:10.1109/LRA.2018.2861569.
- [2] Y. Kong, Y. Fu, Human action recognition and prediction: A survey, arXiv preprint arXiv:1806.11230.
- [3] B. Liu, E. Adeli, Z. Cao, K.-H. Lee, A. Shenoi, A. Gaidon, J. C. Niebles, Spatiotemporal relationship reasoning for pedestrian intent prediction, *IEEE Robotics and Automation Letters* 5 (2) (2020) 3485–3492. doi:10.1109/LRA.2020.2976305.
- [4] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. doi:10.1109/CVPR.2016.90.
- [5] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105. doi:10.1145/3065386.
- [6] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.

- [7] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019, pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in neural information processing systems, 2017, pp. 5998–6008.
- [9] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, L. Fei-Fei, Large-scale video classification with convolutional neural networks, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2014, pp. 1725–1732. doi:10.1109/cvpr.2014.223.
- [10] K. Simonyan, A. Zisserman, Two-stream convolutional networks for action recognition in videos, in: Advances in neural information processing systems, 2014, pp. 568–576.
- [11] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model and the kinetics dataset, in: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 6299–6308. doi:10.1109/cvpr.2017.502.
- [12] H. Bilen, B. Fernando, E. Gavves, A. Vedaldi, Action recognition with dynamic image networks, IEEE transactions on pattern analysis and machine intelligence 40 (12) (2017) 2799–2813. doi:10.1109/tpami.2017.2769085.
- [13] C. Rodriguez, B. Fernando, H. Li, Action anticipation by predicting future dynamic images, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 0–0. doi:10.1007/978-3-030-11015-4_10.
- [14] V. Choutas, P. Weinzaepfel, J. Revaud, C. Schmid, Potion: Pose motion representation for action recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7024–7033. doi:10.1109/cvpr.2018.00734.
- [15] H. Kwon, Y. Kim, J. S. Lee, M. Cho, First person action recognition via two-stream convnet with long-term fusion pooling, Pattern Recognition Letters 112 (2018) 161–167. doi:10.1016/j.patrec.2018.07.011.
- [16] F. Baradel, N. Neverova, C. Wolf, J. Mille, G. Mori, Object level visual reasoning in videos, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 105–121.
- [17] Y. LeCun, Y. Bengio, et al., Convolutional networks for images, speech, and time series, The handbook of brain theory and neural networks 3361 (10) (1995) 1995.
- [18] P. Schydlo, M. Rakovic, L. Jamone, J. Santos-Victor, Anticipation in human-robot cooperation: A recurrent neural network approach for multiple action sequences prediction, in: 2018 IEEE International Conference on Robotics and Automation (ICRA), IEEE, 2018, pp. 1–6. doi:10.1109/ICRA.2018.8460924.
- [19] A. K. Dey, G. D. Abowd, Towards a Better Understanding of Context and Context-Awareness, Computing Systems 40 (3) (1999) 304–307. doi:10.1007/3-540-48157-5_29.
- [20] H. Wang, J. Feng, Delving into 3d action anticipation from streaming videos, arXiv preprint arXiv:1906.06521.
- [21] F. Pirri, L. Mauro, E. Alati, V. Ntouskos, M. Izadpanahkakhk, E. Omrani, Anticipation and next action forecasting in video: an end-to-end model with memory, arXiv preprint arXiv:1901.03728.
- [22] S. Agethen, H.-C. Lee, W. H. Hsu, Anticipation of human actions with pose-based fine-grained representations, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0–0.

- [23] Y. Shi, B. Fernando, R. Hartley, Action anticipation with rbf kernelized feature mapping rnn, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 301–317. doi:10.1007/978-3-030-01249-6_19.
- [24] J.-F. Hu, W.-S. Zheng, L. Ma, G. Wang, J.-H. Lai, J. Zhang, Early action prediction by soft regression, IEEE transactions on pattern analysis and machine intelligence doi:10.1109/TPAMI.2018.2863279.
- [25] M. Sadegh Aliakbarian, F. Sadat Saleh, M. Salzmann, B. Fernando, L. Petersson, L. Andersson, Encouraging lstms to anticipate actions very early, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 280–289. doi:10.1109/iccv.2017.39.
- [26] H. Bilen, B. Fernando, E. Gavves, A. Vedaldi, S. Gould, Dynamic image networks for action recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 3034–3042. doi:10.1109/cvpr.2016.331.
- [27] M. Sadegh Aliakbarian, F. Sadat Saleh, M. Salzmann, B. Fernando, L. Petersson, L. Andersson, Encouraging lstms to anticipate actions very early, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 280–289. doi:10.1109/ICCV.2017.39.
- [28] M. S. Aliakbarian, F. Saleh, B. Fernando, M. Salzmann, L. Petersson, L. Andersson, Deep action-and context-aware sequence learning for activity recognition and anticipation, arXiv preprint arXiv:1611.05520.
- [29] S. Gite, H. Agrawal, K. Kotecha, Early anticipation of drivers maneuver in semiautonomous vehicles using deep learning, Progress in Artificial Intelligence 8 (3) (2019) 293–305. doi:10.1007/s13748-019-00177-z.
- [30] S. Gite, H. Agrawal, Early prediction of driver’s action using deep neural networks, International Journal of Information Retrieval Research (IJIRR) 9 (2) (2019) 11–27. doi:10.4018/IJIRR.2019040102.
- [31] D. Wang, Y. Yuan, Q. Wang, Early action prediction with generative adversarial networks, IEEE Access 7 (2019) 35795–35804. doi:10.1109/ACCESS.2019.2904857.
- [32] J. Liu, A. Shahroudy, G. Wang, L.-Y. Duan, A. K. Chichung, Skeleton-based online action prediction using scale selection network, IEEE transactions on pattern analysis and machine intelligence doi:10.1109/TPAMI.2019.2898954.
- [33] A. Jain, H. S. Koppula, B. Raghavan, S. Soh, A. Saxena, Car that knows before you do: Anticipating maneuvers via learning temporal driving models, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 3182–3190. doi:10.1109/ICCV.2015.364.
- [34] A. Jain, A. Singh, H. S. Koppula, S. Soh, A. Saxena, Recurrent neural networks for driver activity anticipation via sensory-fusion architecture, in: 2016 IEEE International Conference on Robotics and Automation (ICRA), IEEE, 2016, pp. 3118–3125. doi:10.1109/ICRA.2016.7487478.
- [35] V. Bloom, V. Argyriou, D. Makris, Linear latent low dimensional space for online early action recognition and prediction, Pattern Recognition 72 (2017) 532–547. doi:10.1016/j.patcog.2017.07.003.
- [36] D. Wang, Y. Yuan, Q. Wang, Early action prediction with generative adversarial networks, IEEE Access 7 (2019) 35795–35804. doi:10.1109/ACCESS.2019.2904857.
- [37] X. Wang, J.-F. Hu, J.-H. Lai, J. Zhang, W.-S. Zheng, Progressive teacher-student learning for early action prediction, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 3556–3565.
- [38] Y. Ji, Y. Yang, X. Xu, H. T. Shen, One-shot learning based pattern transition map for action early recognition, Signal Processing 143

- (2018) 364–370. doi:[10.1016/j.sigpro.2017.06.001](https://doi.org/10.1016/j.sigpro.2017.06.001).
- [39] P. Felsen, P. Agrawal, J. Malik, What will happen next? forecasting player moves in sports videos, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 3342–3351. doi:[10.1109/iccv.2017.362](https://doi.org/10.1109/iccv.2017.362).
- [40] L. Neumann, A. Zisserman, A. Vedaldi, Future event prediction: If and when, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0–0.
- [41] M. Kassner, W. Patera, A. Bulling, Pupil: an open source platform for pervasive eye tracking and mobile gaze-based interaction, in: Proceedings of the 2014 ACM international joint conference on pervasive and ubiquitous computing: Adjunct publication, ACM, 2014, pp. 1151–1160.
- [42] D. M. Blei, A. Kucukelbir, J. D. McAuliffe, Variational inference: A review for statisticians, Journal of the American Statistical Association 112 (518) (2017) 859–877. doi:[10.1080/01621459.2017.1285773](https://doi.org/10.1080/01621459.2017.1285773).
- [43] A. Graves, Practical variational inference for neural networks, in: Advances in neural information processing systems, 2011, pp. 2348–2356.
- [44] D. P. Kingma, M. Welling, Auto-encoding variational bayes, arXiv preprint arXiv:1312.6114.
- [45] C. Blundell, J. Cornebise, K. Kavukcuoglu, D. Wierstra, Weight uncertainty in neural network, in: International Conference on Machine Learning, 2015, pp. 1613–1622.
- [46] Y. Gal, Z. Ghahramani, Dropout as a bayesian approximation, in: 33rd International Conference on Machine Learning, ICML 2016, Vol. 3, 2016, pp. 1661–1680.
- [47] D. P. Kingma, T. Salimans, M. Welling, Variational dropout and the local reparameterization trick, in: Advances in Neural Information Processing Systems, 2015, pp. 2575–2583.
- [48] A. Kendall, Y. Gal, What uncertainties do we need in bayesian deep learning for computer vision?, in: Advances in neural information processing systems, 2017, pp. 5574–5584.
- [49] D. Hafner, D. Tran, A. Irpan, T. Lillicrap, J. Davidson, Reliable uncertainty estimates in deep neural networks using noise contrastive priors, arXiv preprint arXiv:1807.09289.
- [50] Y. Gal, Uncertainty in deep learning, Ph.D. thesis, PhD thesis, University of Cambridge (2016).
- [51] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, Y. A. Sheikh, Openpose: Realtime multi-person 2d pose estimation using part affinity fields, IEEE Transactions on Pattern Analysis and Machine Intelligence (2019) 1–1 doi:[10.1109/TPAMI.2019.2929257](https://doi.org/10.1109/TPAMI.2019.2929257).
- [52] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural computation 9 (8) (1997) 1735–1780. doi:[10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- [53] Y. Gal, Z. Ghahramani, A theoretically grounded application of dropout in recurrent neural networks, in: Advances in neural information processing systems, 2016, pp. 1019–1027.