

Predicting decision-making in the future: Human versus Machine

Hoe Sung Ryu¹[0000–0002–9515–4402]
Uijong Ju²[0000–0002–9391–3938]
Christian Wallraven^{3,*}[0000–0002–2604–9115]

Department of Artificial Intelligence, Korea University, Seoul, Korea¹
hoesungryu@korea.ac.kr

Department of Information Display, Kyung Hee University, Seoul, Korea²
juuijong@khu.ac.kr

Department of Artificial Intelligence & Department of Brain and Cognitive
Engineering, Korea University, Seoul, Korea³
wallraven@korea.ac.kr

Abstract. Deep neural networks (DNNs) have become remarkably successful in data prediction, and have even been used to predict future actions based on limited input. This raises the question: do these systems actually “understand” the event similar to humans? Here, we address this issue using videos taken from an accident situation in a driving simulation. In this situation, drivers had to choose between crashing into a suddenly-appeared obstacle or steering their car off a previously indicated cliff. We compared how well humans and a DNN predicted this decision as a function of time before the event. The DNN outperformed humans for early time-points, but had an equal performance for later time-points. Interestingly, spatio-temporal image manipulations and Grad-CAM visualizations uncovered some expected behavior, but also highlighted potential differences in temporal processing for the DNN.

Keywords: Deep Learning · Video Prediction · Humans versus Machines · Decision-Making · Video Analysis

1 Introduction

The ability to predict, anticipate and reason about future events is the essence of intelligence [17] and one of the main goals of decision-making systems [11]. In general, predicting human behavior in future situations is of course an extremely hard task in an unconstrained setting. Given additional information about the context and the type of behavior to be predicted, however, behavior forecasting becomes a more tractable problem.

Recently, deep neural networks (DNNs) have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection, and many other domains [15]. Although DNNs yield excellent performance in such applications, it is often difficult to get insights into how and why a certain classification result has been made. To address this issue, Explainable Artificial Intelligence (XAI) proposes to make a shift towards more transparent AI [1]. As part of this shift, research has focused on comparing DNN performance in a task with human performance to better understand the underlying decision-making capacities of DNNs [31]. Here it becomes important to look closely at

the metric with which performance is measured - for example, if humans and DNNs have the same, high accuracy in identifying COVID-19 chest radiographs, this does not mean that DNNs use the same image-related cues to solve the task [9]. Therefore, [23] proposed to use metrics *beyond accuracy* to understand more deeply how DNNs and humans differ. Methods from XAI, such as feature visualization, for example, can help to understand which visual input is important for the network in making a certain decision [9].

Examples of behavior prediction that have been tackled with DNNs recently include, for example, predicting a future action by observing only a few portions of an action in progress [27], anticipating the next word in textual data given context [21], or predicting human driving behavior before dangerous situations occur [22]. Hence, DNNs seem to be capable of analyzing the spatio-temporal contents of an event to predict an outcome - the important question, then, becomes, do these networks actually “understand” the event similar to humans, or do they use spurious, high-dimensional correlations to form their prediction [9]?

In the present work, we take an accident situation during driving as a challenging context for studying this question. To explain the situation, imagine you are driving and there is a fork ahead; a prior warning sign alerted you that one of the directions of the fork will lead to a cliff, which will fatally crash the car. The other direction seems safe to drive, until, suddenly an obstacle appears on this direction of the fork, blocking the safe path. How do you react to this sudden change of circumstance? Understanding and modeling human behavior and its underlying factors in such situations can teach us a lot about decision-making under pressure and with high-risk stakes and has many important application areas.

Since it is impossible to study such a situation in the real world, a recent study by [18] employed virtual reality (VR) to investigate exactly this event. In this experiment, participants were trained to navigate a driving course that contained multiple, warning-indicated forks. The aforementioned accident situation was inserted only during the final test-run to see how participants would react to the sudden appearance of the obstacle (turn left and crash into the obstacle, or turn right and crash the car fatally off the cliff). From this study, here, we take the in-car videos that lead up to that final decision and segment them into time periods for predicting the turn direction. Importantly, the resulting short video segments were analyzed by *both* human participants and a DNN to predict the final decision. Using this strategy, we can compare human and DNN performance in predicting decision-making, but also look closer into the decision features of the DNN, using the tools of explainable AI.

2 Related Work

2.1 Predicting the “future”

There has been growing interest in predicting the future through DNNs where machines have to react to human actions as early as possible such as autonomous

driving, human-robotic interaction. In general, anticipating actions before these begin is a challenging problem, as this requires extensive contextual knowledge.

To solve this problem, approaches have used predefined target classes and a few portions of an action from short video segments leading up to an event (e.g., [27] for predicting the future motion of a person). Similarly, in the field of natural language processing, [21] proposed a method of predicting what a person will say next from given contextual data.

Other methods for human trajectory prediction, include, for example, [34, 4] trying to predict pedestrian trajectories from 3D coordinates obtained from stereo cameras and LIDAR sensors using deep learning-based models. Although adding information from various sensors improved the prediction performance, obtaining such data is still a challenging problem in general, leading to image-only approaches such as [25].

2.2 Comparisons of humans and DNNs

In order to better peek into the black box of DNNs, comparing performance between humans and DNNs has become an important research topic. Focusing on the human visual system, research has discussed human-machine comparisons at a conceptual level [20, 16, 7]. Indeed, these works show that deep learning systems not only solve image classification, but also reproduce certain aspects of human perception and cognition. This even goes as far as being able to reproduce phenomena of human vision, such as illusions [14] or crowding [33].

In addition, several studies have been conducted trying to impart more higher-level decision-making skills into DNNs: in [3], for example, ResNet models were able to solve abstract visual reasoning problems (IQ test) significantly better than humans - it is highly unlikely, however, that the model's knowledge representation matches that of humans, as the resulting networks were not able to learn visual relationships efficiently and robustly [12] - see also the many examples of adversarial attacks [2].

In other studies, input stimuli are manipulated or degraded to determine important visual features for human or machine decision-making. In [13], for example, twelve different types of image degradation were tested to compare error-patterns between humans and DNNs in classification. The authors found that the human visual system seemed to be more robust for image manipulations, but it DNNs trained directly on the degraded images were able to outperform humans. Similarly, based on adversarial perturbations, [7, 26] indicate that DNNs may have the potential ability to provide predictions and explanations of cognitive phenomena.

Overall, the field therefore seems to have quite heterogenous results when trying to compare human and machine performance with some studies finding commonalities and others finding critical differences. Here, we try to add to this general topic by comparing human and DNN performance in predicting a decision in a critical accident situation and analyzing the critical DNN features for prediction using similar tools.

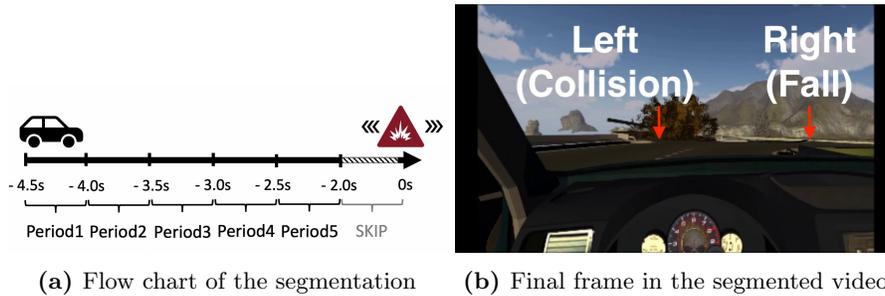


Fig. 1. (a) Segmentation setup and (b) screenshot from experiment video

3 Method

3.1 Experiment setup

Data: To compare the prediction of decision-making between humans and DNNs, we used the original, in-car videos from [18] - both humans and DNNs received the same data and the same task for a fair comparison. First, we trimmed the total of $n_{\text{total}} = 74$ videos from all participants to 2 to 4.5 seconds before the final decision - hence, the actual decision in which the car turned left or right was not part of the video. Let D denote our set of trimmed videos. Second, we partitioned D into five non-overlapping subsets $D_p \subset \{D_1, D_2, D_3, D_4, D_5\}$ with each D_p being a 0.5s-long video, containing 16 frames at $480\text{px} \times 720\text{px}$ (width \times height). For each video, this yields a total of 5 time segments such that, for example, D_5 indicates the final segment running from -2.5s to -2.0s , where 0s would indicate the actual time of the decision. In addition, each of the D_p is labelled either as fall ($n_{\text{fall}} = 23$) and collision ($n_{\text{collide}} = 51$). Given the human decision proportions in this accident situation (most people chose to collide with the obstacle, rather than to crash their car off the cliff), this results in an imbalanced, binary dataset with label ratios of $\approx 31\%$ versus $\approx 69\%$, respectively.

3.2 Behavioral experiment

Participants: A total of 29 participants (18 females, mean age 24.69 ± 3.11 (SD)) were recruited from the student population of Korea University. All participants had normal or corrected-to-normal vision and possessed a driver’s license. The experiment adhered to the tenets of the Declaration of Helsinki and was approved by the Institutional Review Board of Korea University with IRB number KUIRB-2018-0096-02s.

Experimental procedure: As for the experimental procedure, participants were tested in a small, enclosed room with no distractions. Upon entering the room, the procedure of the experiment was explained to them - in particular, that the video clips were part of a longer video sequence leading up to a final decision by a driver whether to collide with the trees or to fall down the cliff. Participants

Layer Name	Filter Shape	Repeats
conv1	$7 \times 7 \times 7$, 64, stride 1	1
conv2_x	$3 \times 3 \times 3$ max pool stride 2 $\left[\begin{array}{l} 3 \times 3 \times 3, 64 \\ 3 \times 3 \times 3, 64 \end{array} \right]$	2
conv3_x	$\left[\begin{array}{l} 3 \times 3 \times 3, 64 \\ 3 \times 3 \times 3, 64 \end{array} \right]$	2
conv4_x	$\left[\begin{array}{l} 3 \times 3 \times 3, 64 \\ 3 \times 3 \times 3, 64 \end{array} \right]$	2
conv5_x	$\left[\begin{array}{l} 3 \times 3 \times 3, 64 \\ 3 \times 3 \times 3, 64 \end{array} \right]$	2
Average Pooling, 512- d FC, Sigmoid		

Table 1. ResNet(2+1)D architecture in our experiments. Convolutional residual blocks are shown illustrated in brackets, next to the number of times each block is repeated in the stack.

were not informed that the decision ratio of the dataset was imbalanced. The whole experiment was conducted on a laptop running at $480\text{px} \times 720\text{px}$ resolution - participants sat a distance of ≈ 60 centimeters, with video segments subtending a visual angle of $\approx 30^\circ$. The behavioral experiment was created in PsychoPy (v3.0) [24].

Each experiment had a sequence of $370 = 5 \times 74$ trials, in which participants were asked to determine whether the car was going to collide or fall. In each trial, a short video segment D_p was randomly chosen and repeatedly shown to the participant until they felt they knew the answer, at which time they were to press the space bar. There was no time limit set by the experimenter, nor were participants explicitly instructed to respond as quickly or as accurately as possible. After the space bar was pressed, the time from start of the segment to the key press was recorded as response time, the video segment stopped looping and disappeared, and a text appeared in the center of the screen: “Which direction will the car go: collide or fall?”. Participants were to press the right arrow button for a collision and the left arrow button for a fall decision. All video segments were pseudo-randomly chosen and shown only once. Dependent variables were response time and response.

3.3 Computational experiment

ResNet(2+1)D architecture: A popular architecture in action recognition consists of a 3D convolutional neural network (CNN), which extends the typical 2D filters of image-based CNNs to 3D convolutional filters. This approach therefore directly extracts spatio-temporal features from videos by creating hierarchical representations and was shown to perform well on large-scale video datasets (e.g., [6]). In our experiments, we use a similar ResNet(2+1)D architecture [32] in which the 3D filters are factorized into 2D spatial convolutions and a 1D temporal convolution, which improves optimization.

In ResNet(2+1)D, the input tensor D_p ($p \in \{1, 2, 3, 4, 5\}$) is $5D$ and has size $B \times F \times C \times W \times H$, where B is the number of mini-batch, F is the number of frames of each video, C is RGB channel, and W and H are the width and height of the frame, respectively.

Our Network takes 16 clips consisting of RGB frames with the size of 112px \times 112px as an input. Each input frame is channel-normalized with mean (0.43216, 0.394666, 0.37645) and SD (0.22803, 0.22145, 0.216989). Down-sampling of the inputs is performed by conv1, conv3_1, conv4_1, and conv5_1. Conv1 is implemented by convolutional striding of $1 \times 2 \times 2$, and the remaining convolutions are implemented by striding of $2 \times 2 \times 2$. Since our overall sample size is small, a fine-tuning strategy with a pre-trained network was used. Its earlier layers remain fixed during training on our data with only the later layers of the network being optimized for prediction. In our experiments, the ResNet(2+1)D model is pre-trained on the Kinetics dataset [32] and fine-tuned on layers conv_5 upwards.

To illustrate the importance of the temporal context, we further experiment with three additional settings: models trained on 8 frame clips (sampled uniformly from each segment) or 2 frame clips (using only the first or the last two consecutive frames in the segment).

During training, batch normalization is applied to all convolutional layers. We deploy the Adamax optimizer with a mini-batch size of 16. Learning rates are updated by using a one-cycle scheduler [30] that starts with a small learning rate of $1e-4$, which is increased after each mini-batch until the maximum learning rate of $8e-3$. All processing was done on an Intel Xeon (Gold 5120 @2.20GHz) CPU and two NVIDIA V100 GPUs using Pytorch version 1.4.0.

Performance measures: For the experiment, we repeated a 5-fold cross-validation process 20 times, which created 100 folds in total. In every batch per epoch, we balanced the label of the training data set to 1 : 1 using sampling with replacement - this was not applied to the test set.

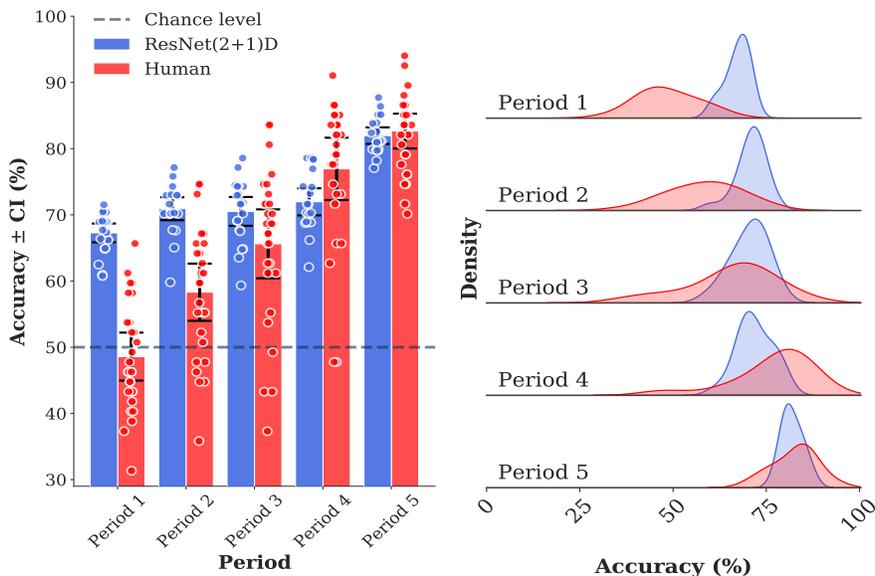
4 Experimental Results

4.1 Comparison of humans and DNN

We compared the ResNet(2+1)D prediction results to human performance on the exact same video sequences. As Figure 2 shows, the model achieves higher recognition rates compared to human participants until Period 3, whereas results seemed more similar to humans in periods 4 and 5.

A two-way analysis of variance (ANOVA) with factors of group (human or DNN) and time periods (5 periods) showed significant differences among groups ($F(1, 48) = 13.8882$, $p < .001$), time periods ($F(4, 192) = 107.2769882$, $p < .001$) and the interaction between group and time periods ($F(4, 192) = 21.6405$, $p < .001$) - see Table 2.

Since the interaction was significant, we next performed multiple pairwise comparisons on all possible combinations. To correct for multiple comparisons,



(a) Prediction accuracy for humans vs DNN (b) Kernel density estimate plot

Fig. 2. Performance comparison of humans and ResNet(2+1)D. The margin of error was calculated at a confidence level of 95%. Gaussian kernels were used to estimate kernel density and bandwidth calculated by Scott’s rule from [28].

Factor	Num DF	Den DF	F-statistic	p-value	η^2
Group	1	47	13.1187	< .001	0.2182
Time	4	188	164.2998	< .001	0.7775
Interaction	4	188	35.9250	< .001	0.4332

Table 2. Two-way ANOVA of performance comparing humans and ResNet

we applied a Bonferroni correction, correcting our alpha-level to 0.005 (= 0.05/10). Results overall showed no significant differences in period 5 ($p = 0.5853$) between human and ResNet(2+1)D, but significant differences from period 1 to 4 - see Table 3.

4.2 Discriminative features - Grad-CAM attention map

In the previous analysis, we compared only the performance of humans versus machines. Here, we employ the popular Grad-CAM method [29] from explainability research to detect the most important spatial (and temporal) information for the prediction. The output of Grad-CAM is a heatmap visualization for a given class label and provides a localization map that highlights important regions in the image.

Contrast	Time	A	B	T	dof	p-adjust	cohen
Group	-	Human	ResNet	-4.2540	34.1198	< .001	-1.0528
Time	-	Period1	Period2	-7.7186	48.0000	< .001	-0.6803
Time	-	Period1	Period3	-7.2275	48.0000	< .001	-1.0843
Time	-	Period1	Period4	-8.6652	48.0000	< .001	-1.8413
Time	-	Period1	Period5	-14.7811	48.0000	< .001	-3.0093
Time	-	Period2	Period3	-3.6096	48.0000	< .001	-0.4241
Time	-	Period2	Period4	-6.8742	48.0000	< .001	-1.2161
Time	-	Period2	Period5	-13.3034	48.0000	< .001	-2.4194
Time	-	Period3	Period4	-5.9944	48.0000	< .001	-0.7870
Time	-	Period3	Period5	-11.3766	48.0000	< .001	-1.9326
Time	-	Period4	Period5	-7.0810	48.0000	< .001	-1.0476
Time * Group	Period1	Human	ResNet	-11.2901	38.9798	< .001	-2.8579
Time * Group	Period2	Human	ResNet	-6.3592	39.4529	< .001	-1.6135
Time * Group	Period3	Human	ResNet	-2.0353	40.1793	0.0484	-0.5183
Time * Group	Period4	Human	ResNet	2.2670	40.9008	0.0287	0.5795
Time * Group	Period5	Human	ResNet	0.5498	42.6486	0.5853	0.1119

Table 3. Results of post-hoc multiple comparisons

As shown in Figure 3 (a), (b) for an example segment of Period 5, both attention maps focus on the central, steering wheel part. Beyond this, however, there are crucial differences in the resulting attention map depending on the condition. In the collision condition (Figure 3(a)), the network fixates areas nearby the tree in every frame of the period. In contrast, the model concentrates on both tree and cliff in the fall condition (Figure 3(b)), but as time progresses, the concentration on the cliff becomes more prominent.

The visualizations for the earliest Period 1 showed that the network focuses on the steering wheel and hill ridges in the collision condition and on similar areas on both trees and cliffs in fall conditions (cf. Figure 3 (c), (d)).

These qualitative differences were highly consistent in all input videos, indicating that the DNN is paying attention to meaningful - and expected - visual features for making the prediction.

4.3 Effect of spatial degradation

We next test the generalizability and robustness of the model by degrading the visual input. We first add spatial Gaussian noise to an image during testing either in the top 60% or the bottom 40% of the image.

Performance in the two conditions for all periods is shown in Figure 5(a). Overall, bottom blurring has virtually no effect on performance, whereas top blurring significantly reduces performance in the initial period, and especially in

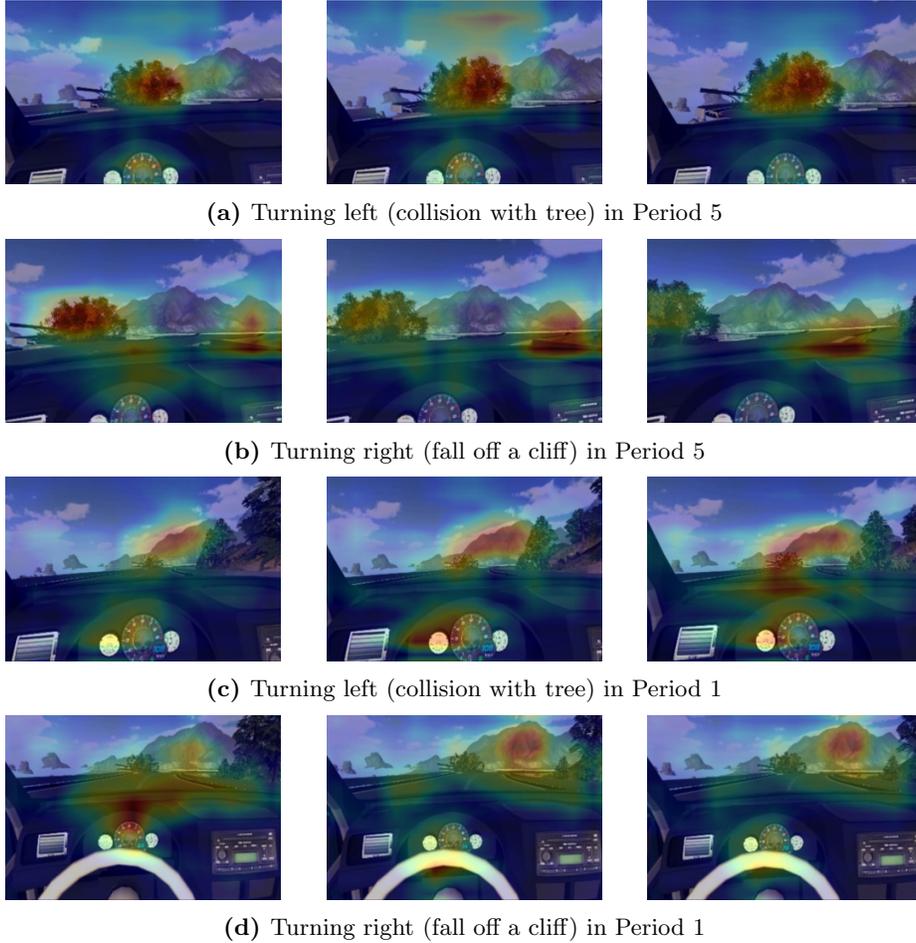


Fig. 3. Attention maps from Grad-CAM at equally sampled time points (left =start; right=end) for Periods 1 and 5. In Period 5, which is the closest period to the final decision, the DNN focuses on tree or cliff. In the earlier Period 1, the DNN puts more focus on the steering wheel or the ridge of the hill.

the final period (performance in Period 5: 72.76% for top-blurring versus 81.85% for bottom blurring and 81.97% for unblurred images). Hence, as one would expect from the attention-map analysis in Figure 3, top-blurring considerably reduces the networks ability for predicting decision-making.

Figure 4 compares the attention maps of the non-blurred with the two blurred conditions to confirm the accuracy results. Indeed, when blurring the top part, almost all activation focuses on the bottom, non-blurred input, virtually disregarding the important features outside the car. Perhaps the remaining focus on the steering wheel may lead to the above-chance prediction performance that could still be observed. As one would expect, blurring the bottom part of the

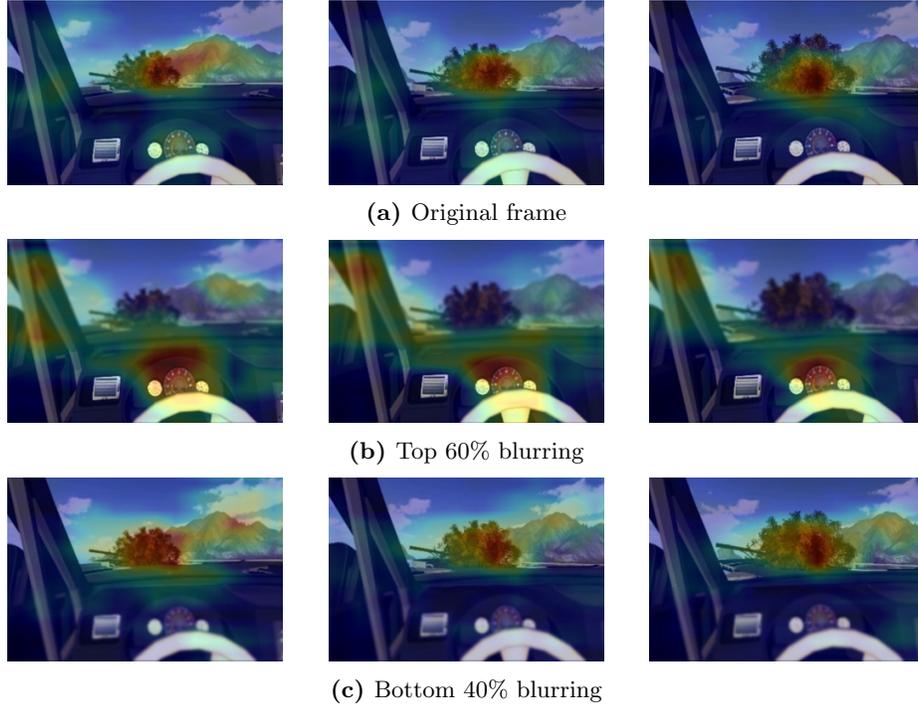


Fig. 4. Effect of spatially-selective blurring on discriminative features. From left to right: start frame to end frame of Period 5. The top row shows discriminative features for the original sequence. Blurring effects are shown applied to the top 60% of the original frame (middle row), or the bottom 40% (bottom row).

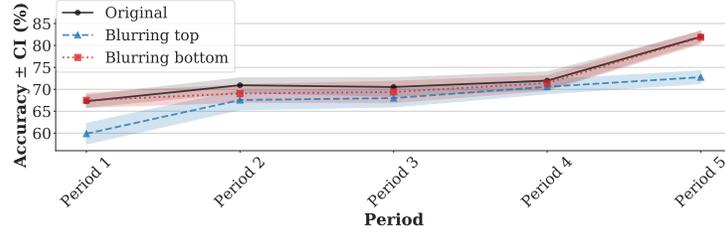
image has virtually no effect on the attention map, barring a slight decrease of focus on the steering wheel (see Figure 4(c)).

4.4 Temporal analysis

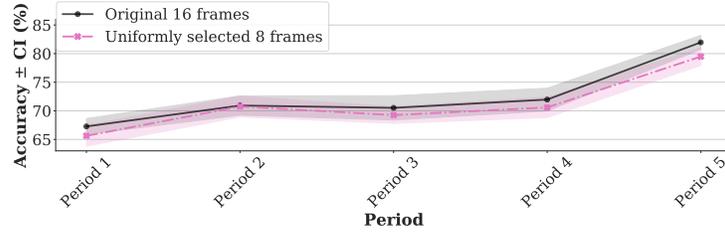
In general, continuous sequences of frames in the video are known to be critical for understanding events (see, for example, [8] for a detailed study on the dynamics of facial expressions). To determine the importance of different temporal aspects, we next present experiments that modified the number of input frames or changed the temporal order of frames.

Changing number of frames: Reducing the number of input frames from 16 frames to 8 frames showed only minor decreases in performance (Figure 5 (b)). A further reduction from the original 16 frames to only 2 frames (Figure 5 (c)) showed varying results: original performance levels could only be obtained with the last 2 frames of the segment, whereas using the first 2 frames of the Period resulted in an overall drop in performance, especially towards the final

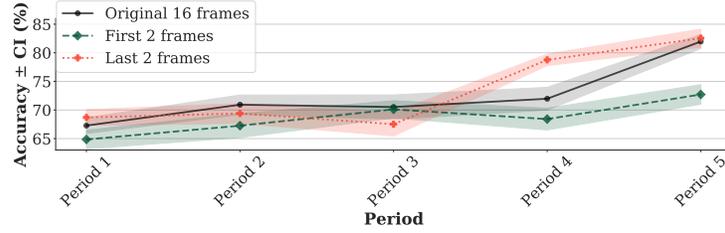
Period. Interestingly, for Period 4, prediction accuracy for the final 2 frames outperformed those obtained by all 16 frames, reaching almost peak accuracy.



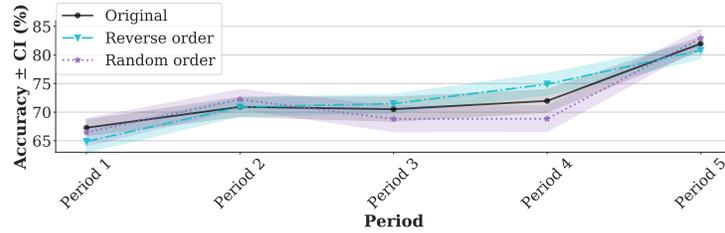
(a) Spatially-selective blurring



(b) 16 frames versus subsampled 8 frames



(c) First 2 frames versus the last 2 frames



(d) Random order versus reversed order

Fig. 5. Prediction accuracy for manipulations in space (blurring (a)) and time (16 vs 8 frames (b), first vs last frames (c), shuffling and time-reversal (d)).

Overall, these results seem to suggest that prediction seems to rely on the final, few frames of each period with limited advantages of adding further frames for temporal context.

Changing temporal order: Given prior results from human studies on the importance of the direction of time and the preservation of temporal structure in general [8], we next reversed time or shuffled the frames. As Figure 5 (d) shows, this has very little effect on performance, indicating that temporal structure itself bears little importance for the model.

5 Conclusion

In this paper, we investigated the difference between humans and a DNN in predicting the final decision in an accident situation. Our results showed that both humans and DNNs increased in accuracy as the time period approached the actual decision - a result in line with expectations, as both the path of the car and the steering movements may “settle” on their final direction at that point. We also found that the DNN made more accurate judgments compared to the human at early time points. Grad-CAM analysis further showed that its attended features are meaningful in terms of semantic content: interestingly, for the early time period, in which the DNN outperforms humans, its focus is on a wider view of the scene (including the ridge of hill) as well as the steering wheel, whereas in later time periods it focuses on closer things (tree or cliff). In future work, we will compare these computational attention maps to those obtained with further human experiments, using, for example, eye-tracking.

Moreover, several analyses were conducted to dive deeper into the underlying spatial and temporal features that may give rise to the prediction accuracy. In terms of spatial features, blurring showed that the outside, top view drove most of the recognition performance, which matches with expectations. In terms of temporal features, we found that the number of input frames does affect performance to some degree, but also showed that even two frames - when properly selected - still yield high performance. Again, it remains to be seen whether human performance would be similar with the same kind of input reduction.

Although most of the analyses so far would match with qualitative expectations about the important features for predictions, our analysis of time reversal and shuffling indicates no adverse effects of these manipulations. This is perhaps somewhat surprising given ample evidence in the human literature that time structure is crucially important in event analysis [8, 19]. Here, human experiments would most likely yield quite reduced performance, indicating that especially the learned *temporal* representation of DNNs may be different to those of humans. Further experiments will be necessary with different pre-training schemes and “deeper” visual hierarchies (vision transformers [10]) to investigate in more detail to what degree human and DNNs representations are similar.

Finally, we want to note that comparing data from human and computational experiments is bound to be complex: the DNN was specifically trained on a decision that was tested later, hence, there was no notion of task generalizability; similarly, humans will attach “meaning” to the outcome of the decision, with the left/right choice carrying consequences that are implicitly understood - such semantic grounding is so far missing from the DNN representation. It will be interesting to see how the current trend of research towards “foundation mod-

els”[5] will create frameworks that are capable of producing more human-level, semantically-rich, and generalizable task solutions.

Acknowledgements: This work was supported by the National Research Foundation of Korea under Grant NRF-2017M3C7A1041824 and by two Institute of Information and Communications Technology Planning and Evaluation (IITP) grants funded by the Korean government (MSIT): Development of BCI based Brain and Cognitive Computing Technology for Recognizing User’s Intentions using Deep Learning (2017-0-00451), and Artificial Intelligence Graduate School Program (Korea University) (2019-0-00079).

References

1. Adadi, A., Berrada, M.: Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access* **6**, 52138–52160 (2018)
2. Akhtar, N., Mian, A.: Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access* **6**, 14410–14430 (2018)
3. Barrett, D., Hill, F., Santoro, A., Morcos, A., Lillicrap, T.: Measuring abstract reasoning in neural networks. In: *International Conference on Machine Learning*. pp. 511–520. PMLR (2018)
4. Bhattacharyya, A., Fritz, M., Schiele, B.: Long-term on-board prediction of people in traffic scenes under uncertainty. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4194–4202 (2018)
5. Bommasani, R., et al.: On the opportunities and risks of foundation models (2021)
6. Carreira, J., Zisserman, A.: Quo vadis, action recognition. A new model and the kinetics dataset. *CoRR*, abs/1705.07750 **2**(3), 1 (2017)
7. Cichy, R.M., Kaiser, D.: Deep neural networks as scientific models. *Trends in cognitive sciences* **23**(4), 305–317 (2019)
8. Cunningham, D.W., Wallraven, C.: Dynamic information for the recognition of conversational expressions. *Journal of Vision* **9**(13), 7–7 (2009)
9. DeGrave, A.J., Janizek, J.D., Lee, S.I.: Ai for radiographic covid-19 detection selects shortcuts over signal. *medRxiv* (2020)
10. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
11. Edwards, W.: The theory of decision making. *Psychological bulletin* **51**(4), 380 (1954)
12. Funke, C.M., Borowski, J., Stosio, K., Brendel, W., Wallis, T.S., Bethge, M.: Five points to check when comparing visual perception in humans and machines. *Journal of Vision* **21**(3), 16–16 (2021)
13. Geirhos, R., Temme, C.R.M., Rauber, J., Schütt, H.H., Bethge, M., Wichmann, F.A.: Generalisation in humans and deep neural networks. *arXiv preprint arXiv:1808.08750* (2018)
14. Gomez-Villa, A., Martin, A., Vazquez-Corral, J., Bertalmío, M.: Convolutional neural networks can be deceived by visual illusions. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 12309–12317 (2019)
15. Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y.: *Deep learning*, vol. 1. MIT press Cambridge (2016)
16. Han, Y., Roig, G., Geiger, G., Poggio, T.: Scale and translation-invariance for novel objects in human vision. *Scientific reports* **10**(1), 1–13 (2020)

17. Hawkins, J., Blakeslee, S.: *On intelligence*. Macmillan (2004)
18. Ju, U., Chuang, L.L., Wallraven, C.: Acoustic cues increase situational awareness in accident situations: A vr car-driving study. *IEEE Transactions on Intelligent Transportation Systems* (2020)
19. Liu, Y., Dolan, R.J., Kurth-Nelson, Z., Behrens, T.E.: Human replay spontaneously reorganizes experience. *Cell* **178**(3), 640–652 (2019)
20. Majaj, N.J., Pelli, D.G.: Deep learning—using machine learning to study biological vision. *Journal of vision* **18**(13), 2–2 (2018)
21. Mikolov, T., Karafiát, M., Burget, L., Černocký, J., Khudanpur, S.: Recurrent neural network based language model. In: *Eleventh annual conference of the international speech communication association* (2010)
22. Ontanón, S., Lee, Y.C., Snodgrass, S., Winston, F.K., Gonzalez, A.J.: Learning to predict driver behavior from observation. In: *2017 AAAI Spring Symposium Series* (2017)
23. Oprea, S., Martinez-Gonzalez, P., Garcia-Garcia, A., Castro-Vargas, J.A., Orts-Escolano, S., Garcia-Rodriguez, J., Argyros, A.: A review on deep learning techniques for video prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020)
24. Peirce, J., Gray, J.R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., Kastman, E., Lindeløv, J.K.: Psychopy2: Experiments in behavior made easy. *Behavior research methods* **51**(1), 195–203 (2019)
25. Poibrenski, A., Klusch, M., Vozniak, I., Müller, C.: Multimodal multi-pedestrian path prediction for autonomous cars. *ACM SIGAPP Applied Computing Review* **20**(4), 5–17 (2021)
26. Ritter, S., Barrett, D.G., Santoro, A., Botvinick, M.M.: Cognitive psychology for deep neural networks: A shape bias case study. In: *International conference on machine learning*. pp. 2940–2949. PMLR (2017)
27. Rodriguez, C., Fernando, B., Li, H.: Action anticipation by predicting future dynamic images. In: *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*. pp. 0–0 (2018)
28. Scott, D.W.: *Multivariate density estimation: theory, practice, and visualization*. John Wiley & Sons (2015)
29. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE international conference on computer vision*. pp. 618–626 (2017)
30. Smith, L.N.: Cyclical learning rates for training neural networks. In: *2017 IEEE winter conference on applications of computer vision (WACV)*. pp. 464–472. IEEE (2017)
31. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International Conference on Machine Learning*. pp. 6105–6114. PMLR (2019)
32. Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M.: A closer look at spatiotemporal convolutions for action recognition. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. pp. 6450–6459 (2018)
33. Volokitin, A., Roig, G., Poggio, T.: Do deep neural networks suffer from crowding? *arXiv preprint arXiv:1706.08616* (2017)
34. Zhang, Z., Gao, J., Mao, J., Liu, Y., Anguelov, D., Li, C.: Stinet: Spatio-temporal-interactive network for pedestrian detection and trajectory prediction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 11346–11355 (2020)