

Interpretable Learning for Self-Driving Cars by Visualizing Causal Attention

Jinkyu Kim and John Canny
EECS, UC Berkeley, Berkeley, CA 94709, USA
{jinkyu.kim, canny}@berkeley.edu

Abstract

Deep neural perception and control networks are likely to be a key component of self-driving vehicles. These models need to be explainable - they should provide easy-to-interpret rationales for their behavior - so that passengers, insurance companies, law enforcement, developers etc., can understand what triggered a particular behavior. Here we explore the use of visual explanations. These explanations take the form of real-time highlighted regions of an image that causally influence the network's output (steering control). Our approach is two-stage. In the first stage, we use a visual attention model to train a convolution network end-to-end from images to steering angle. The attention model highlights image regions that potentially influence the network's output. Some of these are true influences, but some are spurious. We then apply a causal filtering step to determine which input regions actually influence the output. This produces more succinct visual explanations and more accurately exposes the network's behavior. We demonstrate the effectiveness of our model on three datasets totaling 16 hours of driving. We first show that training with attention does not degrade the performance of the end-to-end network. Then we show that the network causally cues on a variety of features that are used by humans while driving.

1. Introduction

Self-driving vehicle control has made dramatic progress in the last several years, and many auto vendors have pledged large-scale commercialization in a 2-3 year time frame. These controllers use a variety of approaches but recent successes [3] suggests that neural networks will be widely used in self-driving vehicles. But neural networks are notoriously cryptic - both network architecture and hidden layer activations may have no obvious relation to the function being estimated by the network. An exception to the rule is visual attention networks [29, 23, 7]. These networks provide spatial attention maps - areas of the image that the network attends to - that can be displayed in a way that is easy for users to interpret. They provide their atten-

tion maps instantly on images that are input to the network, and in this case on the stream of images from automobile video. As we show from our examples later, visual attention maps lie over image areas that have intuitive influence on the vehicle's control signal.

However, attention maps are only part of the story. Attention is a mechanism for filtering out non-salient image content. But attention networks need to find all *potentially* salient image areas and pass them to the main recognition network (a CNN here) for a final verdict. For instance, the attention network will attend to trees and bushes in areas of an image where road signs commonly occur. Just as a human will use peripheral vision to determine that "there is something there", and then visually fixate on the item to determine what it actually is. That is, the attention model must not mask out regions that might be important for driving control, but must look at foliage or other stimulæ to determine that they are not street signs or other vehicles. We therefore post-process the attention network's output, clustering it into attention "blobs" and then mask (set the attention weights to zero) each blob to determine the effect on the end-to-end network output. Blobs that have a causal effect on network output are retained while those that do not are removed from the visual map presented to the user.

Figure 1 shows an overview of our model. Our approach can be divided into three steps: (1) Encoder: convolutional feature extraction, (2) Coarse-grained decoder by visual attention mechanism, and (3) Fine-grained decoder: causal visual saliency detection and refinement of attention map. Our contributions are as follows:

- We show that visual attention heat maps are suitable "explanations" for the behavior of a deep neural vehicle controller, and do not degrade control accuracy.
- We show that attention maps comprise "blobs" that can be segmented and filtered to produce simpler and more accurate maps of visual saliency.
- We demonstrate the effectiveness of using our model with three large real-world driving datasets that contain over 1,200,000 video frames (*approx.* 16 hours).
- We illustrate typical spurious attention sources in driving video and quantify the reduction in explanation complexity from causal filtering.

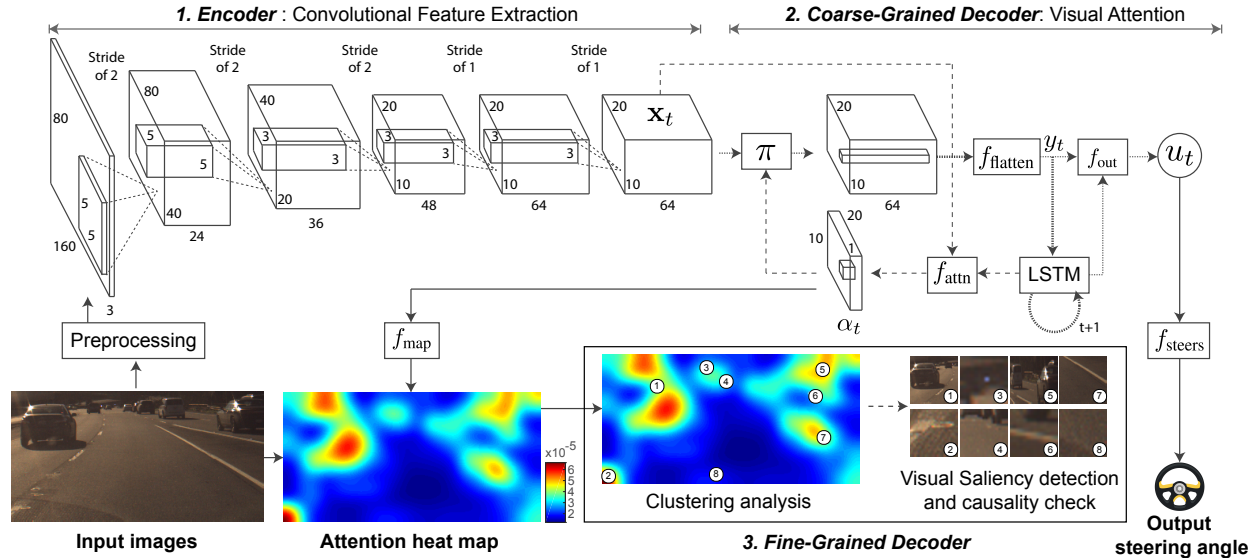


Figure 1. Our model predicts steering angle commands from an input raw image stream in an end-to-end manner. In addition, our model generates a heat map of attention, which can visualize where and what the model sees. To this end, we first encode images with a CNN and decode this feature into a heat map of attention, which is also used to control a vehicle. We test its causality by scrutinizing each cluster of attention blobs and produce a refined attention heat map of causal visual saliency.

2. Related Works

2.1. End-to-End Learning for Self-driving Cars

Self-driving vehicle controllers can be classified as: mediated perception approaches and end-to-end learning approaches. The mediated perception approach depends on recognizing human-designated features (*i.e.*, lane markings and cars) in a controller with if-then-else rules. Some examples include Urmsion *et al.* [27], Buehler *et al.* [4], and Levinson *et al.* [20].

Recently there is growing interest in end-to-end learning vehicle control. Most of these approaches learn a controller by supervised regression to recordings from human drivers. The training data comprise video from one or more vehicle cameras, and the control outputs (steering and possible acceleration and braking) from the driver. ALVINN (Autonomous Land Vehicle In a Neural Network) [21] was the first attempt to use neural network for directly mapping images to navigate the direction of the vehicle. More recently Bojarski *et al.* [3] demonstrated good performance with convolutional neural networks (CNNs) to directly map images from a front-view camera to steering controls. Xu *et al.* [28] proposed an end-to-end egomotion prediction approach that takes raw pixels and prior vehicle state signals as inputs and predicts a sequence of discretized actions (*i.e.*, straight, stop, left-turn, and right-turn). These models show good performance but their behavior is opaque and uninterpretable.

An intermediate approach was explored in Chen *et al.* [6]

who defined human-designated interpretable intermediate features such as the curvature of lane, distances to neighboring lanes, and distances from the front-located vehicles. A CNN is trained to produce these features, and a simple controller maps them to steering angle. They also generated deconvolution maps to show image areas that affected network output. However, there were several difficulties with that work: (i) use of the intermediate layer may cause degradation of control accuracy (ii) the intermediate feature descriptors may provide a limited and ad-hoc vocabulary for explanations and (iii) the authors noted the presence of spurious input features but there was no attempt to remove them. By contrast, our work shows that state-of-the-art driving models can be made interpretable without sacrificing accuracy, that attention models provide more robust image annotation, and causal analysis further improves explanation saliency.

2.2. Visual Explanation

In a landmark work, Zeiler and Fergus [31] used “deconvolution” to visualize layer activations of convolutional networks. LeCun *et al.* [18] provides textual explanations of images as automatically-generated captions. Building on this work, Bojarski *et al.* [2] developed a richer notion of “contribution” of a pixel to the output. However a difficulty with deconvolution-style approaches is the lack of formal measures of how the network output is affected by spatially-extended features (rather than pixels). Attention-based approaches like ours directly extract areas of the im-

age that *did not affect* network output (because they were masked out by the attention model), and causal filtering further removes spurious image areas. Hendricks *et al.* [13] trains a deep network to generate specific explanation without explicitly identifying semantic features. Also, Justin Johnson [16] proposes DenseCap which uses fully convolutional localization networks for dense captioning, their paper achieves both localizing objects and describing salient regions in images using natural language. In reinforcement learning, Zrihem *et al.* [30] proposes a visualization method to interpret the agents action by describing Markov Decision Process model as a directed graph on a t-SNE map.

3. Method

3.1. Preprocessing

Our model predicts continuous steering angle commands from input raw pixels in an end-to-end manner. As discussed by Bojarski *et al.* [3], our model predicts the inverse turning radius \hat{u}_t ($= r_t^{-1}$, where r_t is the turning radius) at every timestep t instead of steering angle commands, which depends on the vehicle’s steering geometry and also result in numerical instability when predicting near zero steering angle commands. The relationship between the inverse turning radius u_t and the steering angle command θ_t can be approximated by Ackermann steering geometry [22] as follows:

$$\theta_t = f_{\text{steers}}(u_t) = u_t d_w K_s (1 + K_{\text{slip}} v_t^2) \quad (1)$$

where θ_t in degrees and v_t (m/s) is a steering angle and a velocity at time t , respectively. K_s , K_{slip} , and d_w are vehicle-specific parameters. K_s is a steering ratio between the turn of the steering and the turn of the wheels. K_{slip} represents the relative motion between a wheel and the surface of road. d_w is the length between the front and rear wheels. Our model therefore needs two measurements for training: timestamped vehicle’s speed and steering angle commands.

To reduce computational cost, each raw input image is down-sampled and resized to $80 \times 160 \times 3$ with nearest-neighbor scaling algorithm. For images with different raw aspect ratios, we cropped the height to match the ratio before down-sampling. A common practice in image classification is to subtract the mean RGB value computed on the training set from each pixel [12, 24]. This is effective to achieve zero-centered inputs which are originally in different scales. Driving datasets, however, do not show that various scales. For instance, the camera gains are (automatically or in advance) calibrated to capture such high-quality images in a certain dynamic range. In our experiment, we could not obtain significant improvement by the use of mean subtraction. Instead, we change the range of pixel intensity values and convert to HSV colorspace, which is

commonly used for its robustness in problems where color description plays an integral role.

We utilize a single exponential smoothing method [15] to reduce the effect of human factors-related performance variation and the effect of measurement noise. Formally, given a smoothing factor $0 \leq \alpha_s \leq 1$, the simple exponential smoothing method is defined as follows:

$$\begin{pmatrix} \hat{\theta}_t \\ \hat{v}_t \end{pmatrix} = \alpha_s \begin{pmatrix} \theta_t \\ v_t \end{pmatrix} + (1 - \alpha_s) \begin{pmatrix} \hat{\theta}_{t-1} \\ \hat{v}_{t-1} \end{pmatrix} \quad (2)$$

where $\hat{\theta}_t$ and \hat{v}_t are the smoothed time-series of θ_t and v_t , respectively. Note that they are same as the original time-series when $\alpha_s = 1$, while values of α_s closer to zero have a greater smoothing effect and are less responsive to recent changes. The effect of applying smoothing methods is summarized in Section 4.4. Note that the use of Kalman filter could be better for fighting measurement noise but suffers from obtaining a long-term composite effect. Our use of the exponentially smoothing methods, therefore, can be justified by obtaining a long-term integral of controls.

3.2. Encoder: Convolutional Feature Extraction

We use a convolutional neural network to extract a set of encoded visual feature vector, which we refer to as a convolutional feature cube \mathbf{x}_t . Each feature vectors may contain high-level object descriptions that allow the attention model to selectively pay attention to certain parts of an input image by choosing a subset of feature vectors.

As depicted in Figure 1, we use a 5-layered convolution network that is utilized by Bojarski *et al.* [3] to learn a model for self-driving cars. As discussed by Lee *et al.* [19], we omit max-pooling layers to prevent spatial locational information loss as the strongest activation propagates through the model. We collect a three-dimensional convolutional feature cube \mathbf{x}_t from the last layer by pushing the preprocessed image through the model, and the output feature cube will be used as an input of the LSTM layers, which we will explain in Section 3.3. Using this convolutional feature cube from the last layer has advantages in generating high-level object descriptions, thus increasing interpretability and reducing computational burdens for a real-time system.

Formally, a convolutional feature cube of size $W \times H \times D$ is created at each timestep t from the last convolutional layer. We then collect \mathbf{x}_t , a set of $L = W \times H$ vectors, each of which is a D -dimensional feature slice for different spatial parts of the given input.

$$\mathbf{x}_t = \{x_{t,1}, x_{t,2}, \dots, x_{t,L}\} \quad (3)$$

where $x_{t,i} \in \mathcal{R}^D$ for $i \in \{1, 2, \dots, L\}$. This allows us to focus selectively on different spatial parts of the given image by choosing a subset of these L feature vectors.

3.3. Coarse-Grained Decoder: Visual Attention

The goal of soft deterministic attention mechanism $\pi(\{x_{t,1}, x_{t,2}, \dots, x_{t,L}\})$ is to search for a good context vector y_t , which is defined as a combination of convolutional feature vectors $x_{t,i}$, while producing better prediction accuracy. We utilize a deterministic soft attention mechanism that is trainable by standard back-propagation methods, which thus has advantages over a hard stochastic attention mechanism that requires reinforcement learning. Our model feeds α weighted context y_t to the system as discussed by several works [23, 29]:

$$\begin{aligned} y_t &= f_{\text{flatten}}(\pi(\{\alpha_{t,i}\}, \{x_{t,i}\})) \\ &= f_{\text{flatten}}(\{\alpha_{t,i}x_{t,i}\}) \end{aligned} \quad (4)$$

where $i = \{1, 2, \dots, L\}$. $\alpha_{t,i}$ is a scalar attention weight value associated with a certain grid of input image in such that $\sum_i \alpha_{t,i} = 1$. These attention weights can be interpreted as the probability over L convolutional feature vectors that the location i is the important part to produce better estimation accuracy. f_{flatten} is a flattening function, which converts the input feature matrix into a 1-D feature vector to be used by the dense layer for LSTM. y_t is thus $D \times L$ -dimensional vector that contains convolutional feature vectors weighted by attention weights. In traditional models [23, 29], $\pi(\{\alpha_{t,i}\}, \{x_{t,i}\})$ maps $L \times D$ convolutional feature cubes to the D dimension only by using the α weighted average context, *i.e.*, $\sum_{i=1}^L \alpha_{t,i}x_{t,i}$, which is but prone to remove spatial information. In our approach, $\pi(\{\alpha_{t,i}\}, \{x_{t,i}\})$ is identity function to preserve spatial information.

As we summarize in Figure 1, we use a long short-term memory (LSTM) network [14] that predicts the inverse turning radius \hat{u}_t and generates attention weights $\{\alpha_{t,i}\}$ at each timestep t conditioned on the previous hidden state h_t and a current convolutional feature cube \mathbf{x}_t . More formally, let us assume a hidden layer $f_{\text{attn}}(x_{t,i}, h_{t-1})$ conditioned on the previous hidden state h_{t-1} and the current feature vectors $\{x_{t,i}\}$. The attention weight $\{\alpha_{t,i}\}$ for each spatial location i is then computed by multinomial logistic regression (*i.e.*, softmax regression) function as follows:

$$\alpha_{t,i} = \frac{\exp(f_{\text{attn}}(x_{t,i}, h_{t-1}))}{\sum_{j=1}^L \exp(f_{\text{attn}}(x_{t,j}, h_{t-1}))} \quad (5)$$

Our network also predicts inverse turning radius \hat{u}_t as an output with additional hidden layer $f_{\text{out}}(y_t, h_t)$ conditioned on the current hidden state h_t and α weighted context y_t .

To initialize memory state c_t and hidden state h_t of LSTM network, we follow Xu *et al.* [29] by averaging of the feature slices $x_{0,i}$ at initial time fed through two addi-

tional hidden layers: $f_{\text{init},c}$ and $f_{\text{init},h}$.

$$c_0 = f_{\text{init},c} \left(\frac{1}{L} \sum_{i=1}^L x_{0,i} \right), \quad h_0 = f_{\text{init},h} \left(\frac{1}{L} \sum_{i=1}^L x_{0,i} \right) \quad (6)$$

As discussed by Xu *et al.* [29], doubly stochastic regularization can encourage the attention model at different parts of the image. At time t , our attention model predicts a scalar $\beta = \text{sigm}(f_{\beta}(h_{t-1}))$ with an additional hidden layer f_{β} conditioned on the previous hidden state h_{t-1} such that

$$y_t = \text{sigm}(f_{\beta}(h_{t-1})) f_{\text{flatten}}(\{\alpha_{t,i}x_{t,i}\}) \quad (7)$$

We use the following penalized loss function \mathcal{L}_1 :

$$\mathcal{L}_1(u_t, \hat{u}_t) = \sum_{t=1}^T |u_t - \hat{u}_t| + \lambda \sum_{i=1}^L \left(1 - \sum_{t=1}^T \alpha_{t,i} \right) \quad (8)$$

where T is the length of time steps, and λ is a penalty coefficient that encourages the attention model to see different parts of the image at each time frame. Section 4.3 describes the effect of using regularization.

3.4. Fine-Grained Decoder: Causality Test

The last step of our pipeline is a fine-grained decoder, in which we refine a map of attention and detect local visual saliencies. Though an attention map from our coarse-grained decoder provides probability of importance over a 2D image space, our model needs to determine specific regions that cause a causal effect on prediction performance. To this end, we assess a decrease in performance when a local visual saliency on an input raw image is masked out.

We first collect a consecutive set of attention weights $\{\alpha_{t,i}\}$ and input raw images $\{\mathcal{I}_t\}$ for a user-specified T timesteps. We then create a map of attention, which we refer \mathcal{M}_t as defined: $\mathcal{M}_t = f_{\text{map}}(\{\alpha_{t,i}\})$. Our 5-layer convolutional neural network uses a stack of 5×5 and 3×3 filters without any pooling layer, and therefore the input image of size 80×160 is processed to produce the output feature cube of size $10 \times 20 \times 64$, while preserving its aspect ratio. Thus, we use $f_{\text{map}}(\{\alpha_{t,i}\})$ as up-sampling function by the factor of eight followed by Gaussian filtering [5] as discussed in [29] (see Figure 2 (A,B)).

To extract a local visual saliency, we first randomly sample 2D N particles with replacement over an input raw image conditioned on the attention map \mathcal{M}_t . Note that, we also use time-axis as the third dimension to consider temporal features of visual saliencies. We thus store spatio-temporal 3D particles $\mathcal{P} \leftarrow \mathcal{P} \cup \{\mathcal{P}_t, t\}$ (see Figure 2 (C)).

We then apply a clustering algorithm to find a local visual saliency by grouping 3D particles \mathcal{P} into clusters $\{\mathcal{C}\}$ (see Figure 2 (D)). In our experiment, we use DBSCAN [10], a density-based clustering algorithm that has

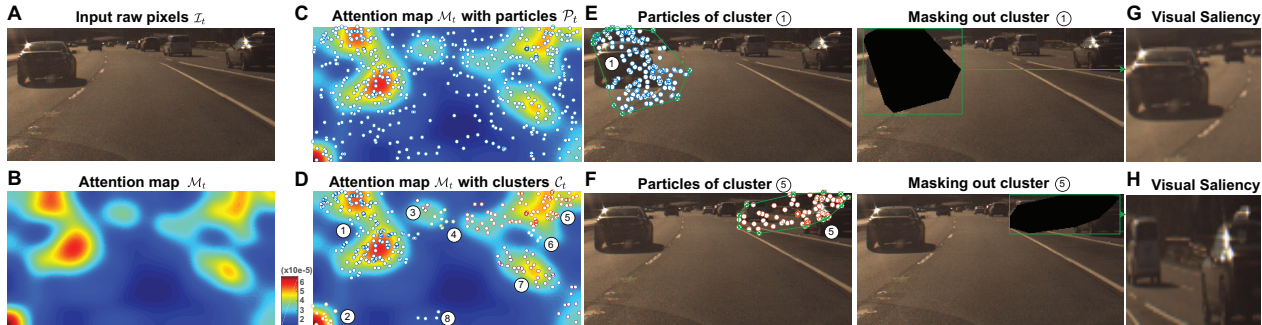


Figure 2. Overview of our fine-grained decoder. Given an input raw pixels \mathcal{I}_t (A), we compute an attention map \mathcal{M}_t with a function f_{map} (B). (C) We randomly sample 3D $N = 500$ particles over the attention map, and (D) we apply a density-based clustering algorithm (DBSCAN [10]) to find a local visual saliency by grouping particles into clusters. (E, F) For each cluster $c \in \mathcal{C}$, we compute a convex hull $\mathcal{H}(c)$ to define its region, and mask out the visual saliency to see causal effects on prediction accuracy (see E, F for clusters 1 and 5, respectively). (G, H) Warped visual saliencies for clusters 1 and 5, respectively.

advantages to deal with a noisy dataset because they group particles together that are closely packed, while marking particles as outliers that lie alone in low-density regions. For points of each cluster c and each time frame t , we compute a convex hull $\mathcal{H}(c)$ to find a local region of each visual saliency detected (see Figure 2 (E, F)).

For points of each cluster c and each time frame t , we iteratively measure a decrease of prediction performance with an input image which we mask out a local visual saliency. We compute a convex hull $\mathcal{H}(c)$ to find a local, and mask out each visual saliency by assigning zero values for all pixels lying inside each convex hull. Each causal visual saliency is generated by warping into a fixed spatial resolution ($=64 \times 64$) as shown in Figure 2 (G, H).

4. Result

4.1. Datasets

As explained in Table 1, we obtain two large-scale datasets that contain over 1,200,000 frames (≈ 16 hours) collected from Comma.ai [8], Udacity [26], and Hyundai Center of Excellence in Integrated Vehicle Safety Systems and Control (HCE) at Berkeley. These three datasets acquired contain video clips captured by a single front-view camera mounted behind the windshield of the vehicle. Alongside the video data, a set of time-stamped sensor measurement is contained, such as vehicle’s velocity, acceleration, steering angle, GPS location and gyroscope angles. Thus, these datasets are ideal for self-driving studies. Note that, for sensor logs unsynced with the time-stamps of video data, we use the estimates of the interpolated measurements. Videos are mostly captured during highway driving in clear weather on daytime, and there included driving on other road types, such as residential roads (with and without lane markings), and contains the whole driver’s activities such as staying in a lane and switching lanes. Note also that, we

exclude frames when the vehicle stops which happens when $\hat{v}_t < 1$ m/s.

4.2. Training and Evaluation Details

To obtain a convolutional feature cube \mathbf{x}_t , we train the 5-layer CNN explained in Section 3.2 by using additional 5-layer fully connected layers (*i.e.*, # hidden variables: 1164, 100, 50, and 10, respectively), of which output predicts the measured inverse turning radius u_t . Incidentally, instead of using addition fully-connected layers, we could also obtain a convolutional feature cube \mathbf{x}_t by training from scratch with the whole network. In our experiment, we obtain the $10 \times 20 \times 64$ -dimensional convolutional feature cube, which is then flattened to 200×64 and is fed through the coarse-grained decoder. Other recent types of more recent expressive networks may give a performance boost over our CNN configuration. However, exploration of other convolutional architectures would be out of our scope.

We experiment with various numbers of LSTM layers (1 to 5) of the soft deterministic visual attention model but did not observe any significant improvements in model performance. Unless otherwise stated, we use a single LSTM layer in this experiment. For training, we use Adam optimization algorithm [17] and also use dropout [25] of 0.5 at hidden state connections and Xavier initialization [11]. We randomly sample a mini-batch of size 128, each of batch contains a set Consecutive frames of length $T = 20$. Our model took less than 24 hours to train on a single NVIDIA Titan X Pascal GPU. Our implementation is based on Tensorflow [1] and code will be publicly available upon publication.

Two datasets (Comma.ai [8] and HCE) we used were available with images captured by a single front-view camera. This makes it hard to use the data augmentation technique proposed by Bojarski *et al.* [3], which generated images with artificial shifts and rotations by using two addi-

	Dataset		
	Comma.ai [8]	HCE	Udacity [26]
# frame	522,434	80,180	650,690
FPS	20Hz	20Hz/30Hz	20Hz
Hours	≈ 7 hrs	≈ 1 hr	≈ 8 hrs
Condition	Highway/Urban	Highway	Urban
Location	CA, USA	CA, USA	CA, USA
Lighting	Day/Night	Day	Day

Table 1. Dataset details. Over 16 hours (>1,200,000 video frames) of driving dataset that contains a front-view video frames and corresponding time-stamped measurements of vehicle dynamics. The data is collected from two public data sources, Comma.ai [8] and Udacity [26], and Hyundai Center of Excellence in Vehicle Dynamic Systems and Control (HCE).

tional off-center images (left-view and right-view) captured by the same vehicle. Data augmentation may give a performance boost, but we report performance without data augmentation.

4.3. Effect of Choosing Penalty Coefficient λ

Our model provides a better way to understand the rationale of the models decision by visualizing where and what the model sees to control a vehicle. Figure 3 shows a consecutive input raw images (with sampling period of 5 seconds) and their corresponding attention maps (*i.e.*, $\mathcal{M}_t = f_{\text{map}}(\{\alpha_{t,i}\})$). We also experiment with three different penalty coefficients $\lambda \in \{0, 10, 20\}$, where the model is encouraged to pay attention to wider parts of the image (see differences between the bottom 3 rows in Figure 3) as we have larger λ . For better visualization, an attention map is overlaid by an input raw image and color-coded; for example, red parts represent where the model pays attention. For quantitative analysis, prediction performance in terms of mean absolute error (MAE) is explained on the bottom of each figure. We observe that our model is indeed able to pay attention on road elements, such as lane markings, guardrails, and vehicles ahead, which is essential for human to drive.

4.4. Effect of Varying Smoothing Factors

Recall from Section 3.1 that the single exponential smoothing method [15] is used to reduce the effect of human factors variation and the effect of measurement noise for two sensor inputs: steering angle and velocity. A robust model for autonomous vehicles would yield consistent performance, even when some measurements are noisy. Figure 4 shows the prediction performance in terms of mean absolute error (MAE) on a comma.ai testing data set. Various smoothing factors $\alpha_s \in \{0.01, 0.05, 0.1, 0.3, 0.5, 1.0\}$ are used to assess the effect of using smoothing methods.

Dataset	Model	MAE in degree [SD]	
		Training	Testing
Comma.ai [8]	CNN+FCN [3]	.421 [0.82]	2.54 [3.19]
	CNN+LSTM	.488 [1.29]	2.58 [3.44]
	Attention ($\lambda=0$)	.497 [1.32]	2.52 [3.25]
	Attention ($\lambda=10$)	.464 [1.29]	2.56 [3.51]
	Attention ($\lambda=20$)	.463 [1.24]	2.44 [3.20]
	CNN+FCN [3]	.246 [.400]	1.27 [1.57]
HCE	CNN+LSTM	.568 [.977]	1.57 [2.27]
	Attention ($\lambda=0$)	.334 [.766]	1.18 [1.66]
	Attention ($\lambda=10$)	.358 [.728]	1.25 [1.79]
	Attention ($\lambda=20$)	.373 [.724]	1.20 [1.66]
	CNN+FCN [3]	.457 [.870]	4.12 [4.83]
	CNN+LSTM	.481 [1.24]	4.15 [4.93]
Udacity [26]	Attention ($\lambda=0$)	.491 [1.20]	4.15 [4.93]
	Attention ($\lambda=10$)	.489 [1.19]	4.17 [4.96]
	Attention ($\lambda=20$)	.489 [1.26]	4.19 [4.93]

Table 2. Control performance comparison in terms of mean absolute error (MAE) in degree and its standard deviation. Control accuracy is not degraded by incorporation of attention compared to an identical base CNN without attention. *Abbreviation:* SD (standard deviation)

With setting $\alpha_s=0.05$, our model for the task of steering estimation performs the best. Unless otherwise stated, we will use α_s as 0.05.

4.5. Quantitative Analysis

In Table 2, we compare the prediction performance with alternatives in terms of MAE. We implement alternatives that include the work by Bojarski *et al.* [3], which used an identical base CNN and a fully-connected network (FCN) without attention. To see the contribution of LSTM, we also test a CNN and LSTM, which is identical to ours but does not use the attention mechanism. For our model, we test with three different values of penalty coefficients $\lambda \in \{0, 10, 20\}$.

Our model shows competitive prediction performance than alternatives. Our model shows 1.18–4.15 in terms of MAE on testing dataset. This confirms that incorporation of attention does not degrade control accuracy. The average run-time for our model and alternatives took less than a day to train each dataset. The Udacity dataset contains more hard-to-predict drivers activities than other two datasets. For instance, this dataset is mostly collected while driving on residential roads with many turns at intersections. This is a challenge for a simple end-to-end controller but our method still shows reasonable performance.

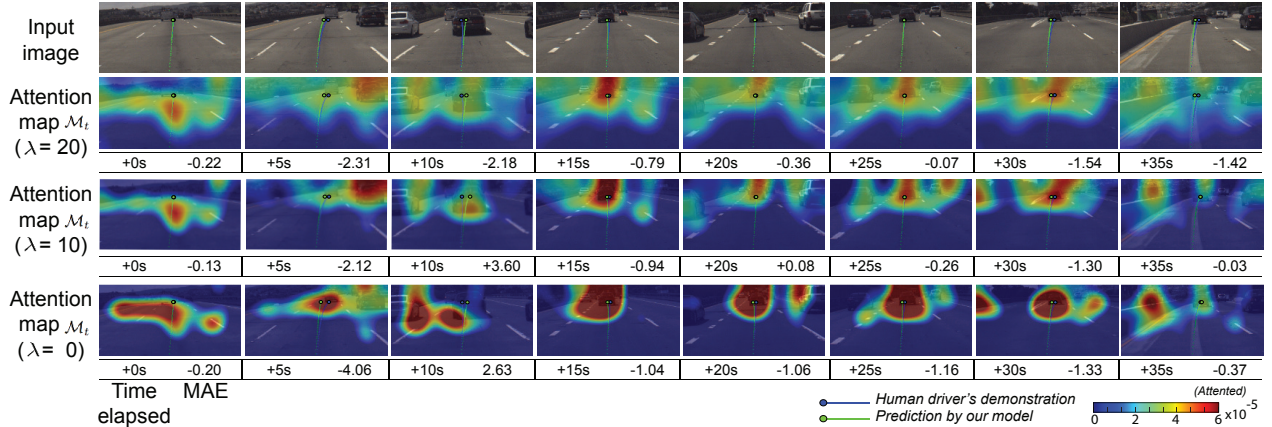


Figure 3. Attention maps over time. Unseen consecutive input image frames are sampled at every 5 seconds (see from left to right). (Top) Input raw images with human drivers demonstrated curvature of path (blue line) and predicted curvature of path (green line). (From the bottom) We illustrate attention maps with three different regularization penalty coefficients $\lambda \in \{0, 10, 20\}$. Each attention map is overlaid by an input raw image and color-coded. Red parts indicate where the model pays attention. *Data*: Comma.ai [8]

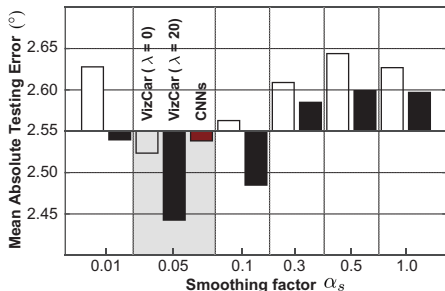


Figure 4. Effect of applying a single exponential smoothing method over various smoothing factors from 0.01 to 1.0. We use two different penalty coefficients $\lambda \in \{0, 20\}$. With setting $\alpha_s = 0.05$, our model performs the best. *Data*: Comma.ai [8]

4.6. Effect of Causal Visual Saliencies

Recall from Section 3.4, we post-process the attention networks output by clustering it into attention blobs and filtering if they have a causal effect on network output. Figure 5 (A) shows typical examples of an input raw image, an attention networks output with spurious attention sources, and our refined attention heat map. We observe our model can produce a simpler and more accurate map of visual saliency by filtering out spurious attention blobs. In our experiment, 62% and 58% out of all attention blobs are indeed spurious attention sources on Comma.ai [8] and HCE datasets (see Figure 5 (B)).

5. Discussion

The proposed method highlights regions that causally influence deep neural perception and control networks for self-driving cars. Thus, it would be worth exploring a po-

tential overlap between the causally salient image areas and what and where human drivers is really paying their attention while driving. Due to the lack of data capturing human drivers' gaze movement, we leave this comparison to a future work. Das *et al.* [9] recently reported that several attention-based visual question answering (VQA) models tend to look at different regions of the image unlike as humans do. However, this comparison is still an active research area for self-driving cars.

Along with devising the fine-grained decoder, we may consider using feature-level masking approach. Masking out convolutional features of attended region can provide a computationally efficient way by removing multiple forward passes of the convolutional network. This approach, however, may suffer from low deconvolutional spatial resolution, which makes challenge to view as a unit apart and divide the whole attention map into distinct attended objects, such as cars or lane markings.

6. Conclusion

We described an interpretable visualization for deep self-driving vehicle controllers. It uses a visual attention model augmented with an additional layer of causal filtering. We tested with three large-scale real driving datasets that contain over 16 hours of video frames. We showed that (i) incorporation of attention does not degrade control accuracy compared to an identical base CNN without attention (ii) raw attention highlights interpretable features in the image and (iii) causal filtering achieves a useful reduction in explanation complexity by removing features which do not significantly affect the output.

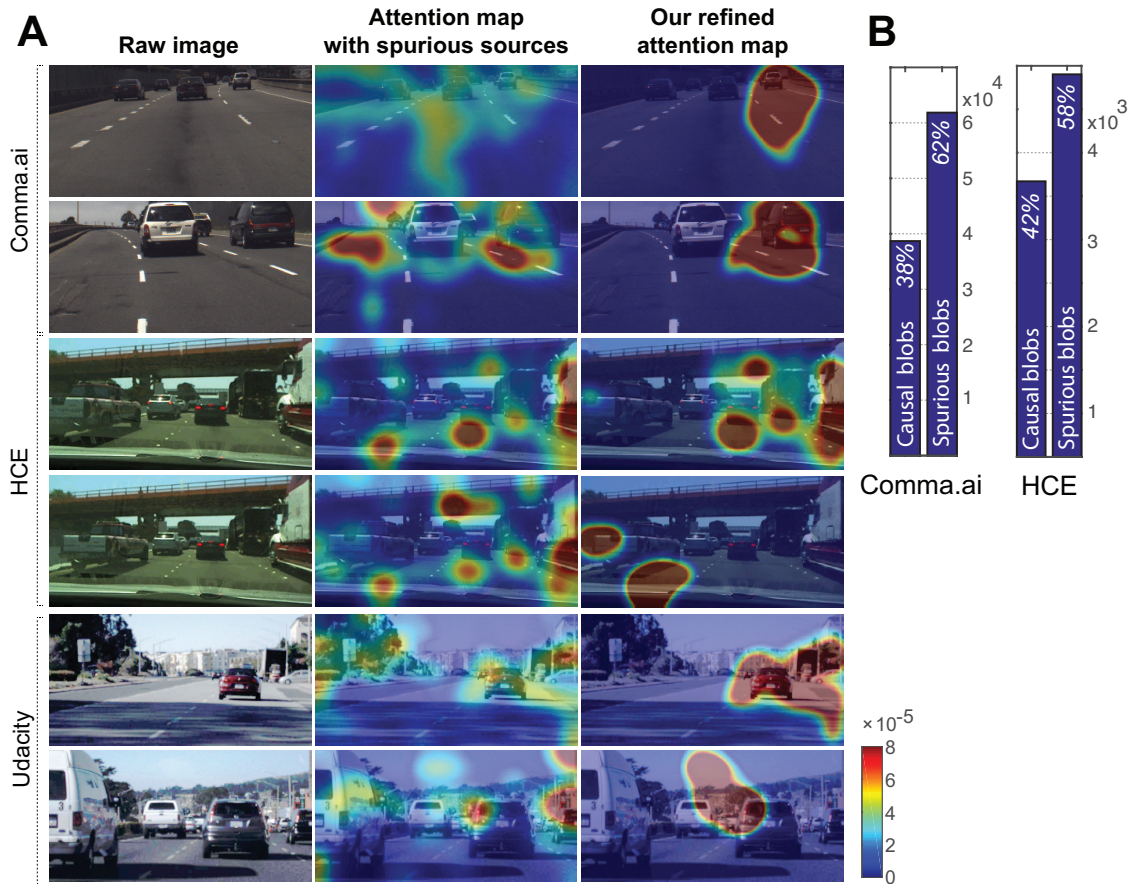


Figure 5. (A) We illustrate examples of (left) raw input images, their (middle) visual attention heat maps with spurious attention sources, and (right) our attention heat maps by filtering out spurious blobs to produce simpler and more accurate attention maps. (B) To measure how much the causal filtering is simplifying attention clusters, we quantify the number of attention blobs before and after causal filtering.

Acknowledgment

The authors would like to thank the anonymous reviewers of this paper and Daniel Seita for their helpful comments. This work was supported by Berkeley DeepDrive and Samsung Scholarship.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous systems, 2015. *Software available from tensorflow.org*, 1, 2015. 5
- [2] M. Bojarski, A. Choromanska, K. Choromanski, B. Firner, L. Jackel, U. Muller, and K. Zieba. Visualbackprop: visualizing cnns for autonomous driving. *arXiv preprint arXiv:1611.05418*, 2016. 2
- [3] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016. 1, 2, 3, 5, 6
- [4] M. Buehler, K. Iagnemma, and S. Singh. *The DARPA urban challenge: autonomous vehicles in city traffic*, volume 56. Springer, 2009. 2
- [5] P. Burt and E. Adelson. The laplacian pyramid as a compact image code. *IEEE Transactions on communications*, 31(4):532–540, 1983. 4
- [6] C. Chen, A. Seff, A. Kornhauser, and J. Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2722–2730, 2015. 2
- [7] S. Chen, S. Zhang, J. Shang, B. Chen, and N. Zheng. Brain inspired cognitive model with attention for self-driving cars. *arXiv preprint arXiv:1702.05596*, 2017. 1
- [8] Comma.ai. Public driving dataset. <https://github.com/commaai/research>, 2017. [Online; accessed 07-Mar-2017]. 5, 6, 7
- [9] A. Das, H. Agrawal, C. L. Zitnick, D. Parikh, and D. Batra. Human attention in visual question answering: Do humans and deep networks look at the same regions? *arXiv preprint arXiv:1606.03556*, 2016. 7
- [10] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A density-based algorithm for discovering clusters in large spatial

- databases with noise. In *KDD*, volume 96, pages 226–231, 1996. 4, 5
- [11] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010. 5
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 3
- [13] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell. Generating visual explanations. In *European Conference on Computer Vision*, pages 3–19. Springer, 2016. 3
- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 4
- [15] R. Hyndman, A. B. Koehler, J. K. Ord, and R. D. Snyder. *Forecasting with exponential smoothing: the state space approach*. Springer Science & Business Media, 2008. 3, 6
- [16] J. Johnson, A. Karpathy, and L. Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4565–4574, 2016. 3
- [17] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [18] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 2
- [19] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609–616. ACM, 2009. 3
- [20] J. Levinson, J. Askeland, J. Becker, J. Dolson, D. Held, S. Kammel, J. Z. Kolter, D. Langer, O. Pink, V. Pratt, et al. Towards fully autonomous driving: Systems and algorithms. In *Intelligent Vehicles Symposium (IV), 2011 IEEE*, pages 163–168. IEEE, 2011. 2
- [21] D. A. Pomerleau. Alvin, an autonomous land vehicle in a neural network. Technical report, Carnegie Mellon University, Computer Science Department, 1989. 2
- [22] R. Rajamani. *Vehicle dynamics and control*. Springer Science & Business Media, 2011. 3
- [23] S. Sharma, R. Kiros, and R. Salakhutdinov. Action recognition using visual attention. *arXiv preprint arXiv:1511.04119*, 2015. 1, 4
- [24] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- [25] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014. 5
- [26] Udacity. Public driving dataset. <https://www.udacity.com/self-driving-car>, 2017. [Online; accessed 07-Mar-2017]. 5, 6
- [27] C. Urmson, J. Anhalt, D. Bagnell, C. Baker, R. Bittner, M. Clark, J. Dolan, D. Duggins, T. Galatali, C. Geyer, et al. Autonomous driving in urban environments: Boss and the urban challenge. *Journal of Field Robotics*, 25(8):425–466, 2008. 2
- [28] H. Xu, Y. Gao, F. Yu, and T. Darrell. End-to-end learning of driving models from large-scale video datasets. *arXiv preprint arXiv:1612.01079*, 2016. 2
- [29] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, volume 14, pages 77–81, 2015. 1, 4
- [30] T. Zahavy, N. B. Zrihem, and S. Mannor. Graying the black box: Understanding dqns. *arXiv preprint arXiv:1602.02658*, 2016. 3
- [31] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014. 2