

Explainable Deep Learning for Video Recognition Tasks: A Framework & Recommendations

Liam Hiley¹

lhiley@cardiff.ac.uk

Alun Preece¹

apreece@cardiff.ac.uk

Julia Hicks²

hicksya@cardiff.ac.uk

¹ Crime and Security Research Institute
Cardiff University
Cardiff, UK

² School of Engineering
Cardiff University
Cardiff, UK

Abstract

The popularity of Deep Learning for real-world applications is ever-growing. With the introduction of high performance hardware, applications are no longer limited to image recognition. With the introduction of more complex problems comes more and more complex solutions, and the increasing need for explainable AI. Deep Neural Networks for Video tasks are amongst the most complex models, with at least twice the parameters of their Image counterparts. However, explanations for these models are often ill-adapted to the video domain. The current work in explainability for video models is still overshadowed by Image techniques, while Video Deep Learning itself is quickly gaining on methods for still images. This paper seeks to highlight the need for explainability methods designed with video deep learning models, and by association spatio-temporal input in mind, by first illustrating the cutting edge for video deep learning, and then noting the scarcity of research into explanations for these methods.

1 Deep Learning for Video

Deep Neural Networks (DNNs) have provided competent solvers for Image Recognition problems for a decade now. With early efforts in pattern recognition [11] paving the way for super-human performance in 2011 [7]. The recent surge in machine learning has brought with it the innovation to apply DNNs to complex problems in other domains such as bioinformatics, and particle physics [27]. One such domain being spatio-temporal data, most commonly video data, in which the motion of objects is depicted in images over time via animation.

1.1 Datasets

The main drive behind development in video and spatio-temporal deep learning has been human activity recognition. Many datasets have been compiled of human action in videos, either for classical classification/recognition tasks [1, 15, 21, 25, 46], sparse localisation/detection

[24, 50, 57] or other. The collection of video data proves difficult, with most datasets up until recently being small [25, 46], 13000 and 7000 samples respectively, very short [50, 57], 150 and 128 hours respectively, or notably short [24], 77 hours total. What's more, while the number of samples would be considered small for an image dataset, the need for data is intuitively greater for video deep learning models, which consistently rely on larger sets of sometimes more complex parameters for inference.

The development of the Kinetics dataset [21] somewhat revived the field by providing a comparatively very large dataset, at 300,000 samples initially. Kinetics boasts a rich variety of human action, mined from youtube videos and labelled using their titles, that has proven useful for pretraining feature extractors for large models that results in a comprehensive feature space for application on smaller more specific tasks. Kinetics has since become an analogue to the benchmark ImageNet [38].

Recently, similar efforts have been made to compile massive-scale video datasets of activity with [1] and [30], although they have not yet found the same success as [21].

Other groups have made efforts solving the question of gathering large enough sets of data for video machine learning tasks. Namely by accessing video surveillance footage [50, 57], a rich and varied source of human activity that is for the most part organic and unrehearsed. This medium for human activity captures more of what's considered the mundane or ordinary activity that might be less common in a source such as YouTube. In [50], researchers collected footage of crimes such as Abuse, Burglary and Assault that can be associated with one, continuous activity, with the hopes of providing data for anomaly detection, for which single-event crimes are perfect examples in a video stream mostly comprised of regular day-to-day activity.

1.2 Convolutional Neural Networks, Recurrent Neural Networks and Long Short-Term Memory Models

Recently the jump has been made from Still Image Processing tasks to streams of images in video input, with a focus on Human Activity Recognition [9, 19, 20, 44, 54, 56]. At first, attempts were made to pool together spatial features extracted from video frames, treated as individual images by applying successful image processing architectures, Convolutional Neural Networks (CNNs) [20]. The power behind CNNs lie in their kernel operators. The convolutional kernel acts as a sliding window, passing over an image and outputting the correlation between that section of the image and the kernels weights, which represent visual features such as lines and corners in the early layers, to scales on fish at much later layers. This technique is simple enough to adapt from image networks, and had the advantage of pretrained weights from massive scale image datasets such as [38]. However, it fails to retain temporal information from the stream, essentially flattening the frame block into a set of unique instances.

To encode temporal information into the model's decision, each model in the frame-by-frame representation could be folded or spooled into a Recurrent Neural Network (RNN) [9], a method popular in natural-language processing and speech recognition that function by spooling together many networks, each for a point in time. By weighting the network outputs, the full model can effectively learn sequence data, and is trained by unfolding into a large DNN. LSTMs improve upon RNNs in that they can retain temporal information over longer sequences via their Long Short-Term Memory cells, which replace the original 'cell' or network body in each time step's network. However, both are costly to train, with back-

propagation flowing through the entire unspooled network. Both also require a large spatial feature extractor to translate the input frame to a feature representation.

1.3 3D Convolutions

A video stream can be viewed as a ordered set of RGB frames, the temporal aspect of the video comes from the relationship between frames. By modelling the input data, not as a set of 2D data, but as a block of 3D data, temporal information is no longer a relation between inputs but is now inherent in the input. Each pixel becomes a voxel, encoding the colour intensity of a point in space and time. To match this, convolutional and pooling kernels need to be inflated to 3D [19]. This method relies on a single-stream spatio-temporal feature extractor to provide a rich enough feature space to contain the extremely varied movements found in human activity (e.g. Skiing vs. Chopping onions).

In [19], they present C3D, a small convolutional network with one fully connected layer for high-level reasoning. The input to the model is a 7 channel video, consisting of R, G, B frames as well as X and Y gradient and optical flow fields (discussed further in 1.4). This 3D representation resulted in the network surpassing baseline methods, including 1.2 on the TRECVID Human Action Surveillance dataset [57]. Since C3D, efforts have been made to retrace the success of popular CNN architectures from Image Recognition. In [13], members of the ResNet [14] family of models were able to outperform C3D and other contemporary methods on the Kinetics [21] and UCF-101 [46] action recognition datasets, showing that 3D Convolutions can be treated similarly to their 2D counterparts. The identifying feature of a ResNet model is it's use of residual skip connections, that connect input to a convolutional block to it's output so that the 'signal' is strengthened in a manner.

1.4 Two-Stream Networks

Rather than learn sequences of features using a recurrent layer, Simonyan and Zisserman [44] implemented a two-stream approach that aggregates the decision from two models in order to classify the video.

The first model, or stream is a regular CNN taking RGB images as input, similar to the early works using pooling from frame-by-frame image models. The second stream operates on the temporal dimension of the video, which is approximated via stacked optical flow frames. These optical flow fields can be taken to represent the motion between intermediate frames in a video, by encoding the apparent velocities between pixels using brightness as a measure of movement. Like the RNN technique, this method has the benefit of considering the temporal aspect of the video, which intuitively must be considered for input that is temporal in nature. Two-Stream Networks are also in some respect recurrent, in that the optical flow is recurrently optimised. Two-stream approaches also have the added benefit over RNNs, of capturing finer low-level motion [6]. Two-stream architectures have since gone on to be improved in various ways [54, 56] and optical flow is still considered beneficial to activity recognition problems to date, with two-stream approaches even taken to 3D CNNs [6]. This two-stream, 3D CNN, dubbed I3D, named for its pretraining method, inflating ImageNet weights from 2D, averages Stacked Optical Flow and RGB trained models at test time. This achieves state-of-the-art performance on [21], [46] and [25]. In [55], it is shown to be possible to 'hallucinate' optical flow fields for still image input using a Hybrid Video Memory machine. The HVM uses a memory of similar still images to the input to generat an optical

flow for the input from noise, based on the relevant optical flow fields given to it with the memorised images. Therefore, for a still image of an action/activity, the model could relatively cheaply and successfully generate temporal features that would then be used as input for the second stream of a two-stream architecture, whereas before the decision would rest solely on the RGB channels.

2 Explaining spatio-temporal models

2.1 Kernel visualisation and Global Explanations

Efforts to explain deep learning models developed for video tasks often seem to be included as an afterthought, an additional justification of the method to supplement the benchmark performance. One logical method is to visualise the kernels, that learn the spatio-temporal features that make these models unique. Kernel visualisation is common in Image tasks as a method of linking deep learning decisions. In [10], images were optimised to maximally activate single filters at a time in one of two deep learning models trained on the MNIST handwritten digit dataset [26]. This was adapted to Convolutional networks in [45]. [33] provides an interactive and comprehensive investigation of feature visualisation for a large deep Image recognition model. In [44], the previous method was adapted for filters trained on optical flow, in an effort to justify the model by way of linking the features it learns to previous hand-crafted methods. In [6] the method is used to visualise the 3D filters, slice by slice rather than as cubes, with the aim of showing the richness of the filters. Again, the aim of these works is not to investigate the features learned by the networks, but to propose the networks themselves. In [54], Deep Draw is applied to a Temporal Segment Network, a Two-Stream Network that adapts sparse sampling to classify using the entire video. The result is an image optimised from noise that maximally activates a particular neuron in the logits, or final, layer, relating to an output class. The authors of DeepDraw [34] note the issues of optimising on noise without constraint, mainly that it generates crowded images that are difficult to decipher. The use of DeepDraw in [54] gives insight into the models understanding of the classes identity in the feature space. Showing intuitive associations between motion and activities, such as ripples in water with the activity of Diving, in the optical flow stream which emphasises the influence of moving bodies of water on the decision. The spatial stream for the same class however activates much more strongly it seems on the presence of humans, topless or in swimsuits. Interestingly, [53], a work on spatio-temporal feature learning with the C3D [19], did not use this method, preferring instead Deconvolution [58] which visualises a filters activations on a given input sample. This local explanation of the learned features, in the [28, 36] sense of the term local, is used to infer how such features have developed. They note that the C3D filters at first focus on the spatial dimension but quickly transfer focus to the motion in the frame, the temporal dimension. Also notable, is that [12], somewhat of a successor of [53], does not attempt to visualise the features it sets out to learn. An application of feature visualisation, in the spirit of explainability can be found in the work of Nerinovsky [32]. Here the kernels at different layers of the popular I3D network [6] are visualised, producing animated videos of swirling features coming in and out of focus. Unlike the recognisable eyes, noses and other such features of animals, it is hard to gain an intuition from these sliced cubes as to what each might detect in a video of human activity, although at earlier layers it is clearer to determine the general type of motion that activates each kernel.

While still explaining the model using its weights, [22] takes a very different approach to visualisation. Concept Activation Vectors (CAV) are partially handcrafted features for model explanations. CAVs are generated by manually segmenting the model’s training dataset into two groups, such that one group contains all samples that represent or at least contain instances of a ‘concept’ as defined by a human. A linear model is then trained on the two subsets, where the input is the activations of the weights of the model on the samples from each subset. The concept activation vector is then the normal of the hyperplane separating the two subsets, towards the concept class. This results in an explainer model that can flag the model as having seen a concept in its input, based on the model’s activations for that input. The authors have so far only applied this method to image data, but suggest video as one of a few possible extensions.

2.2 Explaining Local Decisions in the Input Space

Local explanations, explanations that explain the decision on a single sample, have found significant success in the Image domain, with most of the techniques developed with that input medium in mind. Sensitivity analysis [4, 35] can be considered one of the earliest attempts at local, visual explanation with gradients of class probability, vectors in the direction of the decision boundary, are displayed on digit classification to illustrate how to change the digit 2, for example, so that the model thinks it’s an 8. While useful and intuitive for simple examples like digit classification, sensitivity analysis often gets noisy for larger more complex problems.

2.2.1 Layer-wise Relevance Propagation and White-box methods

Increasingly popular are the Layer-wise Relevance Propagation (LRP) family of techniques, first defined in [3] which take a white-box approach to local explanations, in that they assume access to the model internals. LRP is proposed as a method for pixel-wise decomposition of relevance to a decision. Although multiple variations on the method are presented, the main theory centres around the following conditions that must be satisfied for the explanation to be considered an LRP explanation.

- The relevance at each layer must sum to the output of the model.
- The relevance at any neuron in the layer other than the output layer is the sum of incoming relevances to that layer.

The authors clarify that this can be satisfied by meaningless explanations and as such the formula for LRP should be taken more as a framework to follow in developing explanation methods. The work has since been adapted and implemented by other researchers [43, 59], and maintained and improved by the original authors [5]. In [31] the original authors propose an implementation of their prior theoretical work based on Taylor expansions of the decision function at a pre-defined root point. The relevance of a neuron is then defined locally in comparison to the root point. This method has since been applied in various fields for explaining deep models [2, 40, 41, 47]. Note: [47] is an application of the LRP method to compressed domain human action recognition. And as such while not a method covered in section 1, is still a very rare example of an investigation of explainability in the spatio-temporal domain.

In [51], the authors note that two models that provide exactly the same output for all inputs, regardless of model internals or implementation, should provide identical explanations, an axiom known as Implementation Invariance. By implementing discrete intermediate gradients, for which the chain rule cannot be applied, [3] and [43] are not effective explanation methods, as relevance in the implementation may be attributed even though it is not relevant in the input. The authors go on to define integrated gradients, a method that instead integrates the gradient of the model function against a root point, at all intermediate points. This method is then shown to satisfy the completeness rule, that the relevance sums to the difference between the input and the root, that is defined as desirable in [3].

LRP methods are inherently applicable to the spatio-temporal domain, since they attribute relevance from the output onto the input space, irrespective of dimensionality. Describing the relevance of the temporal dimension on the decision at a voxel is however still only implied through difference in similar spatial features between frames. That is, by reconstructing the relevance field into a video, it is difficult to see whether an area in a frame is attributed high relevance because of its spatial information, or because of its temporal information. In [47] they instead opt to show the amount of relevance over time, including the frames as reference at key points. This better shows the effect of motion on the decision, which suggests that explanations for video, cannot be digested in the same way that the inputs are.

Other white-box methods have been developed that do not use the LRP framework. CAM [60] notes the use of Global Average Pooling, originally developed for training regularisation, for localisation of feature map activations. This produces a heatmap of the areas in the image that activated strongly to a particular class. The authors note the general loss of spatial information in hidden layers, and thus this technique is not compatible with models that make use of such layers. Grad-CAM [42] addresses this by first propagating the gradient of the class output node back through the hidden layers, with respect to the last convolutional layers activations. From then on CAM can be applied. CAM and Grad-CAM produce class discriminative heatmaps that perform very well on object localisation without any additional training. This has the added benefit of showing objects in the scene that the model recognises as relevant to the decision. These, again, are applicable to video data, much moreso in the case of [42] since the majority of models for video tasks use at least one fully connected layer.

The SHAP (SHapley Additive exPlanations) framework for local explanations [29] finds similarities in [3, 36, 43] and proposes a framework based on methods from game theory, to improve all methods. SHAP explanations are based on multiple explainers for different models. The deep explainers are gradient based and white-box. The authors suggest that SHAP values are the only possible consistent, and locally accurate additive feature attribution method.

Perhaps the least efficient but most model-faithful explanation technique is proposed in [37], where models are trained with a term for explainability included in the loss. The authors constrain input gradients on features considered to be irrelevant to the decision. This discourages the model from learning on those features, which should result in a model that focuses more on objects in the scene that a human deems significant.

To our knowledge, the only explanation method in this format developed solely for video can be found in [48], and [16]. In [48] the authors map the activations to the input in a way that is similar theoretically and visually to [42]. This method, known as Saliency Tubes, provides cylindrical heatmaps that visualise the focus of attention in the input video, through each frame. As an alternate visualisation, the stack of frames is staggered so that the path

of the tube through the video can be seen at once. This method for visualising the motion is useful in better translating the temporal aspect of the explanation without animating the frames. In [16], the authors suggest separating the relevance generated by the deep Taylor method, for 3D models, into its spatial and temporal components. They also provide a novel yet simple way of approximating this that shows that models do seem to attribute relevance to some objects in a scene because of their motion, much more so than their appearance.

2.2.2 Black-box methods

Other efforts have been made to develop explanations for models without access to the model internals. The most popular of these techniques, known as Local Interpretable Model Explanations or LIME [36], seeks to approximate the decision function by many closely sampled input points, which all center around the input point to be explained. It can then attribute positive or negative influence on the decision function to the differences in the sampled inputs, and overlay this on the original input. LIME has found much success and support in the field, and is implemented for text, image, and tabular data. In [49] they note the instability of LIME: Since it uses random seeding to draw samples, via segmentation in images, multiple explanations of the same input for the same model can result in very different attributions. The authors suggest aggregating explanations, to improve stability.

The theory behind LIME itself is extendable to all popular machine learning input domains. The issue lies in how the samples are drawn from the input. For text processing, words can be considered atoms to be used for sampling, the same can be said for items in tabular data. LIME's solution to sampling images is based on segmentation. This again results in instability, and would for any 3D segmentation adapted for video inputs.

2.2.3 Explanation by example

In [23], a different approach to local explanations is taken. Instead of attributing relevance to features in the input sample, the implementation of which all previous methods [3, 29, 31, 36, 43] have argued over, Influence Functions instead explains the model's decision on an input in terms of the training samples that most influenced that decision. They do this by discretely deriving the loss function at the explained sample, with respect to each input in the training data. This is computationally massive task to achieve, especially in the chosen example domain of images, with large datasets like ImageNet. However, the authors also show that even approximations to the derivative can provide insight into what influenced the decision.

This method does not attempt to alter the input, and only relies on being able to provide that input as an example. Therefore, it is the most naturally applicable to the video domain. The explanation can be presented in the same manner as the model input without any processing involved.

3 The Interpretability of Explanations

While providing an explanation that is faithful to the model, i.e. accurately represents the models entire decision process, is a quality sought after by many in the literature [10, 29, 31, 34, 58] ([37] by far the most so), the resulting explanations are often noisy and sometimes indecipherable to a human, showing that transparency is not necessarily the only quality to

optimise for when generating explanations. Other methods seek instead to be more digestible for humans [22, 36]. Many works have sought to provide guidelines on more interpretable explanations, in recent years. In [36], they note the tradeoff between an explanation’s fidelity to a model and its interpretability to humans. Lipton [28] defines some qualities of an Interpretable model, namely transparency, and post-hoc interpretability, which covers the explanation methods noted above (Section 2). He goes on to note a few key cautions to take when optimising interpretability: Linear models are not necessarily more interpretable than deep networks, Transparency can contend with performance if parameters are limited so that they are more understandable. Also, that a plausible explanation, i.e. one that perfectly highlights a salient object, might not be a true explanation of the model, and so should not be the only factor when grading that explanation. It is suggested in [8] that interpretability is not necessarily only human understanding. This work also introduces δ -interpretability, or the quality of a model being able to improve performance on a task by way of explanation. Interpretability and Explainability are first distinguished in [52], where Explainability is the capability of the model to provide an explanation that is faithful to the causal factors of that decision. Interpretability is defined here as the information on the decision an agent can gain by use of the explanation, given the transparency of the model. This sets up Interpretability as a qualitative measurement for Explanations. The authors go on to suggest that since the role an agent plays in the ecosystem formed around a supporting AI, Interpretability is subjective to the agent. Context such as field knowledge, motivation and time constraints must factor into the explanation provided, as different information is valuable based on the use the agent intends to find from the explanation. There is still disagreement on how to quantify the explainability and interpretability of an explanation, [17] suggests guidelines for measuring explainability, that would also assure interpretability in the above sense. In [39] attribution methods are quantitatively assessed by ordering the pixels of the relevance map and then repeatedly removing the first most relevant pixel. By degrading the image in this way, the authors theorise that explanations that best represent the decision will cause the greatest drop in performance for that sample when passed again through the model. This method confirms that [3] outperforms [35] and [58] in model fidelity. This method’s accuracy is called into question in [18], where it is suggested that removing information from an image, results in artifacts that essentially shifts the distribution of input features to one that the model was not originally trained on. Therefore, it is suggested instead to retrain the model after removal of the feature, and test the models performance then. It was shown through this that models are surprisingly robust to feature removal, and can still reliably classify images with very little information left from the original sample.

4 Discussion

It seems the main issues with interpretability and explainability of Deep Learning models for video tasks is first and foremost a lack of attention. The scope of explainable AI is still very much centered on the Image domain, and efforts in explaining spatio-temporal models do not attempt to adapt techniques to the new modality. Work to highlight the hidden temporal nature of explanations past animating the frames, in [47, 48] currently suggests unraveling the video into its frames and their explanations, but this is then either cluttered (as in [48]) or loses the videos sense of sequence. Future work on explaining motion in video will likely use separate explanations of temporal saliency on optical flow fields, as these networks are still widely used. However, developing a system in which spatial and

temporal relevance do not interfere with one another when displayed together in one image, would provide a visually efficient mode of explanation that does not require referencing between multiple images, and would allow the explanation to be viewed in the same way as its input. Explanation could also add insight into the influence of motion on model decisions. The main approaches to video tasks are thought to capture temporal information at different abstractions [6], explanations that are heavily motion focused, with little relevance attributed to spatial features, would support this statement.

5 Framework for Future Research

In this paper we have addressed some gaps in the current works for explainable activity recognition with video deep learning models. To suggest some possible solutions for future efforts, we propose a framework for improving this field, from its current state:

1. Many of the current approaches to explainability for video models are adapted from methods for still image models. To faithfully explain video models, it will be necessary to design new methods intended primarily for these models.
2. Applications of video models in a deployment setting will often feature real-time processing, for example in a surveillance scenario. This is necessary to take into account when implementing explanation methods, since heavier-weight methods would be infeasible in such a scenario.
3. Furthermore, considering applications must be taken into account, to ensure the success of the method. Explaining models with complex representations of motion such as 3D CNNs and C-RNNs can result in hard to interpret results. Explanations that aim for transparency and faithfulness might fall short for a user with no background in machine learning or computer vision. This problem is addressed in [52], where the authors design a framework for the user roles in a explainable deep learning system. They highlight the need to consider a user's background and their motivation for requesting the explanation, when providing that explanation. This should be the case for future video deep learning explanation methods as well.

6 Conclusion

In summary, video deep learning is currently at the forefront of machine vision, with a rich and varied body of work committed to comprehensively and compactly representing time in imagery. However this success and interest has not attracted much work in explaining these models, in contrast to image recognition for which Explainable AI is likely now most popular. Visualising the features learned by these spatio-temporal models has been thought of by some when designing new architectures, but mostly to justify that the models are learning *something*, which is an observation oriented towards a researcher and perhaps less useful when justifying a model for deployment. Local explanation methods have also been extended to these models, but it is clear that out of the box these methods are very much anchored to their origins in the image domain and as such seem an ill fit to the video input. Current work in adapting popular explanation methods such as visualisation [32], deep Taylor [47] and Grad-CAM [48] show promise and possibility for understanding these models, and there are clear directions to take in the future such as exploring the use of optical

flow, in the explanation as much as it is already in the decision, as well as a method to highlight temporally-salient regions in an explanation in order to distinguish the models decision making process from that of a similar model learned on images.

We provide a framework for addressing these problems, notably the necessity to develop native video explanation methods, that these methods should be lightweight enough to run near real-time, and that the developers take into consideration the intended users and their usecases for these techniques.

References

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016.
- [2] Leila Arras, Franziska Horn, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. "what is relevant in a text document?": An interpretable machine learning approach. *PLoS one*, 12(8):e0181142, 2017.
- [3] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- [4] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(Jun):1803–1831, 2010.
- [5] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *International Conference on Artificial Neural Networks*, pages 63–71. Springer, 2016.
- [6] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.
- [7] Dan CireşAn, Ueli Meier, Jonathan Masci, and Jürgen Schmidhuber. Multi-column deep neural network for traffic sign classification. *Neural networks*, 32:333–338, 2012.
- [8] Amit Dhurandhar, Vijay Iyengar, Ronny Luss, and Karthikeyan Shanmugam. Tip: Typifying the interpretability of procedures. *arXiv preprint arXiv:1706.02952*, 2017.
- [9] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634, 2015.
- [10] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. 2009.
- [11] Kuniyuki Fukushima. Neocognitron. *Scholarpedia*, 2007. doi: 10.4249/scholarpedia.1717.
- [12] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3154–3160, 2017.
- [13] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. *arXiv:1512.03385 [cs]*, 2015. arXiv: 1512.03385.
- [15] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Nieves. ActivityNet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [16] Liam Hiley, Alun Preece, Yulia Hicks, Harrison Taylor, and David Marshall. Discriminating spatial and temporal relevance in deep Taylor decompositions for explainable activity recognition. In *Workshop on Explainable Artificial Intelligence (XAI), IJCAI*, 2019.
- [17] Robert R Hoffman, Shane T Mueller, Gary Klein, and Jordan Litman. Metrics for explainable ai: Challenges and prospects. *arXiv preprint arXiv:1812.04608*, 2018.
- [18] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. Evaluating feature importance estimates. *arXiv preprint arXiv:1806.10758*, 2018.
- [19] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
- [20] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [21] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [22] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viégas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). *arXiv preprint arXiv:1711.11279*, 2017.
- [23] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, pages 1885–1894. JMLR. org, 2017.
- [24] Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 780–787, 2014.
- [25] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563. IEEE, 2011.
- [26] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- [27] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553): 436, 2015.
- [28] Zachary C. Lipton. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*, 2016.
- [29] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [30] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Tom Yan, Lisa Brown, Quanfu Fan, Dan Gutfrueud, Carl Vondrick, et al. Moments in time dataset: one million videos for event understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–8, 2019. ISSN 0162-8828. doi: 10.1109/TPAMI.2019.2901464.
- [31] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
- [32] Arseny Nerinovsky. Visualizations of the i3d network. URL <https://github.com/Arseny-N/c3d-vis>.
- [33] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The Building Blocks of Interpretability. *distill.pub*, March 2018.
- [34] Audun M. Øygaard. Visualizing GoogLeNet classes. 2015. URL <https://www.auduno.com/2015/07/29/visualizing-googlenet-classes>
- [35] Peter M. Rasmussen, Tanya Schmah, Kristoffer H. Madsen, Torben E. Lund, Stephen C. Strother, and Lars K. Hansen. Visualization of nonlinear classification models in neuroimaging, 2012.
- [36] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should I trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. ACM, 2016.
- [37] Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. *arXiv preprint arXiv:1703.03717*, 2017.
- [38] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *arXiv:1409.0575 [cs]*, 2014. URL <http://arxiv.org/abs/1409.0575>. arXiv: 1409.0575.
- [39] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems*, 28(11):2660–2673, 2017.

- [40] Robin Tibor Schirrmeyer, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggenberger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for eeg decoding and visualization. *Human brain mapping*, 38(11):5391–5420, 2017.
- [41] Kristof T Schütt, Huziel E Saucedo, P-J Kindermans, Alexandre Tkatchenko, and K-R Müller. Schnet—a deep learning architecture for molecules and materials. *The Journal of Chemical Physics*, 148(24):241722, 2018.
- [42] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017.
- [43] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3145–3153. JMLR. org, 2017.
- [44] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014.
- [45] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [46] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [47] Vignesh Srinivasan, Sebastian Lapuschkin, Cornelius Hellge, Klaus-Robert Müller, and Wojciech Samek. Interpretable human action recognition in compressed domain. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1692–1696. IEEE, 2017.
- [48] Alexandros Stergiou, Georgios Kapidis, Grigorios Kalliatakis, Christos Chrysoulas, Remco Veltkamp, and Ronald Poppe. Saliency tubes: Visual explanations for spatio-temporal convolutions, 2019.
- [49] Mitchell Stiffler, Adam Hudler, Eunjin Lee, Dave Braines, David Mott, and Daniel Harborne. An analysis of reliability using lime with deep learning models. 2018.
- [50] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6479–6488, 2018.
- [51] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR. org, 2017.
- [52] Richard Tomsett, Dave Braines, Dan Harborne, Alun Preece, and Supriyo Chakraborty. Interpretable to whom? a role-based model for analyzing interpretable machine learning systems. *arXiv preprint arXiv:1806.07552*, 2018.

- [53] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015.
- [54] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*, pages 20–36. Springer, 2016.
- [55] Yali Wang, Lei Zhou, and Yu Qiao. Temporal hallucinating for action recognition with few still images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5314–5322, 2018.
- [56] Yifan Wang, Jie Song, Limin Wang, Luc Van Gool, and Otmar Hilliges. Two-stream sr-cnns for action recognition in videos. In *BMVC*, 2016.
- [57] Ming Yang, Shuiwang Ji, Wei Xu, Jinjun Wang, Fengjun Lv, Kai Yu, Yihong Gong, Mert Dikmen, Dennis J Lin, and Thomas S Huang. Detecting human actions in surveillance videos. In *TREC video retrieval evaluation workshop*, 2009.
- [58] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [59] Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation backprop. *International Journal of Computer Vision*, 126(10):1084–1102, 2018.
- [60] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.