# Interpretable Deepfake Detection via Dynamic Prototypes

**Loc Trinh, Michael Tsang, Sirisha Rambhatla, Yan Liu**
Department of Computer Science
University of Southern California
Los Angeles, CA 90089
{loctrinh, tsangm, sirishar, yanliu.cs}@usc.edu

## Abstract

Deepfake is one notorious application of deep learning research, leading to massive amounts of video content on social media ridden with malicious intent. Therefore detecting deepfake videos has emerged as one of the most pressing challenges in AI research. Most state-of-the-art deepfake solutions are based on black-box models that process videos frame-by-frame for inference, and they do not consider temporal dynamics, which are key for detecting and explaining deepfake videos by humans. To this end, we propose Dynamic Prototype Network (DPNet) – a simple, interpretable, yet effective solution that leverages dynamic representations (i.e., *prototypes*) to explain deepfake visual dynamics. Experiment results show that the explanations of DPNet provide better overlap with the ground truth than state-of-the-art methods with comparable prediction performance. Furthermore, we formulate temporal logic specifications based on these prototypes to check the compliance of our model to desired temporal behaviors.

## 1 Introduction

While artificial intelligence (AI) plays a major role in revolutionizing many industries, it has also been used to generate and spread malicious misinformation. In this context, *Deepfake* videos – which can be utilized to alter the identity of a person in a video – have emerged as perhaps the most sinister form of misinformation, posing a significant threat to communities around the world [49, 8, 19, 21, 46, 48]. As deepfakes become pervasive, ascertaining the trustworthiness of a video and making a determination of its authenticity becomes critical.

To address this challenge, a series of excellent work has been conducted on detecting deepfakes [47, 31, 1, 38]. While they have achieved good progress towards the prediction task to a certain extent, there is still significant room for improvement. First, even though existing work focus on the authentication problem, very few of them address the interpretability issue. That is, explaining *why* a model predicts a certain video as real or fake, which can be crucial for maintaining trustworthy content dissemination. Second, humans detect a deepfake video by examining the unnatural dynamics caused by the distortions induced by the generation model [24, 37, 59]. Yet most state-of-the-art deepfake detection techniques analyze a potential video frame-by-frame, and have not explored these temporal dynamics [35, 38]. As a result, there is a need for an *interpretable* deepfake detection technique which considers temporal dynamics and at the same time gives insight into the temporal inconsistencies in deepfake videos.

To this end, we propose DPNet – an interpretable prototype-based neural network that captures dynamic features, such as unnatural movements and temporal artifacts, and uses them to explain *why* a particular prediction was made. Specifically, DPNet works by first learning the prototypical representations of the temporal inconsistencies within the latent space, utilizing spatial-temporal
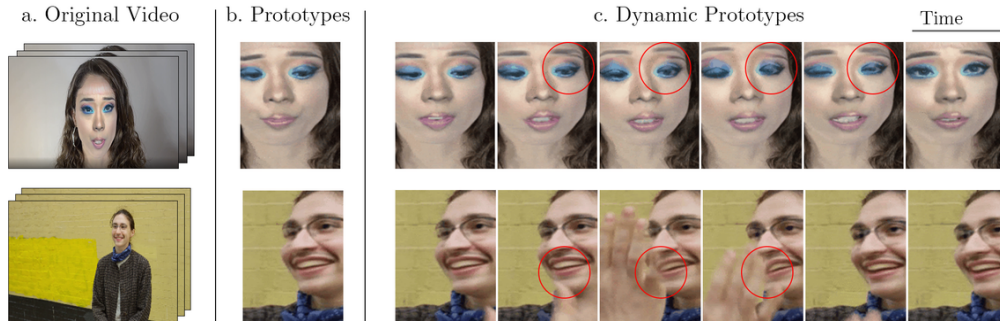
Figure 1: Examples of static vs. dynamic explanations for deepfake videos. Qualitatively, seeing temporal artifacts allow a human judge to quickly determine whether a video is real or fake. Red circles indicate regions of interest. Best view as GIFs (See Appendix C).

information from the inputs. Then, it makes predictions based on the similarities between the dynamics of a test video and the small set of learned dynamic prototypes. Lastly, the dynamic prototypes are then intermittently projected to the closest representative video patch from the training dataset, yielding an immediate human-understandable interpretation of the learned prototypes.

The primary advantages of `DPNet` are as follows:

- **Faithful explanations via case-based reasoning**: `DPNet` follows a case-based reasoning approach that utilizes previously learned dynamics - as evidence (i.e *cases*) - to tackle a new unseen testing video. This helps the model to explain why a certain prediction was made, in a way that is reflective of the network's underlying reasoning process.

- **Visual dynamic explanations**: `DPNet` provides explanations in the form of visual dynamics, via the learned prototypes, that are principled, accessible, and easy for humans to understand.

- **Temporal logic specifications**: `DPNet` automatically learns the dynamic prototypes, which can be used to formulate temporal logic specifications that check the robustness of the model and verify whether it conforms to desired temporal behaviors.

## 2 Related Work

**Deepfake detection**: A prominent line of work for deepfake detection focuses on hand-crafting facial features from the video, such as eye color and missing reflections [31], 3D head poses [54], and facial movements [2, 6]. However, these approaches do not scale well to larger and more sophisticated deepfakes. To address this problem, researchers propose to leverage advances in image detection via convolutional neural networks (CNNs) and process the deepfakes frame-by-frame. Examples include spatial pyramid pooling module to detect resolution-inconsistent facial artifacts [27], applying mesoscopic features [1], ImageNet-based [12] model XceptionNet [38] (which achieved state-of-the-art results by fine-tuning on deepfake datasets), and variants of Capsule Networks [35, 34]. All these methods process deepfake videos frame-by-frame and do not consider temporal dynamics in the videos. However, it is known that humans detect a deepfake video by examining the unnatural dynamics caused by the distortions induced by the generation model [24, 37, 59]. Even though several attempts have been made to explore multi-modal and temporal information in deepfakes (e.g, two-stream CNNs [58], recurrent neural networks [39], and inter-frame dissimilarities using optical flow [4]), their performance is inferior to that of frame-based methods.

**Interpretable neural networks**: One popular approach to explaining deep neural networks is *posthoc* analysis via gradient and perturbation-based methods [42, 55, 5, 41, 44, 40, 33, 13, 53]; however, they do not simplify the inherently complex underlying architectures. Instead, another line of research tries to build networks that are interpretable by design, with a built-in to self-explain. The advantage of this approach is that the interpretability is presented via *units* of explanation — general concepts and not necessarily raw inputs. This can be seen in the work of [3] for basis concept learning and [23, 26, 32] for case-based reasoning and prototype learning. Chen et al. [11] proposed learning prototypes for image classification to make predictions based on similarity to class-specific image patches. However, this type of interpretability has not been brought to video classification to leverage the rich temporal information within videos.
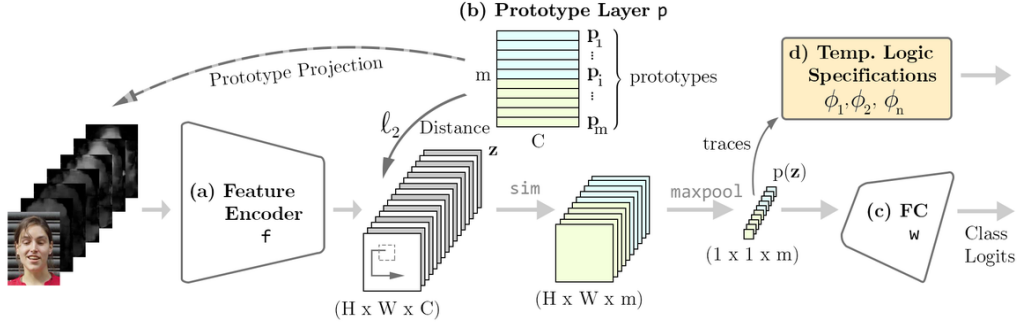
Figure 2: `DPNet` architecture for deepfake detection.

Although deep learning-based video classification models have been developed for video understanding [57, 29, 57, 51, 16], there is much to be desired in terms of interpretability, especially when compared to intrinsically interpretable models. In contrast, our proposed approach directly captures temporal artifacts in deepfakes as dynamic prototypes and uses them as visual explanations to explain predictions, which is specifically important for deepfake detection.

**Safety verification**: Neural networks (NN), in general, cannot provide any guarantees regarding model safety, leading to critical failures, one being adversarial examples [45, 17]. The importance of safety verification, especially in critical domains such as healthcare or autonomous driving, was highlighted when the discovery of such attacks prompted many lines of work in robust verification in ML [25, 22, 18, 20, 52, 10, 30]. To further reason about safety and robustness in time, temporal logic has been broadly in previous work [7, 36, 43, 15], and recent work in using temporal logic to verify time-series and NN-based perception system have shown promises [50, 13]. The dynamic prototypes from our interpretable models provide a convenient vehicle for us to formulate temporal logic specifications for videos and reveal valuable insights into the model's compliance to desired temporal behaviors.

## 3 Dynamic Prototype Network

We introduce our Dynamical Prototype Network (`DPNet`), the loss function, and the training procedure in this section. In addition, we highlight the steps that our network took to predict a new video, and how those exact steps can be interpreted in a human-friendly way.

### 3.1 `DPNet` architecture

The proposed architecture is shown in Figure 2. Formally, let $\mathcal{V} = \{(\mathbf{v}_i, y_i)\}_{i=1}^N$ be the video dataset, where $\mathbf{v}_i$ be a deepfake video sequence of length $T_{\mathbf{v}_i}$, and $y_i \in \{0, 1\}$ is the label for fake and real. As shown here, `DPNet` consists of four components: the feature encoder $f$, the prototype layer $p$, the fully-connected layer $w$, and the robustness temporal logic checker.

**Feature encoder** $f(\cdot)$: The feature encoder $f$ encodes a processed video input $\mathbf{x}_i \in \mathbb{R}^{224 \times 224 \times S}$ into a hidden representation $\mathbf{z} \in \mathbb{R}^{H \times W \times C}$. Here the input $\mathbf{x}_i$ to the encoder $f$ is formed by stacking one RGB frame with precomputed optical flow fields between several consecutive frames, yielding $S$ channels (Figure 2a). We let the input to the `DPNet` be a fixed-length $T < T_{\mathbf{v}_i}$, and randomly selected the initial starting frame for $\mathbf{x}_i$. This allows us to explicitly describe the motion of facial features between video frames, while *simultaneously* presenting the RGB pixel information to the network. Furthermore, since the optical flow fields can also be viewed as image channels, they are well suited for image-based convolutional networks. The feature encoder $f$ outputs a convolutional tensor $\mathbf{z} = f(\mathbf{x}_i)$ of shape $(H, W, C)$ that is forwarded to the prototype layer.

**Prototype layer** $p(\cdot)$: The network learns $m$ prototype vectors $\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_m$ of shape $(1, 1, C)$ in the latent space, each corresponds to a *dynamic prototype* in the architecture; see section 3.2.1 for the learning procedure. The prototype layer $p$ computes the squared $\ell_2$ distance between the prototype vectors $\mathbf{p}_j$ and each *patch* (of shape $(1, 1, C)$) of the encoded input $\mathbf{z}$. Thus generating $m$ distance maps of shape $(H, W)$. The shape of the prototype vectors is chosen to represent the smallest dynamic patch within the encoded input $\mathbf{z}$.

Formally, in Figure 2b, the prototype layer $p$ computes $m$ similarity scores:

$$p(\mathbf{z}) = [\, p_1(\mathbf{z}), \ p_2(\mathbf{z}), \ \ldots, \ p_m(\mathbf{z}) \,]^T \tag{1}$$

3

where, the similarity score between a prototype $\mathbf{p}_j$ and $\mathbf{z}$, denoted as $p_j(\cdot)$ is given by

$$p_j(\mathbf{z}) = \max_{\mathbf{z}' \in \text{patches}(\mathbf{z})} \frac{1}{1 + \|\mathbf{z}' - \mathbf{p}_j\|_2^2} \tag{2}$$

**Fully-connected layer** $w(\cdot)$: This layer computes weighted sums of similarity scores, $\mathbf{a} = \mathbf{W}\, p(\mathbf{z})$, where $\mathbf{W} \in \mathbb{R}^{K \times m}$ are the weights, and $K$ denotes the number of classes ($K = 2$ for DPNet). We then use a softmax layer to compute the predicted probability as follows,

$$\hat{y}_i = \frac{\exp(a_i)}{\sum_{j=1}^{K} \exp(a_j)}. \tag{3}$$

Note that, we allocate $m_k$ prototypes for each class $k \in \{0, 1\}$ s.t. $\sum_k m_k = m$. In other words, every class is represented by $m_k$ prototypes in the final model.

**Verification of dynamic prototypes via temporal logic**: Our DPNet architecture allows for the usage of formal methods to verify the robustness of our model. Given the direct computational path from $f$ to $p$ to $w$, we can verify whether the learned prototypes satisfy some desired temporal behaviors. Here, we used Timed Quality Temporal Logic (TQTL) similar to [13], but instead, we consider each video as a *data stream* with each frame as a *time-step*. We hereby give a brief summary of the TQTL language, and the specifications we used to verify our model.

*Timed Quality Temporal Logic.* The set of TQTL formulas $\phi$ over a finite set of Boolean-value predicates $Q$ over attributes of prototypes, a finite set of time variables ($V_t$), and a finite set of prototype indexes ($V_p$) is inductively defined according to the following grammar:

$$\phi ::= \texttt{true} \,|\, \pi \,|\, \neg\phi \,|\, \phi_1 \vee \phi_2 \,|\, \phi_1 \,\mathbf{U}\, \phi_2 \,|\, x \leq y + n \,|\, x.\phi \,|\, \exists p_i @ x, \phi \,| \tag{4}$$

where $\pi \in Q$, and $\phi_1$ and $\phi_2$ are valid TQTL formulas. $\pi$ has the form $\pi \equiv f_\pi(t_{1...n}, p_{1...m}) \sim C$, where $\sim$ is a comparison operator, i.e. $\sim \in \{<, \leq, >, \geq, =, \neq\}$, and $C \in \mathbb{R}$. For example, the predicate for "the similarity of prototype $\mathbf{p}_1$ to an input at time step 2 is greater than 0.9" is $f(t_2, id_1) > 0.9$. We hereby use $S(\cdot)$ to denote the prototype similarity score.

In the grammar above, $x, y \in V_t, n \in \mathbb{N}, p_i \in V_p$, and $\mathbf{U}$ is the "until" operator. The time constraints of TQTL are represented in the form of $x \leq y + n$. The freeze time quantifier $x.\phi$ assigns the current time to a variable variable x before processing the subformula $\phi$. The quantifiers $\exists$ and $\forall$ respectively existentially and universally quantify over the prototype IDs in a given frame. In addition, we use three additional operators: ($\psi$ Implies $\phi$) $\psi \rightarrow \phi \equiv \neg\psi \vee \phi$, (Eventually $\psi$) $\Diamond\psi \equiv \texttt{true}\,\mathbf{U}\,\psi$, and (Always $\psi$) $\Box\psi \equiv \neg\Diamond\neg\psi$. The semantics of TQTL can be find in the Appendix B.

*Specifications.* We verify that if our model predicts $fake$ for a testing video $V$, throughout the video, there exists a clip at time $t$ where a prototype $\mathbf{p}_i \in \mathbf{P}_{fake}$ is most similar to it compared to all prototypes of the $real$ class $\mathbf{p}_k \in \mathbf{P}_{real}$ for all time $0 \leq t' \leq T_V$. This verifies the existence of a key frame that our prototypes 'see' $fake$ strongly. The formula $\phi_{1,key\_frame}$ denotes this specification:

$$\begin{aligned}
\phi_{1,key\_frame} = \;& \Diamond(t.\exists p_k @ t, Class(V) = fake(real) \wedge p_k \in P_{fake(real)} \\
& \rightarrow \Box(t'.((0 < t' \wedge t' < T_V) \\
& \rightarrow \forall p_j @ t', p_j \in P_{real(fake)} \wedge S(t, p_k) > S(t', p_j) )))
\end{aligned}$$

Moreover, we specify that if a prototype is *non-relevant*, its similarity to a frame should be consistently low across. An example of this safety specification is that a prototype representing a $fake$ temporal artifact should not be highly activated in a $real$ video at any time. We specify this notion below, where numerical values are *user-specified thresholds* that control the strictness of the specification:

$$\begin{aligned}
\phi_{2,non\_relevance} = \;& \Box(t.\forall p_i @ t, Class(V) = fake(real) \wedge p_i \in P_{real(fake)} \\
& \rightarrow S(t, p_i) < 0.4 \wedge \Box(t'.(t \leq t' \wedge t' \leq t + 5) \\
& \rightarrow |S(t', p_i) - S(t, p_i)| < 0.1))
\end{aligned}$$

## 3.2 Training procedure

We aim to learn meaningful latent representation which ensures that the prototype vectors a) are close to the input video patches (*fidelity*), b) of different classes are well-separated (*Separability*), and c) are interpretable by humans (*grounded*). We leverage previous works to incorporate appropriate loss functions to enforce these [3, 27]. Furthermore, we also introduce a *diversity* loss term to ensure that prototypes of the same class are non-overlapping. Specifically, we jointly optimize the feature encoder $f$ along with the prototype vectors $\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_m$ in the prototype layer $p$ to minimize the the cross-entropy loss on training set while regularizing for the desiderata.

### 3.2.1 Loss function

Let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be our training dataset, where $\mathbf{x}_i$ is our stacked input extracted from video $\mathbf{v}_i$.

**Full objective.** For hyperparameters $\lambda_c, \lambda_s,$ and $\lambda_d$, our overall objective function that we wish to minimize is given by

$$\mathcal{L}(\mathcal{D}; \theta) = CE(\mathcal{D}; \theta) + \lambda_c R_{clus}(\mathcal{D}; \theta) + \lambda_s R_{sep}(\mathcal{D}; \theta) + \lambda_d R_{div}(\mathcal{D}; \theta) \tag{5}$$

where $CE(\cdot)$, $R_{clus}(\cdot)$, and $R_{sep}(\cdot)$ are the cross-entropy loss, clustering loss, and separation loss, respectively. Here, $\theta$ are the trainable parameters for the feature encoder $f$ and the prototype layer $p$.

The cross-entropy loss here imposes prediction accuracy and is given by

$$CE(\mathcal{D}; \theta) = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K -\mathbb{1}\big[y_i = k\big] \log(\hat{y}_k). \tag{6}$$

The clustering loss $R_{clus}$ minimizes the squared $\ell_2$ distance between some latent patch within a training image and its closest prototype vector from that class, and is given by

$$R_{clus}(\mathcal{D}; \mathbf{p_1}, \mathbf{p_2}, \ldots, \mathbf{p_m}) = \frac{1}{N} \sum_{i=1}^N \min_{\mathbf{p}_j \in \mathbf{P}_{y_i}} \min_{\mathbf{z} \in \text{patches}(\mathbf{x}_i)} \|\mathbf{z} - \mathbf{p}_j\|_2^2, \tag{7}$$

where $\mathbf{P}_{y_i}$ is the set of prototype vectors allocated to the class $y_i$. The separation loss $R_{sep}$ encourages every patch of a training image to stay away from the prototypes *not* of its own class.

$$R_{sep}(\mathcal{D}; \mathbf{p_1}, \mathbf{p_2}, \ldots, \mathbf{p_m}) = -\frac{1}{N} \sum_{i=1}^N \min_{\mathbf{p}_j \notin \mathbf{P}_{y_i}} \min_{\mathbf{z} \in \text{patches}(\mathbf{x}_i)} \|\mathbf{z} - \mathbf{p}_j\|_2^2. \tag{8}$$

Similar loss functions have also been used in [11, 32].

**Diversity loss**: We propose a cosine similarity-based regularization term which penalizes prototype vectors of the same class for overlapping with each other, given by

$$R_{div}(\mathbf{p_1}, \mathbf{p_2}, \ldots, \mathbf{p_m}) = \sum_{k=1}^K \sum_{\substack{i \neq j \\ \mathbf{p}_i, \mathbf{p}_j \in \mathbf{P}_k}} \max(0, \cos(\mathbf{p}_i, \mathbf{p}_j) - s_{max}), \tag{9}$$

where $\cos(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a}^\top \mathbf{b}}{\|\mathbf{a}\| \cdot \|\mathbf{b}\|}$ for vectors $\mathbf{a}$ and $\mathbf{b}$ is the cosine similarity, and $s_{max}$ is a hyperparameter for the maximum similarity allowed. This cosine similarity-based loss considers the angle between the prototype vectors regardless of their length. Hence, it allows us to penalize the similarity between the prototypes up to a threshold, leading to diverse and more expressive representations.

### 3.2.2 Prototype projection and grounding

To achieve grounding, while training we intersperse the following projection step after every few epochs. Specifically, we project the prototype vectors to actual video patches from training videos that contain those dynamics as follows,

$$\mathbf{p}_j \leftarrow \text{argmin}_{\mathbf{z}' \in \text{patches}(f(\mathbf{x}_i))} \|\mathbf{z}' - \mathbf{p}_j\|_2^2 \ \forall i \ \text{s.t.} \ y_i = k, \tag{10}$$

for all prototype vectors of class $k$, i.e. $\mathbf{p}_j \in \mathbf{P}_k$. This step projects each prototype vector of a given class to the closest latent representation of a training video patch that is also from that same class. As a result, the prediction of a test video is made based on the similarities it has with the learned prototype vectors. Consequently, the test prediction is grounded on the training videos.

## 4 Evaluation

We discuss our evaluations of `DPNet` on benchmark deepfake datasets for prediction performance, model interpretation, and robustness to temporal logic specifications.

### 4.1 Experimental settings

Table 1: Basic information for both FaceForensics++ and DeepFakeDetection datasets. [28]

| Dataset | Real | | Fake | | Avg. FPS | Release Date |
| --- | --- | --- | --- | --- | --- | --- |
| | Video | Frame | Video | Frame | | |
| FaceForensics++ (FF++) | 1000 | 509.9k | 1000 | 509.9k | 26.7 | 2019.01 |
| DeepFakeDetection (DFD) | 363 | 315.4k | 3068 | 2242.7k | 24.0 | 2019.09 |

We use two benchmark deepfake datasets: FaceForensics++ (FF++) [38] and the Google/Jigsaw Deepfake Detection dataset (DFD) [14]. Both FF++ and DFD are large-scale and contain manipulated

Figure 3: The reasoning process of the `DPNet`. The prediction for each class is based on the evidence between the dynamics of the input and a small set of dynamic prototypes. Best view as GIFs (See C).



Figure 4: Left column depicts the different classes of temporal artifacts and unnatural movements found by `DPNet` across both DFD and FF++ datasets. Best view as GIFs (See C). The right column shows the effect of diversity regularization (9) on prototype similarity scores across test videos.

and natural images of human faces (Table 1). Both datasets provide ground truth manipulation masks showing what part of the face was manipulated. We use a standard preprocessing technique that extracts frames and crops out facial areas from each video [56]. During `DPNet` training phases, the input is formed by stacking 1 RGB frame followed by 8 pre-computed optical flow fields that are uniformly separated. The input registers a temporal signature of roughly 1s for each dataset, which should be sufficient to capture salient deepfake artifacts. `DPNet` uses a pre-trained ResNet101 as a backbone network. Since our input contains temporal frames, we also perform cross-modality pre-training. The encoder $f$ and prototypes $\mathbf{p}$ were trained with learning rates of $1e^{-4}$ and $1e^{-3}$ respectively. From cross-validation, $\lambda_c, \lambda_s, \lambda_d$ are set to $(0.8, -0.8, 5.0)$. $m_k = 20$, and prototype vectors are randomly initialized. Further details of experiment settings are provided in Appendix A.

### 4.2 Comparison with state-of-the-art methods

In Table 2, we present the performance of `DPNet` on FF++ and DFD, as well as baselines that do not have interpretable interpretations or temporal aspects. For binary classification, we report both the accuracy, and more importantly the AUC. `DPNet` achieves 0.988 AUC on FF++ and 95.00% accuracy,

6

Table 2: Test Accuracy and AUC on DFD and FF++ of `DPNet` and other baselines. Prototype overlap percentage with ground truth manipulation masks are also provided in column 5. The - entries denote numbers not reported in prior work, or non-applicability to the metric.

| Dataset | Method | Binary Accuracy (%) | AUC (%) | Overlap (%) |
|---------|--------|---------------------|---------|-------------|
| FF++ | Two-Stream (Zhou e al. 2017, [58]) | - | $\simeq 70.01$ | - |
| | Multi-Task Learning (Nguyen et al. 2019, [34]) | - | 76.30 | - |
| | MesoNet (Afcha et al. 2018, [1]) | $\simeq 93.00$ | - | - |
| | Capsule (Nguyen et al. 2019, [35]) | 93.11 | 96.60 | - |
| | Xception (Rossler et al. 2019, [38]) | $\simeq \mathbf{97.00}$ | **99.70** | - |
| | ProtoPNet (Chen et al. 2019, [11]) | 92.00 | 96.22 | 83.67 |
| | DPNet, 40p (Ours) | 95.00 | **98.75** | **94.82** |
| | DPNet, 80p (Ours) | **96.25** | 98.20 | 92.99 |
| DFD | ProtoPNet (Chen et al. 2019, [11]) | **92.99** | **95.78** | 64.60 |
| | DPNet, 40p (Ours) | 91.83 | 91.33 | 72.88 |
| | DPNet, 80p (Ours) | 90.37 | 92.22 | **79.29** |

which is on par with existing state-of-the-art of 0.996 and 95.73% AUC achieved by ExceptionNet. Note that ExceptionNet does not provide any form of intrinsic interpretability. Importantly, `DPNet` outperforms ProtoPNet [11] by 0.025 in AUC and 3% in accuracy, while the interpretation of `DPNet`'s prototypes aligns much closer with ground truths. We provide two settings with doubling number of prototypes ($m = 40, 80$) per class. We also note that `DPNet` underperforms ProtoPNet in predictive performance for DFD, but the gain in ground truth overlap percentage indicates that perhaps ProtoPNet utilized background features that correlate with training labels but are not truly discriminating.

**Quantitative interpretation metric.** We measure the interpretations against available ground truth to determine what percentage each of the prototype's prototypical patch overlaps with the ground truth masks. Intuitively, each prototype looks at a prototypical patch within the input region and uses those patches to make a prediction. Hence, we want to see how much each prototype actually looks at our region of interest (the manipulated region). We use a state-of-the-art attribution method, Integrated Gradients [44], to retrieve the pixel importance of each input with respect to the similarity score of a prototype vector $\mathbf{p}_j$. We compute the overlap of the 95th percentile importance map $I_j$, where values less than the threshold are set as 0, and 1 otherwise, with the binary mask of the manipulated region $M$. Formally, each prototype computes a prototype percentage overlap (PPO):
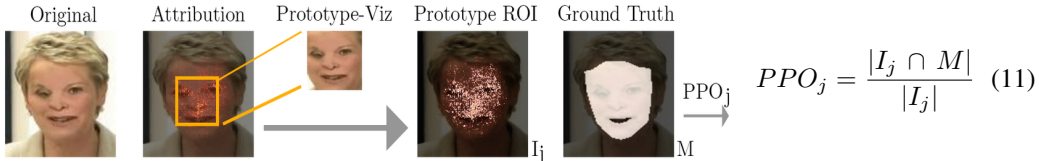


$$PPO_j = \frac{|I_j \cap M|}{|I_j|} \quad (11)$$

Figure 5: Evaluation of the network's explanations over ground truths. The $PPO_j$ score is averaged for the top 10 most similar prototypes across the entire validation set, yielding our final $PPO$ score for each of the models.

**Analysis.** As seen in Table 2, `DPNet` outperforms ProtoPNet by at least 8 points across both datasets. The added temporal information allows `DPNet` to capture important dynamics that cannot be learned in other baseline methods that only process videos frame-by-frame. Not only do the prototypes in our proposed approach focus more on the regions of interest, but the interpretations that we get from these prototypes are more accessible and easier for humans to understand. Figure 4 visualize some of the learned prototypes demonstrating unnatural movements and behaviors captured by `DPNet`.

Moreover, given the case-based reasoning nature of the network, it is easy to understand how each similar prototype vector contributes to the final prediction. The final layer allows users to read off the importance of each prototype to a testing image. Attention models, on the other hand, do not follow case-based reasoning. It is possible to use *posthoc* analysis to get 'prototypes' like ours from convolutional filters, but these are not faithful to the computation. `DPNet`'s prototype vectors actually participate in calculating the similarity scores combined in the final layer, which makes it more understandable. Overall, although there exists some accuracy gap between `DPNet` and state-of-the-arts, the added interpretability over videos instead of images brings richer explanations that are easy for a layperson to understand, especially in this case of deepfakes.

### 4.3 Visualization and interpretation of learned prototypes

We visualize the prototypes by taking the attribution with respect to each prototype similarity score. During the training steps, while projecting, we kept track of the original training clip that is closest in latent space to each prototype. Hence, during testing, we can visualize the prototypes that are most similar to the testing videos.

Figure 3 visually represents how our network makes predictions based on examining the dynamics withing the video and comparing that against the learned dynamics prototype. Figure 4 presents classes of temporal artifacts captured by DPNet within the learned prototypes. One class of artifacts features heavy discoloration Fig.4a, which is already interpretable as images, but changes of discoloration over time added interpretability. Other classes which feature more subtle discoloration or movement Fig. 4b,c are harder to interpret with just one image, especially ones that feature subtle facial jitterings and unnatural oscillations Fig. 4d. Hence, providing explanations that are dynamic overtime stays faithful to the learned dynamic prototypes while increasing interpretability for end users. Combined with the interpretable structure of the network, one can easily understand which temporal artifacts are under scrutiny within an input deepfake video.

### 4.4 Temporal specifications over learned prototypes

We checked the robustness of DPNet and interpretable baselines against our temporal specifications specified in Table 3. Overall, both approaches satisfy the $\phi_{1,key\_frame}$ specification up to high-percentage, with DPNet performing slightly better, especially with the $fake$ traces. This indicates that with high-probability, we can find a discriminative prototype within the video that is most relevant in explaining why a video is real or fake. On the other hand, the models perform poorly with the stringent specification $\phi_{2,non\_relevance}$, which requires non-relevant prototypes to both stay low and not change at all over time. Since DPNet utilized temporal information via optical flows, this is a stricter specification to enforce as flows can change drastically

Table 3: Percentage of traces satisfying temporal specifications. Rows in each block represent the percentage over positive, negative, and all traces.

|  | $\phi_1$ | $\phi_2$ | $\phi_3$ |
|---|---|---|---|
| ProtoPNet | 95.23 | 49.73 | 70.89 |
|  | 89.09 | 42.18 | 58.76 |
|  | 92.00 | 45.75 | 64.50 |
| DPNet, 40p | 93.65 | 28.04 | 74.07 |
|  | 91.46 | 16.11 | 56.39 |
|  | 92.50 | 21.75 | 64.75 |

across consecutive frames, hence the lower percentage of satisfying traces. We experiment by relaxing the consecutive "next 5 frames" non-changing constraint, $\phi_{3,non\_relevant,relaxed}$, which now only enforce non-relevant prototype similarities to be low. The flexibility of temporal logic along with interpretable model allows end-users to specify and enforce certain desiderata in their detection framework. This further increases fairness, interpretability, and utilities of these frameworks.

## 5 Discussion and Conclusion

**Summary** We introduce DPNet – an *interpretable* deepfake detection technique which leverages (and gives insight into) the temporal dynamics to determine if a given video is real or a deepfake. In addition to the prototype-based interpretations, we draw motivations from the temporal logic specifications in formal methods to design robustness metrics to analyze the decision-making process of the underlying model. Our prototype layer along with the specifications can be used with other deep fake detection models to interpret them and can be of independent interest.

**Limitations and Future Work** Incorporating temporal logic specifications provides a rich avenue for future work on interpretability. To this end, it will be advantageous to learn the specifications instead of designing them. This, however, is non-trivial and remains an open challenge in the literature as prior work revolves around rule-mining via reducing to the boolean satisfiability problem [9]. As a result, the present work shows one way to add logic-based interpretability via formal methods.

**Conclusion** The primary goal of transparency is to provide users with the tools and context to enable them to make an informed decision. Another essential piece is to build trust into the decisions. To this end, this present work aims to explain, in a human interpretable form and by logic specifications, as to why and a determination was made, providing a safeguard against the misinformation spread via deepfakes.

## Broader Impact

There is a growing need to use automatic deepfake detection models to detect and combat deepfakes. However, a reliance on deepfake detectors places enormous trust on these models. This work aims to help justify this trust by improving the interpretability of deepfake detection. In addition to model interpretability, this work offers insights into what parts of deepfake videos can be used to discern deepfakes, which may inform people how to detect deepfakes themselves. The risk of this work is that its insights may ultimately be used to improve deepfake generators. Although improvements in deepfake detection and generation may become a vicious cycle, this should not hinder research on explaining deepfake detectors. For example, are already discussions on using deepfake detectors to protect the videos of world leaders [2]. Model interpretations strengthen the accountability of deepfake detectors. Our work takes a step towards exposing this problem and encourages future research in explaining deepfake detectors.

## References

[1] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: a compact facial video forgery detection network. *CoRR*, abs/1809.00888, 2018.

[2] Shruti Agarwal, Hany Farid, Yuming Gu, Mingming He, Koki Nagano, and Hao Li. Protecting world leaders against deep fakes. In *CVPR Workshops*, 2019.

[3] David Alvarez-Melis and Tommi S. Jaakkola. Towards robust interpretability with self-explaining neural networks. *CoRR*, abs/1806.07538, 2018.

[4] Irene Amerini, Leonardo Galteri, Roberto Caldelli, and Alberto Bimbo. Deepfake video detection through optical flow based cnn. In *CVPR Workshops*, pages 1205–1207, 10 2019.

[5] 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, 07 2015.

[6] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L. Morency. Openface 2.0: Facial behavior analysis toolkit. In *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*, pages 59–66, 2018.

[7] Patricia Bouyer. Model-checking timed temporal logics. *Electronic Notes in Theoretical Computer Science*, 231:323 – 341, 2009. Proceedings of the 5th Workshop on Methods for Modalities (M4M5 2007).

[8] Sarah Cahlan. How misinformation helped spark an attempted coup in gabon. *The Washington Post*, 2020.

[9] Alberto Camacho and Sheila A McIlraith. Learning interpretable models expressed in linear temporal logic. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 29, pages 621–630, 2019.

[10] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017.

[11] Chaofan Chen, Oscar Li, Alina Barnett, Jonathan Su, and Cynthia Rudin. This looks like that: deep learning for interpretable image recognition. *CoRR*, abs/1806.10574, 2018.

[12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.

[13] Adel Dokhanchi, Hani Ben Amor, Jyotirmoy V. Deshmukh, and Georgios Fainekos. Evaluating perception systems for autonomous vehicles using quality temporal logic. In Martin Leucker and Christian Colombo, editors, *Runtime Verification- 18th International Conference, RV 2018, Proceedings*, pages 409–416. Springer Verlag, January 2019.

[14] Nicholas Dufour, Andrew Gully, Per Karlsson, Alexey Victor Vorbyov, Thomas Leung, Jeremiah Childs, and Christoph Bregler. Deepfakes detection dataset by google & jigsaw.

[15] Matthew B. Dwyer, George S. Avrunin, and James C. Corbett. Property specification patterns for finite-state verification. In *Proceedings of the Second Workshop on Formal Methods in Software Practice*, page 7–15, New York, NY, USA, 1998. Association for Computing Machinery.

[16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. *CoRR*, abs/1812.03982, 2018.

[17] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2015.

[18] Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Qin, Jonathan Uesato, Relja Arandjelovic, Timothy A. Mann, and Pushmeet Kohli. On the effectiveness of interval bound propagation for training verifiably robust models. *ArXiv*, abs/1810.12715, 2018.

[19] Karen Hao. An ai app that "undressed" women shows how deepfakes harm the most vulnerable. *MIT Technology Review*, 2019.

[20] Xiaowei Huang, Marta Z. Kwiatkowska, Sen Wang, and Min Wu. Safety verification of deep neural networks. In *CAV*, 2017.

[21] David Ingram. A face-swapping app takes off in china, making ai-powered deepfakes for everyone. *NBC*, 2019.

[22] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. *ArXiv*, abs/1702.01135, 2017.

[23] Been Kim, Cynthia Rudin, and Julie A Shah. The bayesian case model: A generative approach for case-based reasoning and prototype classification. In *Advances in Neural Information Processing Systems*, pages 1952–1960, 2014.

[24] Pavel Korshunov and Sébastien Marcel. Deepfakes: a new threat to face recognition? assessment and detection. *ArXiv*, abs/1812.08685, 2018.

[25] Mathias Lécuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. *2019 IEEE Symposium on Security and Privacy (SP)*, pages 656–672, 2018.

[26] Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. *CoRR*, abs/1710.04806, 2017.

[27] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. *CoRR*, abs/1811.00656, 2018.

[28] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. *arXiv: Cryptography and Security*, 2019.

[29] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019.

[30] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *ArXiv*, abs/1706.06083, 2017.

[31] F. Matern, C. Riess, and M. Stamminger. Exploiting visual artifacts to expose deepfakes and face manipulations. In *2019 IEEE Winter Applications of Computer Vision Workshops (WACVW)*, pages 83–92, 2019.

[32] Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. Interpretable and steerable sequence learning via prototypes. *CoRR*, abs/1907.09728, 2019.

[33] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *CoRR*, abs/1706.07979, 2017.

[34] Huy H. Nguyen, Fuming Fang, Junichi Yamagishi, and Isao Echizen. Multi-task learning for detecting and segmenting manipulated facial images and videos. *CoRR*, abs/1906.06876, 2019.

[35] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using capsule networks to detect forged images and videos. *CoRR*, abs/1810.11215, 2018.

[36] Oded Padon, Jochen Hoenicke, Giuliano Losa, Andreas Podelski, Mooly Sagiv, and Sharon Shoham. Reducing liveness to safety in first-order logic. *Proc. ACM Program. Lang.*, 2(POPL), December 2017.

[37] Ivan Petrov, Daiheng Gao, Nikolay Chervoniy, Kunlin Liu, Sugasa Marangonda, Chris Um'e, Mr. Dpfks, RP Luis, Jian Jiang, Sheng Zhang, Pingyu Wu, Bo Zhou, and Weiming Zhang. Deepfacelab: A simple, flexible and extensible face swapping framework. *ArXiv*, abs/2005.05535, 2020.

[38] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. *CoRR*, abs/1901.08971, 2019.

[39] Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. Recurrent convolutional strategies for face manipulation detection in videos. *CoRR*, abs/1905.00582, 2019.

[40] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. Grad-cam: Why did you say that? visual explanations from deep networks via gradient-based localization. *CoRR*, abs/1610.02391, 2016.

[41] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. *CoRR*, abs/1704.02685, 2017.

[42] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *preprint*, 12 2013.

[43] A Prasad Sistla. Safety, liveness and fairness in temporal logic. *Formal Aspects of Computing*, 6(5):495–511, 1994.

[44] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *CoRR*, abs/1703.01365, 2017.

[45] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *CoRR*, abs/1312.6199, 2014.

[46] Rob Toews. Deepfakes are going to wreak havoc on society. we are not prepared. *Forbes*, 2020.

[47] Rubén Tolosana, Rubén Vera-Rodríguez, Julian Fiérrez, Aythami Morales, and Javier Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *ArXiv*, abs/2001.00179, 2020.

[48] William Turton and Andrew Martin. How deepfakes make disinformation more real than ever. *Bloomberg News*, 2020.

[49] Cristian Vaccari and Andrew Chadwick. Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Social Media + Society*, 6(1):2056305120903408, 2020.

[50] Marcell Vazquez-Chanlatte, Shromona Ghosh, Jyotirmoy V. Deshmukh, Alberto L. Sangiovanni-Vincentelli, and Sanjit A. Seshia. Time series learning using monotonic logical properties. *CoRR*, abs/1802.08924, 2018.

[51] 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. *CoRR*, abs/1608.00859, 2016.

[52] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang, and Suman Jana. Efficient formal safety analysis of neural networks. In *Advances in Neural Information Processing Systems*, pages 6367–6377, 2018.

[53] Lily Weng, Pin-Yu Chen, Lam Nguyen, Mark Squillante, Akhilan Boopathy, Ivan Oseledets, and Luca Daniel. PROVEN: Verifying robustness of neural networks with a probabilistic approach. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6727–6736, Long Beach, California, USA, 09–15 Jun 2019. PMLR.

[54] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses. *CoRR*, abs/1811.00661, 2018.

[55] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. *CoRR*, abs/1311.2901, 2013.

[56] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503, 2016.

[57] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. *European Conference on Computer Vision*, 2018.

[58] Peng Zhou, Xintong Han, Vlad I. Morariu, and Larry S. Davis. Two-stream neural networks for tampered face detection. *CoRR*, abs/1803.11276, 2018.

[59] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2242–2251, 2017.

# Appendix

## A   Experimental Settings

In this section, we discuss the experimental settings for our evaluations of `DPNet` on large-scale deepfake datasets such as FaceForensics++ and Google's DeepFakeDetection.

**Datasets.** The FaceForensics++ (FF++) dataset [38] is a subset of deepfake videos from the Technical University of Munich and University Federico II of Naples' FaceForensics benchmark consisting of 1000 original video sequences sourced from `YouTube` and 1000 synthetic videos generated using `faceswap`. All videos in FF++ mostly contain a trackable frontal face and manipulation masks are also provided for synthetically generated videos.

The Google/Jigsaw DeepFake detection (DFD) [14] dataset consists of 3,068 deepfake videos generated based on 363 original videos of 28 consented individuals of various ages, genders, and ethnic groups. Pairs of actors were selected randomly and deep neural networks swapped the face of one actor onto the head of another. The specific deep neural network synthesis model is not disclosed, but manipulation masks are also provided in this dataset. The statistics for these datasets are provided in the paper and below.

| Dataset | Real | | Fake | | Avg. FPS | Release Date |
|---|---|---|---|---|---|---|
| | Video | Frame | Video | Frame | | |
| FaceForensics++ | 1000 | 509.9k | 1000 | 509.9k | 26.7 | 2019.01 |
| DeepFakeDetection | 363 | 315.4k | 3068 | 2242.7k | 24.0 | 2019.09 |

Table 4: Basic information for both FaceForensics++ and DeepFakeDetection datasets. [28]

**Pre-processing.** We use the standard data pre-processing techniques for deepfakes and videos, i.e. separating the frames from videos and cropping out facial areas using a face detection algorithm MCTNN [56] instead of using the full-frame. When cropping faces, we use a face margin of 0.5-0.7 to get a full cropping of the head, and also use the detected boundaries to crop the mask videos. Consecutive cropped frames are processed in OpenCV to calculate the dense optical flows.

**Training details.** During training, each input is formed by stacking 1 RGB frame with 8 pre-computed optical flow fields starting from that RGB frame, stacking every 3 frames apart. Given the FPS of the two datasets, FF++ and DFD, this registers a temporal signature of roughly 1.0s. The initial input frame is randomly selected from the video each iteration. The selection process gives our input dataloader a pseudo-augmentation aspect so the network does not overfit to any particular segment of a video, but besides this, we do not use any other form of video data augmentations.

We used a pre-trained ResNet101 as the backbone of our architecture. However, given the shape of our input, we also need to perform c*ross modality pre-training* to initialize the weights of the first convolutional layer, i.e we average the weights across the RGB channels and replicate this average by the channel number of temporal network input [51]. Our network parameters $\theta$ are trained using Adam, with a learning rate of $1e^{-4}$ for $f_\omega$ and $1e^{-3}$ for $\mathbf{p}_i$. We use a small number of prototypes to represent facial features, hence we chose $m_k = 20$. The learning rate decays by a factor of 0.5 every 10 epochs and the networks were trained for 150 epochs. We chosen our lambdas during cross-validation and set $\lambda_c, \lambda_s, \lambda_d$ to be (0.8, -0.8, 5.0). The values $\lambda_c, \lambda_s$ are chosen from (.4, -.4), (.6, -.6), (.8, -.8) and (1,-1). The value of $\lambda_d$ is chosen from (1, 10).

**Validation details.** During validation, our model can perform frame-wise evaluation of the frames in a video and aggregate the results. We follow the testing scheme of previous video detection work [51] and sample 25 RGB frames and optical flows stack evenly spaced across deepfake videos. We combine the logits using an aggregation function (`sum` or `avg`). This can be understood as combining the *similarity evidence* between the learned prototypes and the testing video across clips in the video.

## B   TQTL Semantics

For completeness, we detailed the semantics for TQTL for evaluating a specification, similar to the work by Dokhanchi et al. [13]. Consider the data stream $\mathcal{D}, i \in N$ is the index of current frame, $\pi \in P, \phi, \phi_1, \phi_2 \in TQTL$ and evaluation function $\epsilon : V_t \cup V_p \rightarrow \mathbb{N}$, which is the environment over

the time and prototype variables. The quality value of formula $\phi$ with respect to $\mathcal{D}$ at frame i with evaluation $\epsilon$ is recursively assigned as follows:

$$\llbracket \top \rrbracket (\mathcal{D}, i, \epsilon) := +\infty$$

$$\llbracket \pi \rrbracket (\mathcal{D}, i, \epsilon) := \llbracket f_\pi(j_1, \ldots, j_n, id_1, \ldots, id_n) \sim c \rrbracket (\mathcal{D}, i, \epsilon)$$

$$\llbracket x.\phi \rrbracket (\mathcal{D}, i, \epsilon) := \llbracket \phi \rrbracket (\mathcal{D}, i, \epsilon[x \Leftarrow i])$$

$$\llbracket \exists id@x, \phi \rrbracket (\mathcal{D}, i, \epsilon) := \max_{k \in \mathcal{P}} \left( \llbracket \phi \rrbracket (\mathcal{D}, i, \epsilon[id \Leftarrow k]) \right)$$

$$\llbracket x \leq y + n \rrbracket (\mathcal{D}, i, \epsilon) := \begin{cases} +\infty & if \ \epsilon(x) \leq \epsilon(y) + n \\ -\infty & otherwise \end{cases}$$

$$\llbracket \neg \phi \rrbracket (\mathcal{D}, i, \epsilon) := -\llbracket \phi \rrbracket (\mathcal{D}, i, \epsilon)$$

$$\llbracket \phi_1 \vee \phi_2 \rrbracket (\mathcal{D}, i, \epsilon) := \max \left( \llbracket \phi_1 \rrbracket (\mathcal{D}, i, \epsilon), \llbracket \phi_2 \rrbracket (\mathcal{D}, i, \epsilon) \right)$$

$$\llbracket \phi_1 U \phi_2 \rrbracket (\mathcal{D}, i, \epsilon) := \max_{i \leq j} \left( \min \left( \llbracket \phi_2 \rrbracket (\mathcal{D}, j, \epsilon), \min_{i \leq k < j} \llbracket \phi_1 \rrbracket (\mathcal{D}, k, \epsilon) \right) \right)$$

We say that $D$ satisfies $\phi$ $(\mathcal{D} \models \phi)$ iff $\llbracket \phi \rrbracket (\mathcal{D}, 0, \epsilon_0) > 0$, where $\epsilon_o$ is the initial environment. On the other hand, a data stream $\mathcal{D}'$ does not satisfy a TQTL formula $\phi$ $(\mathcal{D}' \not\models \phi)$, iff $\llbracket \phi \rrbracket (\mathcal{D}, 0, \epsilon_0) \leq 0$. The quantifier $\exists id@x$ is the maximum operation on the quality values of formula $\llbracket \phi \rrbracket$ corresponding to the prototypes IDs at frame $x$.

## C   GIFs & Code

For the following visualizations and figures in the main paper, we have included the corresponding gifs in folders within the supplementary materials. The authors are committed tno publicly release code as open source for reproducibility for the camera-ready deadline.

## D   Additional Visualizations

### D.1   More examples of prototypes and classes of temporal artifacts learned



Figure 6: Different classes of temporal artifacts and unnatural movements found by DPNet across both DFD and FF++ datasets. Top block are fake prototypes and bottom are real prototypes. a) heavy discolouration, b) subtle discolouration, c) subtle disappearance, d) unnatural movement, e) combined eye-mouth movement, and f) head movement. Best view as GIFs.

14

## D.2    More examples of how `DPNet` classify a video

(Test Clip)(Test Prototype Attr.)          (Prototype Similarity)     (Train Prototype Attr.Clip)(Train Clip)

Top 1 activated fake prototype for this image: 6 (1.853004813194275)

(a)

Top 2 activated fake prototype for this image: 10 (1.0717833042144775)

Top 3 activated fake prototype for this image: 8 (0.8912854194641113)

Cumulative score: 8.131853103637695

Top 1 activated real prototype for this image: 38 (0.003509079571813345)

Top 2 activated real prototype for this image: 33 (0.00348912226036191)

Top 3 activated real prototype for this image: 30 (0.0034866277128458023)

Cumulative score: 0.06459416449069977

Top 1 activated fake prototype for this image: 6 (1.9246346950531006)

(b)

Top 2 activated fake prototype for this image: 8 (1.8038361072540283)

Top 3 activated fake prototype for this image: 16 (1.1743850708007812)

Cumulative score: 9.638761520385742



15

(Test Clip) (Test Prototype Attr.)     (Prototype Similarity)   (Train Prototype Attr.Clip) (Train Clip)

Figure 7: More examples of how DPNet classify a video. (a) and (b) are deepfakes, and (c) is genuine. Best view as GIFs. The prediction for each class is based on the evidence between the dynamics of the input and a small set of dynamic prototypes.