

---

# Supplementary Materials: Improving Deep Learning Interpretability by Saliency Guided Training

---

Aya Abdelsalam Ismail, Héctor Corrada Bravo\*, Soheil Feizi\*  
{asalam, sfeizi}@cs.umd.edu, corradah@gene.com  
Department of Computer Science, University of Maryland

## Experiments

**Computational resources:** All experiments were conducted on a single 12GB NVIDIA RTX2080Ti GPU. **Saliency Methods:** Captum [7] implementation was used for different saliency methods.

### Saliency Guided Training for Images

#### Datasets and Classifiers

- **MNIST** [10]: a database of handwritten digits. The classifier consists of two CNN layers with kernel size 3 and stride of 1 followed by two fully connected layers, two dropout layers with  $p = 0.25$  and  $p = 0.5$ , and the 10 output neurons.
- **CIFAR10** [8]: a low-resolution classification dataset with 10 different classes representing airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks. ResNet18 [3] was used as a classifier, ResNet18 is a deep CNN with “identity shortcut connection,” i.e., skip connections, that skip one or more layers to solve the vanishing gradient problem faced by deep networks.
- **BIRD** [2]: A kaggle datasets of 260 bird species. Images were gathered from internet searches by species name. VGG16 [12] was used as a classifier, the last few dense layers and the output layer were modified to accommodate the number of classes in this dataset.

Dataset	# Training	# Testing	# Classes	Features	Test Accuracy		$\lambda$	$k$ (as a % of feature)
					Traditional	Sal. Guided		
MNIST	60000	10000	10	$1 \times 28 \times 28$	99.4	99.3	1	50%
CIFAR10	50000	10000	10	$3 \times 32 \times 32$	92.0	91.5	1	50%
BIRD	38518	1350	260	$3 \times 224 \times 224$	96.6	96.9	1	50%

Table 1: Datasets used for Image experiments.  $k$  is the percentage of overall features masked during saliency guided training. For example, in MNIST number of features masked  $\lceil 0.5 \times 28 \times 28 \rceil = 392$ .

**Masking** For images, low salient features are replaced by a random variable within the color channel input range. For example, in an RGB image, if pixel  $2 \times 3$  is to be masked  $1 \times 2 \times 3$  would be replaced with a random variable within R channel range, similarly  $2 \times 2 \times 3$  and  $3 \times 2 \times 3$  would be replaced with a random variable within G and B channel range respectively.

#### Saliency Map Quality for Images

The examples shown in Figure 1, Figure 2, and Figure 3 were correctly classified by both models. Gradients are scaled per sample to have values between -1 and 1. Overall, saliency maps produced

---

\* Authors contributed equally

by saliency guided training are less noisy than those produced by traditional training and tend to highlight the object itself rather than the background. The distributions of gradient values per sample show that most features have small values (near zero) with a higher separation of high saliency features away from zero for saliency guided training.

### Model Accuracy Drop

We compare interpretable and traditional training for different saliency methods with modification-based evaluation. Each experiment is repeated five times. Figure 4 shows the mean and standard error for model degradation on different gradient-based methods.

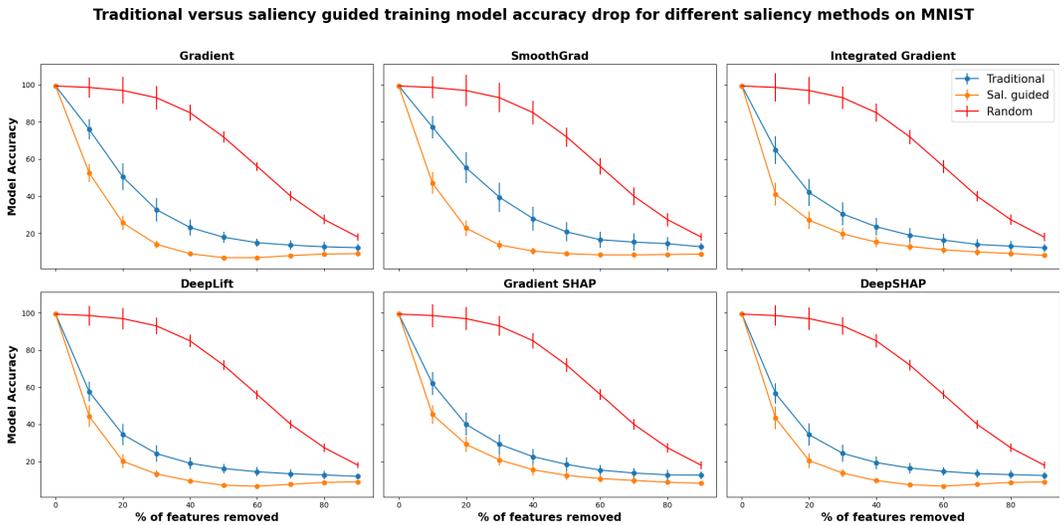


Figure 4: The mean and standard error for model accuracy drop when removing features with high saliency using traditional and saliency guided training for different gradient-based methods against a random baseline. A steeper drop indicates better performance. We find that regardless of the saliency method used, the performance improves by saliency guided training.

### Fine-tuning with Saliency guided Training

We investigate the effect of training traditionally and fine-tuning with saliency guided training. This would be particularly useful for large datasets like imagenet. Table 2 shows the area under accuracy drop curve (AUC) on MNIST Figure 4 for gradient when training traditionally, training using saliency guided procedure and fine-tuning (smaller AUC indicates better performance). We find that fine-tuning improves the performance over traditionally trained networks.

Training Procedures	AUC
Traditional	3360.4
Saliency Guided	1817.6
Fine-tuned	2258.8

Table 2: Area under accuracy drop curve on MNIST for different training procedures

Note that, there is not much gain in training performance when training from scratch versus fine-tuning for small datasets like MNIST. However, for larger datasets like CIFAR10, we observed a clear decrease in the number of epochs when fine-tuning the network. The number of epochs for traditional training CIFAR10 is on average 118, saliency training is 124 while fine-tuning takes only 70 epochs.

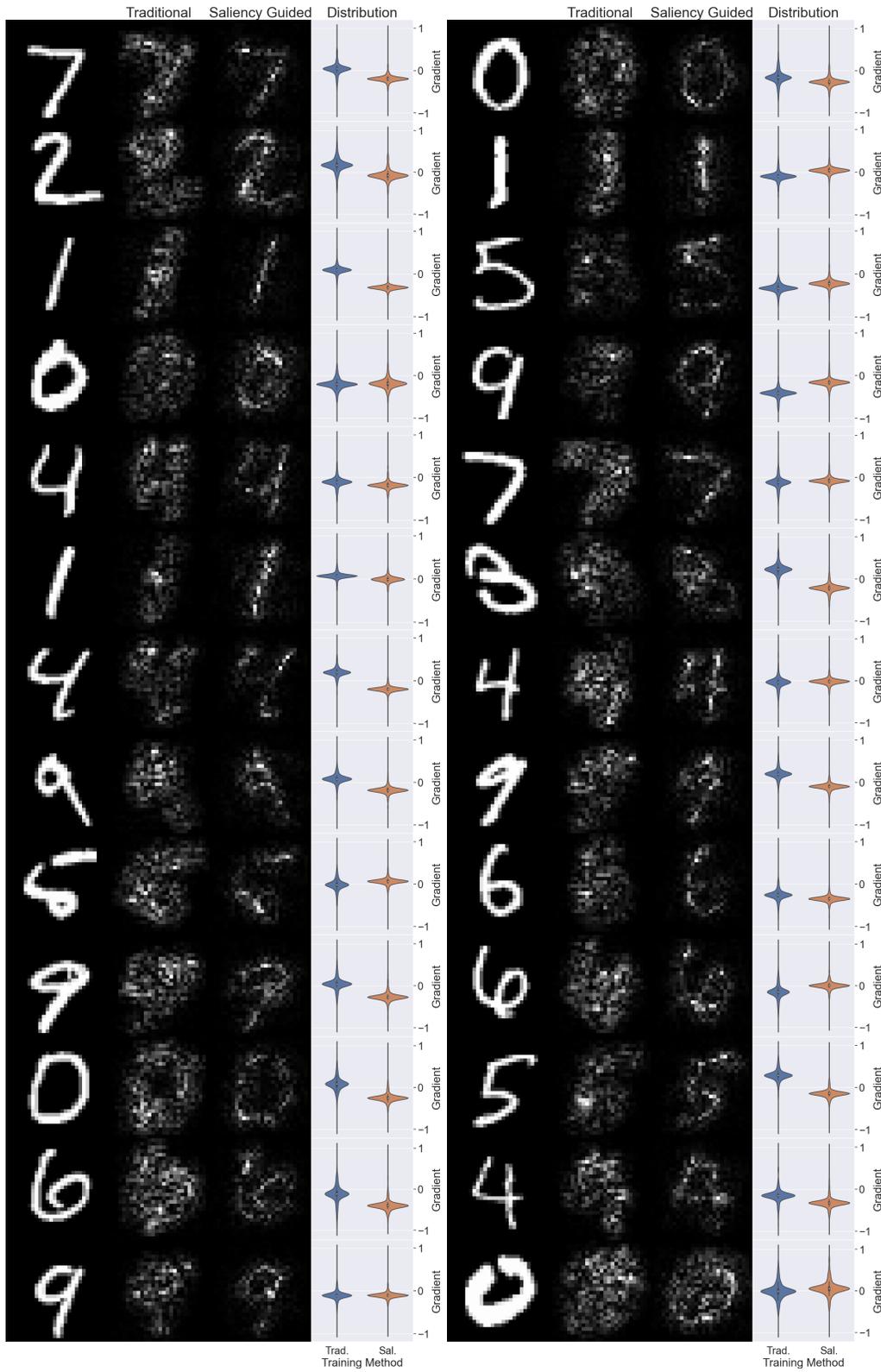


Figure 1: Saliency maps and saliency distribution for Traditional and Saliency Guided Training on MNIST

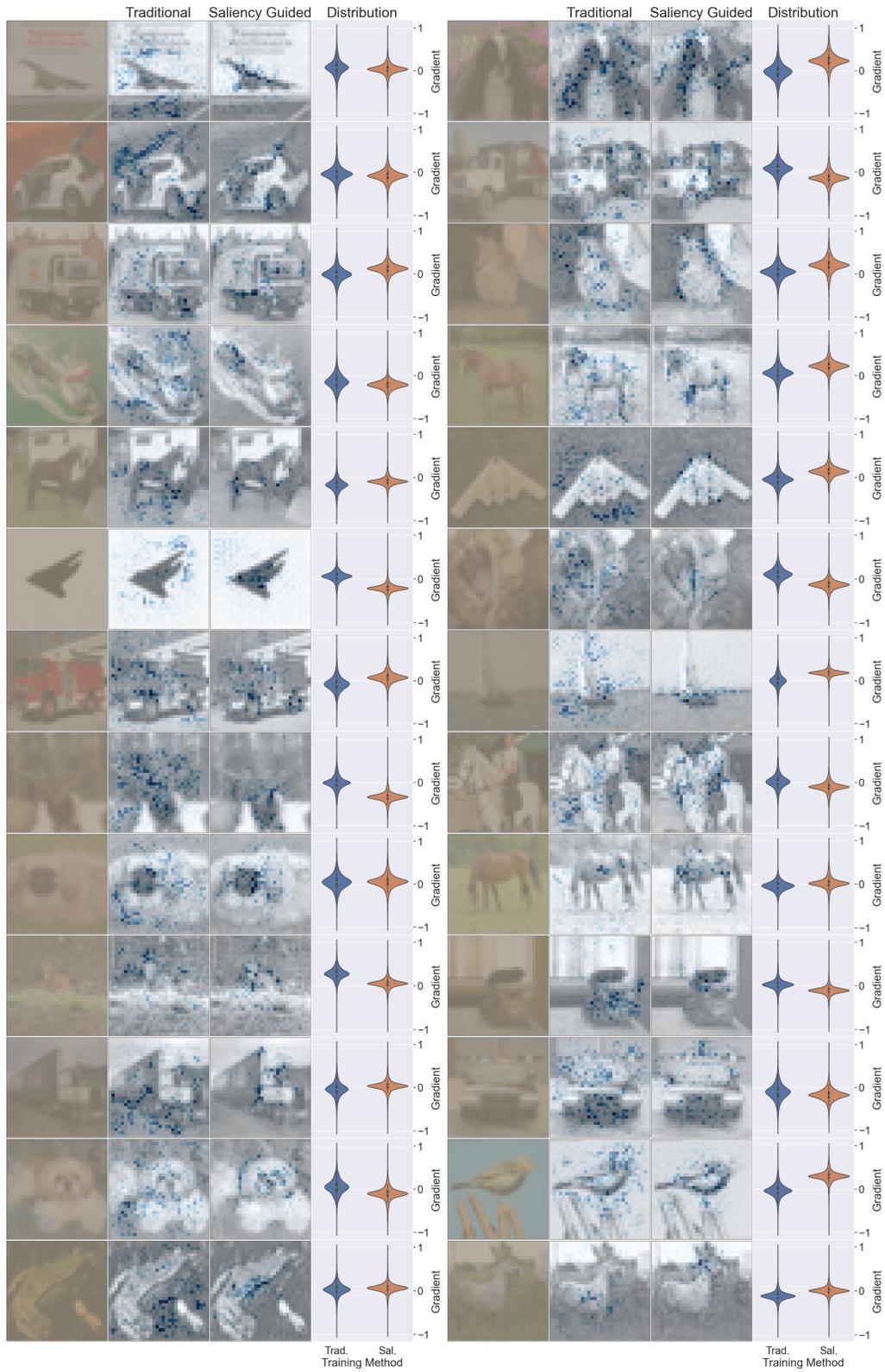


Figure 2: Saliency maps and saliency distribution for Traditional and Saliency Guided Training on CIFAR10

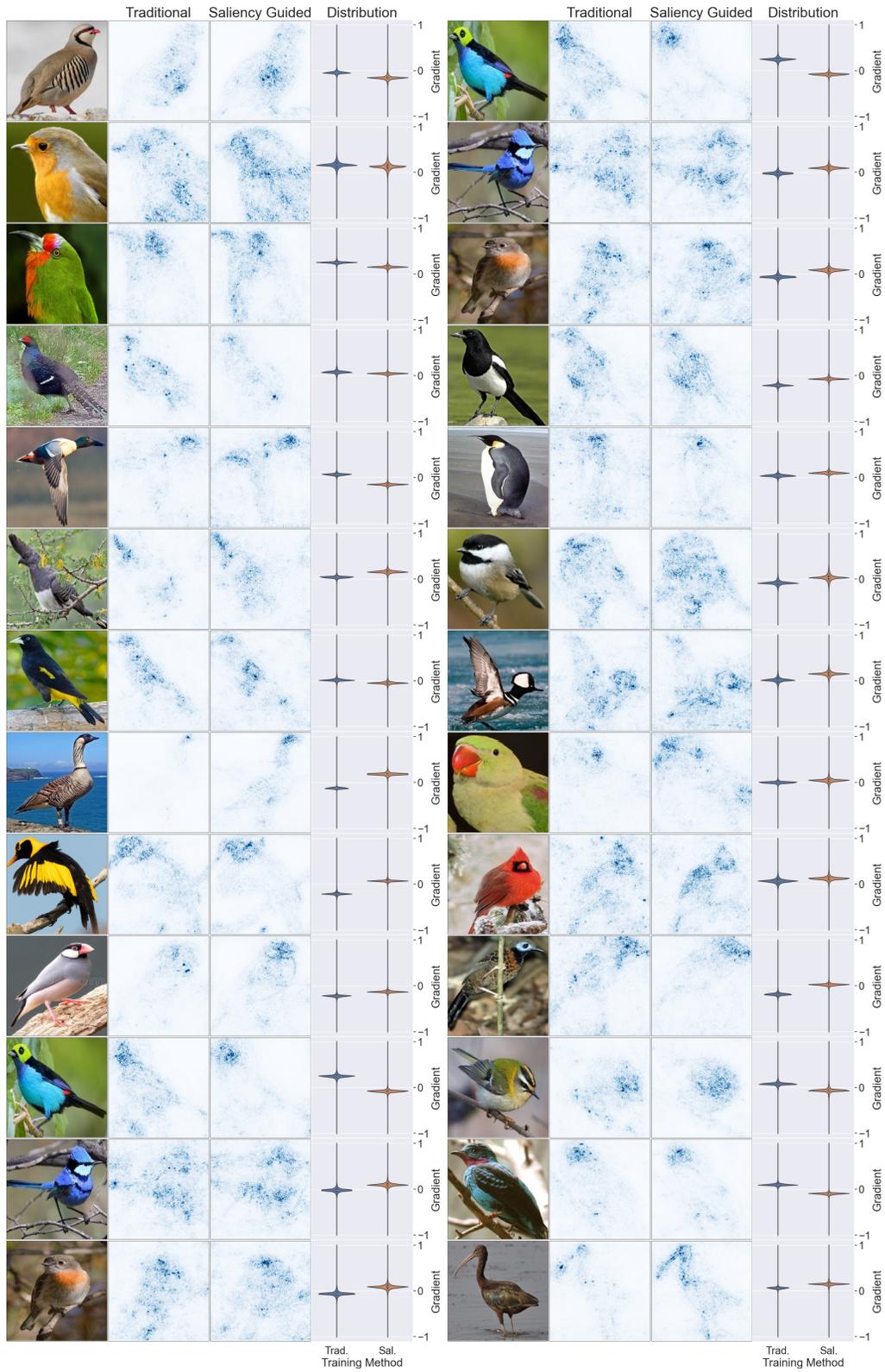


Figure 3: Saliency maps and saliency distribution for Traditional and Saliency Guided Training on BIRD

## Saliency Guided Training for Language

We compare the interpretability of different models trained on language tasks using the ERASER [1] benchmark.

**Datasets** For all datasets, words embedding were generated from Glove [11] and a bidirectional LSTM [4] was used for classifications. Details about each dataset is available in Table 3

Dataset	# Training	# Testing	# Classes	Tokens	Sentences	Test Accuracy		$\lambda$	$k$ (as a % of tokens)
						Traditional	Sal. Guided		
Movie Review	1600	200	2	774	36.8	0.8890	0.8980	1	60%
FEVER	97957	6111	2	327	12.1	0.7234	0.7255	1	80%
e-SNLI	911928	16429	3	16	1.7	0.9026	0.9068	1	70%

Table 3: Overview of datasets in the ERASER benchmark. Number of labels, dataset size, and average numbers of sentences and tokens in each document.  $k$  is the percentage of overall tokens within a particular document.

**Masking** For language tasks, masking is a bit more tricky. We tried multiple masking function, including:

- **Removing** the masking function creates new input such that  $\tilde{X}$  contains only high salient word from the original input  $X$ .
- **Replace with token “[UNK]”** the masking function replaces the low salient word with the token “[UNK]” i.e., unknown.
- **Replace with token “[SEP]”** the masking function replaces the low salient word with the token “[SEP]” i.e., white space.
- **Replace with random word** the masking function replaces the low salient with a random word from vocabulary.
- **Replace with last high salient word** the masking function replaces the low salient word with the previous high salient word.

Over the three datasets, we found that the last masking function (replace with last high salient word) gave the best results. We believe that the masking function can also be dataset-dependent. This particular experiment aims to prove that saliency guided training improves interpretability on language tasks. We will consider finding the optimal masking function for different language tasks in our future work.

**Metrics** ERASER provides two metrics to measure interpretability. *Comprehensiveness* evaluates if all features needed to make a prediction are selected. To calculate an explanation comprehensiveness, a new input  $\bar{X}_i$  is created such that  $\bar{X}_i = X_i - R_i$  where  $R_i$  is predicted rationales. Let  $f_\theta(X_i)_j$  be the prediction of model for class  $j$ . The model comprehensiveness is calculated as:

$$\text{Comprehensiveness} = f_\theta(X_i)_j - f_\theta(\bar{X}_i)_j$$

A high score here implies that the explanation removed was influential in the predictions. The second metric is *Sufficiency* that evaluates if the extracted explanations contain enough signal to make a prediction. The following equation gives the explanation sufficiency:

$$\text{Sufficiency} = f_\theta(X_i)_j - f_\theta(R_i)_j$$

A lower score implies that the explanations are adequate for a model prediction.

To evaluate the faithfulness of continuous importance scores assigned to tokens by models, the soft score over features provided by the model is converted into discrete rationales  $R_i$  by taking the top- $k_d$  values, where  $k_d$  is a threshold for dataset  $d$ . Denoting the tokens up to and including bin  $k$ , for instance,  $i$  by  $R_{ik}$ , an aggregate *comprehensiveness* measure is defined as:

$$\frac{1}{|\mathcal{B}| + 1} \left( \sum_{k=0}^{|\mathcal{B}|} f_\theta(X_i)_j - f_\theta(\bar{X}_{ik})_j \right)$$

*Sufficiency* is defined similarly. Here tokens are grouped into  $k = 5$  bins by grouping them into the top 1%, 5%, 10%, 20% and 50% of tokens, with respect to the corresponding importance score. This metric is referred to as Area Over the Perturbation Curve (AOPC). For reference, we report these when random scores are assigned to tokens. Results are shown in the main paper Table 1.

### Saliency Guided Training for Time Series

We evaluated saliency guided training on a multivariate time series, both quality on multivariate time series MNIST and quantitatively through synthetic data.

### Saliency Maps Quality for Multivariate Time Series

We compare the saliency maps produced on MNIST treated as a multivariate time series with 28 time steps each having 28 features. Figure 5, Figure 6, and Figure 7 shows the saliency maps produced by different saliency methods for Temporal Convolutional Network (TCN), LSTM with Input-Cell Attention and, Transformers respectively. There is a visible improvement in saliency quality across different networks when saliency guided training is used. The most significant improvement was found in TCNs.

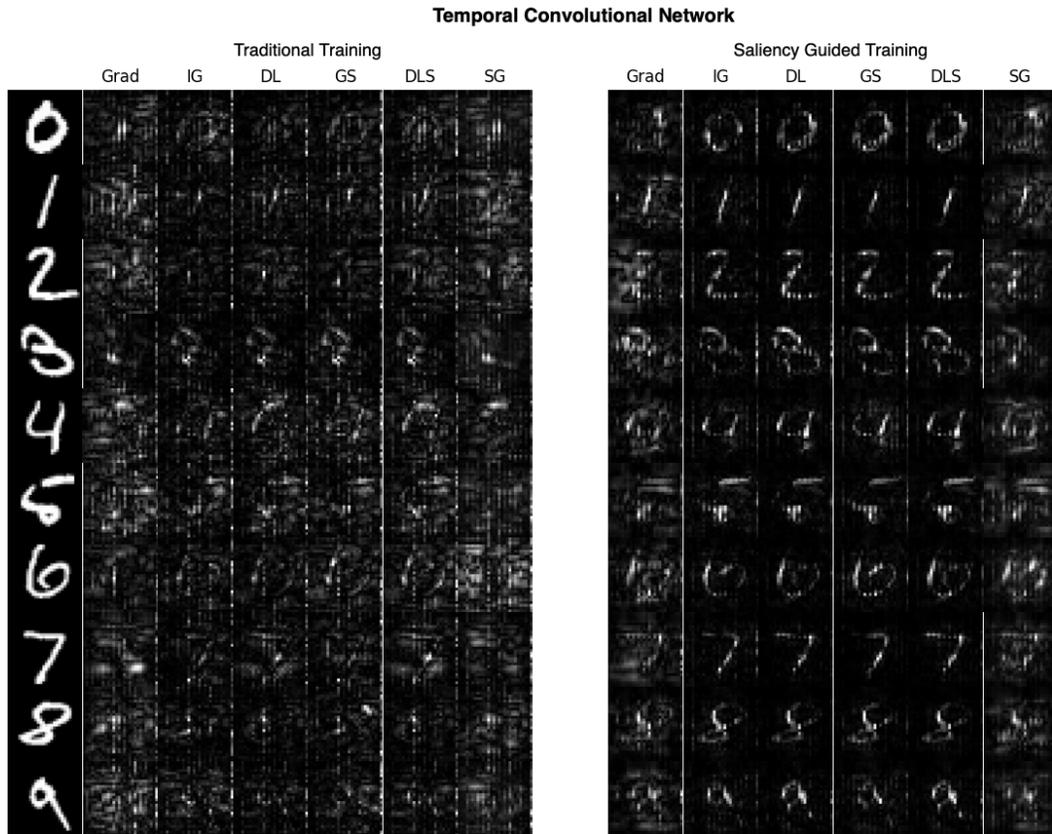


Figure 5: Saliency maps produced for (TCN, saliency method) pairs.

### Quantitative Analysis on Synthetic Data

We evaluated saliency guided training on a multivariate time series benchmark proposed by Ismail et al. [6]. The benchmark consists of 10 synthetic datasets, each examining different design aspects in typical time series datasets. Properties of each dataset is shown in Figure 8. Informative features are highlighted by the addition of a constant  $\mu$  to the positive class and subtraction of  $\mu$  from the negative class. For the following experiments  $\mu = 1$ . Details of each dataset is available in table 4.

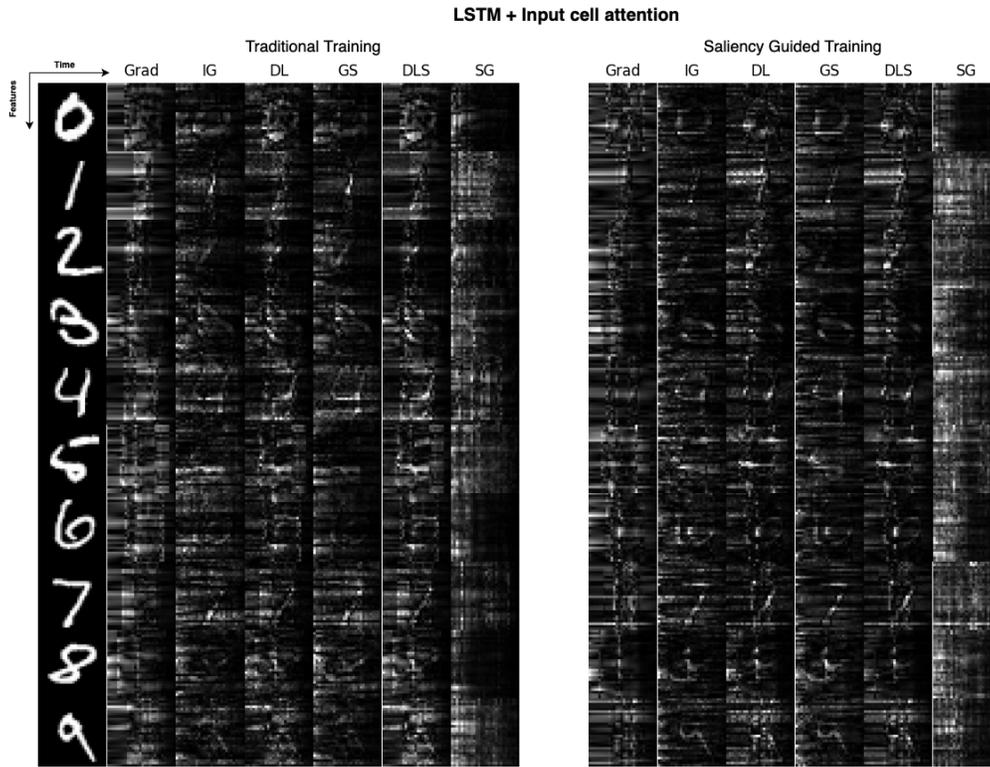


Figure 6: Saliency maps produced for (*LSTM with Input-Cell Attention, saliency method*) pairs.

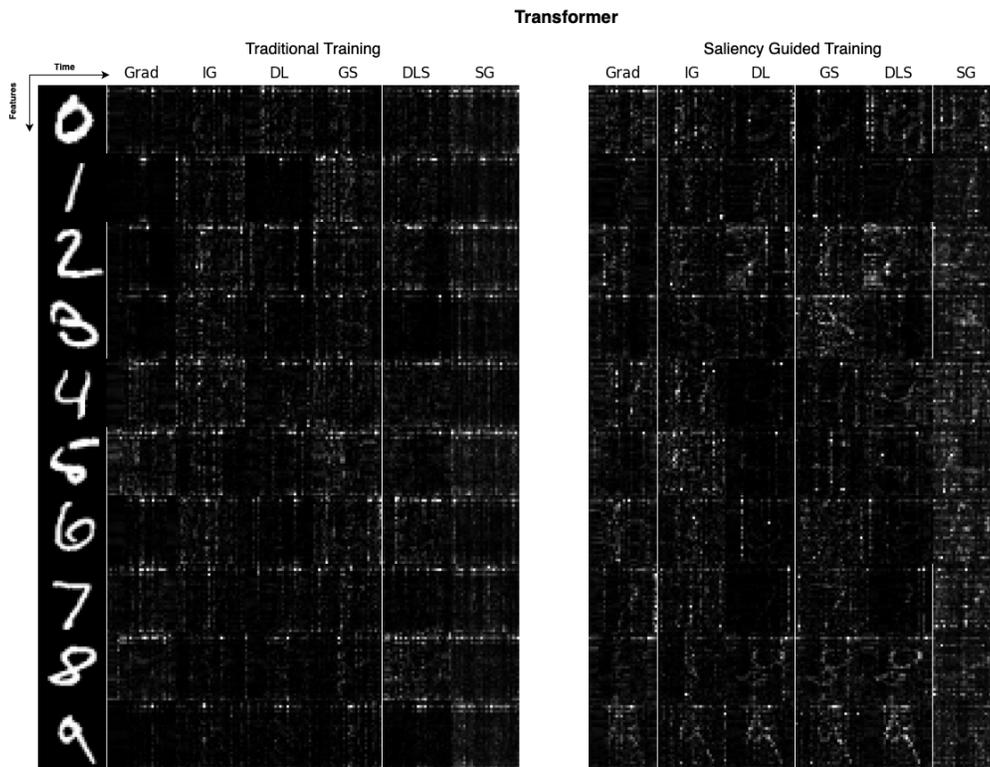


Figure 7: Saliency maps produced for (*Transformers, saliency method*) pairs.

Design Aspect	Levels	Middle Box		Moving Box		Positional Box		Rare Time		Rare Feature	
		N	S	N	S	T	F	N	M	N	M
Signal Position over Time	Same	•	•					•			•
	Different			•	•	•			•		
Signal Position over Features	Same	•	•			•		•	•	•	
	Different			•	•		•				•
Signal Difference	Value	•	•	•	•			•	•	•	•
	Position					•	•				
	Shape										
Signal Abundance over Time	Abundant	•		•		•	•			•	•
	Rare		•		•			•	•		
Signal Abundance over Features	Abundant	•		•		•	•	•	•		
	Rare		•		•					•	•

Figure 8: Figure from Ismail et al. [6]: Different evaluation datasets used for benchmarking saliency methods. Some datasets have multiple variations shown as sub-levels. N/S: normal and small shapes, T/F: temporal and feature positions, M: moving shape. All datasets are trained for binary classification. Examples are shown above each dataset, where dark red/blue shapes represent informative features.

Dataset	# Training	# Testing	# Time Steps	# Feature	# Informative Time steps	# Informative Features
Middle	1000	100	50	50	30	30
Small Middle	1000	100	50	50	15	15
Moving Middle	1000	100	50	50	30	30
Moving Small Middle	1000	100	50	50	15	15
Rare Time	1000	100	50	50	6	40
Moving Rare Time	1000	100	50	50	6	40
Rare Features	1000	100	50	50	40	6
Moving Rare Features	1000	100	50	50	40	6
Postional Time	1000	100	50	50	20	20
Postional Feature	1000	100	50	50	20	20

Table 4: Synthetic dataset details: Number of training samples, number of testing samples, number of time steps per sample, number of features per time step, number of time steps with informative features, and number of informative features in an informative time step.

Following Ismail et al. [6], we compare 4 neural architectures: LSTM [4], LSTM with Input-Cell Attention [5], Temporal Convolutional Network (TCN) [9] and, Transformers [13]. Each (*neural architecture, dataset*) pair was trained both traditionally and using saliency guided training. Test accuracy is reported in Table 5

Datasets	LSTM		LSTM+ Input-Cell		TCN		Transformer	
	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.
Middle	99	100	100	100	100	100	100	100
Small Middle	100	100	99	100	100	100	100	100
Moving Middle	100	100	100	100	100	100	99	100
Moving Small Middle	100	100	99	100	100	100	99	100
Rare Time	100	100	99	100	100	100	100	100
Moving Rare Time	100	100	100	100	100	100	99	100
Rare Features	100	100	100	100	100	100	100	100
Moving Rare Features	99	100	99	99	100	100	99	100
Positional Time	100	100	100	100	100	100	100	100
Positional Feature	100	100	99	100	99	100	100	100

Table 5: Test accuracy of different (*neural architecture, dataset*) pairs.

Quantitatively measuring the interpretability of a (*neural architecture, saliency method*) pair involves applying the saliency method, ranking features according to the saliency values, replacing high salient features with uninformative features from the original distribution at different percentages. Finally, we measure the model accuracy drop, weighted precision, and recall.

The area under precision curve (AUP) and the area under the recall curve (AUR) are calculated by the precision/recall values at different levels of degradation. Similar to Ismail et al. [6], we compare the AUP and AUR with a random baseline; since the baseline might be different for different models, we reported the difference between metrics values generated using the saliency method and the baseline. All experiments were ran 5 times the mean  $Diff(AUP)$ , and  $Diff(AUR)$  is shown in Tables [6-9].

The results in Tables [6-9] show the follows: **LSTM**: Saliency guided training along with Integrated Gradient has the best precision and recall. **LSTM with Input Cell Attention**: Saliency guided training improves the performance of different saliency methods and datasets. DeepSHAP gives the best precision, while DeepSHAP gives the best recall. **TCN**: overall, saliency guided training improves the performance of different saliency methods and datasets. Integrated Gradient, Gradient SHAP, and DeepSHAP are best performing saliency methods. **Transformers**: have the worst interpretability. Using saliency guided training improved recall but not precision.

Metric	Datasets	$\lambda$	$k$	Gradient		Integrated Gradient		DeepLIFT		Gradient SHAP		DeepSHAP		SmoothGrad		
				Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	
$Diff(AUP)$	Middle	1	30%	-0.280	-0.280	-0.261	-0.036	-0.267	-0.270	-0.263	-0.124	-0.267	-0.271	-0.283	-0.269	
	Small Middle	1	60%	-0.071	-0.066	-0.053	<b>0.052</b>	-0.070	-0.055	-0.056	-0.033	-0.070	-0.055	-0.072	-0.044	
	Moving Middle	1	60%	-0.265	-0.277	-0.218	-0.169	-0.237	-0.264	-0.222	-0.237	-0.239	-0.264	-0.259	-0.263	
	Moving Small Middle	1	5%	-0.059	-0.060	-0.035	<b>0.051</b>	-0.043	-0.045	-0.042	-0.013	-0.044	-0.046	-0.056	-0.037	
	Rare Time	1	30%	-0.076	-0.076	-0.075	-0.065	-0.076	-0.076	-0.075	-0.071	-0.076	-0.076	-0.076	-0.068	
	Moving Rare Time	1	50%	-0.067	-0.058	-0.042	<b>0.016</b>	-0.053	-0.042	-0.045	-0.010	-0.054	-0.043	-0.061	-0.032	
	Rare Feature	1	30%	-0.063	-0.075	-0.039	<b>0.006</b>	-0.047	-0.073	-0.038	-0.027	-0.048	-0.073	-0.069	-0.059	
	Moving Rare Feature	1	10%	-0.062	-0.069	-0.021	<b>0.012</b>	-0.040	-0.056	-0.032	-0.029	-0.041	-0.056	-0.059	-0.044	
	Postional Time	1	30%	-0.116	-0.119	-0.040	-0.006	-0.107	-0.112	-0.058	-0.046	-0.108	-0.113	-0.111	-0.102	
	Postional Feature	1	2%	-0.064	-0.104	-0.042	-0.104	-0.028	-0.089	-0.043	-0.105	-0.031	-0.091	-0.055	-0.053	
	$Diff(AUR)$	Middle	1	30%	0.072	<b>0.076</b>	0.128	<b>0.153</b>	0.125	<b>0.135</b>	0.122	<b>0.132</b>	0.114	<b>0.126</b>	<b>0.070</b>	0.031
		Small Middle	1	60%	-0.043	<b>0.037</b>	0.048	<b>0.157</b>	0.029	<b>0.116</b>	0.038	<b>0.129</b>	0.007	<b>0.102</b>	-0.032	<b>0.009</b>
Moving Middle		1	60%	0.060	<b>0.073</b>	0.119	<b>0.124</b>	0.110	<b>0.124</b>	<b>0.119</b>	0.117	0.099	<b>0.115</b>	<b>0.061</b>	0.042	
Moving Small Middle		1	5%	-0.032	-0.004	0.046	<b>0.135</b>	0.043	<b>0.073</b>	0.042	<b>0.093</b>	0.025	<b>0.060</b>	-0.023	-0.025	
Rare Time		1	30%	-0.244	-0.137	-0.132	<b>0.043</b>	-0.169	-0.021	-0.116	<b>0.005</b>	-0.189	-0.043	-0.145	-0.108	
Moving Rare Time		1	50%	-0.222	-0.070	-0.092	<b>0.075</b>	-0.103	<b>0.018</b>	-0.065	<b>0.060</b>	-0.131	0.002	-0.144	-0.035	
RareFeature		1	30%	0.182	<b>0.197</b>	<b>0.219</b>	0.180	0.217	<b>0.223</b>	<b>0.216</b>	<b>0.216</b>	0.211	<b>0.219</b>	<b>0.191</b>	0.166	
Moving Rare Feature		1	10%	0.143	<b>0.162</b>	0.191	<b>0.196</b>	0.191	<b>0.202</b>	0.194	<b>0.196</b>	0.183	<b>0.197</b>	<b>0.162</b>	0.107	
Postional Time		1	30%	-0.032	-0.073	0.072	<b>0.119</b>	<b>0.029</b>	0.021	0.046	<b>0.082</b>	<b>0.012</b>	0.001	-0.019	-0.064	
Postional Feature		1	2%	-0.053	-0.070	<b>0.016</b>	-0.005	-0.002	-0.005	<b>0.004</b>	-0.009	-0.018	-0.025	-0.056	-0.083	

Table 6: Difference in weighted AUP and AUR for (*LSTM, saliency method*) pairs. Overall, the best preference was achieved when using Integrated Gradients as a saliency method and saliency guided training as a training procedure.

Metric	Datasets	$\lambda$	$k$	Gradient		Integrated Gradient		DeepLIFT		Gradient SHAP		DeepSHAP		SmoothGrad		
				Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	
$Diff(AUP)$	Middle	1	40%	0.014	<b>0.046</b>	0.233	<b>0.252</b>	0.218	<b>0.237</b>	0.244	<b>0.261</b>	0.232	<b>0.247</b>	-0.006	<b>0.026</b>	
	Small Middle	1	60%	0.049	<b>0.150</b>	0.161	<b>0.273</b>	0.169	<b>0.305</b>	0.170	<b>0.273</b>	0.180	<b>0.312</b>	0.038	<b>0.091</b>	
	Moving Middle	1	80%	0.044	<b>0.082</b>	<b>0.262</b>	0.251	<b>0.260</b>	0.256	<b>0.276</b>	0.256	<b>0.277</b>	0.261	0.010	<b>0.044</b>	
	Moving Small Middle	1	5%	0.044	<b>0.055</b>	0.181	<b>0.201</b>	0.179	<b>0.196</b>	0.187	<b>0.200</b>	0.190	<b>0.204</b>	0.022	<b>0.029</b>	
	Rare Time	1	40%	0.186	<b>0.278</b>	0.271	<b>0.378</b>	0.323	<b>0.412</b>	0.279	<b>0.373</b>	0.338	<b>0.424</b>	0.133	<b>0.209</b>	
	Moving Rare Time	1	80%	0.144	<b>0.276</b>	0.233	<b>0.388</b>	0.269	<b>0.417</b>	0.238	<b>0.381</b>	0.282	<b>0.429</b>	0.103	<b>0.167</b>	
	Rare Feature	1	30%	0.032	<b>0.101</b>	0.163	<b>0.270</b>	0.166	<b>0.266</b>	0.174	<b>0.278</b>	0.180	<b>0.274</b>	0.039	<b>0.105</b>	
	Moving Rare Feature	1	5%	-0.002	-0.004	0.120	<b>0.124</b>	<b>0.116</b>	<b>0.116</b>	0.124	<b>0.127</b>	<b>0.126</b>	<b>0.126</b>	-0.003	-0.004	
	Postional Time	1	40%	0.117	<b>0.186</b>	0.184	<b>0.225</b>	0.236	<b>0.314</b>	0.197	<b>0.252</b>	0.248	<b>0.316</b>	0.093	<b>0.187</b>	
	Postional Feature	1	5%	-0.021	<b>0.007</b>	0.072	<b>0.083</b>	0.080	<b>0.113</b>	0.089	<b>0.101</b>	0.088	<b>0.122</b>	-0.031	-0.012	
	$Diff(AUR)$	Middle	1	40%	0.028	<b>0.084</b>	0.163	<b>0.176</b>	0.160	<b>0.180</b>	0.162	<b>0.173</b>	0.157	<b>0.177</b>	-0.001	<b>0.044</b>
		Small Middle	1	60%	0.064	<b>0.176</b>	0.186	<b>0.217</b>	0.189	<b>0.217</b>	0.182	<b>0.212</b>	0.183	<b>0.213</b>	0.031	<b>0.159</b>
Moving Middle		1	80%	0.060	<b>0.117</b>	0.174	<b>0.180</b>	0.175	<b>0.187</b>	0.173	<b>0.177</b>	0.175	<b>0.183</b>	0.021	<b>0.072</b>	
Moving Small Middle		1	5%	0.079	<b>0.101</b>	<b>0.202</b>	0.201	<b>0.199</b>	0.194	<b>0.198</b>	0.196	<b>0.194</b>	0.186	0.029	<b>0.052</b>	
Rare Time		1	40%	0.139	<b>0.203</b>	0.214	<b>0.225</b>	0.214	<b>0.233</b>	0.211	<b>0.223</b>	0.211	<b>0.233</b>	0.103	<b>0.191</b>	
Moving Rare Time		1	80%	0.118	<b>0.213</b>	0.198	<b>0.226</b>	0.200	<b>0.233</b>	0.193	<b>0.224</b>	0.194	<b>0.232</b>	0.070	<b>0.197</b>	
RareFeature		1	30%	0.077	<b>0.181</b>	0.196	<b>0.223</b>	0.197	<b>0.224</b>	0.193	<b>0.222</b>	0.193	<b>0.223</b>	0.074	<b>0.172</b>	
Moving Rare Feature		1	5%	<b>0.059</b>	0.039	0.188	<b>0.189</b>	<b>0.191</b>	0.186	0.182	<b>0.183</b>	<b>0.186</b>	0.180	<b>0.038</b>	0.028	
Postional Time		1	40%	0.140	<b>0.201</b>	0.188	<b>0.200</b>	0.203	<b>0.225</b>	0.185	<b>0.202</b>	0.201	<b>0.224</b>	0.109	<b>0.188</b>	
Postional Feature		1	5%	-0.017	<b>0.043</b>	0.141	<b>0.146</b>	0.145	<b>0.166</b>	0.141	<b>0.150</b>	0.132	<b>0.157</b>	-0.041	<b>0.005</b>	

Table 7: The difference in weighted AUP and AUR for different (*LSTM with Input-Cell Attention, saliency method*) pairs. The use of saliency guided training improved the performance of most saliency methods. Overall, DeepSHAP and DeepLIFT produced the best precision and recall, respectively, when combined with saliency guided training.

Metric	Datasets	$\lambda$	$k$	Gradient		Integrated Gradient		DeepLIFT		Gradient SHAP		DeepSHAP		SmoothGrad	
				Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.
Diff(AUP)	Middle	1	50%	0.127	<b>0.217</b>	0.283	<b>0.393</b>	0.350	<b>0.398</b>	0.290	<b>0.384</b>	0.365	<b>0.416</b>	0.090	<b>0.194</b>
	Small Middle	1	40%	0.164	<b>0.260</b>	0.299	<b>0.433</b>	0.312	<b>0.419</b>	0.302	<b>0.418</b>	0.328	<b>0.442</b>	0.156	<b>0.253</b>
	Moving Middle	1	70%	0.122	<b>0.197</b>	0.287	<b>0.342</b>	0.332	<b>0.367</b>	0.286	<b>0.329</b>	0.345	<b>0.387</b>	0.047	<b>0.182</b>
	Moving Small Middle	1	80%	<b>0.065</b>	0.043	<b>0.194</b>	0.151	0.169	<b>0.191</b>	<b>0.190</b>	0.152	0.183	<b>0.200</b>	<b>0.037</b>	0.023
	Rare Time	1	50%	0.184	<b>0.290</b>	0.314	<b>0.363</b>	<b>0.324</b>	0.309	0.314	<b>0.360</b>	<b>0.352</b>	0.319	0.177	<b>0.226</b>
	Moving Rare Time	1	50%	0.142	<b>0.182</b>	0.260	<b>0.333</b>	<b>0.257</b>	0.243	0.258	<b>0.330</b>	<b>0.275</b>	0.251	0.122	<b>0.179</b>
	Rare Feature	1	30%	0.058	<b>0.244</b>	0.246	<b>0.451</b>	0.252	<b>0.422</b>	0.249	<b>0.453</b>	0.286	<b>0.450</b>	0.085	<b>0.259</b>
	Moving Rare Feature	1	5%	-0.003	<b>0.004</b>	0.116	<b>0.134</b>	0.112	<b>0.114</b>	0.122	<b>0.129</b>	0.122	<b>0.123</b>	<b>0.007</b>	0.005
	Postional Time	1	70%	<b>0.115</b>	0.072	<b>0.180</b>	0.114	<b>0.233</b>	0.069	<b>0.187</b>	0.117	<b>0.237</b>	0.035	<b>0.106</b>	0.066
	Postional Feature	1	10%	0.082	<b>0.176</b>	0.151	<b>0.199</b>	0.136	<b>0.162</b>	0.155	<b>0.203</b>	0.137	<b>0.175</b>	0.058	<b>0.159</b>
Diff(AUR)	Middle	1	50%	0.133	<b>0.161</b>	0.190	<b>0.207</b>	0.202	<b>0.209</b>	0.188	<b>0.205</b>	0.201	<b>0.205</b>	0.054	<b>0.128</b>
	Small Middle	1	40%	0.086	<b>0.230</b>	0.194	<b>0.240</b>	0.202	<b>0.240</b>	0.189	<b>0.239</b>	0.196	<b>0.241</b>	0.039	<b>0.230</b>
	Moving Middle	1	70%	0.134	<b>0.144</b>	0.191	<b>0.195</b>	0.201	<b>0.208</b>	0.186	<b>0.194</b>	0.201	<b>0.203</b>	0.036	<b>0.121</b>
	Moving Small Middle	1	80%	<b>0.118</b>	0.117	<b>0.204</b>	0.199	0.193	<b>0.195</b>	<b>0.196</b>	0.196	0.186	<b>0.190</b>	-0.001	<b>0.065</b>
	Rare Time	1	50%	0.173	<b>0.215</b>	0.199	<b>0.233</b>	<b>0.225</b>	0.221	0.193	<b>0.230</b>	<b>0.226</b>	0.204	0.125	<b>0.151</b>
	Moving Rare Time	1	50%	0.106	<b>0.198</b>	0.177	<b>0.220</b>	<b>0.195</b>	0.189	0.167	<b>0.224</b>	<b>0.191</b>	0.179	-0.057	<b>0.149</b>
	Rare Feature	1	30%	0.152	<b>0.222</b>	0.222	<b>0.239</b>	0.223	<b>0.239</b>	0.219	<b>0.239</b>	0.224	<b>0.239</b>	0.130	<b>0.217</b>
	Moving Rare Feature	1	5%	0.101	<b>0.122</b>	0.198	<b>0.204</b>	<b>0.206</b>	0.205	0.195	<b>0.201</b>	0.196	<b>0.198</b>	0.048	<b>0.055</b>
	Postional Time	1	70%	0.126	<b>0.128</b>	0.156	<b>0.172</b>	<b>0.194</b>	0.160	0.147	<b>0.165</b>	<b>0.181</b>	0.110	0.039	<b>0.102</b>
	Postional Feature	1	10%	0.126	<b>0.174</b>	0.177	<b>0.196</b>	<b>0.180</b>	0.175	0.172	<b>0.194</b>	0.164	<b>0.154</b>	0.049	<b>0.160</b>

Table 8: The difference in weighted AUP and AUR for different (*TCN, saliency method*) pairs. The use of saliency guided training improved the performance of most saliency methods. Overall, when combined with saliency guided training, Integrated Gradients and DeepSHAP produced the best precision. For recall, Integrated Gradients, DeepLift, Gradient SHAP, and DeepSHAP seem to perform similarly, again, the best performance was achieved when saliency guided training is used.

Metric	Datasets	$\lambda$	$k$	Gradient		Integrated Gradient		DeepLIFT		Gradient SHAP		DeepSHAP		SmoothGrad	
				Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.	Trad.	Sal.
Diff(AUP)	Middle	1	30%	-0.179	-0.213	<b>0.051</b>	-0.004	-0.116	-0.176	<b>0.067</b>	-0.064	-0.062	-0.222	-0.069	-0.150
	Small Middle	1	60%	-0.034	-0.057	<b>0.054</b>	0.042	-0.018	-0.034	<b>0.066</b>	0.024	<b>0.009</b>	-0.060	<b>0.006</b>	-0.022
	Moving Middle	1	90%	-0.188	-0.146	<b>0.062</b>	0.018	-0.142	-0.067	<b>0.065</b>	-0.011	-0.130	-0.143	-0.091	-0.157
	Moving Small Middle	1	70%	-0.002	-0.008	0.031	<b>0.039</b>	0.017	<b>0.029</b>	0.026	<b>0.037</b>	<b>0.036</b>	0.016	-0.021	-0.035
	Rare Time	1	50%	<b>0.038</b>	-0.006	<b>0.118</b>	0.057	-0.019	<b>0.017</b>	<b>0.132</b>	0.049	<b>0.014</b>	-0.010	-0.029	-0.031
	Moving Rare Time	1	50%	<b>0.066</b>	0.062	<b>0.110</b>	0.049	-0.021	<b>0.046</b>	<b>0.117</b>	0.055	-0.009	<b>0.033</b>	-0.026	-0.027
	Rare Feature	1	30%	-0.049	-0.045	0.029	<b>0.139</b>	-0.002	-0.004	0.033	<b>0.088</b>	<b>0.008</b>	-0.015	0.005	<b>0.028</b>
	Moving Rare Feature	1	10%	-0.034	-0.031	0.041	<b>0.055</b>	0.008	<b>0.014</b>	0.038	<b>0.049</b>	0.008	<b>0.022</b>	-0.002	-0.013
	Postional Time	1	60%	-0.060	-0.078	<b>0.084</b>	0.029	-0.048	-0.047	0.102	-0.001	<b>0.013</b>	-0.072	<b>0.026</b>	-0.057
	Postional Feature	1	10%	-0.094	-0.099	<b>0.032</b>	0.012	-0.061	-0.097	<b>0.046</b>	0.008	-0.029	-0.098	<b>0.019</b>	0.006
Diff(AUR)	Middle	1	30%	<b>0.087</b>	0.053	<b>0.167</b>	0.146	0.155	0.112	<b>0.157</b>	0.111	<b>0.119</b>	-0.051	<b>0.040</b>	-0.025
	Small Middle	1	60%	<b>0.085</b>	0.030	0.186	<b>0.189</b>	<b>0.128</b>	0.113	<b>0.173</b>	0.171	<b>0.077</b>	-0.025	<b>0.060</b>	0.036
	Moving Middle	1	90%	<b>0.071</b>	0.134	0.164	<b>0.185</b>	0.136	<b>0.181</b>	0.150	<b>0.183</b>	0.057	<b>0.130</b>	0.019	<b>0.040</b>
	Moving Small Middle	1	70%	0.118	<b>0.137</b>	0.171	<b>0.177</b>	0.157	<b>0.171</b>	0.160	<b>0.171</b>	0.098	<b>0.117</b>	-0.004	-0.019
	Rare Time	1	50%	<b>0.152</b>	0.139	<b>0.206</b>	0.172	0.116	<b>0.135</b>	<b>0.199</b>	0.157	0.077	<b>0.088</b>	-0.027	-0.073
	Moving Rare Time	1	50%	0.184	<b>0.185</b>	<b>0.198</b>	0.170	0.124	<b>0.175</b>	<b>0.186</b>	0.168	0.059	<b>0.127</b>	-0.033	-0.013
	Rare Feature	1	30%	0.115	<b>0.152</b>	0.184	<b>0.217</b>	0.187	<b>0.188</b>	0.173	<b>0.203</b>	<b>0.144</b>	0.129	0.087	<b>0.135</b>
	Moving Rare Feature	1	10%	0.101	<b>0.122</b>	0.179	<b>0.183</b>	0.174	<b>0.180</b>	0.165	<b>0.176</b>	0.125	<b>0.149</b>	<b>0.060</b>	0.034
	Postional Time	1	60%	<b>0.091</b>	0.071	<b>0.193</b>	0.172	<b>0.154</b>	0.150	<b>0.184</b>	0.151	<b>0.145</b>	0.059	<b>0.103</b>	-0.004
	Postional Feature	1	10%	<b>0.017</b>	0.013	0.170	<b>0.149</b>	<b>0.123</b>	0.057	<b>0.160</b>	0.131	<b>0.105</b>	-0.072	<b>0.094</b>	0.073

Table 9: The difference in weighted AUP and AUR for different (*Transformers, saliency method*) pairs. In this benchmark, Transformers seem to have the worst interpretability. Using saliency guided training improved recall but not precision. Overall best precision was achieved when combining traditional training with Gradient SHAP. While best recall was achieved when using saliency guided training and Integrated Gradients.

## References

- [1] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. Eraser: A benchmark to evaluate rationalized nlp models. *arXiv preprint arXiv:1911.03429*, 2019.
- [2] Gerry. 265 bird species, 2021. URL <https://www.kaggle.com/gpiosenska/100-bird-species/>.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [4] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. In *Neural computation*, 1997.
- [5] Aya Abdelsalam Ismail, Mohamed Gunady, Luiz Pessoa, Hector Corrada Bravo, and Soheil Feizi. Input-cell attention reduces vanishing saliency of recurrent neural networks. In *Advances in Neural Information Processing Systems*, 2019.

- [6] Aya Abdelsalam Ismail, Mohamed Gunady, Héctor Corrada Bravo, and Soheil Feizi. Benchmarking deep learning interpretability in time series predictions. *arXiv preprint arXiv:2010.13924*, 2020.
- [7] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqu Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.
- [8] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master's thesis, Department of Computer Science, University of Toronto*, 2009.
- [9] Colin Lea, Michael Flynn, Rene Vidal, Austin Reiter, and Gregory Hager. Temporal convolutional networks for action segmentation and detection. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [10] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database, 2010.
- [11] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [12] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.