

Common Attribution Methods

Method	Paper	Dataset	Info
GradCam	<a href="#">Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization</a>	general public images	Image Captioning and VQA achieved using CNN + LSTM on locally activated regions
Integrated Gradients	<a href="#">Axiomatic Attribution for Deep Networks</a>	ImageNet	GoogleNet arch used
ViT Shapley	<a href="#">LEARNING TO ESTIMATE SHAPLEY VALUES WITH VISION TRANSFORMERS</a>	ImageNette, MURA, musculoskeletal radiographs, Oxford IIIT-Pets	expensive (not in production)
LRP	<a href="#">Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers</a>	CIFAR-10, ImageNET, MIT Places	based on first order Taylor expansion

Attention Based Methods

Method	Paper	Dataset	Info
Raw Attention	<a href="#">Self-attention for raw optical Satellite Time Series Classification</a>	Crop Satellite Data	LSTM-RNN, MS-ResNet, TempCNN, hyper para selection - HyperOpt, classification task
Attention Rollout & Flow	<a href="#">Quantifying Attention Flow in Transformers</a>	subject-verb agreement	information flow in the network with a DAG, maximum flow algorithms, pretrained Bert
Partial LRP	<a href="#">AttnLRP: Attention-Aware Layer-Wise Relevance Propagation for Transformers</a>	SQuAD v2 (QA), Wikipedia summary	handles attribution flow through non linear components better
Beyond Attention	<a href="#">Transformer Interpretability Beyond Attention Visualization</a>	ImageNet, Movie Review, ERASER	relevancy propagation rule applicable to both positive and negative attributions, integrate the attention and the relevancy scores, and combine the integrated results for multiple attention blocks.

Method	Paper	Dataset	Info
<b>GradSam</b>	<a href="#">Grad-SAM: Explaining Transformers via Gradient Self-Attention Maps</a>	Stanford Sentiment Tree (SST), AgNews, IMDB, MultiRC	
Beyond Intuition	<a href="#">Beyond Intuition: Rethinking Token Attributions inside Transformers</a>	ImageNet, Movie Review, Newgroups Review	attention perception and reasoning feedback, head-wise and token-wise approximations, single & bi modality

## Pruning-based Methods

Method	Paper	Dataset	Info
<b>IA-Red<sup>2</sup></b>	<a href="#">IA-RED2: Interpretability-Aware Redundancy Reduction for Vision Transformers</a>	ImageNet, Kinetics-400	multitask interpreter, both model-agnostic and task-agnostic
Vision DiffMask	<a href="#">VISION DIFFMASK: Faithful Interpretation of Vision Transformers with Differentiable Patch Masking</a>	CIFAR-10, ImageNet	activations of the model's hidden layers
X Pruner	<a href="#">X-Pruner: eXplainable Pruning for Vision Transformers</a>	CIFAR-10, ILSVRC-12	novel explainability-aware mask, DeiT, Swin transformer
<b>EViT</b>	<a href="#">NOT ALL PATCHES ARE WHAT YOU NEED: EXPEDITING VISION TRANSFORMERS VIA TOKEN REORGANIZATIONS</a>	ImageNet, JFT-300M	attentive token identification, inattentive token fusion

## Inherently Explainable Methods

Method	Paper	Dataset	Info
ViT-CX	<a href="#">ViT-CX: Causal Explanation of Vision Transformers</a>	ImageNet, MURA	vit feature maps -> mask, high degree of redundancy, clustering of masks
<b>ViT-Net</b>	<a href="#">ViT-Net: Interpretable Vision Transformers with Neural Tree Decoder</a>	CUB-200-2011, Stanford Cars, Stanford Dogs	neural tree encoder
<b>R-Cut</b>	<a href="#">R-Cut: Enhancing Explainability in Vision Transformers with Relationship Weighted Out and Cut</a>	ImageNet, LRN	

Other

Method	Paper	Dataset	Info
eX-ViT	<a href="#">eX-ViT: A Novel explainable vision transformer for weakly supervised semantic segmentation</a>	PASCAL VOC 2012, MS COCO 2014	Explainable Multi-Head Attention (E-MHA) module and Attribute-guided Explainer (AttE) module