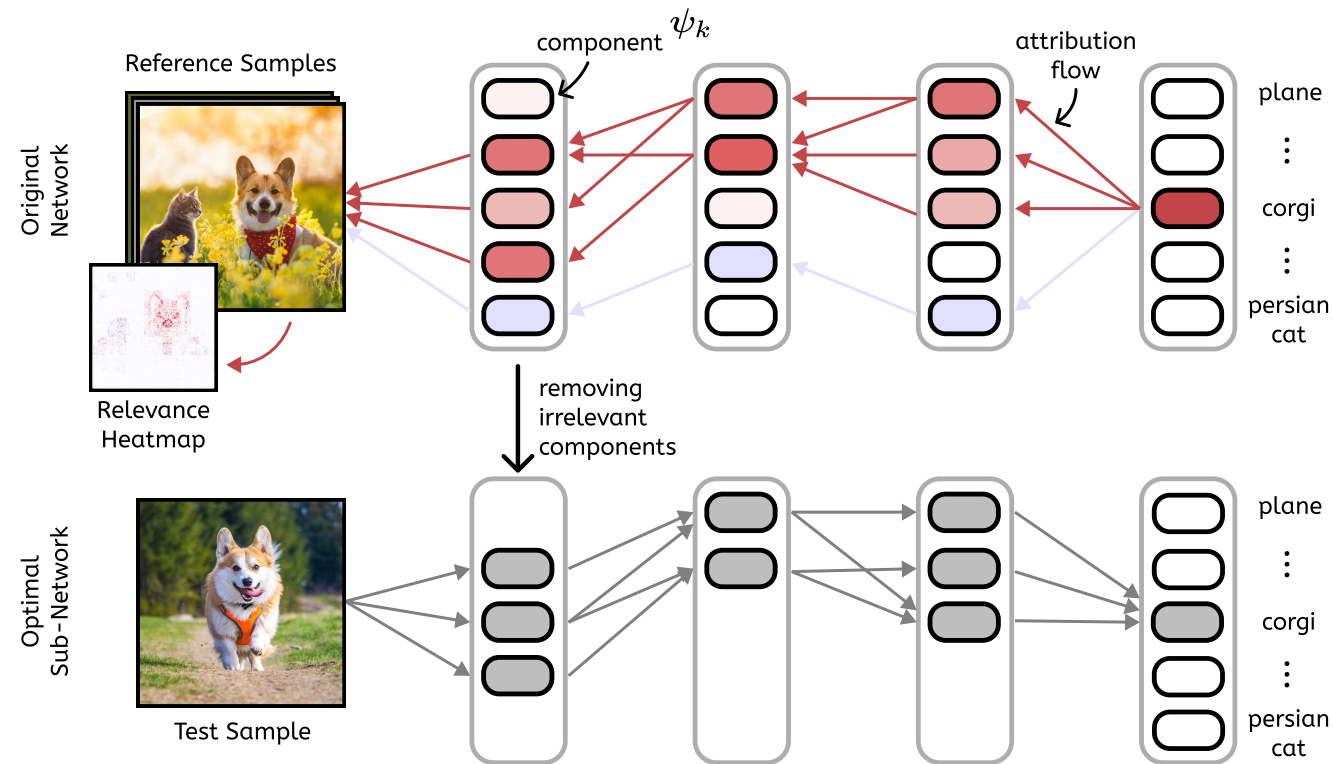# When XAI meets Compression & Sub-graph Discovery

## Pruning By Explaining Revisited: Optimizing Attribution Methods to Prune CNNs & Transformers

Sayed Mohammad Vakilzadeh Hatefi, Maximilian Dreyer, Reduan Achtibat, Thomas Wiegand, Wojciech Samek, Sebastian Lapuschkin

ECCV

Fraunhofer HHI · BIFOLD · TECHNISCHE UNIVERSITÄT BERLIN · iTOBOS · TEMA

PXP Pruning by eXplaining in PyTorch

---

## Pruning by Explaining



Reference Samples · Relevance Heatmap · Original Network · component $\psi_k$ · attribution flow · plane · corgi · persian cat · removing irrelevant components · Optimal Sub-Network · Test Sample · plane · corgi · persian cat

## Optimize XAI for Pruning



Performance-Efficiency Tradeoff · performance · original · optimized (ours) · naive · random · better · sparsity

Original Network · input layer · intermediate layers · output layer · over-parameterized DNN

Naive Attribution-based Pruning · +XAI · sparse DNN with mostly relevant structures remaining

Optimized Attribution-based Pruning (ours) · +XAI optimization · optimally sparse DNN with only relevant structures remaining

A Corgi · Confidence: 86% · Faithful Explanation · Deviation from Original Explanation · Cosine Similarity · Pruning Rate · Drop of Confidence · Prediction Confidence · Pruning Rate

Pruning Rate: 10% · 25% · 35% · 50% · 70% · 80% · 95%

Optimized LRP Pruner · Confidence: 85% · 87% · 85% · 78% · 27% · 0% · 0%

Naive LRP Pruner · Confidence: 87% · 85% · 82% · 36% · 0% · 0% · 0%

Random Pruner · Confidence: 85% · 54% · 25% · 3% · 0% · 0% · 0%

---

### → Our Pruning Framework

Given **a set of reference samples** $\mathcal{X}_{\text{ref}}$ defined as:

$$\mathcal{X}_{\text{ref}} = \{x_1, x_2, \ldots, x_{n_{\{\text{ref}\}}}\}$$

**Importance score** of a component $\psi_k$ can be computed by:

$$\bar{R}_{\psi_k} = \frac{1}{n_{\text{ref}}} \sum_{i=1}^{n_{\text{ref}}} R_{\psi_k}(x_i)$$

But, how should we compute $R_{\psi_k}$? In other words, **what is a reliable pruning criterion?**

**+** Use relevance scores of **Layer-wise Relevance Propagation**:

$$R_{i \leftarrow j}^{(l-1,l)} = \frac{z_{ij}}{z_j} R_j^l$$

What is an **advantage** of this criterion?

**+ LRP**'s relevance scores are intrinsically **normalized** due to their conservation property across layers.

How large should be the set of reference samples $\mathcal{X}_{\text{ref}}$?

**+** The more samples used for attribution, the more stable the pruning is. However, for **CNNs**, the work of [1] has shown that **10 reference samples** per class is sufficient.

**+** For **Transformers** on the other hand, our experiments conveyed that **only 1 reference sample** generates robust relevance scores for pruning.

### → Optimization of XAI Methods

Typically takes place to generate **faithful explanations**, but solutions are **not necessarily optimal for pruning**. So, why don't we **optimize XAI for pruning** directly?
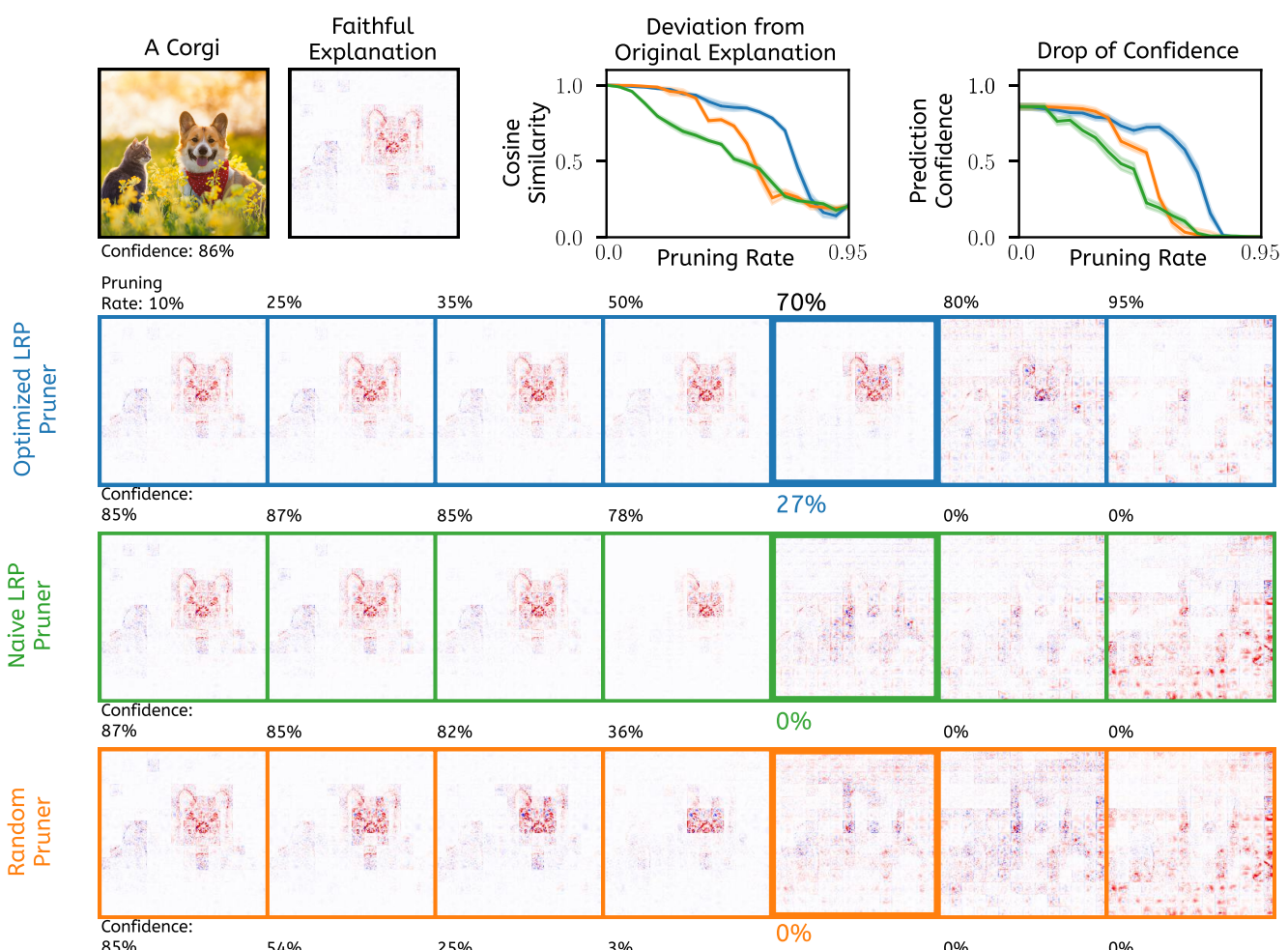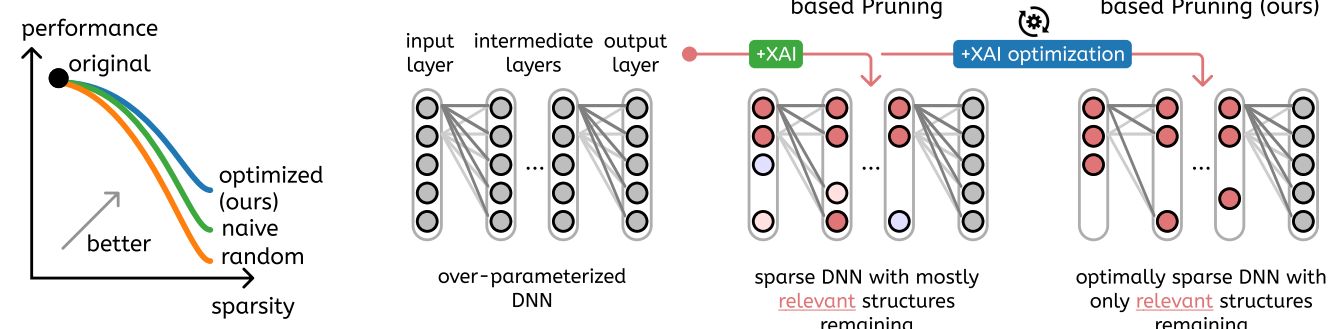
### → CNNs and Transformers in Pruning

**CNNs** have quite **sufficient amount of parameters** and thus not much can be pruned from them based on the default task (i.e., ImageNet 1000-class classification). **Explanations that faithfully attribute CNNs**, perform well on pruning as well.

**Transformers** are typically more **overparameterized** than CNNs, which induce more pruning rates while keeping high performance given the default task. Unlike CNNs, **a faithful explainer** of Transformers **does not guarantee stable pruning**, thus **encouraging extra optimization** of explainer.

Overall, **LRP-Epsilon** [2, 3, 4] is a **promising explainer for pruning** across different architectures.
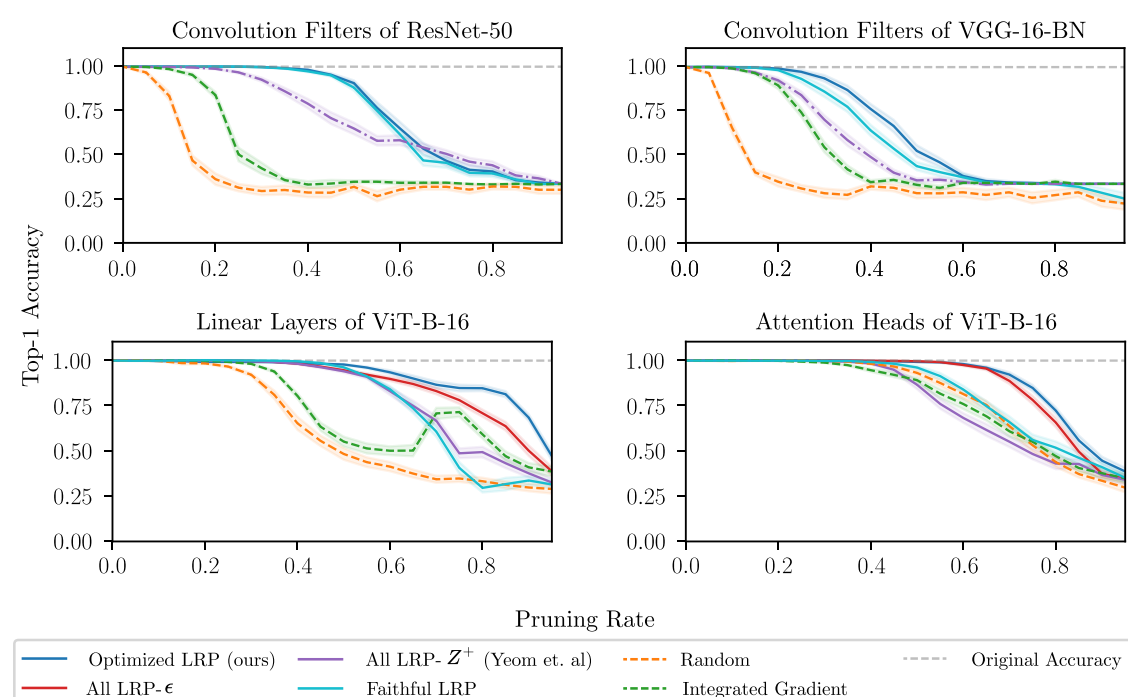
---

## Experiments

### Comparison of Different Pruning Criteria



Convolution Filters of ResNet-50 · Convolution Filters of VGG-16-BN · Linear Layers of ViT-B-16 · Attention Heads of ViT-B-16 · Top-1 Accuracy · Pruning Rate

Legend: Optimized LRP (ours) · All LRP-$\epsilon$ · All LRP-$Z^+$ (Yeom et. al) · Faithful LRP · Random · Integrated Gradient · Original Accuracy

### Overparametrization of CNNs vs Transformers



Convolution Filters of ResNet-18 · Convolution Filters of VGG-16-BN · Attention Heads of ViT-B-16 · Linear Layers of ViT-B-16 · Top-1 Accuracy · Pruning Rate

Legend: LRP-opt · Integrated Gradient · Random · Original Accuracy · 3-classes Domain Restriction · 100-classes Domain Restriction · Full Model (1000 classes)

---

## References

[1] Yeom et al. **Pruning by explaining: A novel criterion for deep neural network pruning**. Pattern Recognition 115, 107899 (2021)

[2] Bach et al. **On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation**. PloS one 10(7), e0130140 (2015)

[3] Montavon et al. **Layer-wise relevance propagation: an overview**. Explainable AI: interpreting, explaining and visualizing deep learning pp. 193–209 (2019)

[4] Achtibat et al. **AttnLRP: Attention-aware layer-wise relevance propagation for transformers**. In: Proceedings of the 41st International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 235, pp. 135–168. PMLR (21–27 Jul 2024)