

Efficient and Flexible Neural Network Training Through Layer-wise Feedback Propagation

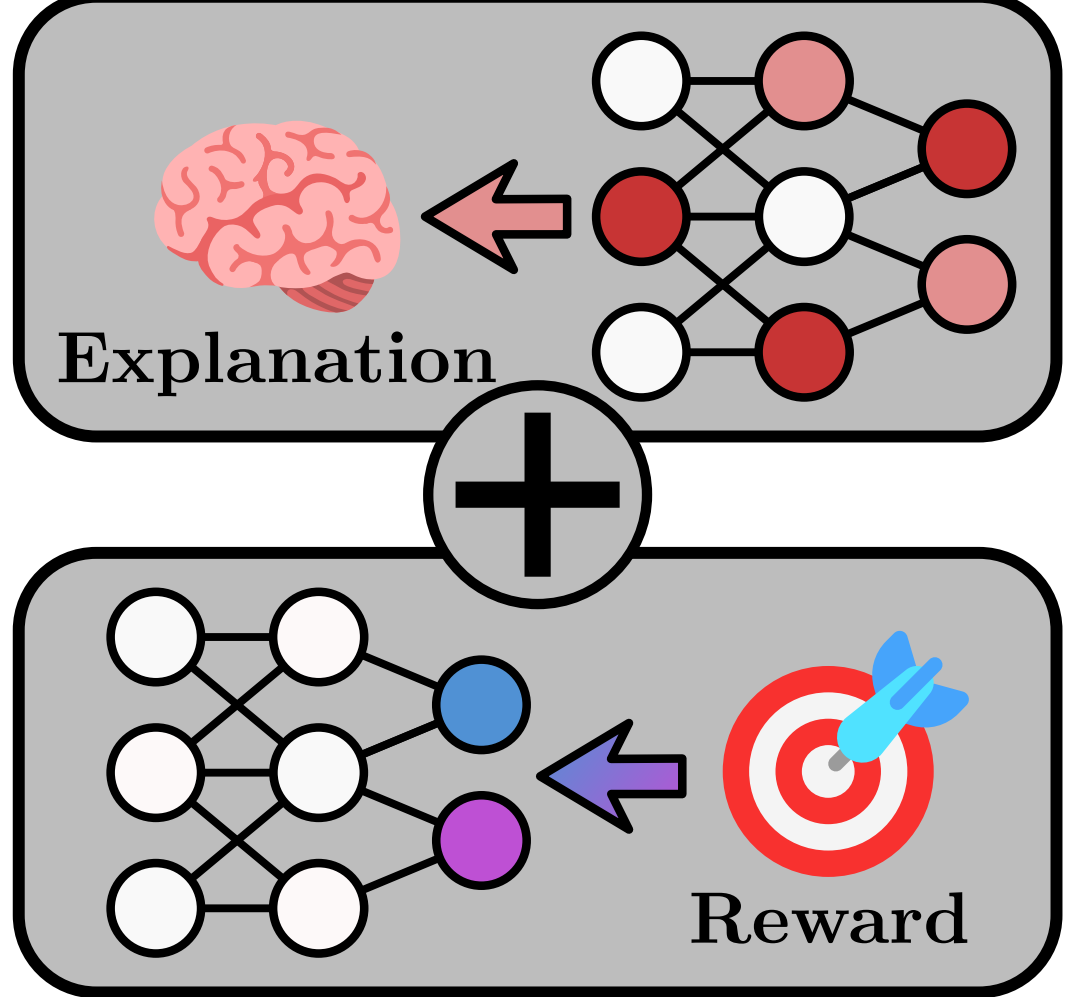
Leander Weber, Jim Berend, Moritz Weckbecker, Alexander Binder, Thomas Wiegand, Wojciech Samek, Sebastian Lapuschkin



TL;DR

Idea

Contribution

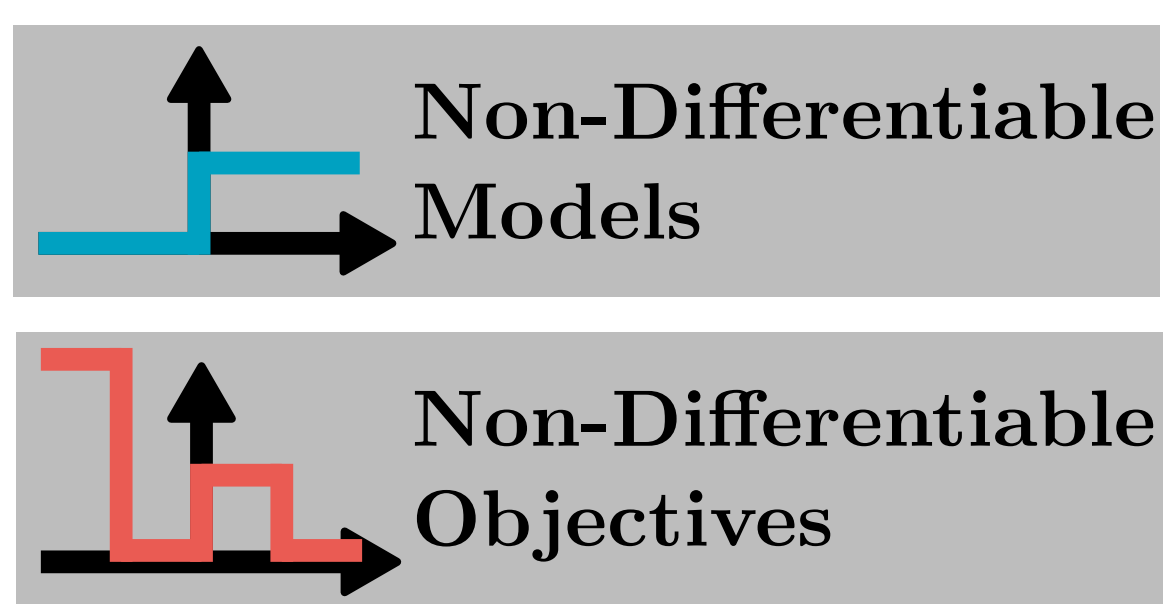


Feedback

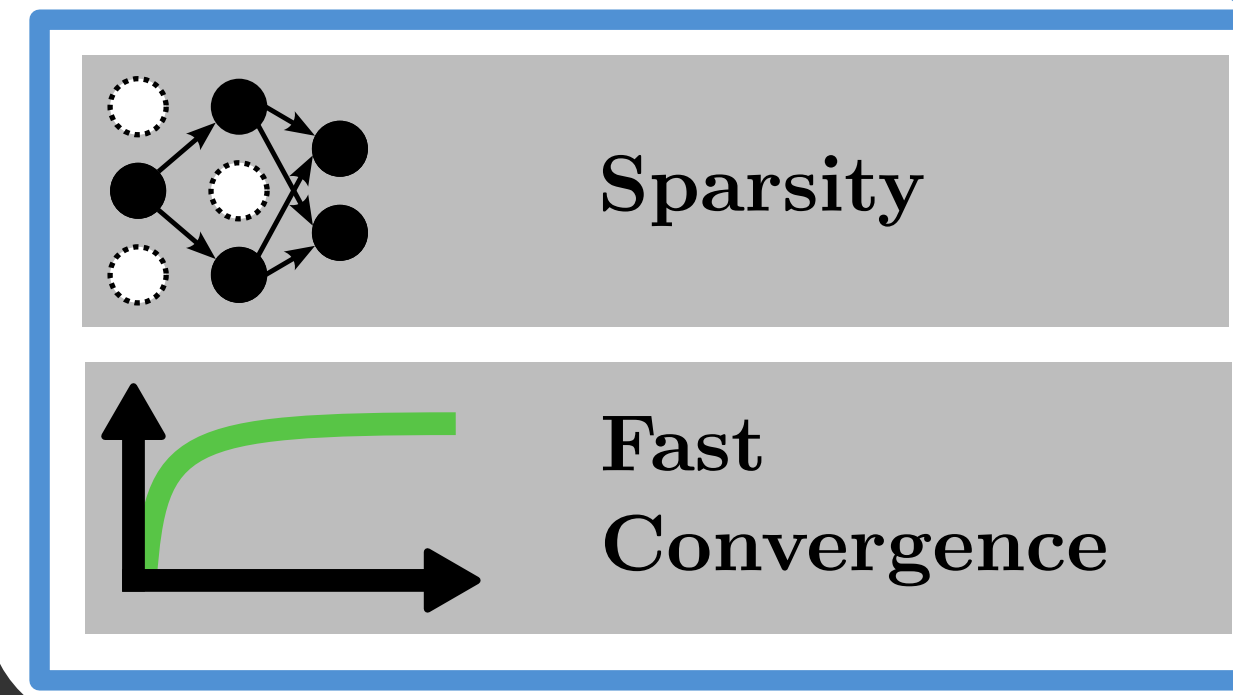
We propose LFP, a novel **training paradigm** that leverages **Local XAI** to distribute a **Feedback Signal**

Strengths

Flexibility



Efficiency

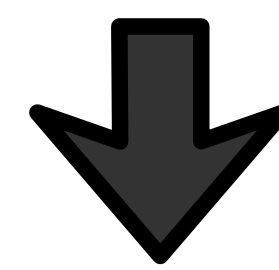


Main Results

Convergence properties similar to SGD

Approximation-free training of non-differentiable models

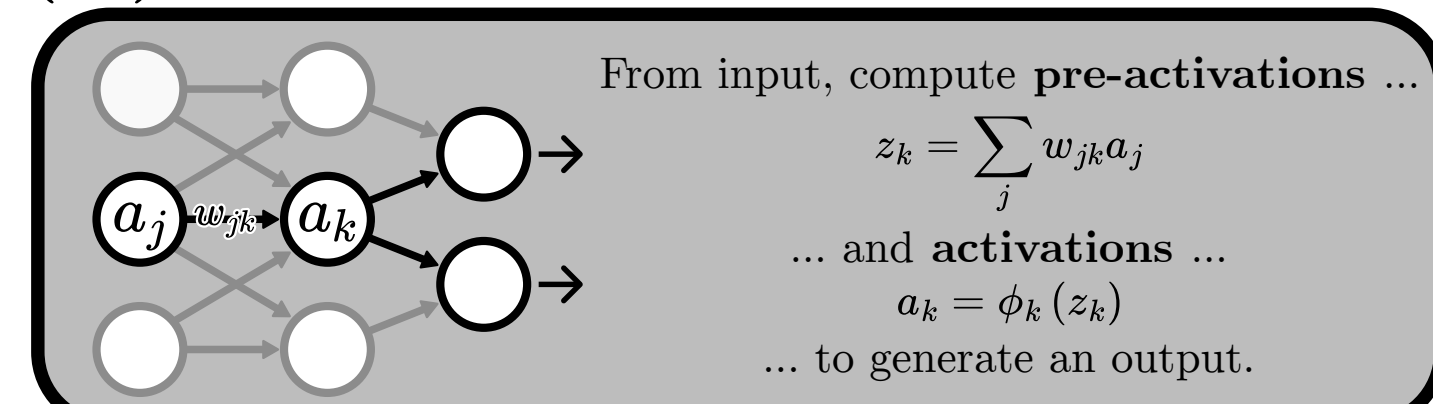
Improved model **sparsity** and **efficient** information representation



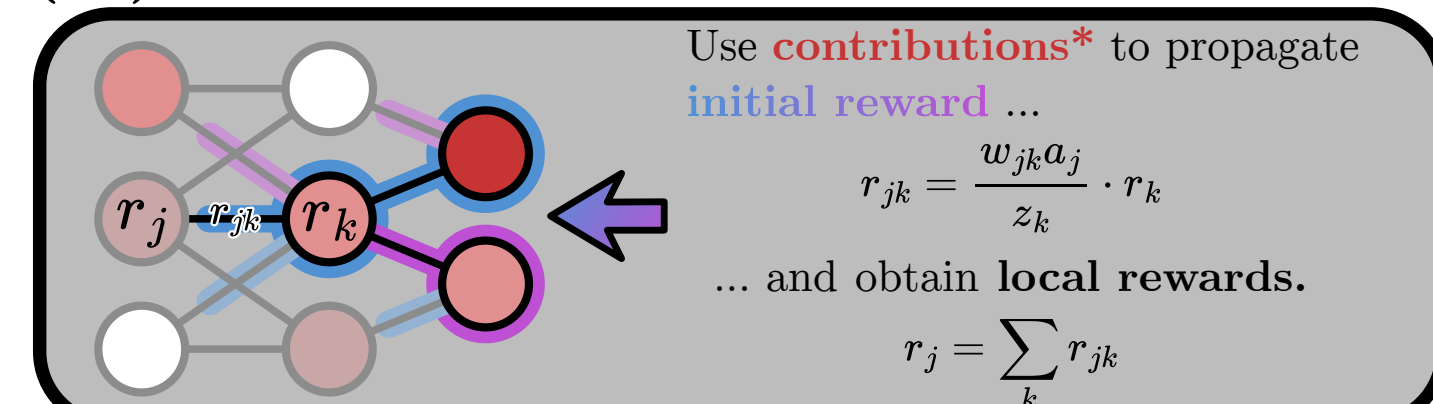
Promising for **energy-restricted applications** such as edge devices

Method

(A) Forward Pass

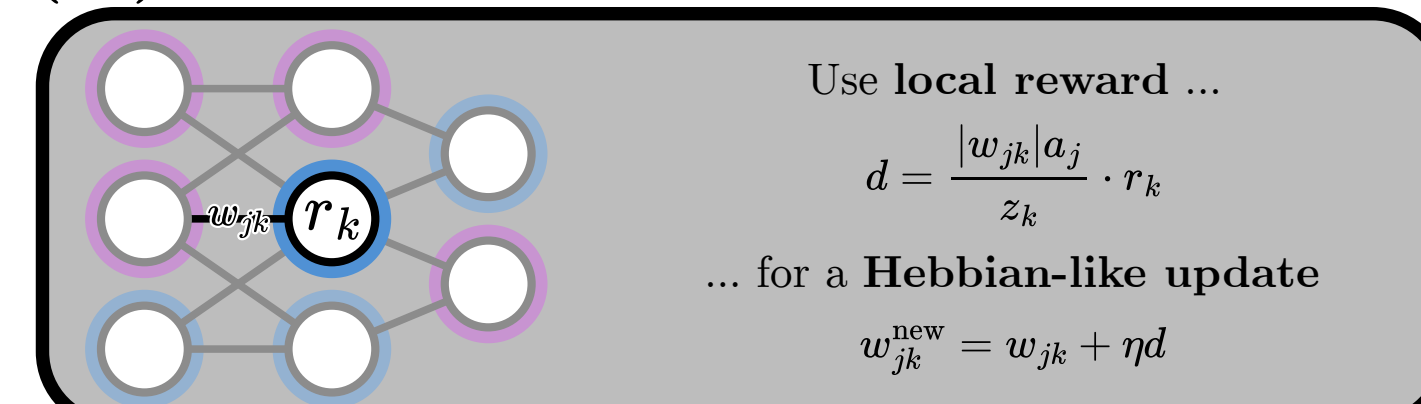


(B) Feedback Propagation



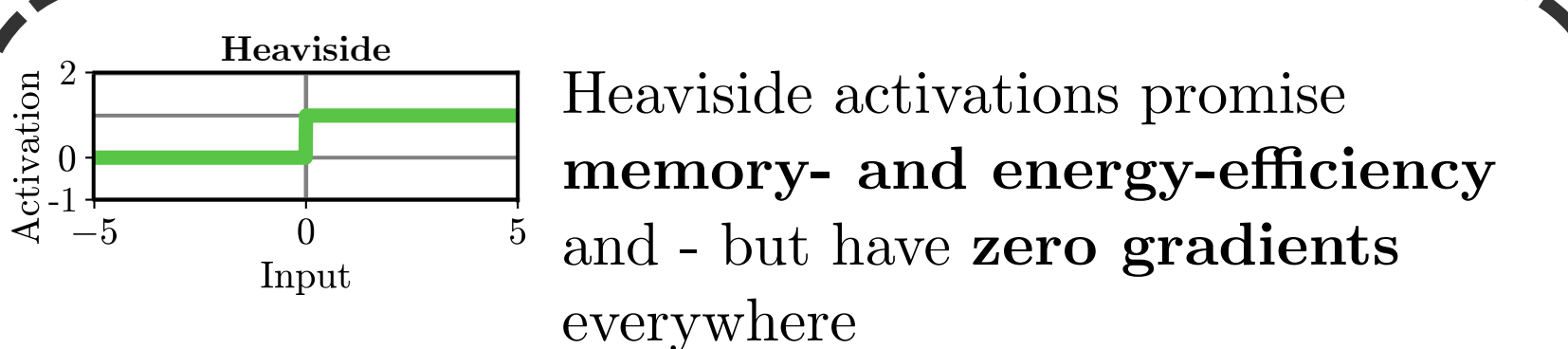
*specifically, we use LRP[l] here, but other local XAI methods may work as well

(C) Parameter Update

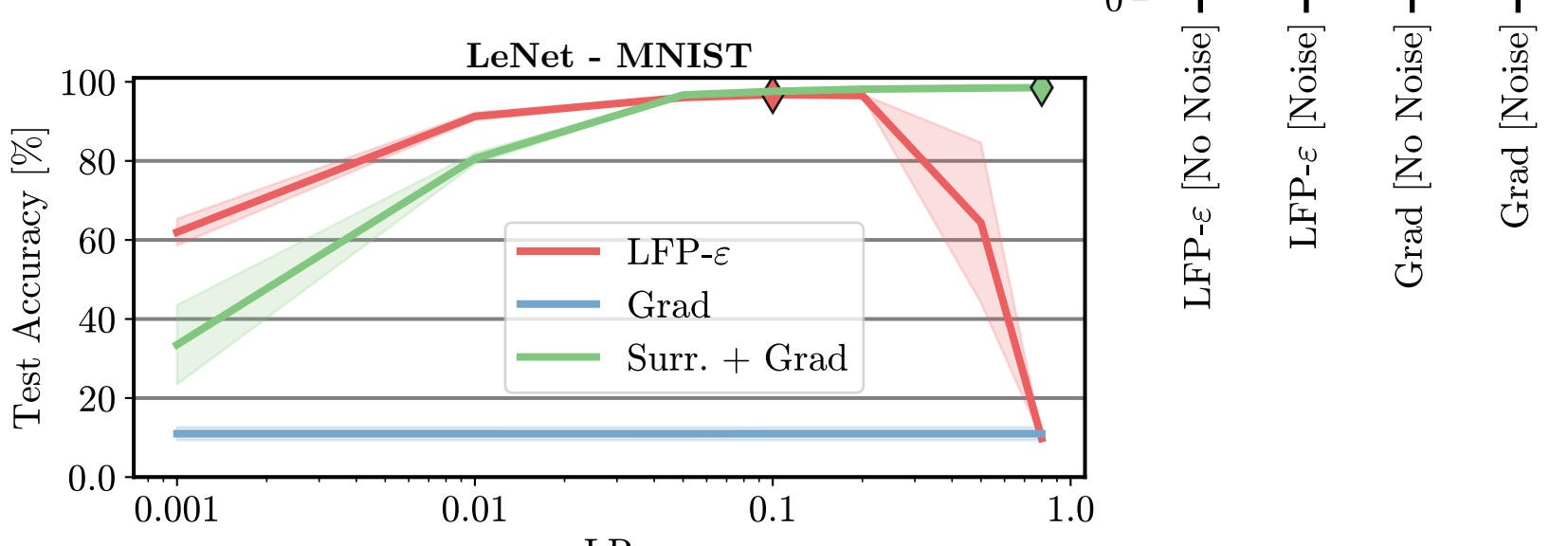


Results

Training Heaviside-activated Models

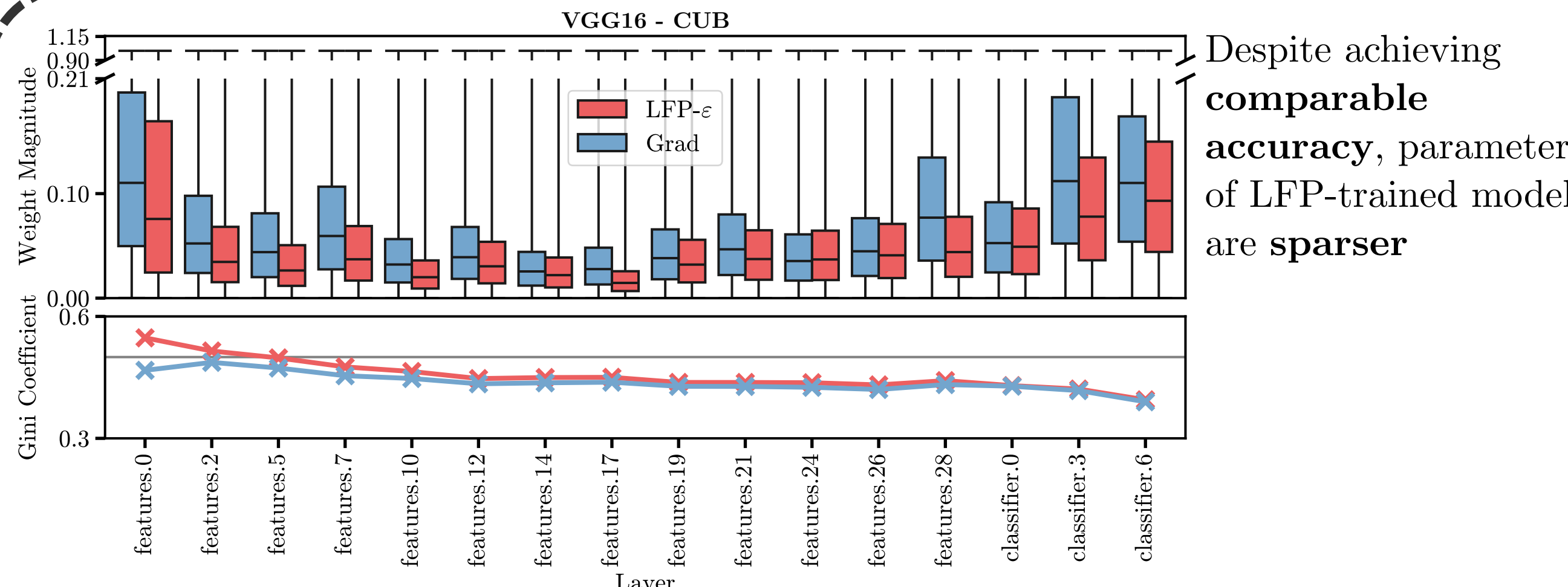


Training Heaviside-activated ANNs is **difficult** for backpropagation-based algorithms. However, adding noise to activations yields a **valid training signal** for LFP.

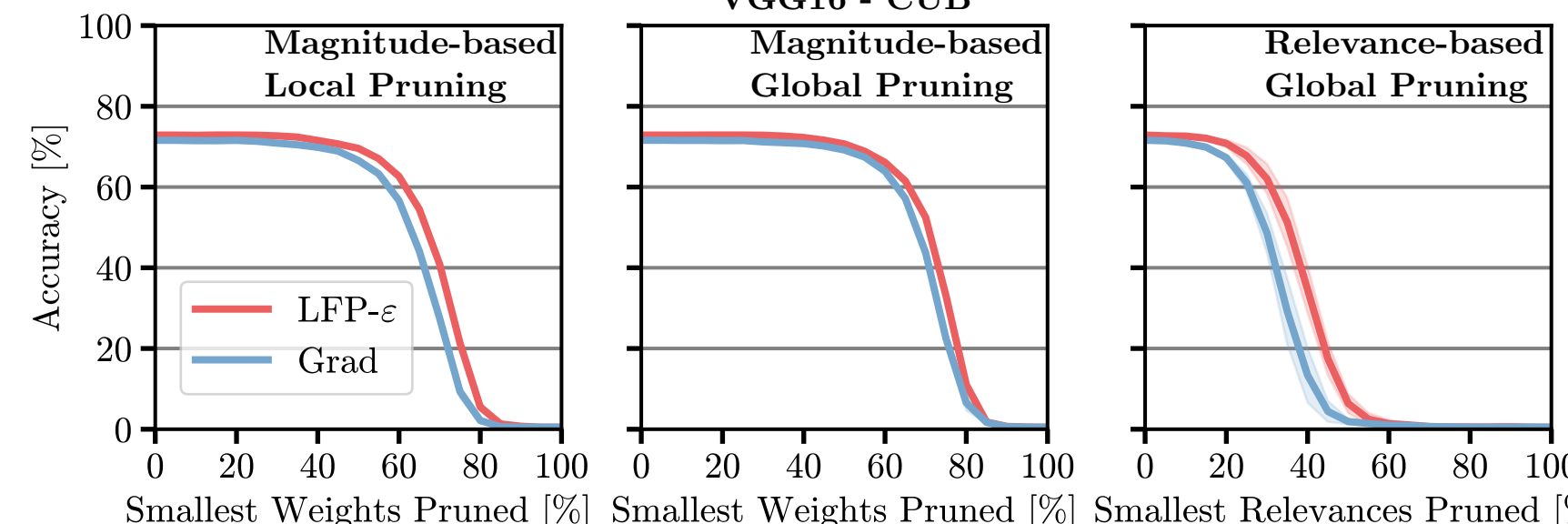


LFP is able to **naturally train SNNs** - **without requiring surrogates** for Heaviside in the backward pass.

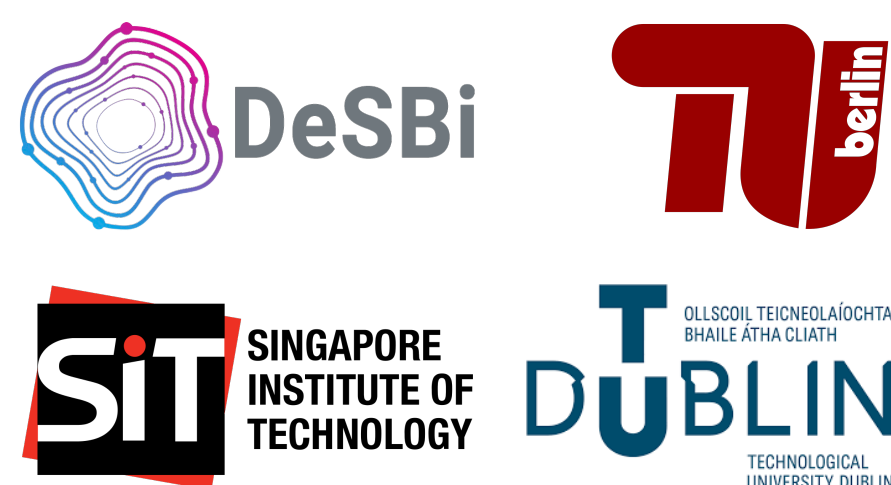
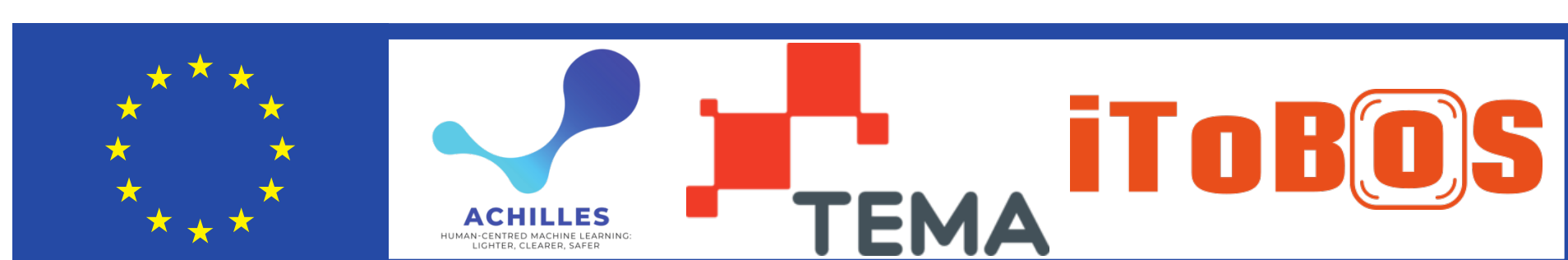
Obtaining Sparse and Efficient Representations



Consequently, LFP-trained models can be **pruned more easily**, indicating more **efficient representation** of the same information



Despite achieving **comparable accuracy**, parameters of LFP-trained models are **sparser**



Check out the paper...



... and the code!



[1] Bach, Sebastian, et al. "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation." PloS one 10.7 (2015): e0130140.

[2] Achitibat, Reduan, et al. "AttnLRP: Attention-Aware Layer-Wise Relevance Propagation for Transformers." Proceedings of the 41st International Conference on Machine Learning, in Proceedings of Machine Learning Research (2024): 235:135-168

