

# VISION DIFFMASK: Faithful Interpretation of Vision Transformers with Differentiable Patch Masking

A. Nalmpantis\*, A. Panagiotopoulos\*, K. Papakostas\*, J. Gkountouras\*, W. Aziz



← Live Demo  
Codebase →

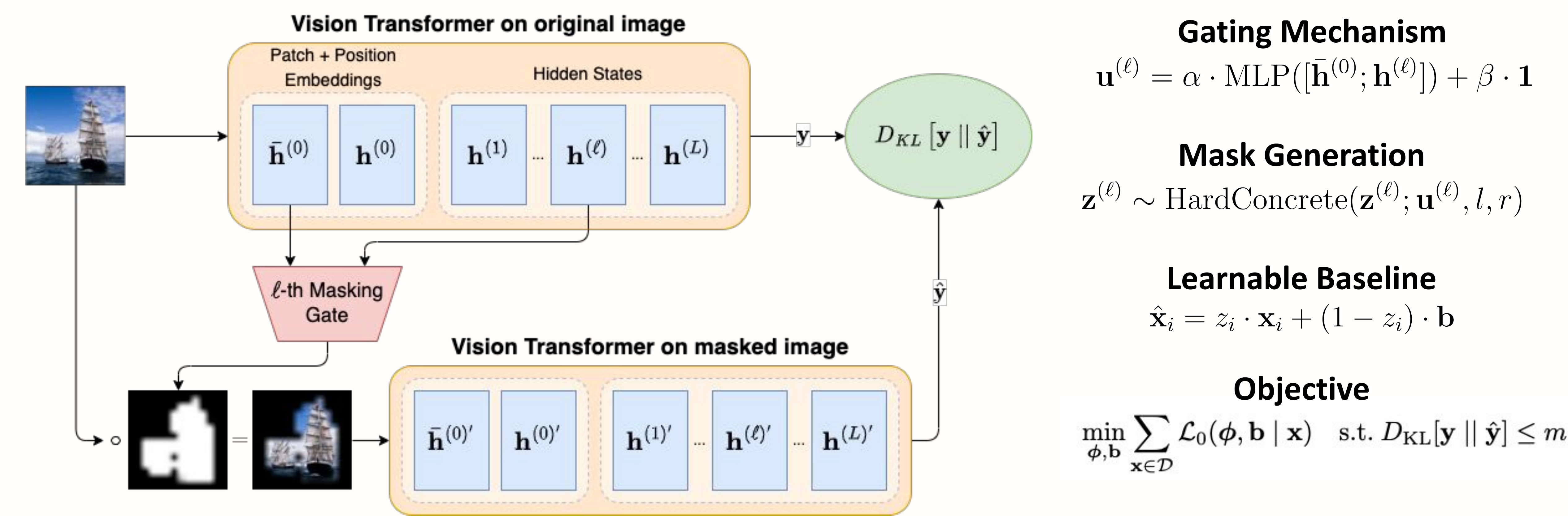


This work received funding from the EU's Horizon Europe RIA (UTTER, contract 101070631).

## Introduction

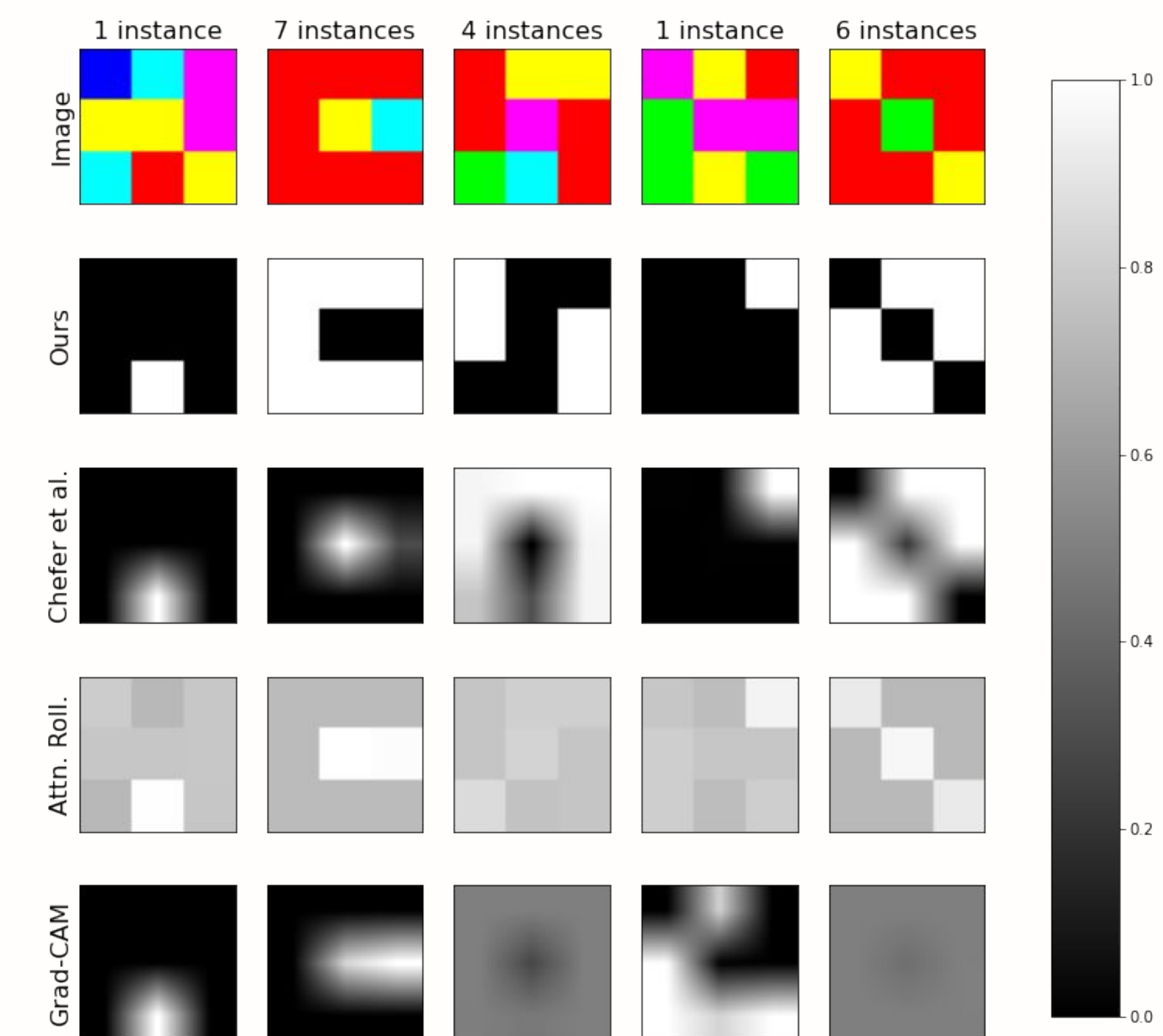
- We introduce a novel post-hoc interpretation method for the Vision Transformer, based on DiffMask [1].
- Our method is *neither* a gradient-based *nor* an attention-based method
- Instead, our method trains a **set of probing models** to find the **minimal subset of an image** that produces the same output distribution with the whole image
- Next to the standard benchmarks, we introduce a **new task** to measure the **faithfulness** of interpretability methods.

## Methodology

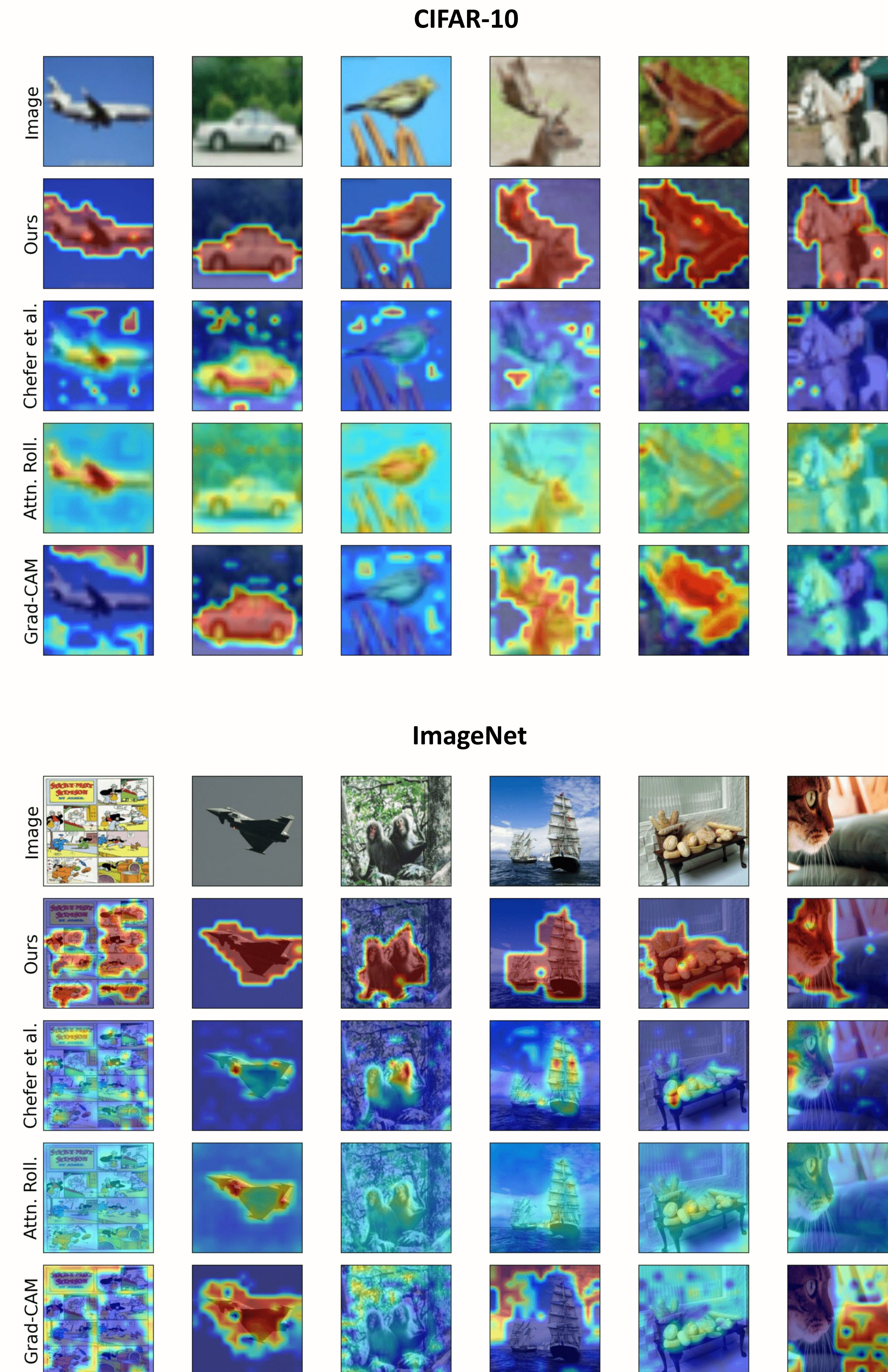


## Faithfulness Test

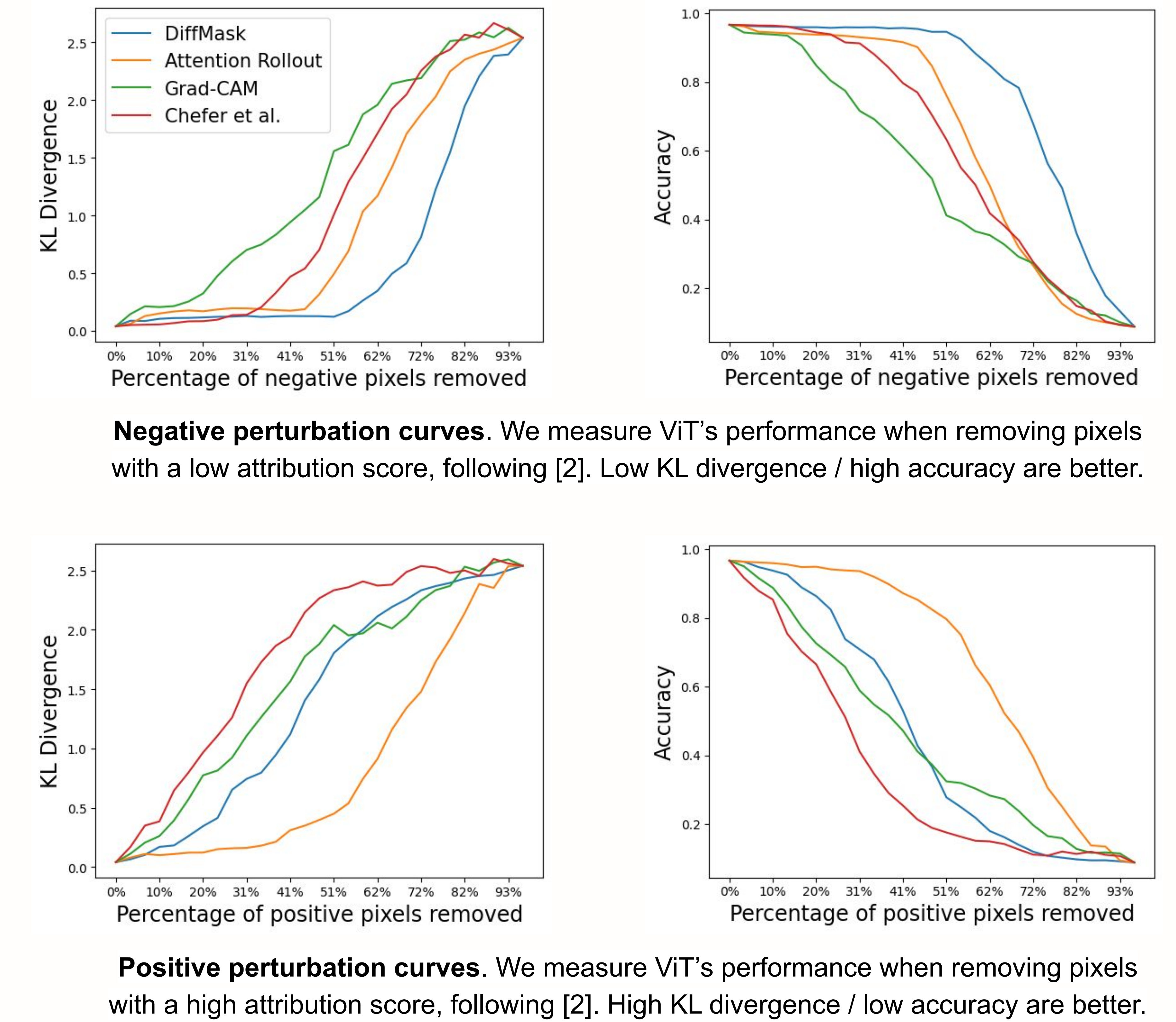
- We train a ViT to *count the number of red patches in an image*
- In this scenario, we **know** what the model should be looking at to make a decision: *either or the red or all the non-red patches*
- Vision DiffMask is the *only* method that is truly **faithful** to the behavior that we anticipate.
- Other methods *both* fail some case *and* are inconsistent between choosing the red patches or their complement



## Qualitative Results



## Quantitative Results



## Conclusion

- We introduced VISION DIFFMASK, a new method for post-hoc interpretability of images.
- VISION DIFFMASK's attributions are experimentally proven to be **consistent** and **faithful**, while also **plausible** to what a human would expect.
- Our method is able to **generalize** well across different datasets (please check our demo 🤖)

## References

- [1] De Cao et al. *How do Decisions Emerge across Layers in Neural Models? Interpretation with Differentiable Masking.* EMNLP 2020
- [2] Chefer et al. *Generic Attention-model Explainability for Interpreting Bi-Modal and Encoder-Decoder Transformers* ICCV 2021