# Interpretable Convolutional Neural Networks

2022.12.09

Ph.D. Student  : Neşe Güneş, M.Sc.
Thesis Advisor : Doc. RNDr. Elena Šikudová, Ph.D.

Interpretable Convolutional Neural Networks

# About article

IEEE Conference on Computer Vision and Pattern Recognition, CVPR

Conferences > 2018 IEEE/CVF Conference on C... ❓

## Interpretable Convolutional Neural Networks

**Publisher: IEEE** | Cite This | 📄 **PDF**

Quanshi Zhang ; Ying Nian Wu ; Song-Chun Zhu   **All Authors**

| **251** Paper Citations | **1270** Full Text Views |

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

ABSTRACT

# Clever Hans



## Right for the wrong reasons

Clever Hans performing in 1904

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).
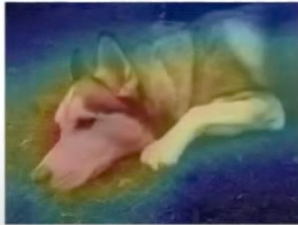
ABSTRACT

# Why interpretability?

- Confounding: Right for the wrong reasons – Clever Hans
- High-stakes decisions: Should patient get a biopsy?
- Responsibility: It's the doctor's responsibility to make a good decision
- Black box models turn computer-aided decisions into automated decisions
    - Doctors won't have the classification results
    - But explanations

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

ABSTRACT

# What about explanations?

- Explaining deep NNs with saliency maps does not work



| | Test Image | Evidence for Animal Being a Siberian Husky | Evidence for Animal Being a Transverse Flute |
|---|---|---|---|
| Explanations Using Attention Maps | | "Explanation" | |

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

# Problem scope

**Tabular Data**
- All features are interpretable
- Features include numerical and categorical data

| | |
|---|---|
| Age | 36 |
| Gender | F |
| Exercise? | Yes |
| Smoking? | No |
| Diabetes? | No |

**Raw data**
- Features are individually uninterpretable
- Pixels, voxels, words, a bit of sound wave



Měla na ruce nejkrásnější náramek , jaký jsem kdy viděl - secesní
víla se na něm proplétala mezi brilianty a smaragdy .
Tenhle náramek měl cenu luxusního auta , jenže jeho krása byla ještě

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

ABSTRACT

# Explainable or Interpretable?

nature > nature machine intelligence > perspectives > article

Perspective | Published: 13 May 2019

## Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead
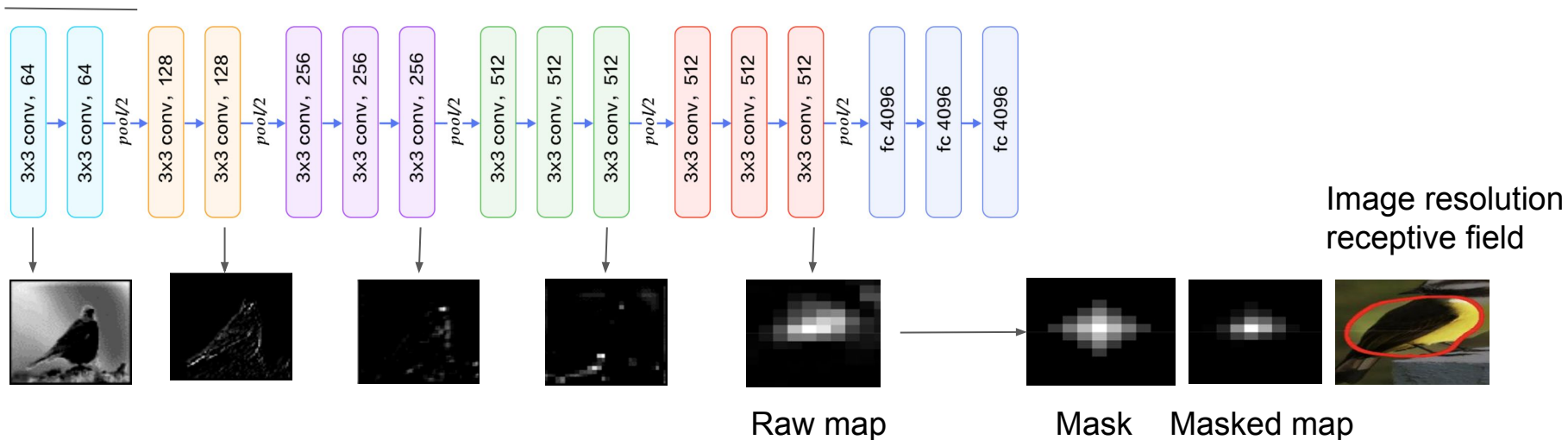
Cynthia Rudin ✉

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

# From feature maps to image regions



Raw map      Mask      Masked map

Image resolution receptive field

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).
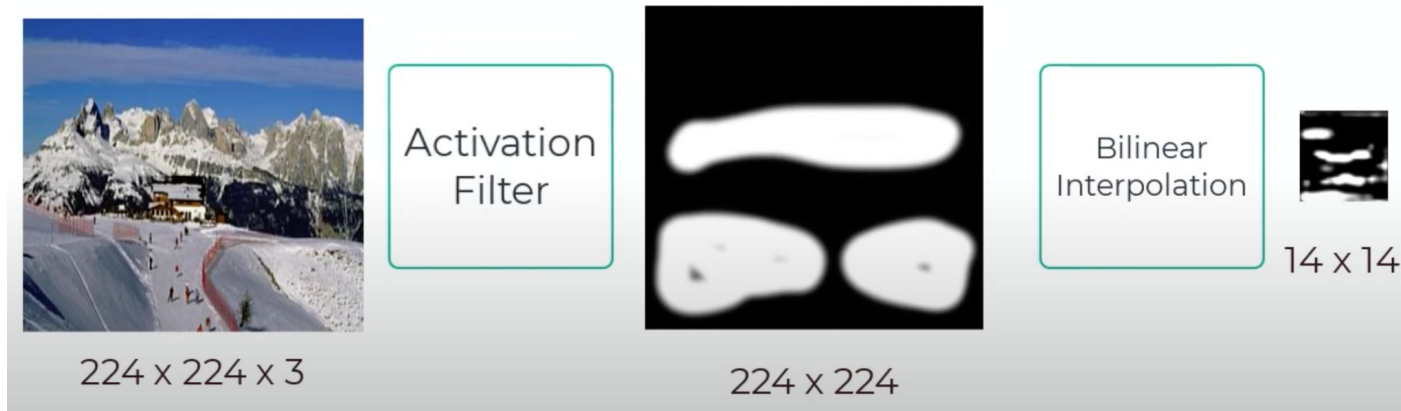
# From feature maps to "segmentation" based on filters



224 x 224 x 3     Activation Filter     224 x 224     Bilinear Interpolation     14 x 14

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).
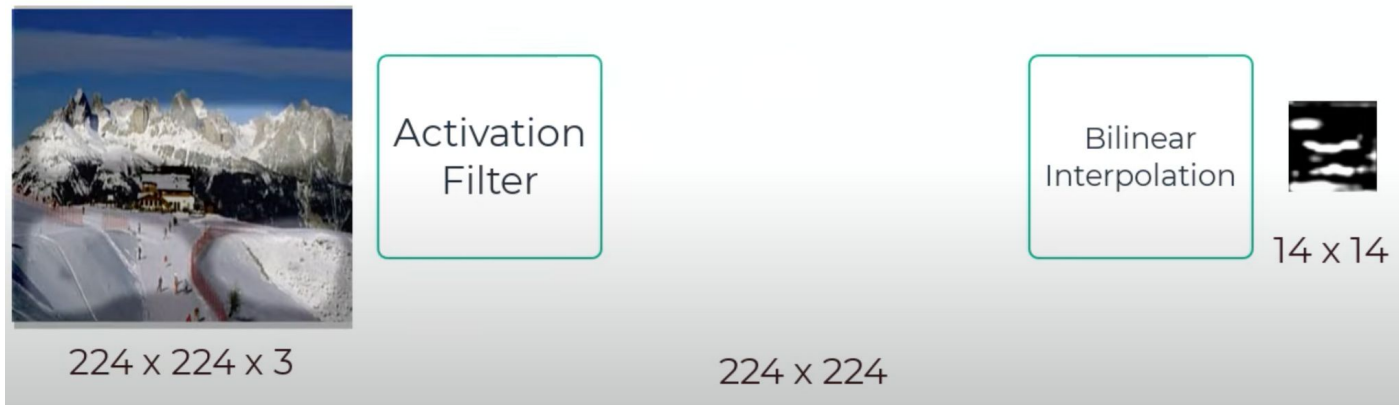
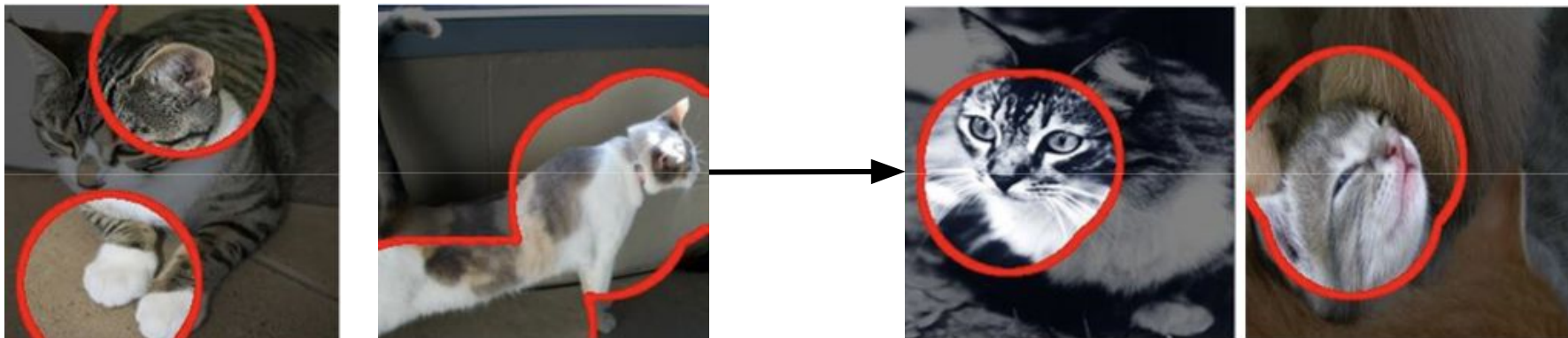# From feature maps to "segmentation" based on filters



[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).
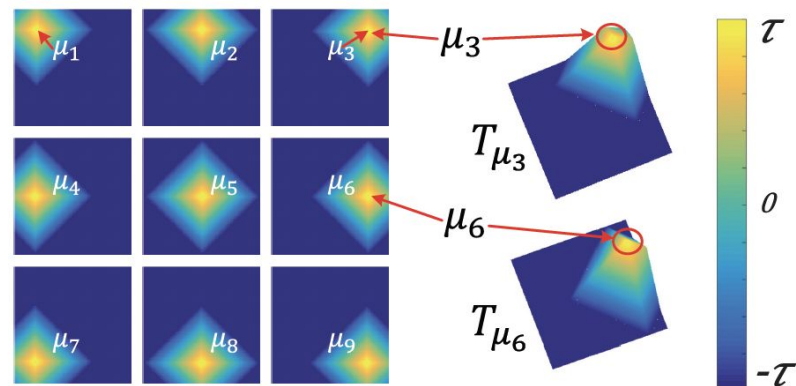
# In a nutshell

❏ From mixture of patterns to object-parts



[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

# How to regularize filters?

- ❏ Understanding the local filter loss
    - ❏ Forgetting irrelevant information
    - ❏ Filter out noisy activations
- ❏ Forward propagation
    - ❏ Part template selection
- ❏ Backpropagation
    - ❏ Determining a target category for each filter



[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).
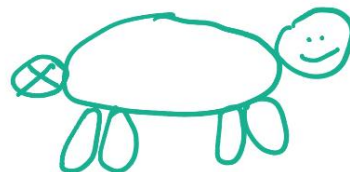
# Time to experiment

❏  Single-category, multi-category classification

  ❏  Back to binary classification!

❏  Metrics: object-part interpretability, location stability: Avoid Picasso-filters

❏  Ground truth annotations for evaluation: GTs are still necessary

❏  Four types of CNNs

  ❏  No ResNet: Who explains the ResNet wins the game

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

# Evaluation metrics

- ❏ Object-part interpretability
- ❏ Stability of object-part locations
    - ❏ Categories to filters
    - ❏ Location deviation of object-parts

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

# Discussion

- ❏ Designing filters only for top conv layer

    - ❏ Future work: designing filters for low conv layers

- ❏ Shared object-parts by different categories

- ❏ Segmentation datasets and annotations

[1] Zhang, Q., Wu, Y. N., & Zhu, S. C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8827-8836).

Neşe Güneş, Ph.D. Student in Ophthalmic Medical Imaging
Charles University, Prague, Czech Republic
E-mail: gunes.nese@matfyz.cuni.cz
Phone: +420 606 318 809

# Na zdraví!