# Classification-by-Components: Probabilistic Modeling of Reasoning over a Set of Components

Sascha Saralajew[1], Lars Holdijk[1], Maike Rees[1], Ebubekir Asan[1], Thomas Villmann[2]

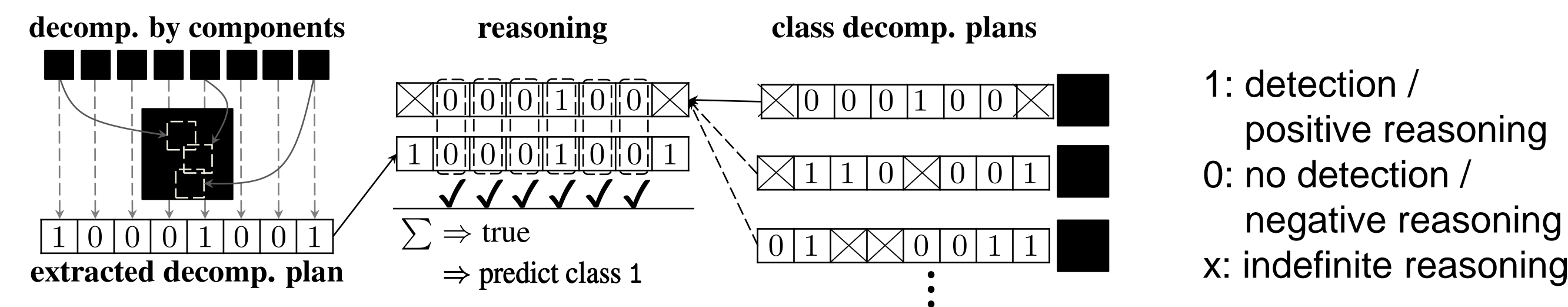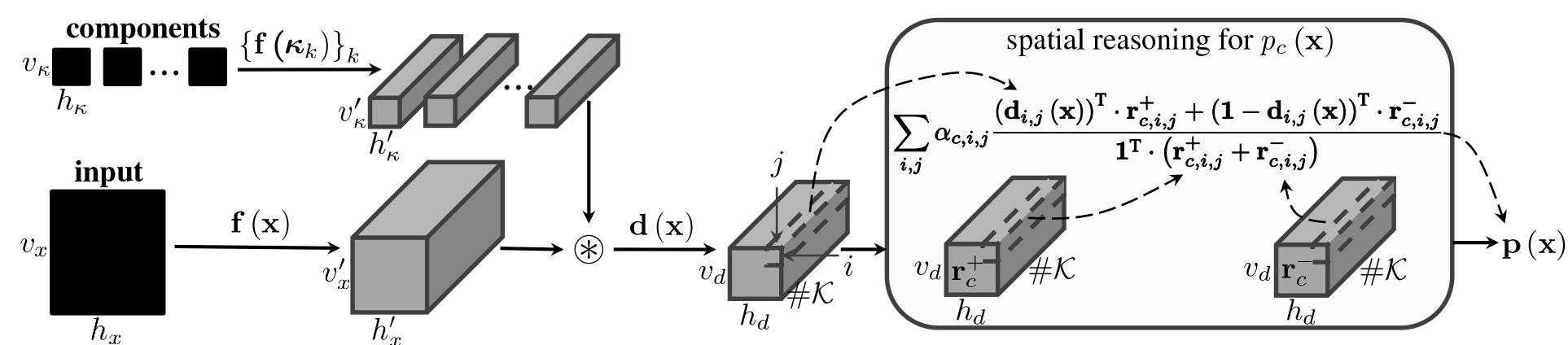[1]Dr. Ing. h.c. F. Porsche AG, [2]University of Applied Sciences Mittweida

## Introduction: Interpretability by Design

- Classification-By-Components networks (CBCs) are interpretable through their architecture
- Inspired by Recognition-by-Components theory [1]: 'Humans recognize complex objects by decomposing them into components'
- Classification output is computed with a probability tree (no softmax)



1: detection / positive reasoning
0: no detection / negative reasoning
x: indefinite reasoning

## The Model: Components and Decomposition Plans



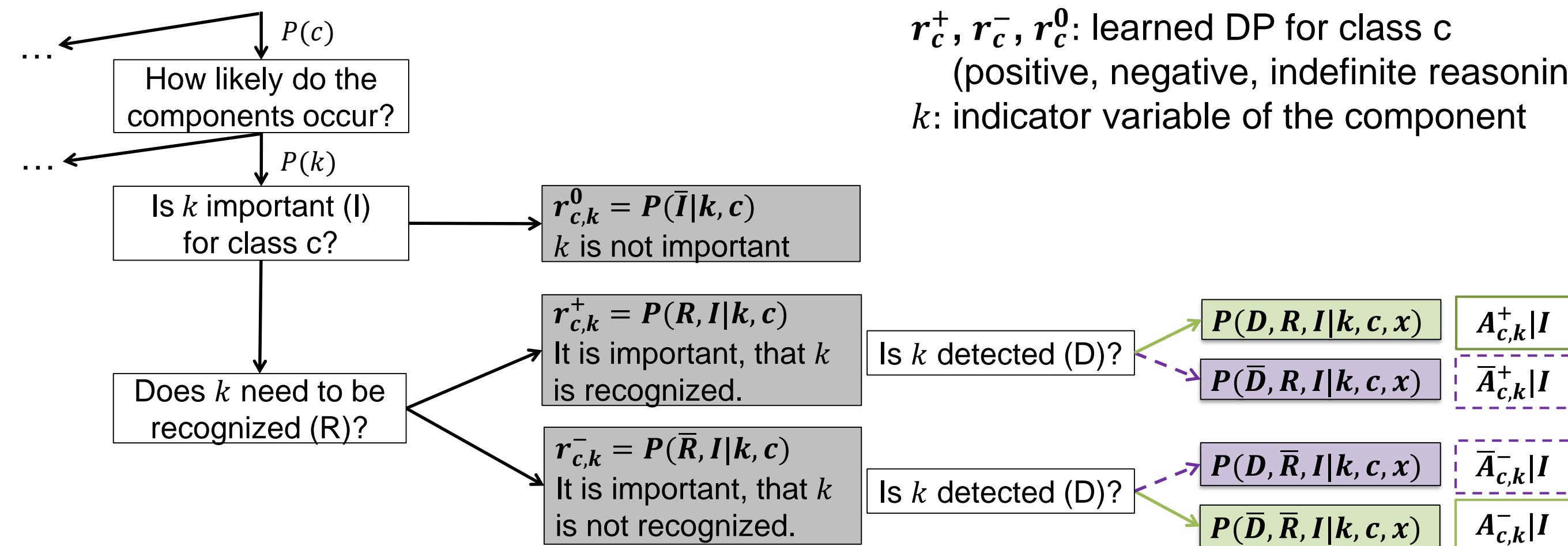**Components:** generic features, either full (input size) or patch (< input size)
**Reasoning:** one Decomposition Plan (DP) per class.
**At inference:** Match extracted DP with learned DPs.
**Training:** With contrastive loss and SGD

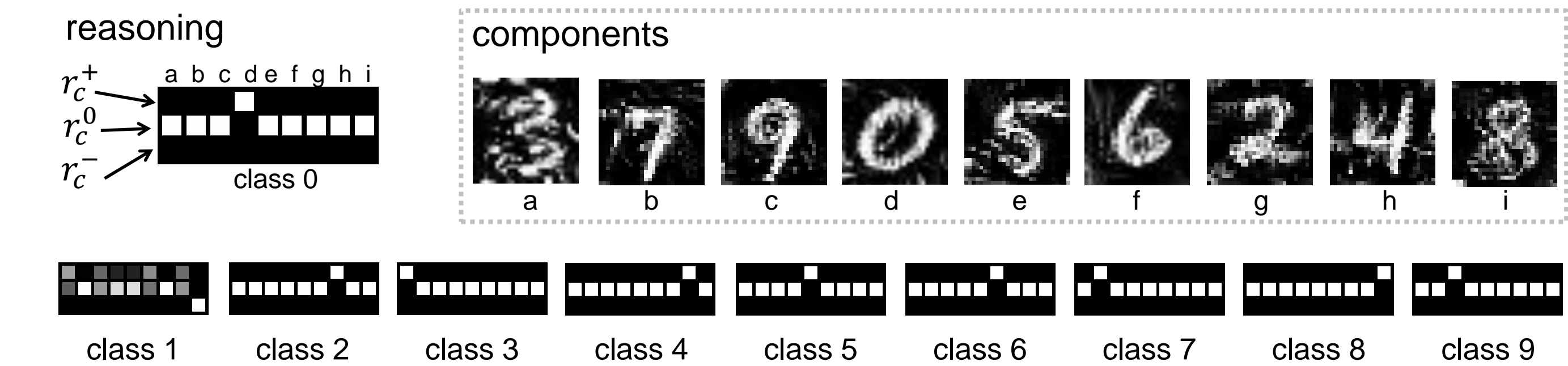$$p_c(x) = \frac{(d(x))^T * r_c^+ + (1 - d(x))^T * r_c^-}{1^T * (1 - r_c^0)}$$

$x$: input
$d(x)$: component detection possibility vector (DP extracted from the input)
$p_c(x)$: class hypothesis probability for class c
$r_c^+, r_c^-, r_c^0$: learned DP for class c (positive, negative, indefinite reasoning)
$k$: indicator variable of the component



## Results: Interpretable Classification

### Components are prototypical and interpretable (MNIST):

To classify a 1, the CBC depends on negative (not look like an 8) and positive reasoning. For the other classes, the CBC depends on strong positive reasoning over one and indefinite reasoning over all the other components.
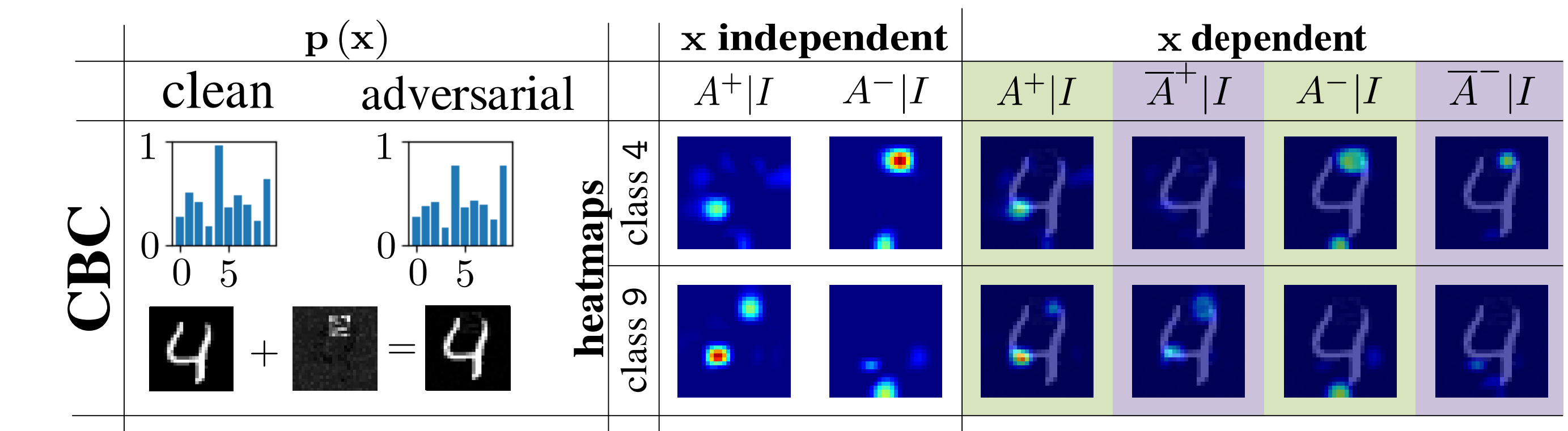


### Components are conceptionally meaningful (ImageNet):

The components with the highest $r_{c,k}^+$ for *dalmatian* (top) and *giant panda* (bottom) are displayed. Components are shared among classes which shows that CBCs are capable of learning complex, class independent structures.



1.00 1.00 1.00 1.00 0.99 0.98 0.88 0.84 0.72 0.70

1.00 1.00 1.00 1.00 1.00 0.98 0.75 0.65 0.61 0.61

### Reasoning explains why an adversarial attack succeeds (MNIST):

The CBC learns to recognize only as few parts as needed to distinguish between classes. For the 4, it checks that there is no stroke at the bottom ($A^-|I$) but a corner at the left ($A^+|I$). The adversary attacks a 4 to being classified as a 9 by inserting exactly enough noise to make the CBC believe that there is a stroke at the top ($\overline{A}^-|I$).



Accuracies: MNIST: 99.32 (CNN: 99.8), ImageNet: 82.4 (AlexNet: 82.8)

## This paper in 30seconds!

- Classification-by-Components networks are probabilistic and interpretable by design.
- They can be stacked onto almost every existing classification NN with almost no computational overhead at inference.
- The concept is to decompose complex objects into components and then compare the extracted with a learned class decomposition plan.
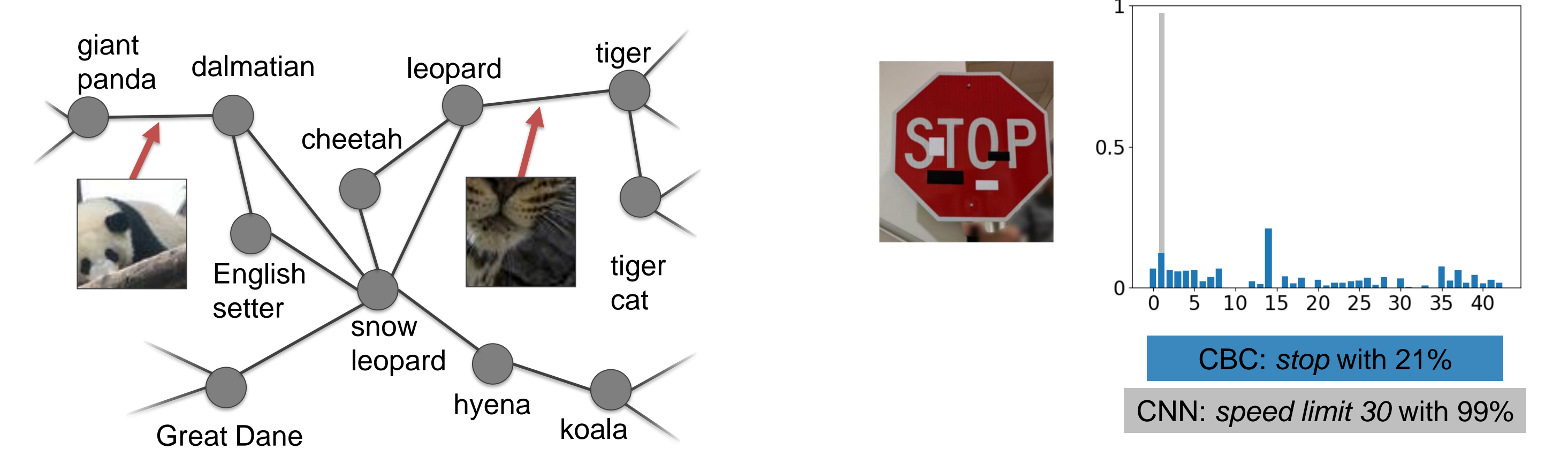
## Conclusion: Interpretability at Scale

- Probabilistic framework for the last and penultimate layer of a NN
- Intuitive visualizations by back projection, heatmaps of DP for agreement ($A$) and disagreement ($\overline{A}$)
- Achieves satisfying accuracies for different datasets with small or large number of classes.

## Current Work: Exciting Improvements in the Making

Investigation of components, that are shared by classes, can reveal a hierarchical structure between classes.

The class hypothesis possibility vector of a CBC might be used to reject adversarial examples.



CBC: *stop* with 21%
CNN: *speed limit 30* with 99%

**Inspirations and Theory**

[1] I. Biederman. Recognition-by-components: A theory of human image understanding. *Psychological review*, 94(2):115, 1987.

[2] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206, 2019