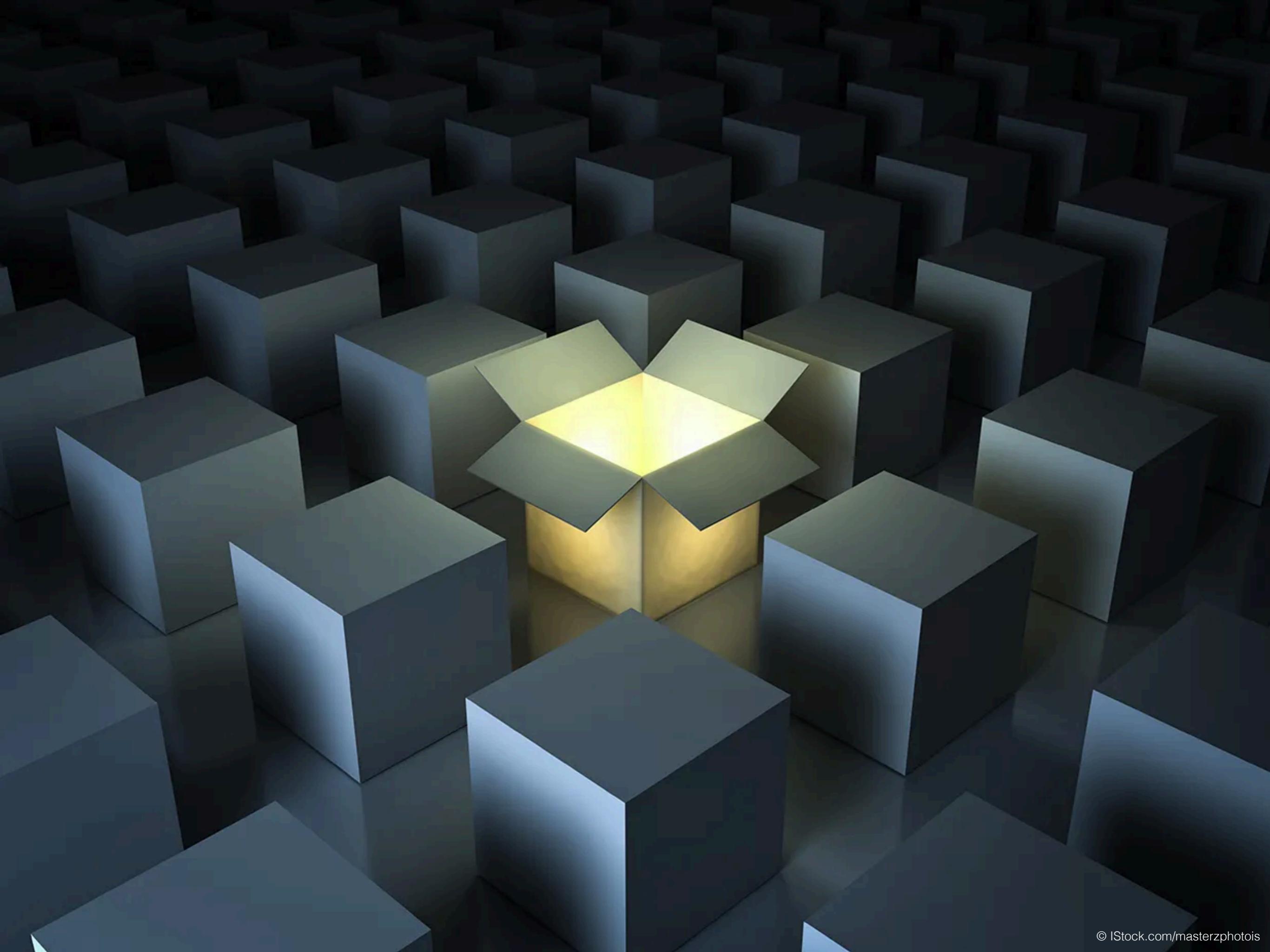




# Explainable AI

## I. Einführung

2. Modell-spezifische Methoden
3. Modell-unspezifische Methoden



# Clever Hans



© Wikipedia

# Explainable AI

## Interpretierbarkeit

Kann ein Mensch den Grund für eine Entscheidung verstehen?

Kann ein Mensch die Entscheidung eines Modells vorhersagen?

## Erklärbarkeit



Ansätze, um interne Funktionen oder Abläufe des Modells zu verdeutlichen

## Warum ist Interpretierbarkeit wichtig?

- ↳ Akzeptanz erhöhen in public mind
- ↳ Debugging (Bias z.B. erkennen)
- ↳ Fairness gewährleisten
  - z.B. KI entscheidet über Credit  
→ muss objektiv und fair sein

# Explainable AI

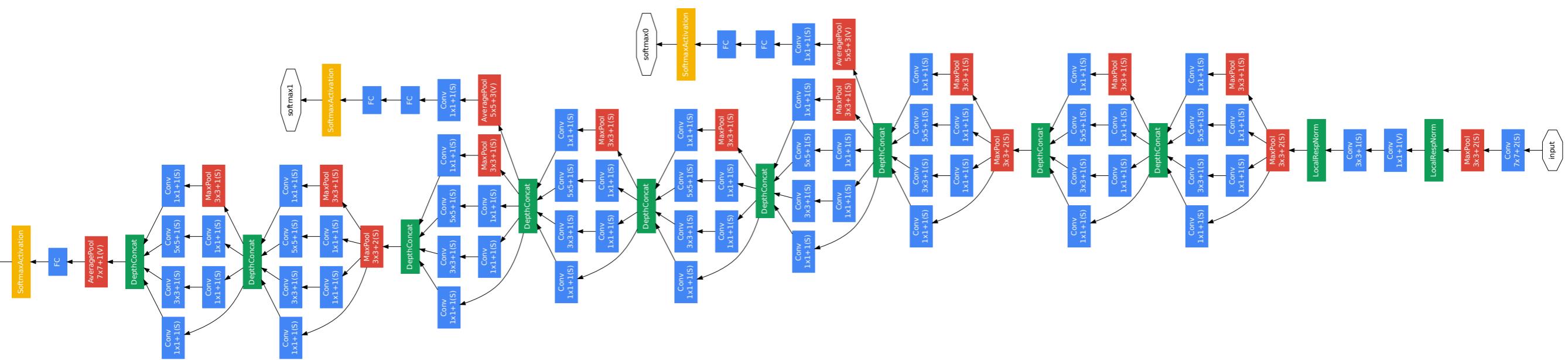
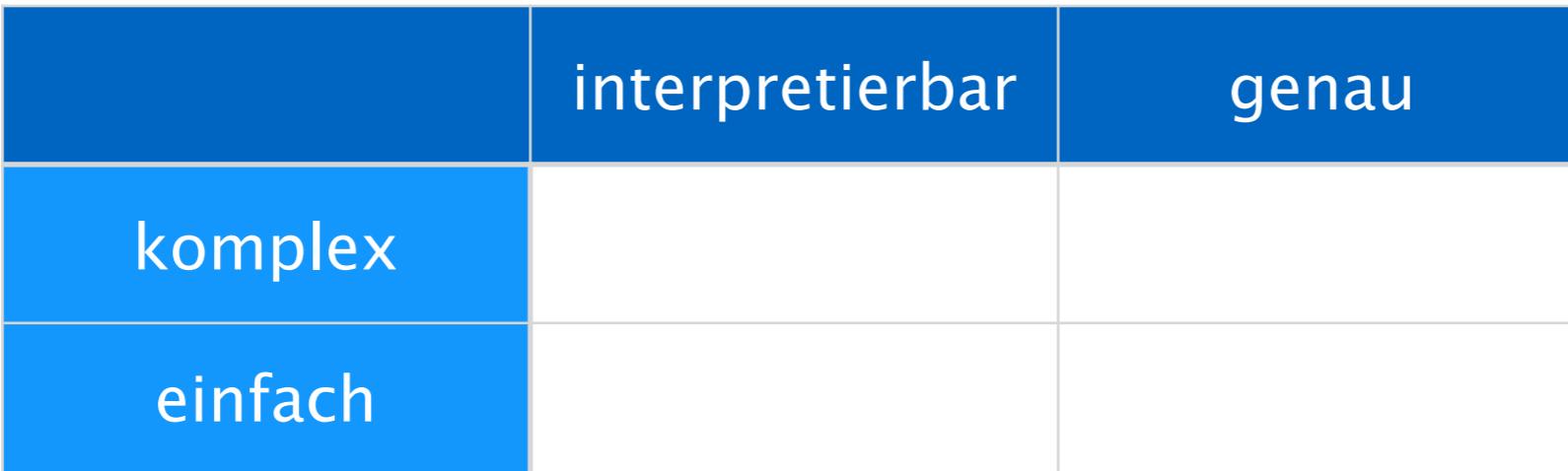
Sehr komplex  $\rightarrow$  sehr genau  $\rightarrow$  nicht einfach interpretierbar

↪ - post-hoc-Explanation

	interpretierbar	genau
komplex	X	✓
einfach	✓	X

Entweder einfache Modelle genau zu machen  
oder komplexe Modelle versuchen zu interpretieren = post-hoc-explanation

# Explainable AI



© Szeged, Liu, Jia, Vermahnet, Reed, Anguelov, Ethan, Vanhoucke, Rabinovich

## Post Hoc Explanation

↳ je tiefer und komplexer ein NN, umso schwieriger ist die Interpretation

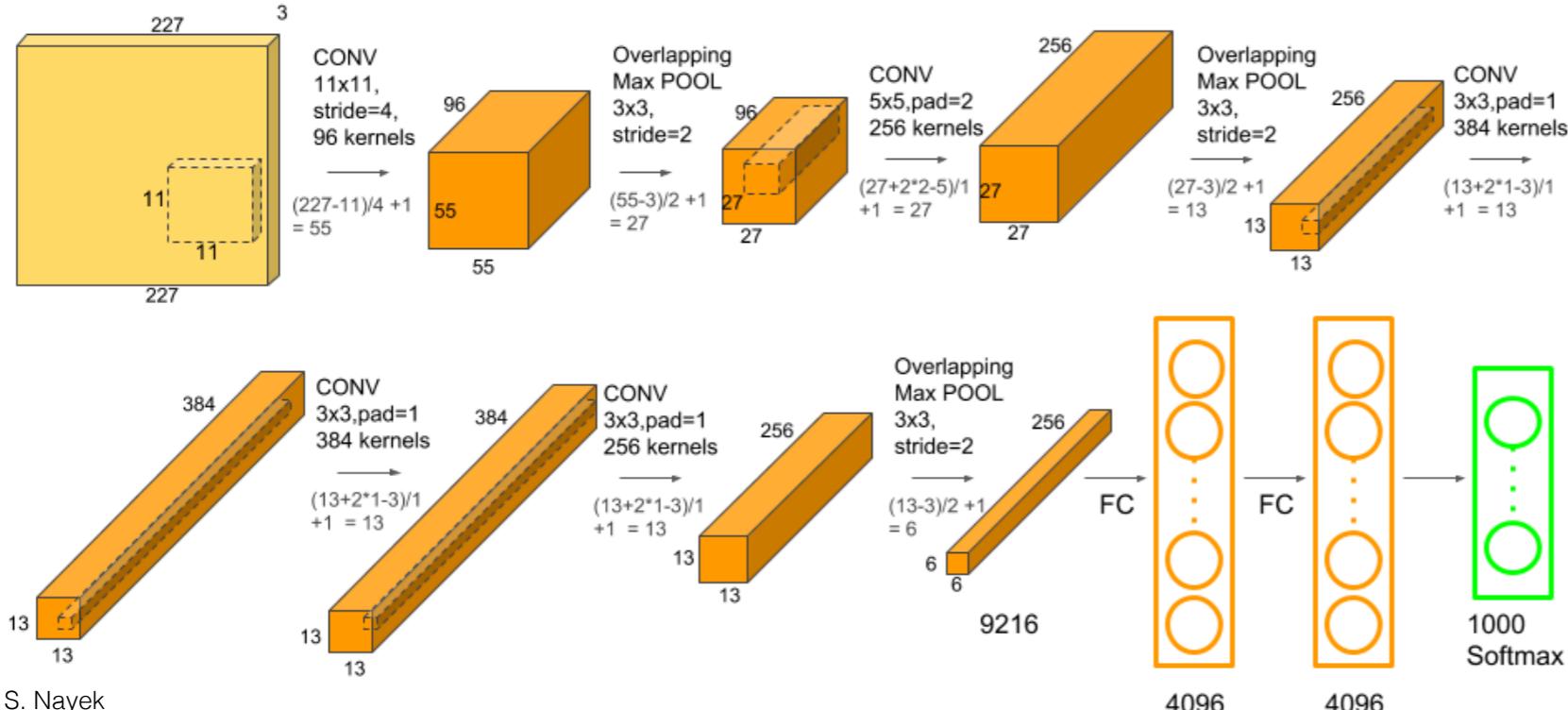


# Explainable AI

1. Einführung
- 2. Modell-spezifische Methoden**
3. Modell-unspezifische Methoden



# Visualisierung der Gewichte



Alex Net

Gewicht

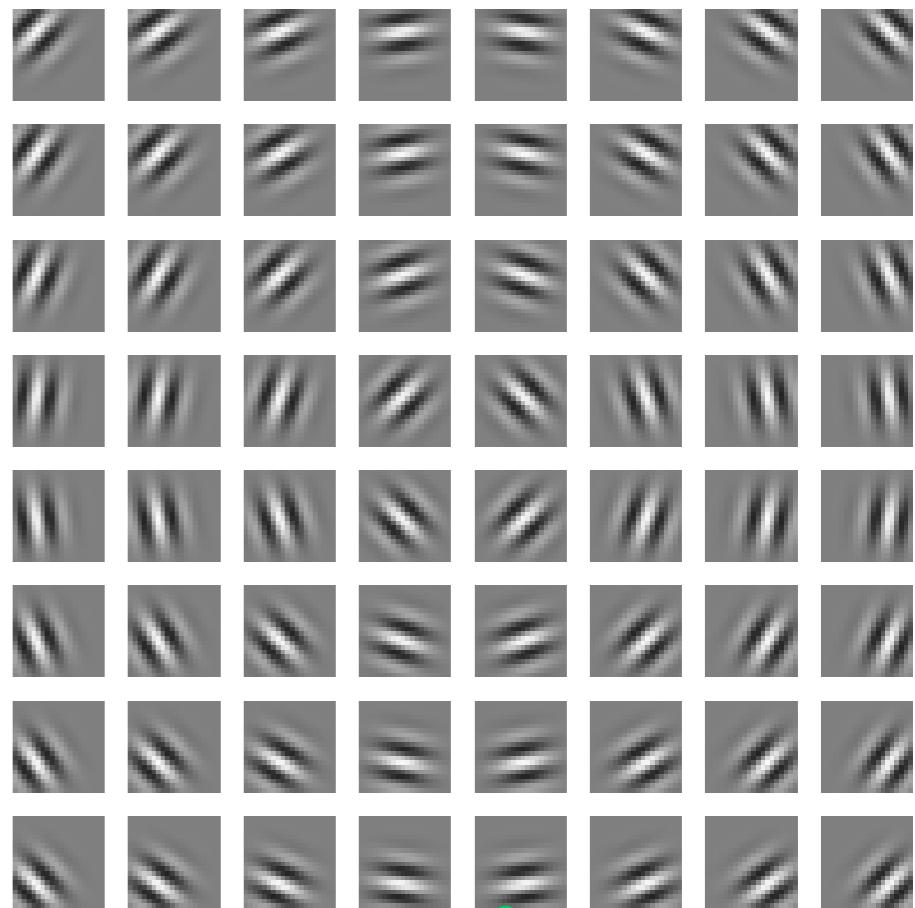
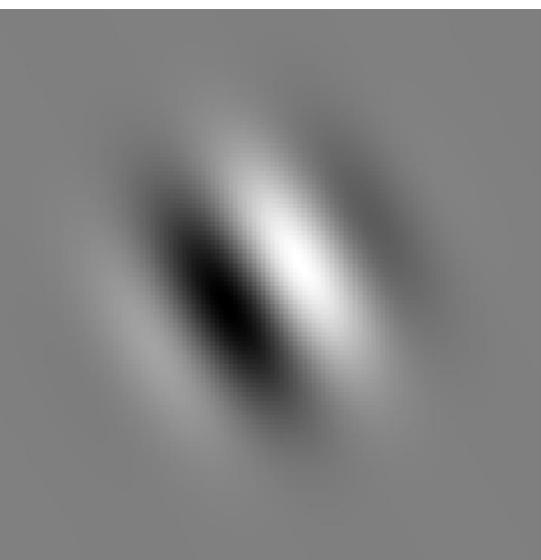
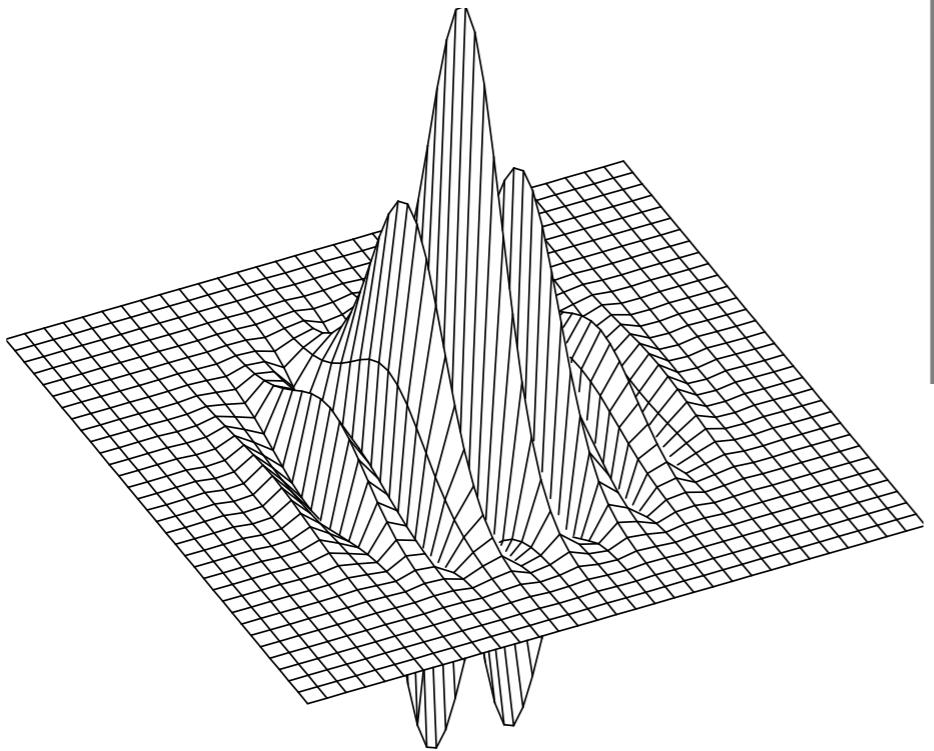
1. Schicht, Faltungsschicht visualisiert



© A. Krizhevsky et al., 2012

# Visualisierung der Gewichte

Gabor - Filter



4 Lekale

Fouriertransformation

gründliche Operation, wo Bild in Frequenzraum überführt wird



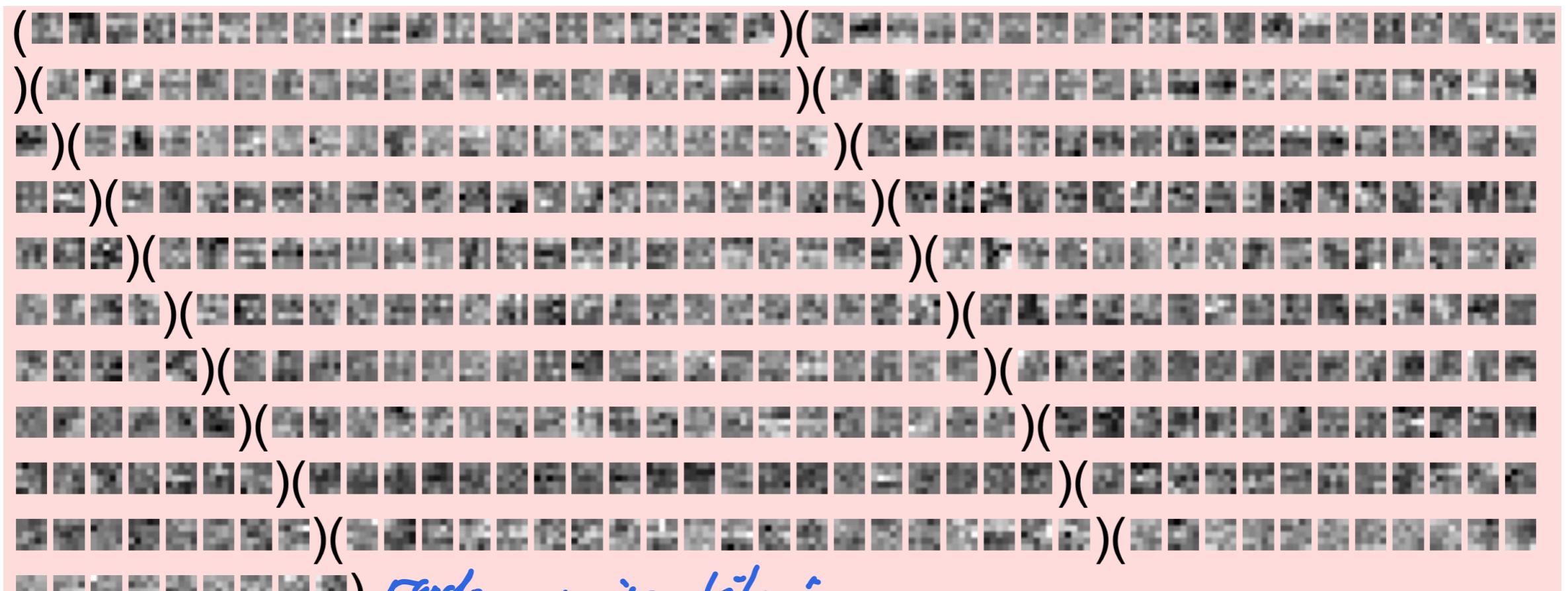
Filterbank

Tiere wie seltsame  
Geben  
© Goodfellow

geleert

# Visualisierung der Gewichte

Faltungshöhe höher Schichten

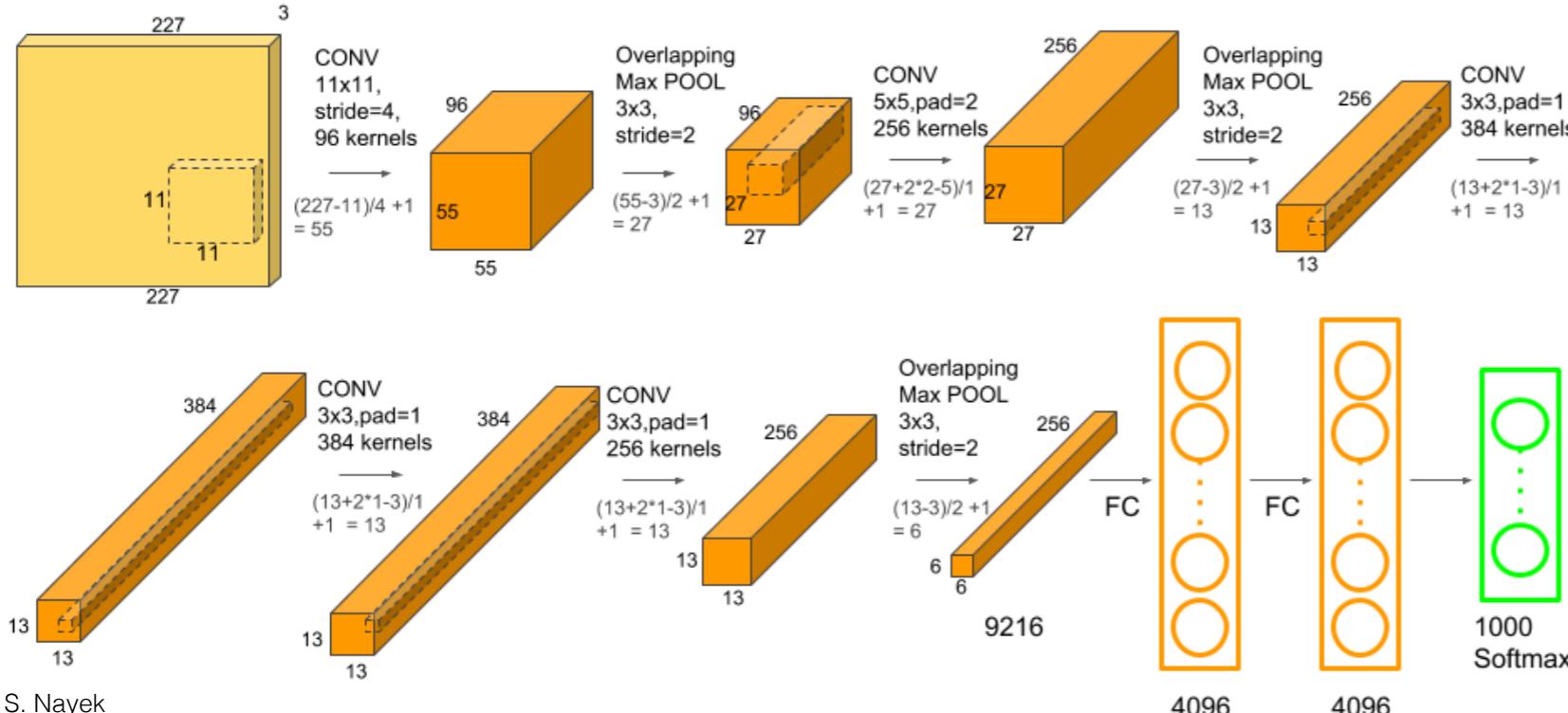


) random noise-like

- ned so gute Idee die zu visualisieren  
muss

© Andrej Karpathy

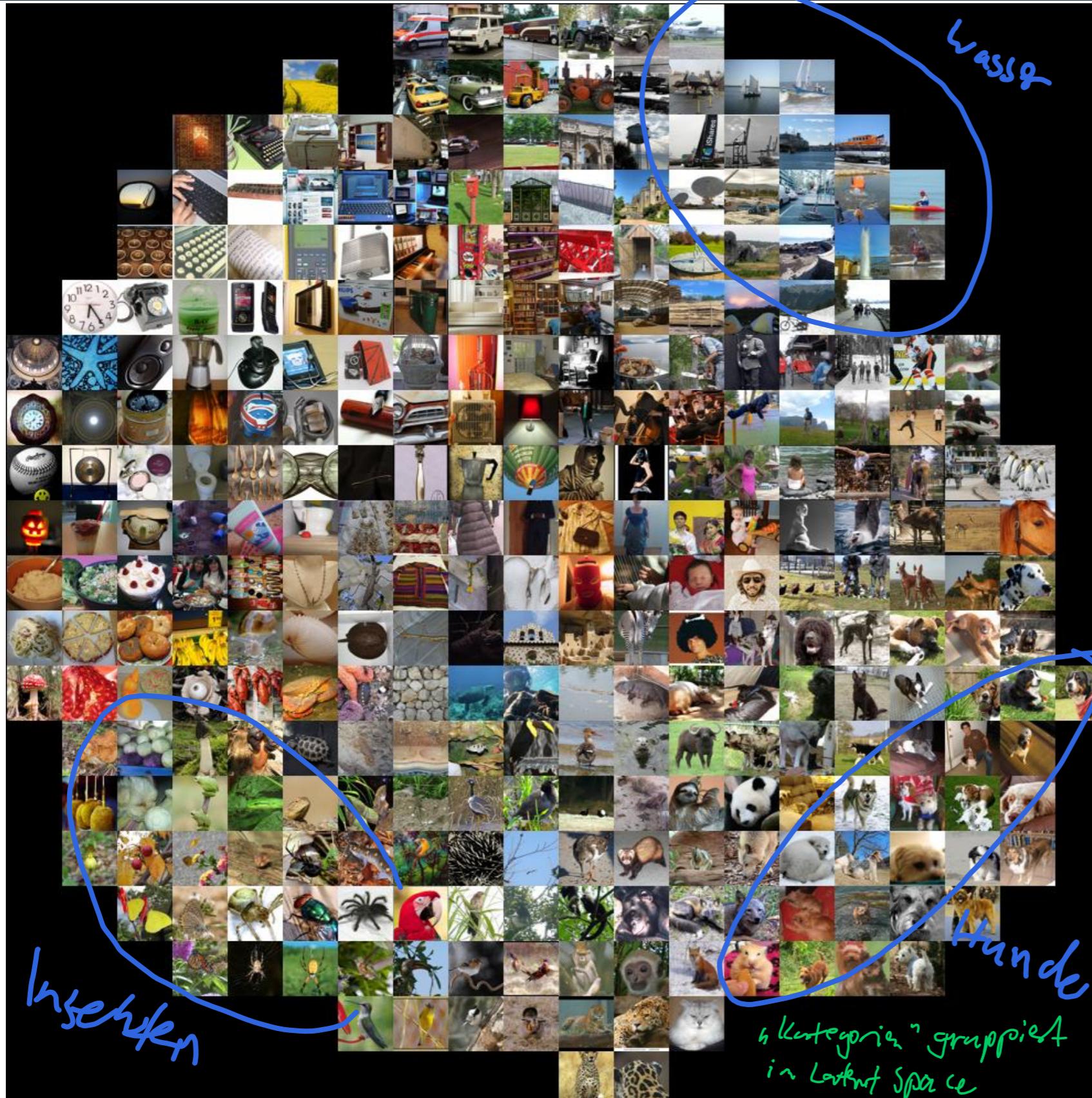
# Visualisierung der Merkmale



© S. Nayek

- ↳ latent space nur visualisieren?
- ↳ Rekurrenz statt Geometrie aussehen
- bild-spezifisch statt nichtspezifisch
- für jedes Bild anders Position im latent space

# Visualisierung der Merkmale



t-SNE-Visualisierung  
des Latent Space des  
Alex-Net  
wo Daten nicht landet wurde das ganze Bild  
darauf geklebt davon  
=> Bild als Deskriptor

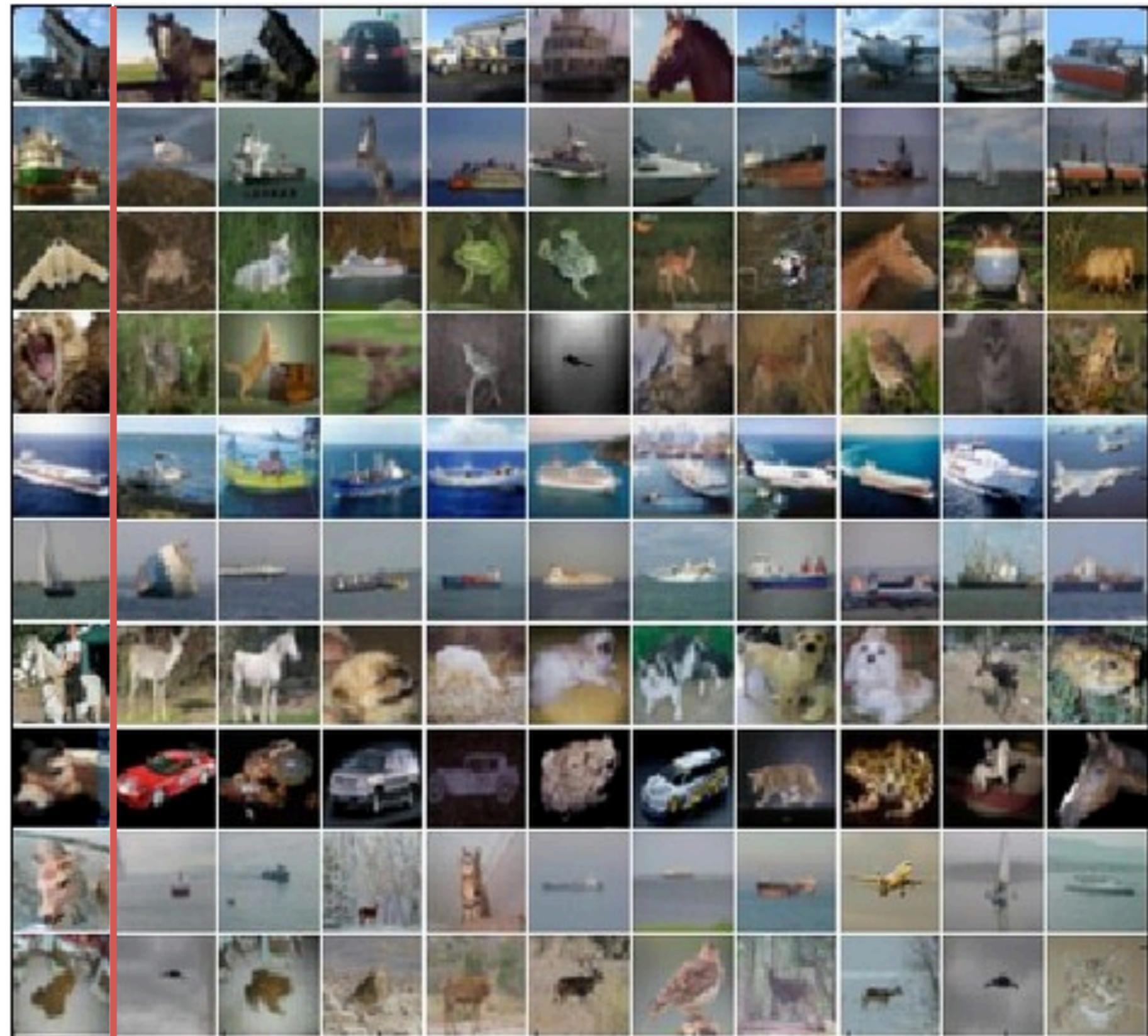
© Andrej Karpathy

# Visualisierung der Merkmale

andere Methode:

Testbilder

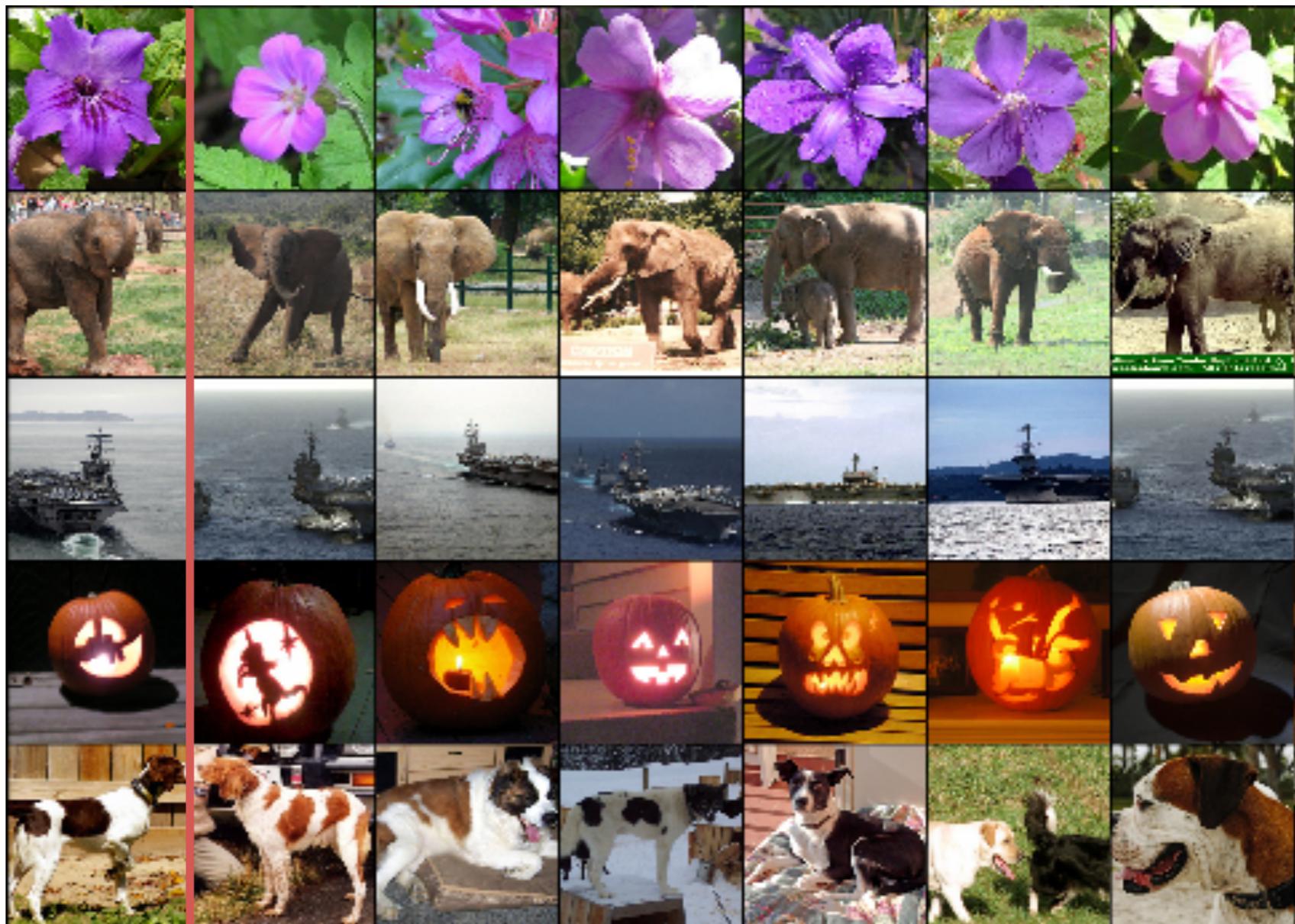
$L_2$ -Norm auf  
Pixelbasis



© Andrej Karpathy

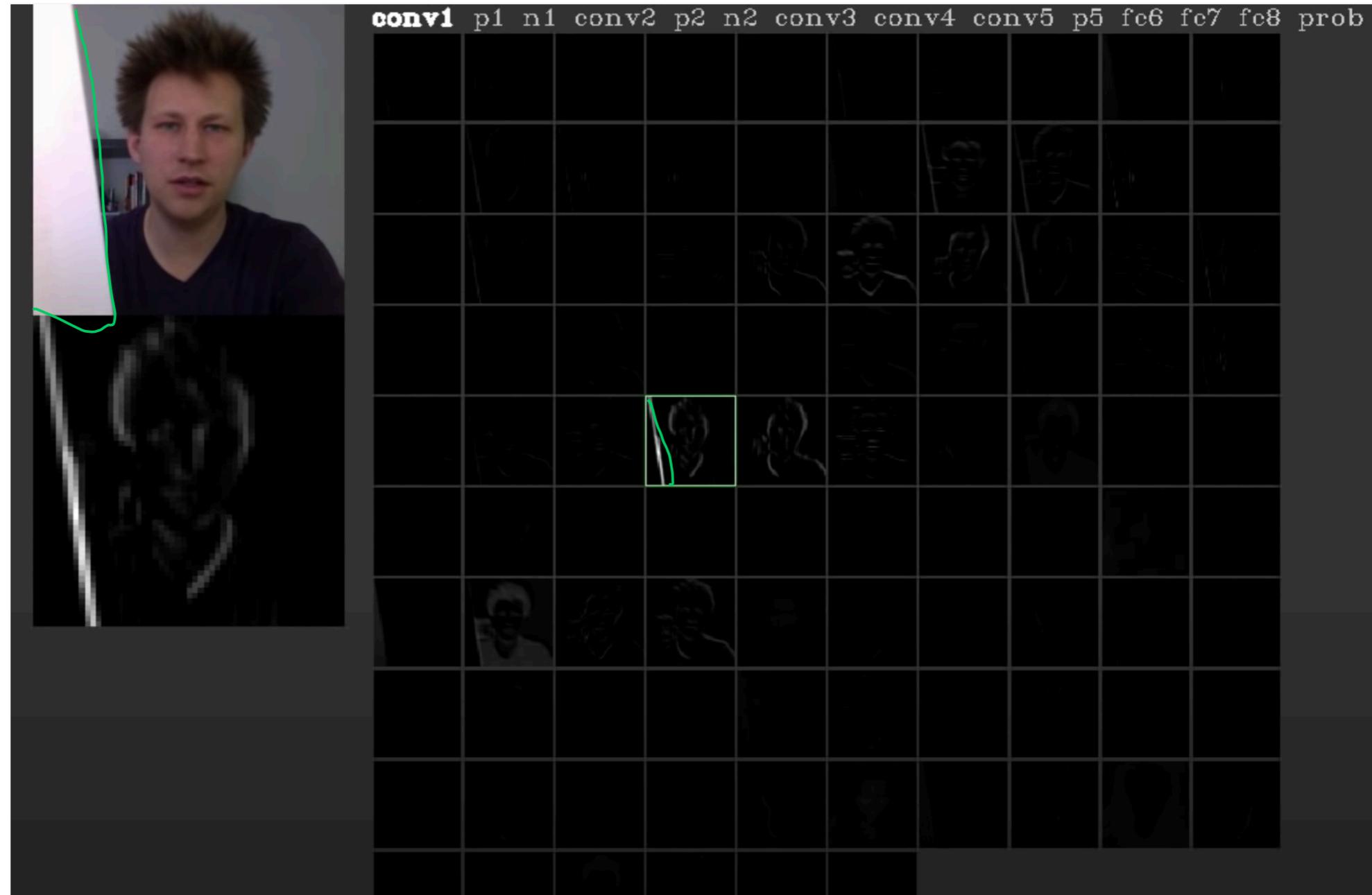
# Visualisierung der Merkmale

$\mathcal{L}_2$ -Norm auf Merkmalsebene des Latent Space



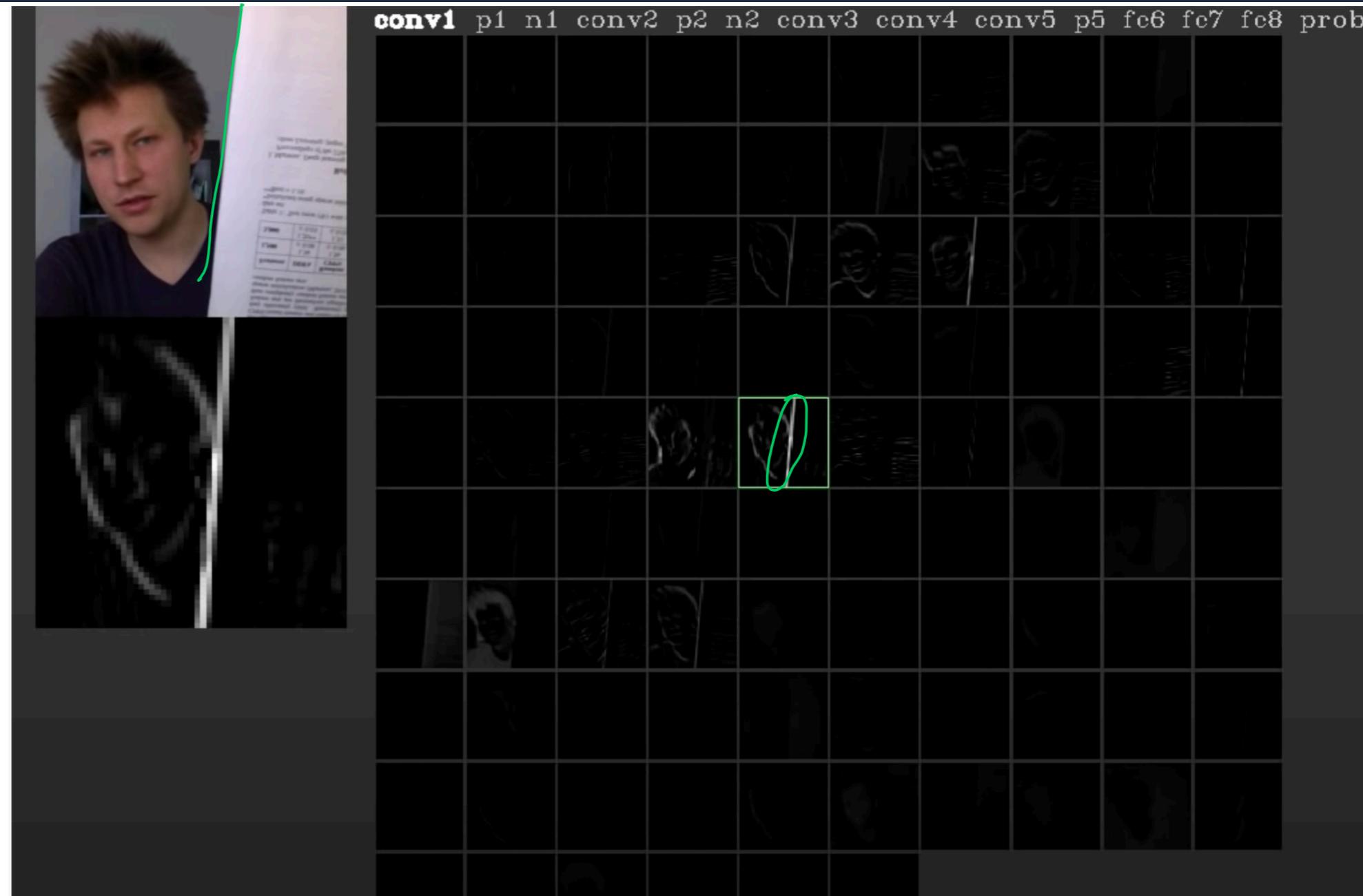
© A. Krizhevsky et al., 2012

# Visualisierung der Merkmale



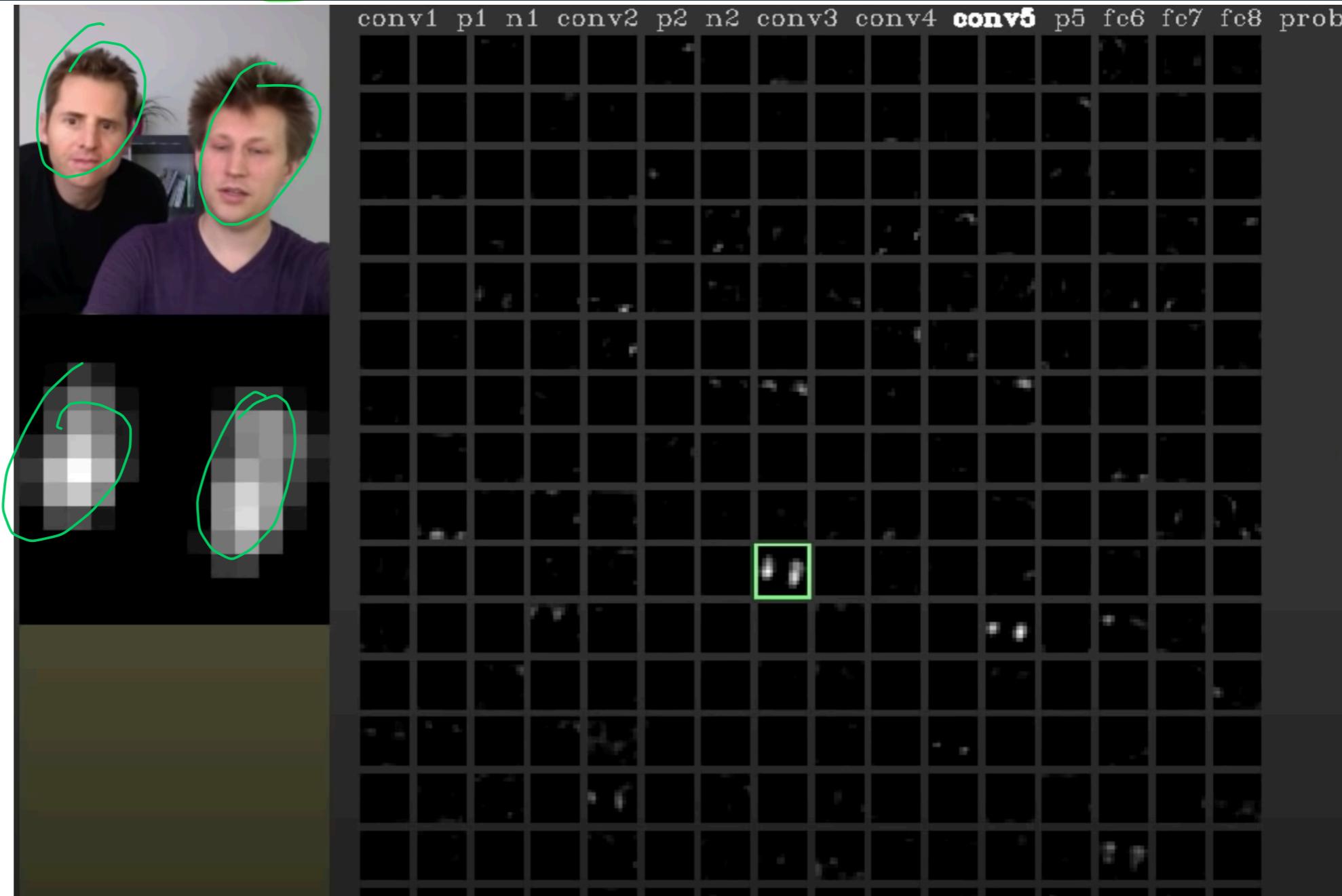
© Yosinski et al.

# Visualisierung der Merkmale



© Yosinski et al.

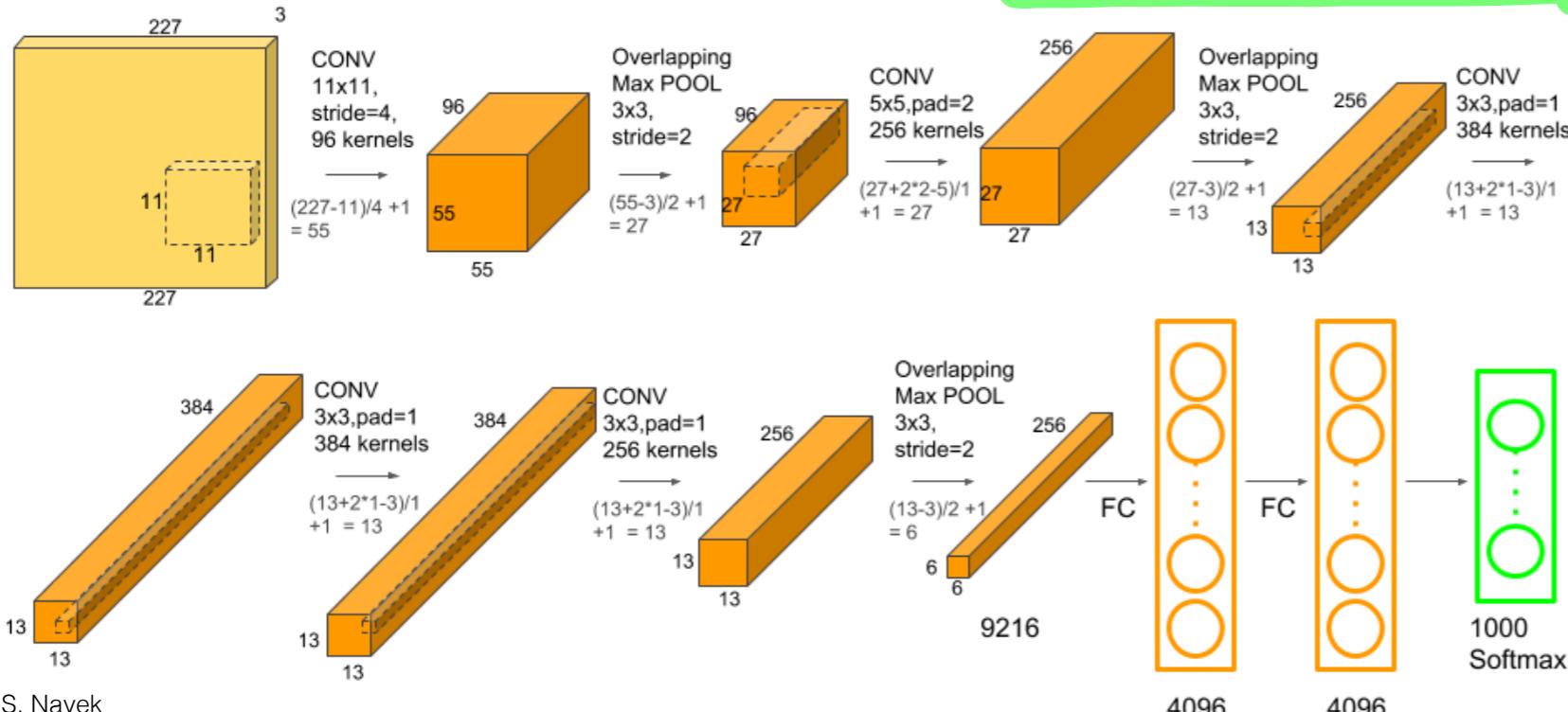
# Visualisierung der Merkmale



↳ Interpretation im Nachhinein experimentell

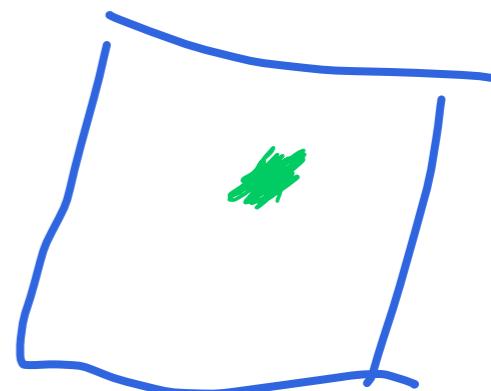
© Yosinski et al.

# Visualisierung von relevanten Bildregionen



„wenn ich das nicht hätte, würde ich dann zu keiner anderen Entscheidung kommen.“

„finden von maximal aktivierenden Bildausschnitten bzgl. einer Mehlaltskarte“



© S. Nayek

→ Maximal aktivierende Bildausschnitte zu korrekter Mehlaltskarte finden

(a) „Mehlaltskarte wählen“

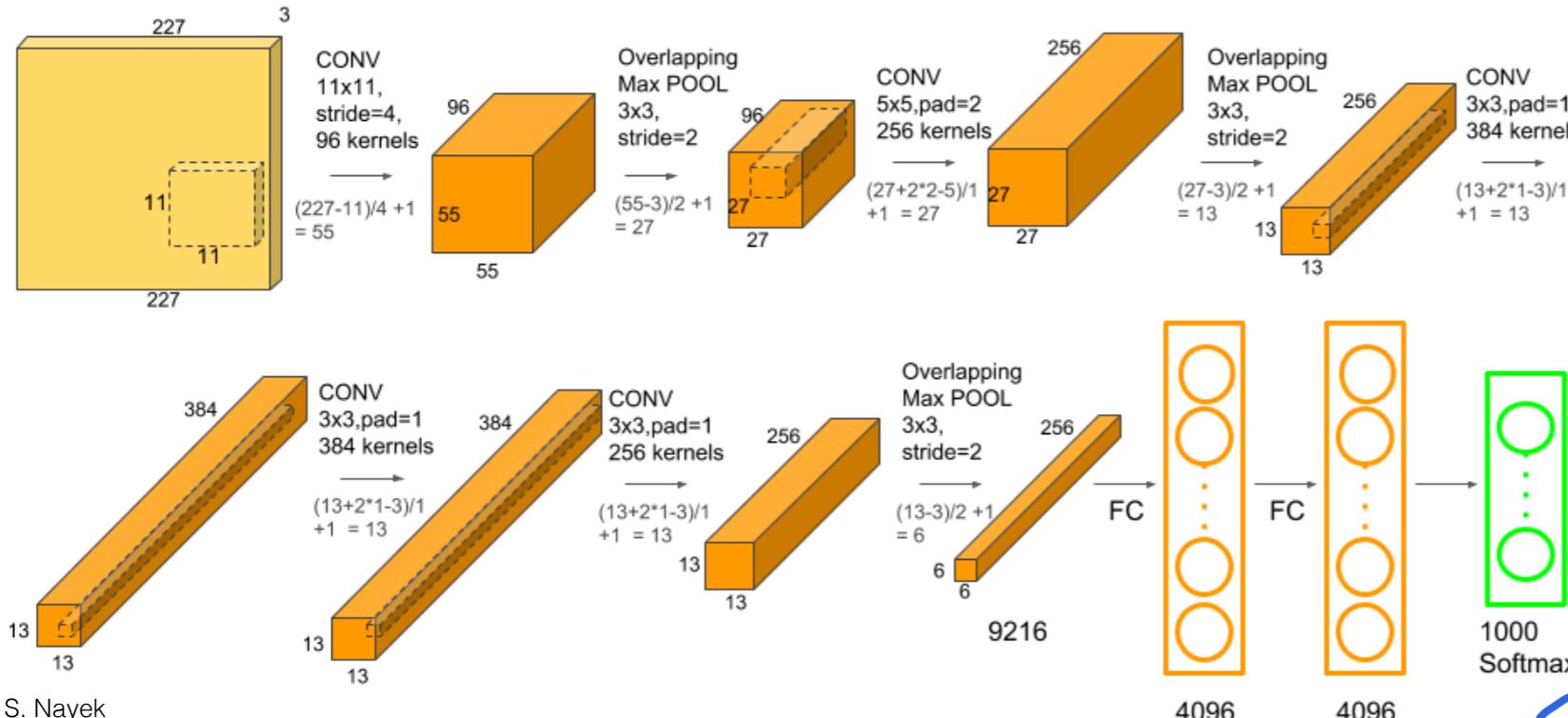
(b) Testbilder feed-forward verarbeiten

(c) Bilder und Position der maximalen Aktivität finden

(d) rückwärts aus Kenntnis der Kernelgrößen Positionen im Bild markieren

© Yosinski et al.

# Visualisierung von relevanten Bildregionen



© S. Nayek

Top 70



© Springenberg et al., 2015

Lineares Modell  $\rightarrow$  leicht erklären

$$y = g(x) = \phi_0 + \sum_i \phi_i x_i$$

Weight      Input

$$1,05 \cdot x_1 + 1,5 \cdot x_2 \approx 120.000$$

Ann:  $x_1, x_2$  stochastisch unabhängig

$\Rightarrow$  1,05 Bedeutung von  $x_1$

$\Rightarrow$  1,5 Bedeutung von  $x_2$

Bestimmung:  $R_i(x) := \frac{dy}{dx_i}$

$\Rightarrow$  Partielle Ableitung von  $y$  nach  $x_i$ :

# Saliency Maps

$$R_i^c(x) = \left| \frac{\partial S_c(x)}{\partial x_i} \right|$$

$S_c(x)$ : Zieloutput des NN  
für Klasse  $c$   
 $x_i$ : Pixel

=> Welche Pixel müssen verändert werden, um den Score einer Klasse maximal zu beeinflussen



© Simonyan et al., 2014

# Saliency Maps



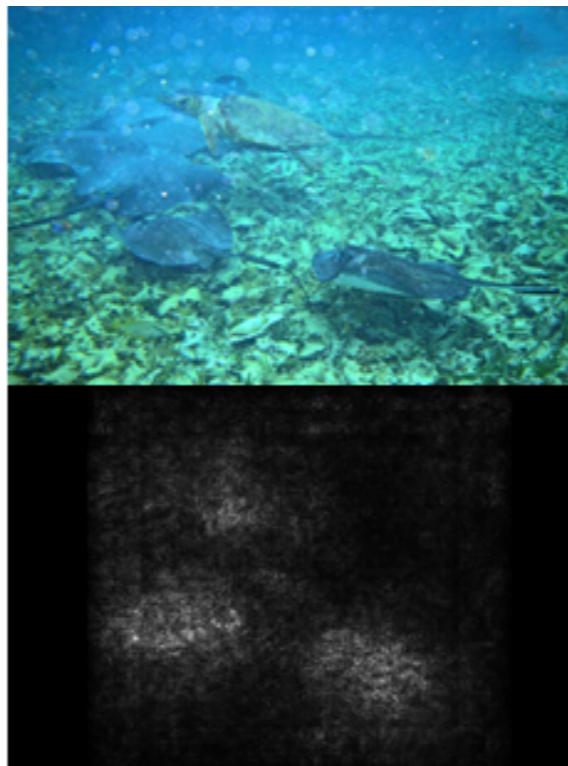
Dad:  
Trommel  
Statt  
in totat  
waschmaschine



Copyright 2005 Ventures, Inc.

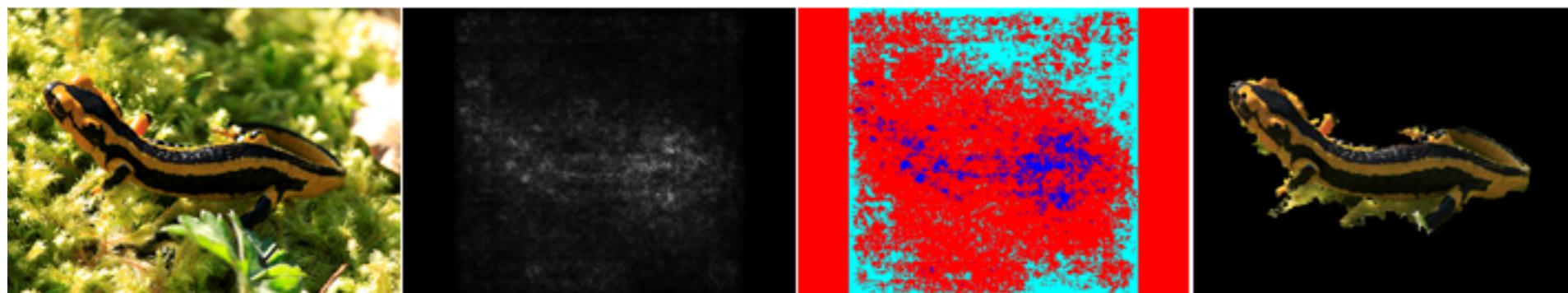
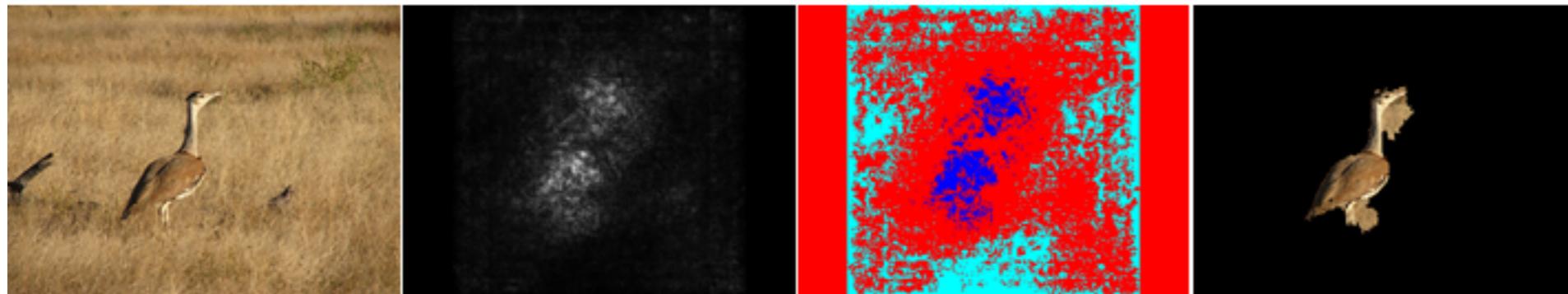
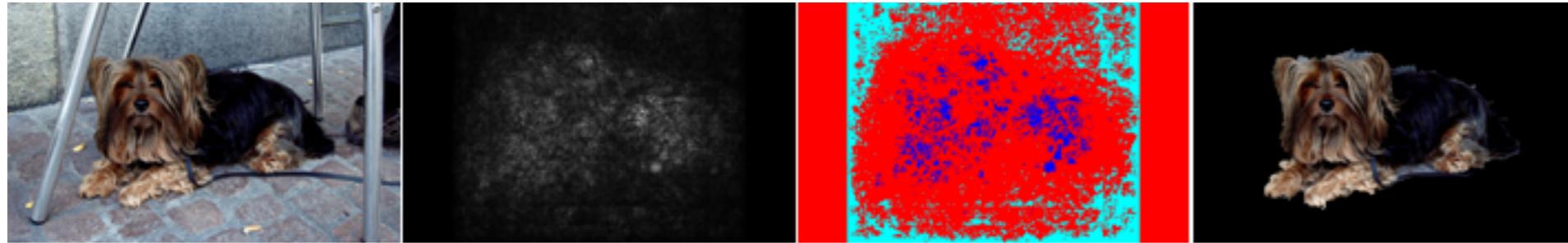


2006-1-8



© Simonyan et al., 2014

# Saliency Maps



Original

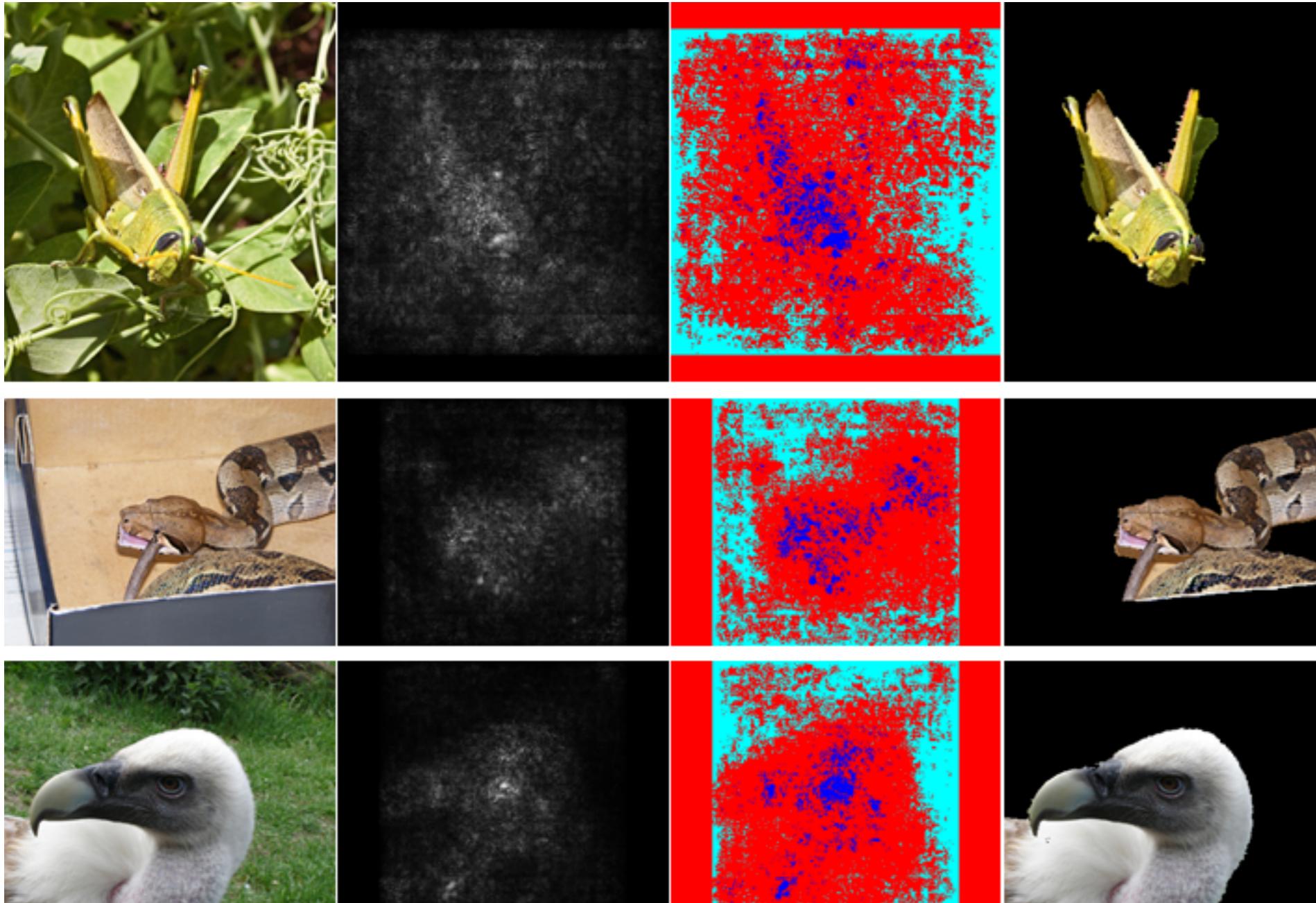
Saliency Map

Thresholding:  
→ blau! Objekt  
→ Cyan: Hintergrund  
→ rot! Hinterhöhlen

Segmentierung  
Graph Cut - Algorithmus

© Simonyan et al., 2014

# Saliency Maps



© Simonyan et al., 2014

$$\phi_1 = 7,05$$

$$\phi_L = 7,5$$

$$x_1 = 700.000$$

$$x_2 = 10.000$$

weicht von der Werte ist fast nur von  $x_1$  bestimmt <sup>weight</sup>

$$\phi_1 \cdot \tilde{x}_1 = 105.000$$

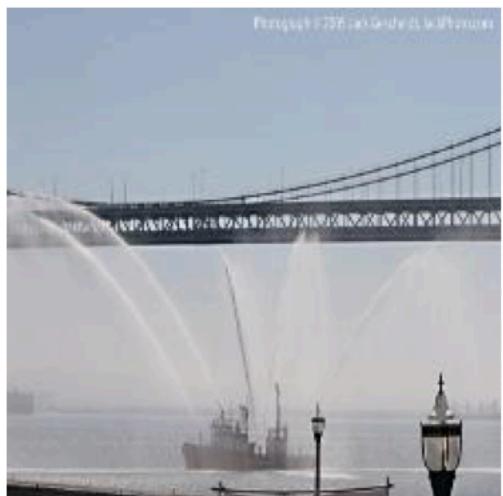
$$\phi_L \cdot \tilde{x}_2 = 15.000$$

$$R_i(x) = \tilde{x}_i \cdot \frac{\partial y_c}{\partial x_i}$$

# Gradient x Input

$$R_i^c(x) = \frac{\partial S_c(x)}{\partial x_i} \cdot x_i$$

*scheint nur  $\cdot x_i$ ; angehängt*



# Gradient x Input

$$R_i^c(\mathbf{x}) = \frac{\partial S_c(\mathbf{x})}{\partial x_i} \cdot x_i$$



hier auch welche

Lösung: Integrated Gradients

# Integrated Gradients

Nicht nur eine Intensität soll das Bild dominieren

Input wird mit  $\alpha$  verändert

$$R_i^c(x) = (x_i - x'_i) \cdot \int_{\alpha=0}^1 \frac{\partial S_c(x' + \alpha(x - x'))}{\partial x_i} d\alpha$$

$$x_i \cdot \int_{\alpha=0}^1 \frac{\partial S_c(\alpha \cdot x)}{\partial x_i} \cdot d\alpha$$

Baseline: Abweichen  $\alpha$  von Nullwerten

$\Rightarrow$  Wie stark erhöht sich der Bildwert, wenn die Präsenz eines Merkmals verstärkt wird

Baseline  $x'$



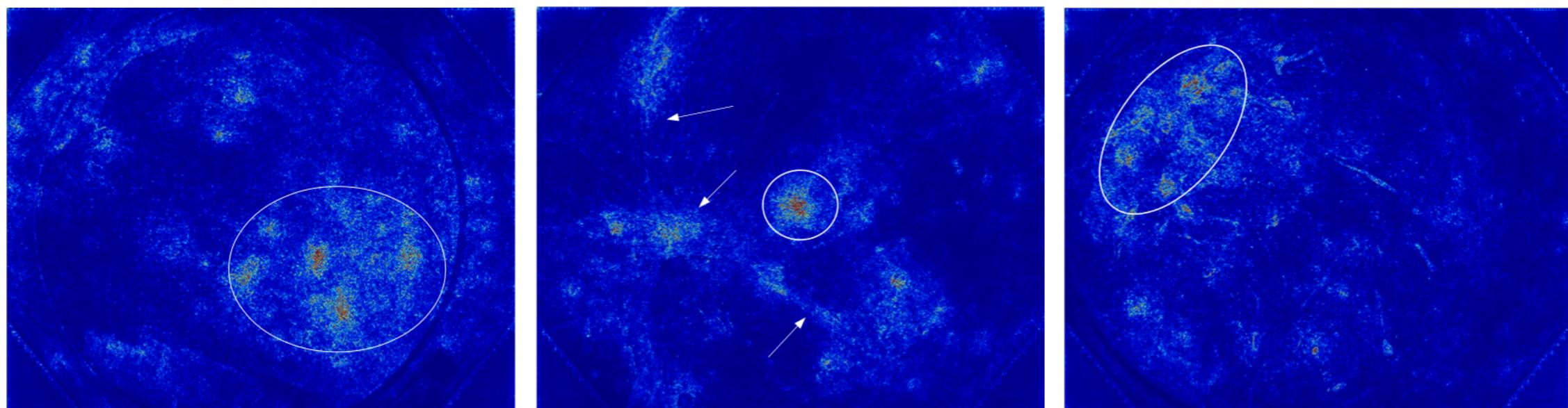
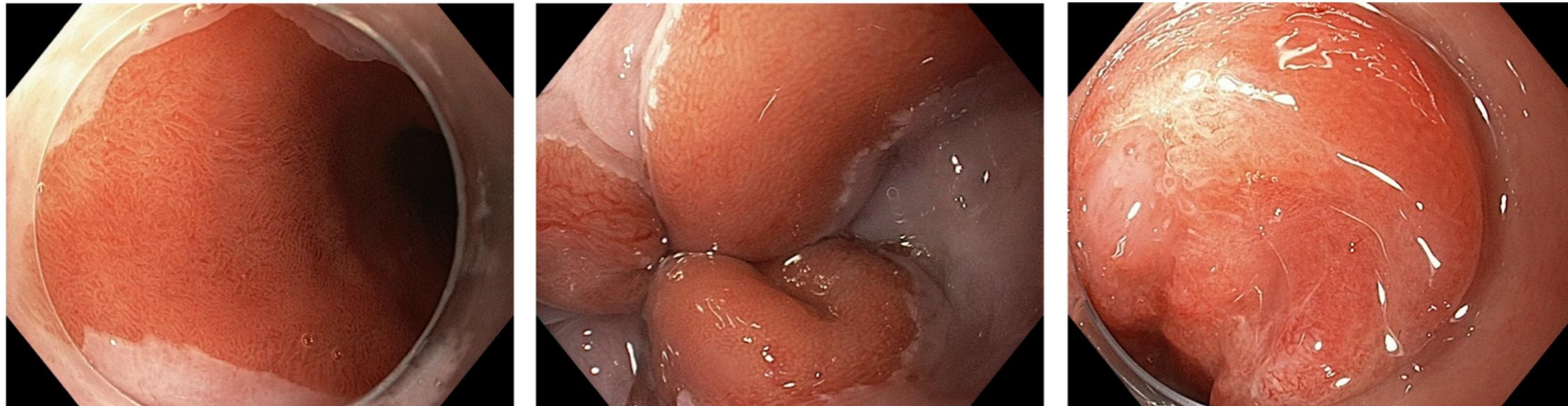
Originalbild  $x$



$\alpha$

© S. Strasser, 2021

# Integrated Gradients



Barrett

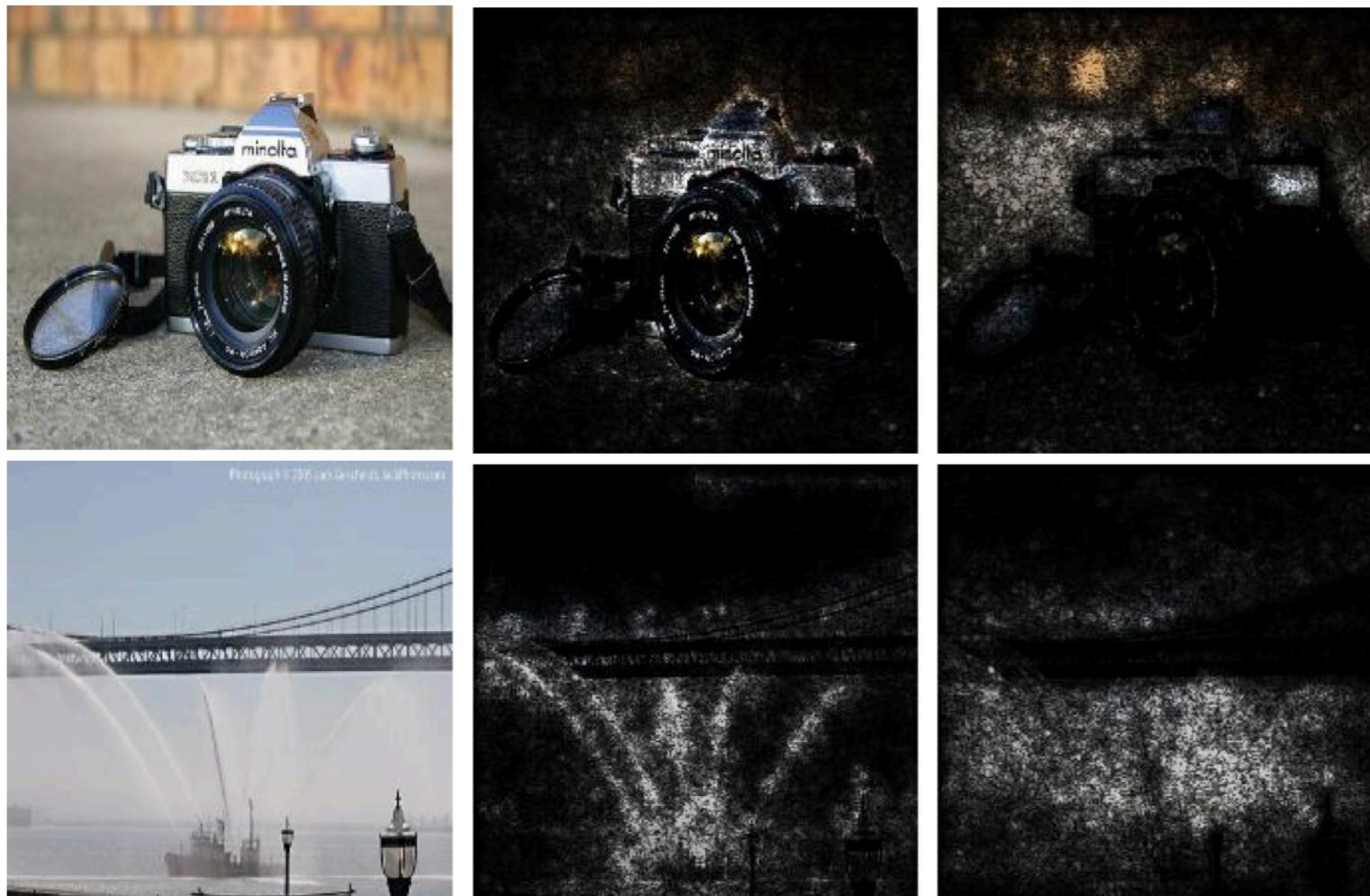
Dyspl

Cancer

Eher nicht so häufig

# Integrated Gradients

besser als vorher

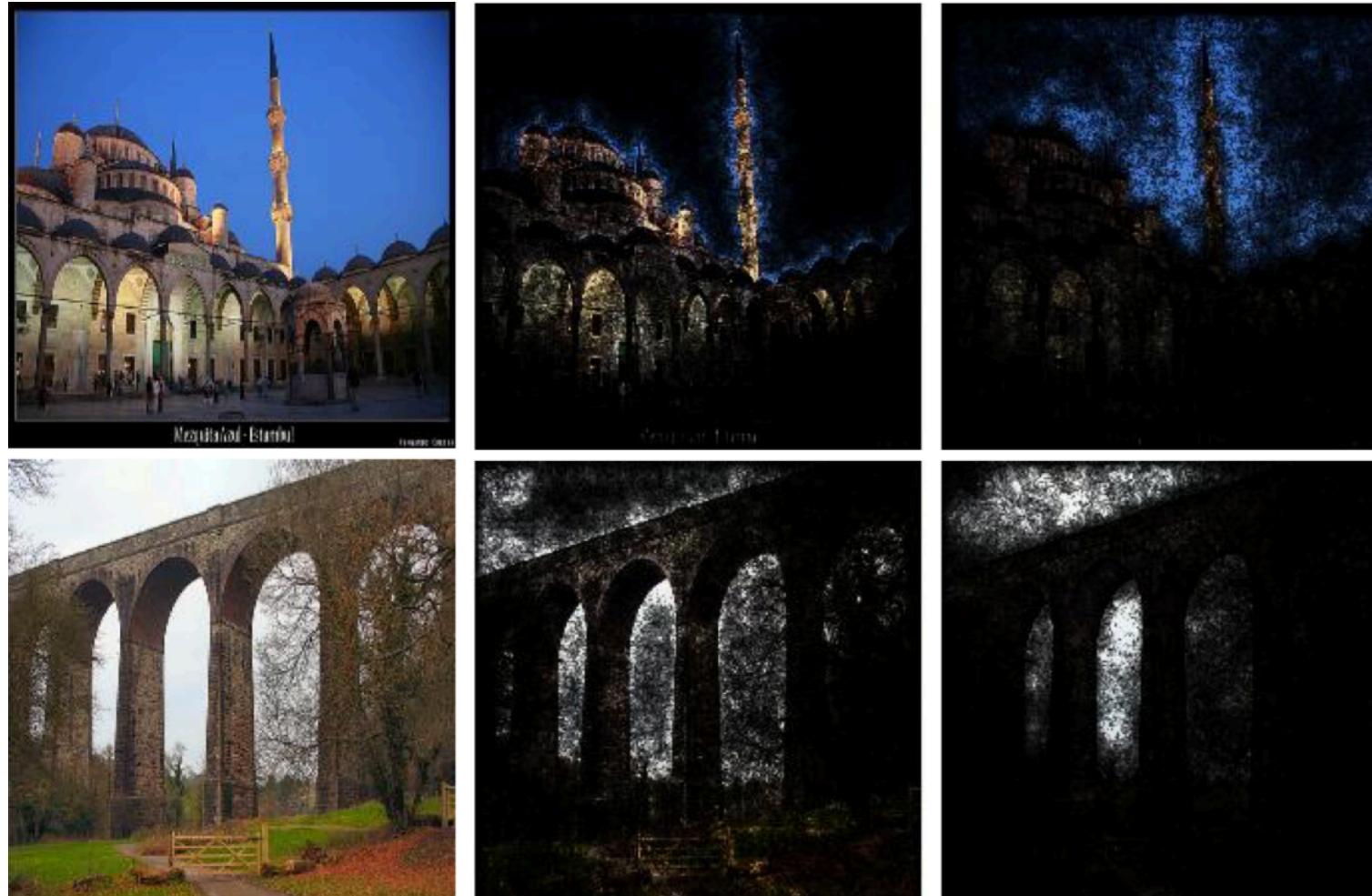


Integrated  
Gradients

Gradient  
+  
Input

© M. Sundararajan et al., 2017

# Integrated Gradients



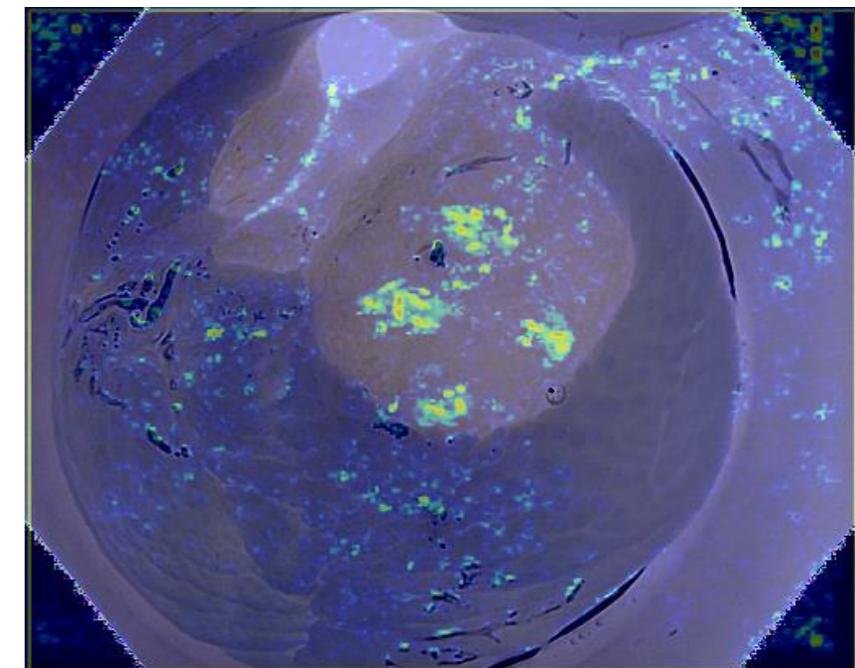
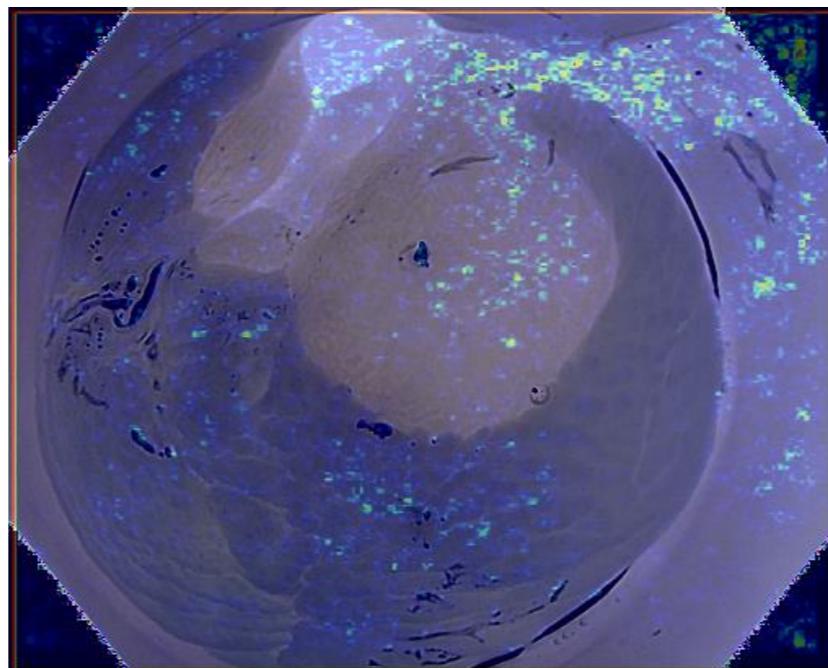
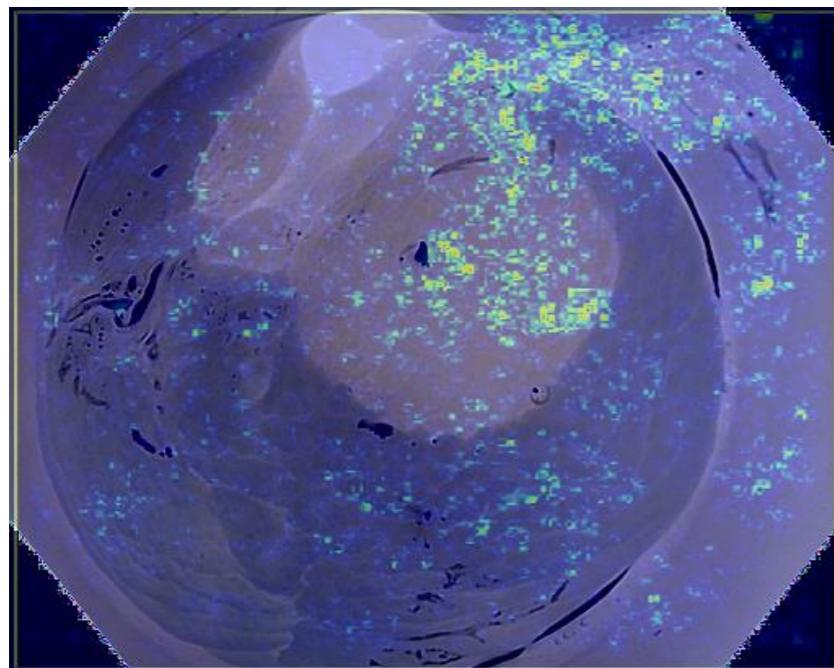
© M. Sundararajan et al., 2017

# Integrated Gradients



© M. Sundararajan et al., 2017

# Überblick





# Explainable AI

1. Einführung
2. Modell-spezifische Methoden
- 3. Modell-unspezifische Methoden**

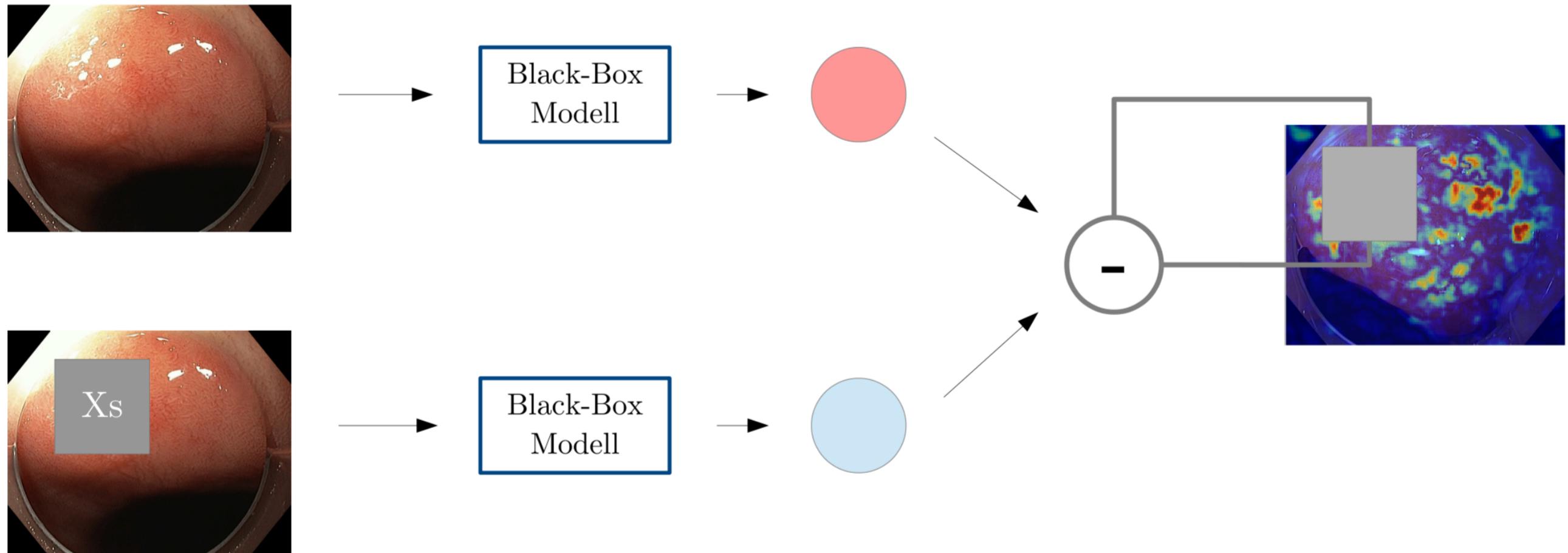
## Fragen & Antworten

- Abflussszenarien nicht vorhanden
- Aktivierungsfunktion: macht aus linearem etwas nichtlineares
- fragt keine Formeln ab, stattdessen stellt er Formeln vor und man soll die dann ableiten/beschreiben
- t-SNE OG (234) nicht erklären, wenn dann die angepasst mit  $\frac{1}{k}$

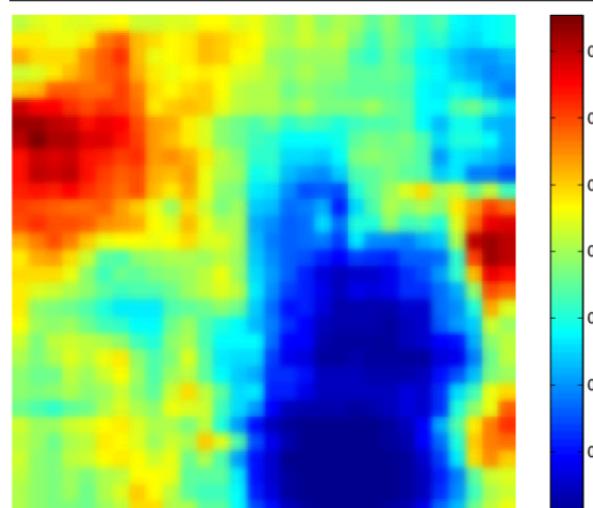
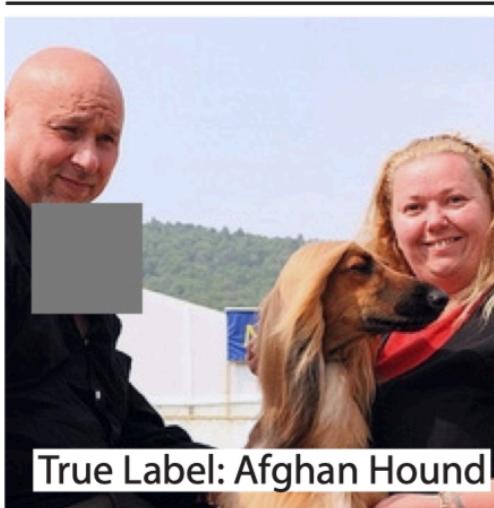
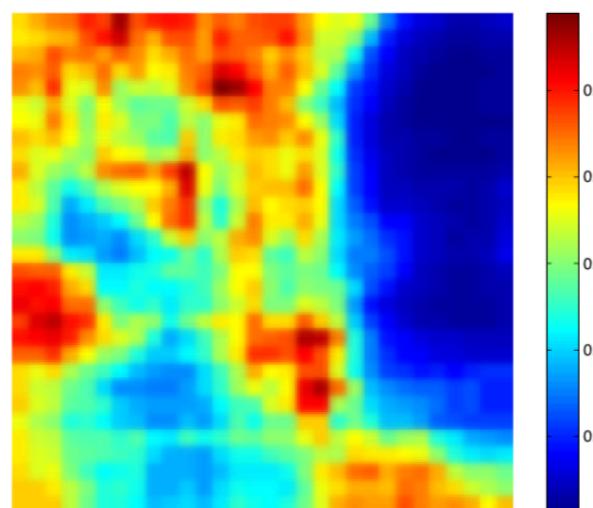
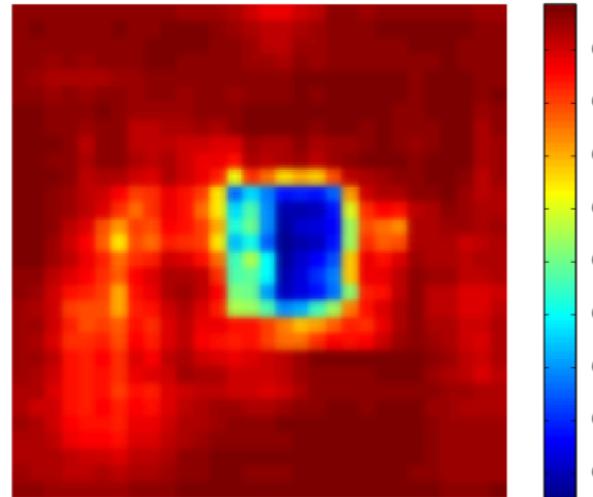
$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{l \neq i} \exp(-\|x_i - x_l\|^2 / 2\sigma_i^2)} \rightarrow p_{ij} = \frac{p_{ij}}{\sum_{l \neq i} p_{lj}} \text{ I approximation}$$

# Okklusion

$$R_i^c(\mathbf{x}) = S_c(\mathbf{x}) - S_c(\mathbf{x} \setminus \mathbf{x}_s)$$

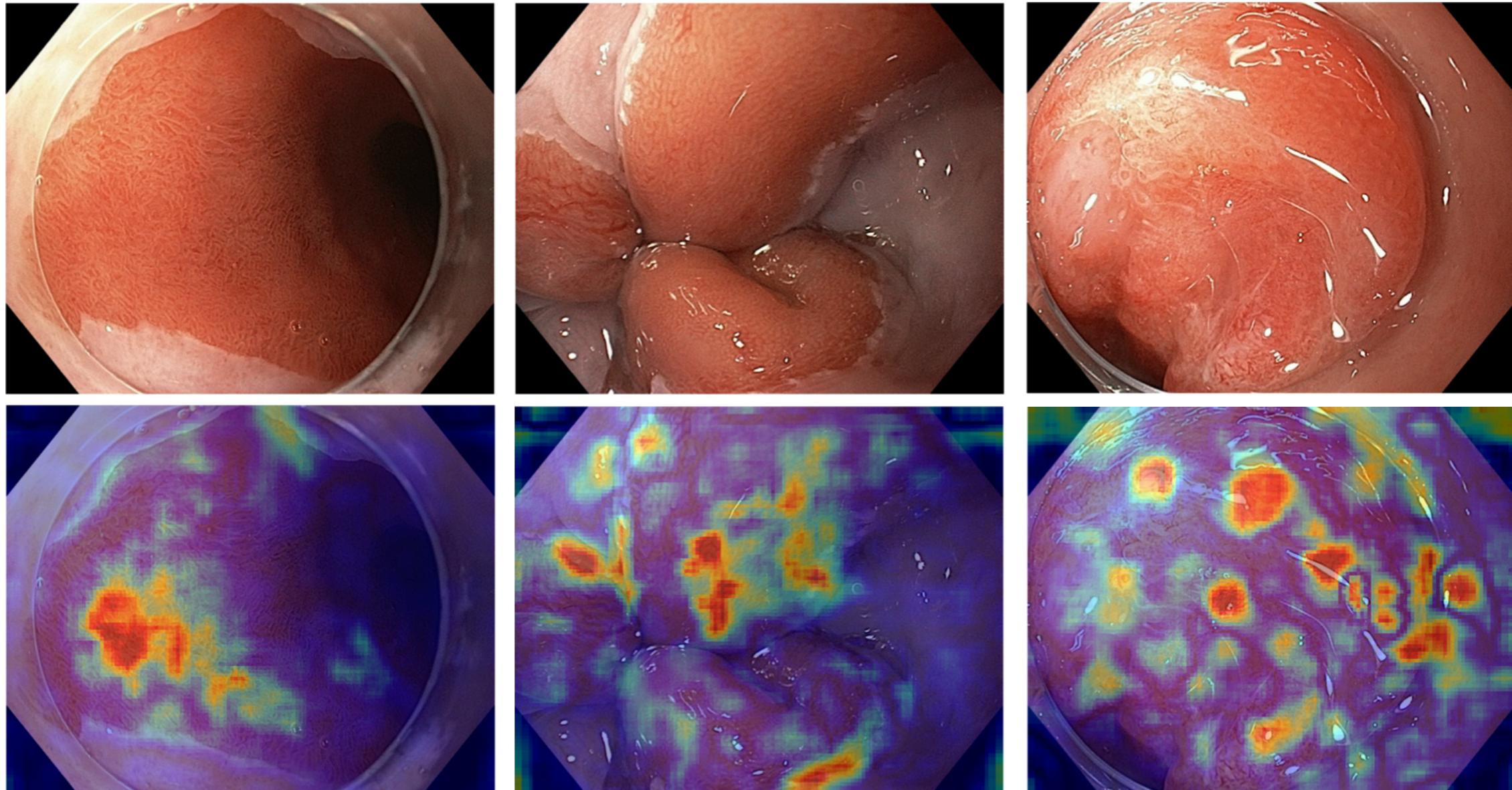


# Okklusion



© Zeiler, Fergus, 2014

# Integrated Gradients



local interpretable model-agnostic explanation

## Lokal

Interpretiert das Black-Box Modell in lokaler Umgebung eines konkreten Beispiels

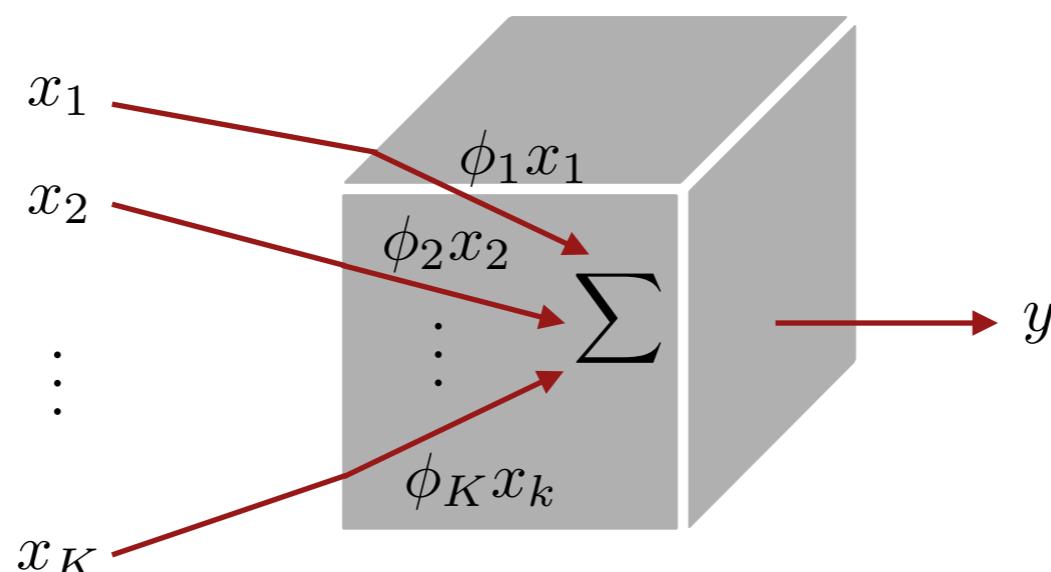
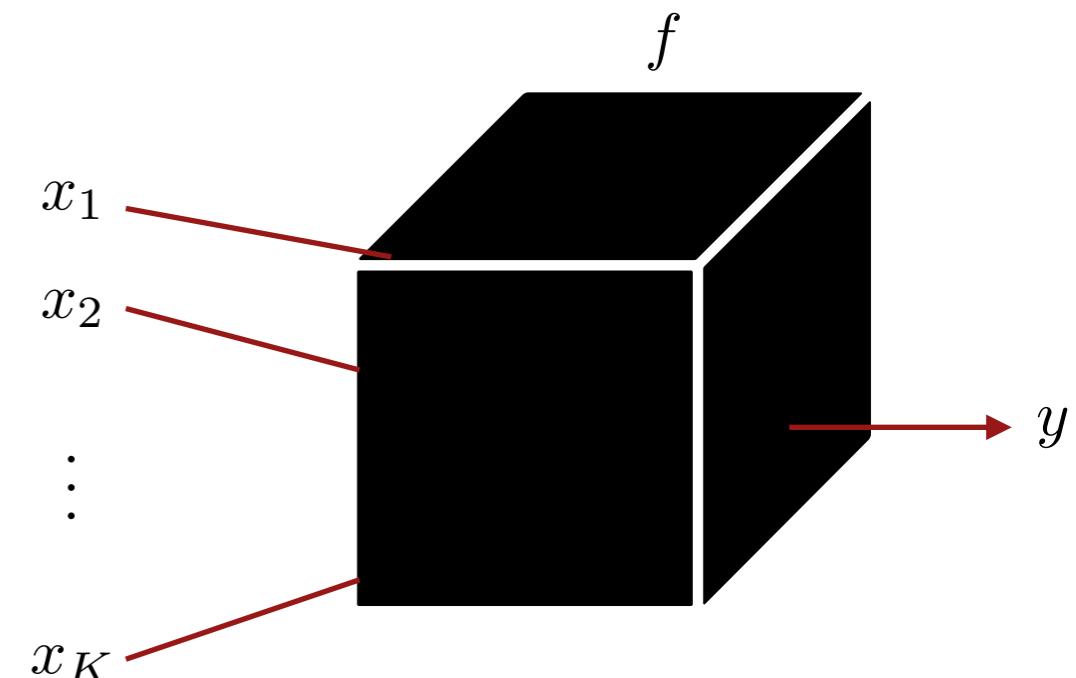
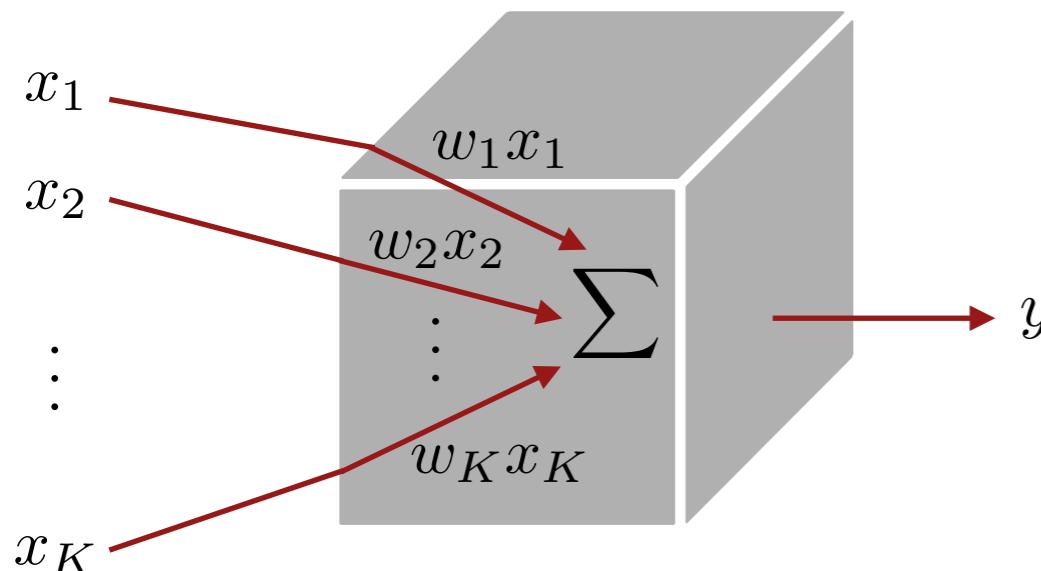
## Interpretierbar

Approximiert das komplexe Black-Box Modell z.B. durch lineare Regression bzw. Logistische Regression

## Model-Agnostic

Das Black-Box Modell bleibt unberührt, keine Details müssen bekannt sein

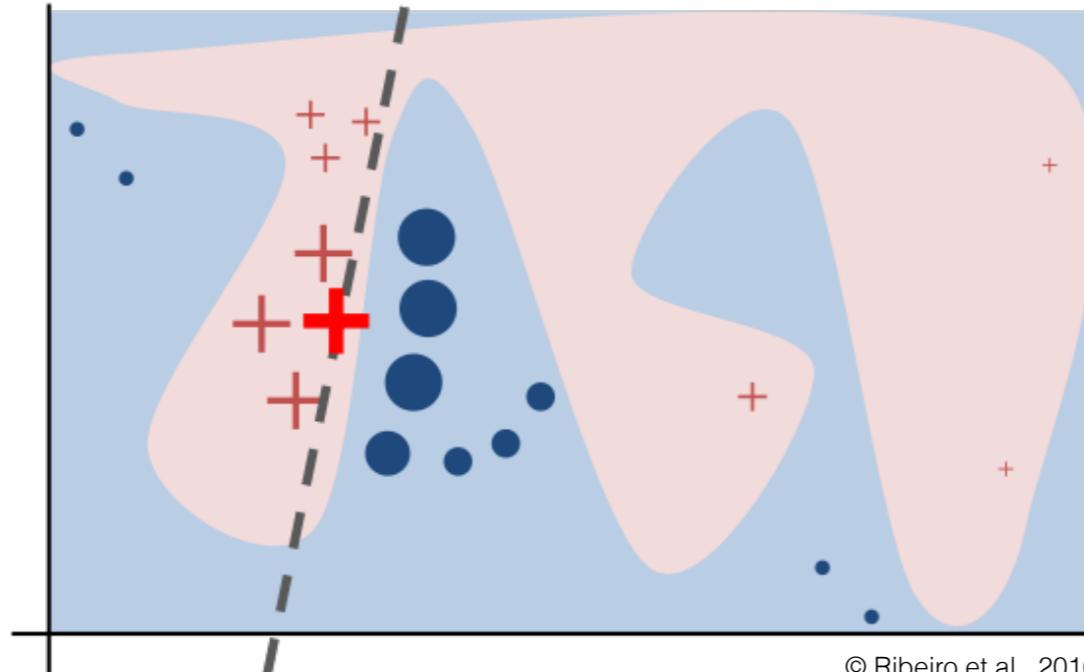
## Idee



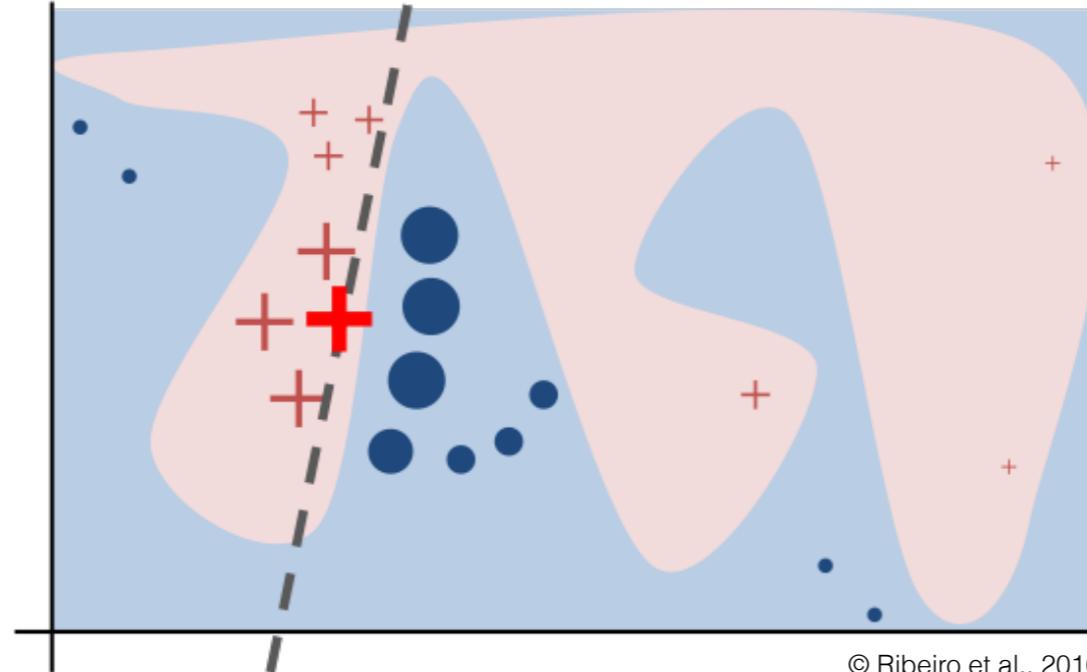
$$\phi_i = \psi_i(f, \mathbf{x}^{(n)})$$

local interpretable model-agnostic explanation

Idee



### Idee



© Ribeiro et al., 2016

Klassifikator  $f$

interpretierbares Modell  $g$

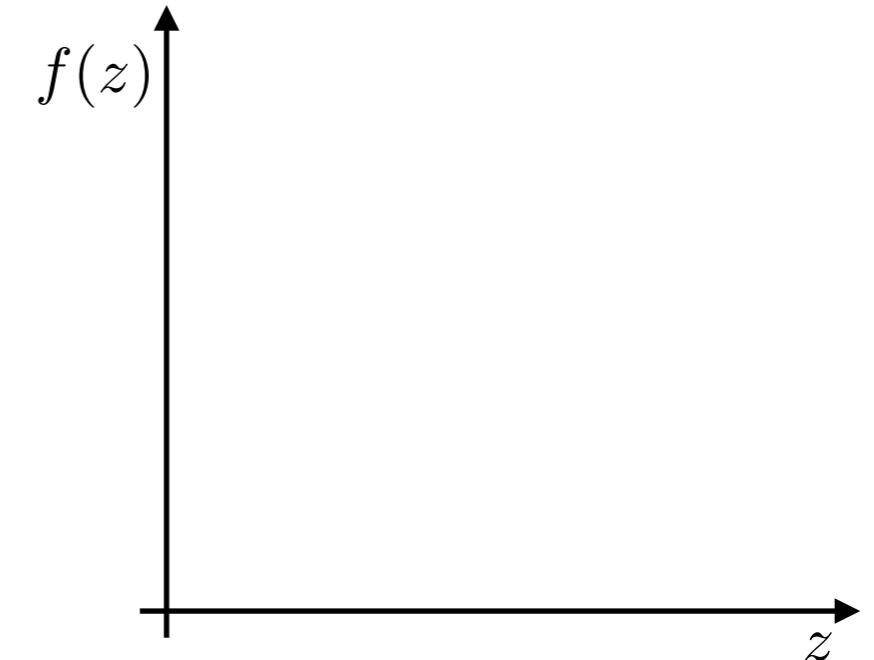
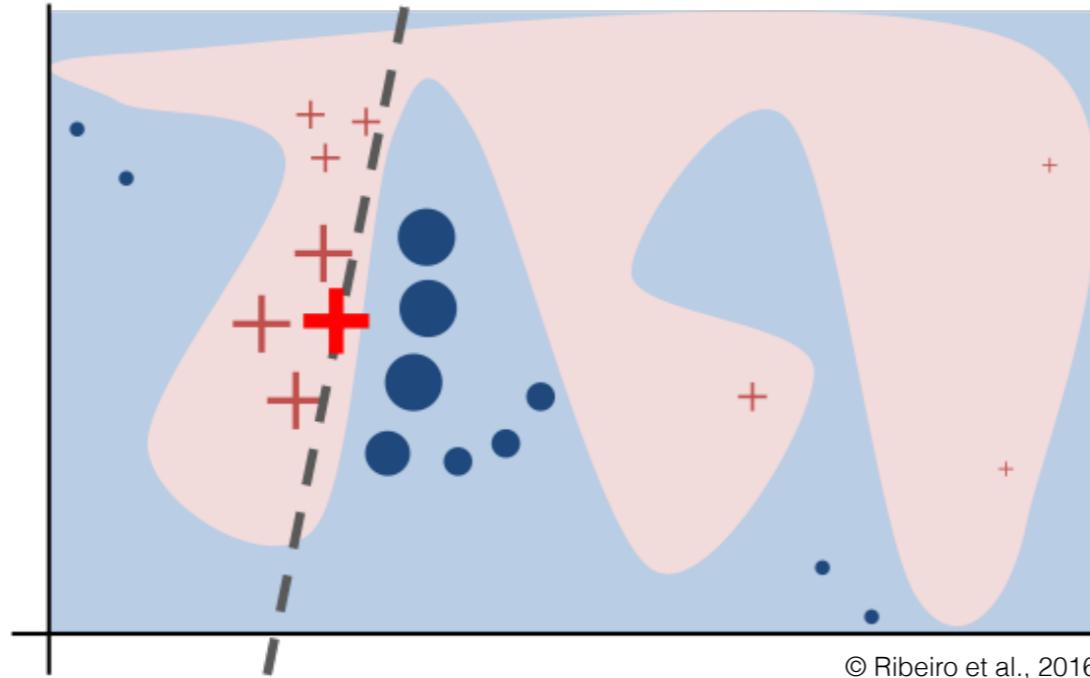
Originalbild  $x$

interpretierbare Repräsentation  $x'$

generierter Datenpunkt  $z$

interpretierbarer generierter Datenpunkt  $z'$

Idee



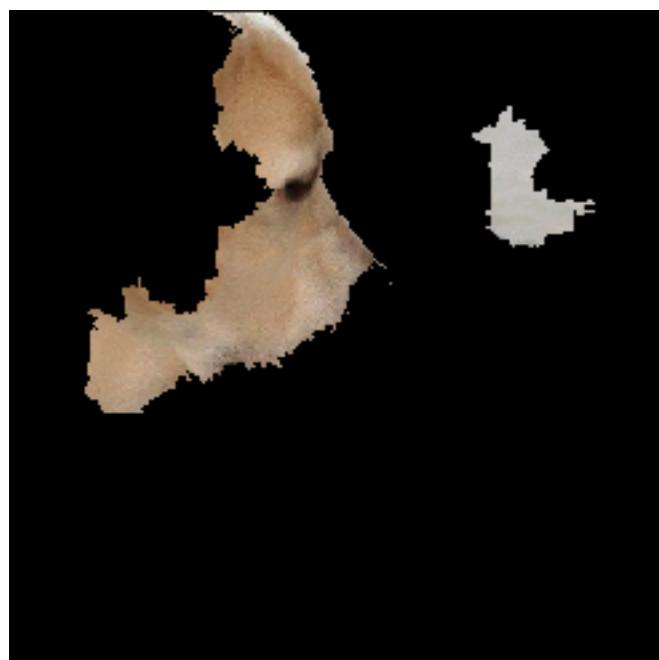
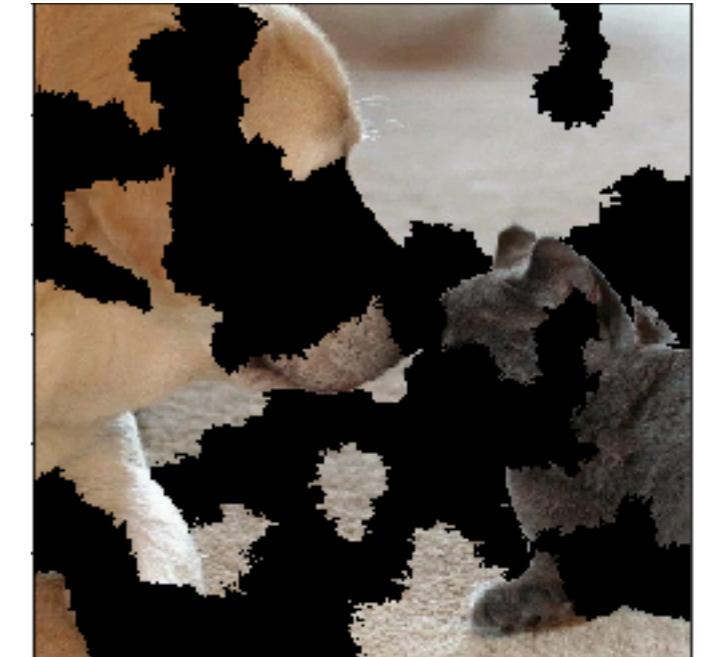
Maß für die Nähe von  $z$  zu  $x$ :  $\pi_x$

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z'} \pi_x(z) (f(z) - g(z'))^2$$

$$g(z') = \phi_0 + \sum_i \phi_i z'_i$$

# LIME

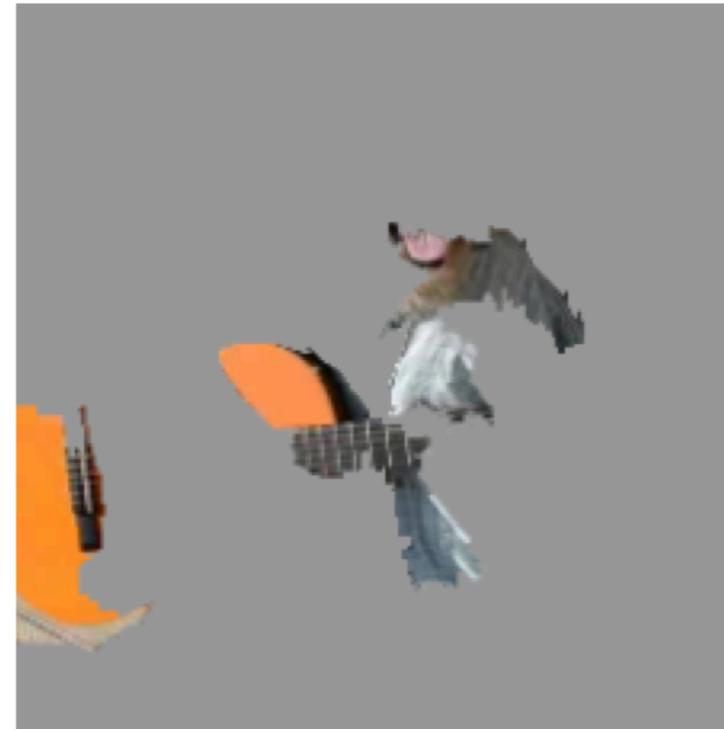
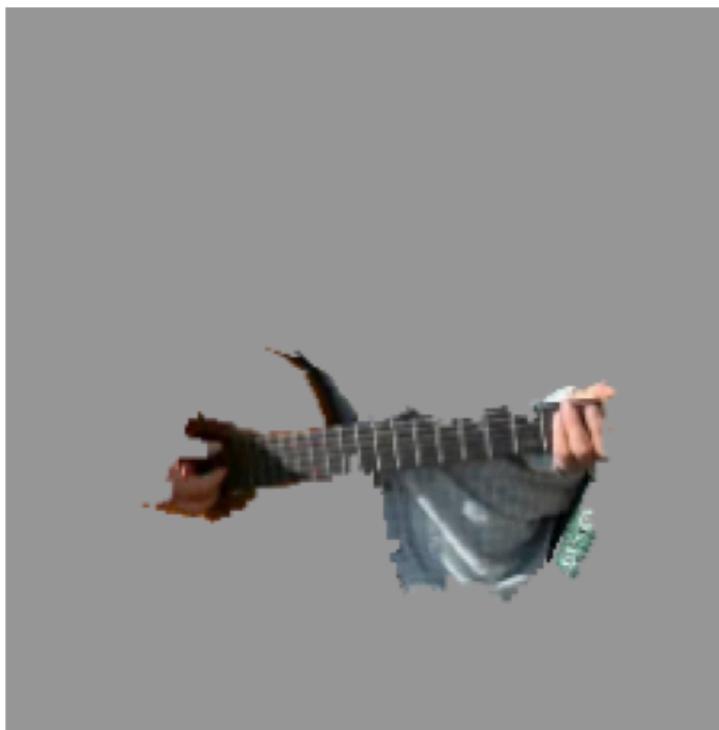
local interpretable model-agnostic explanation



© C. Arteaga

# LIME

local interpretable model-agnostic explanation



© Ribeiro et al., 2016

local interpretable model-agnostic explanation

Offen

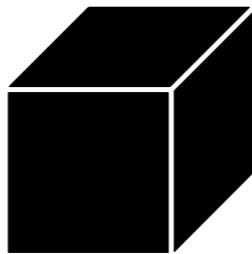
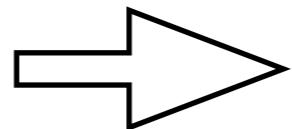
Wie kommen wir zu Superpixeln?



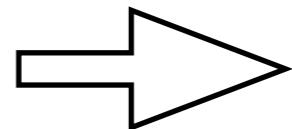
Kann man bestimmte Merkmale als insgesamt wichtiger als andere markieren?



Projektgabe



Projektbewertung



Note  
2,3



## SHapley Additive exPlanations

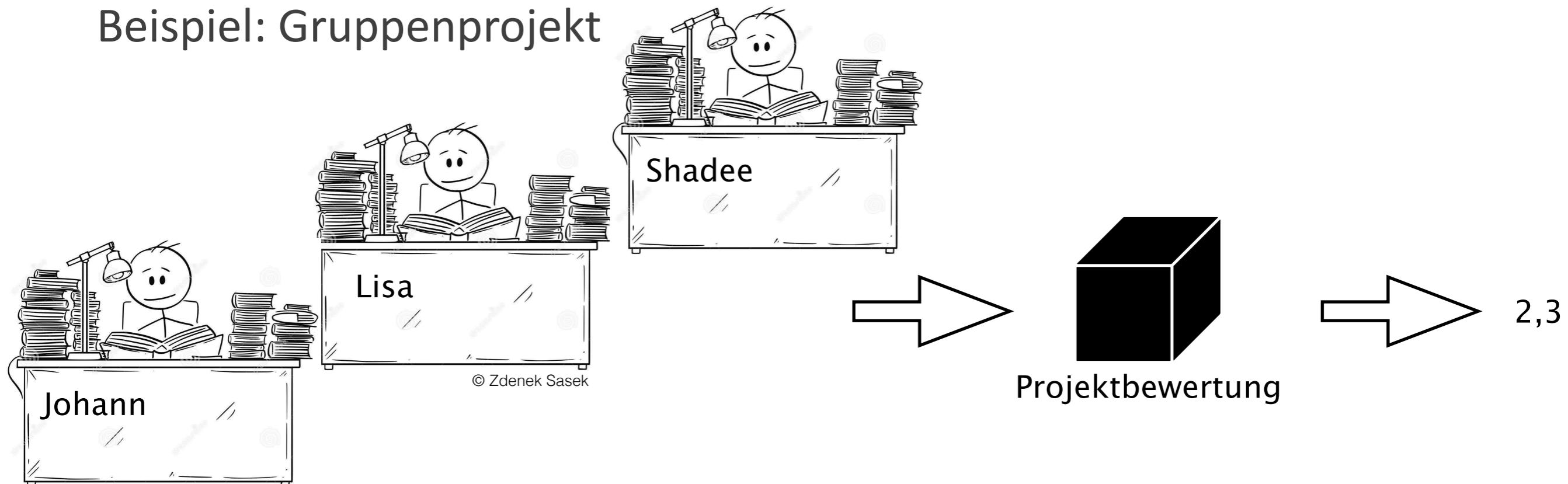
### Shapley Werte

aus der Spieltheorie

liefern die Beiträge einzelner Aspekte zum Gesamtergebnis

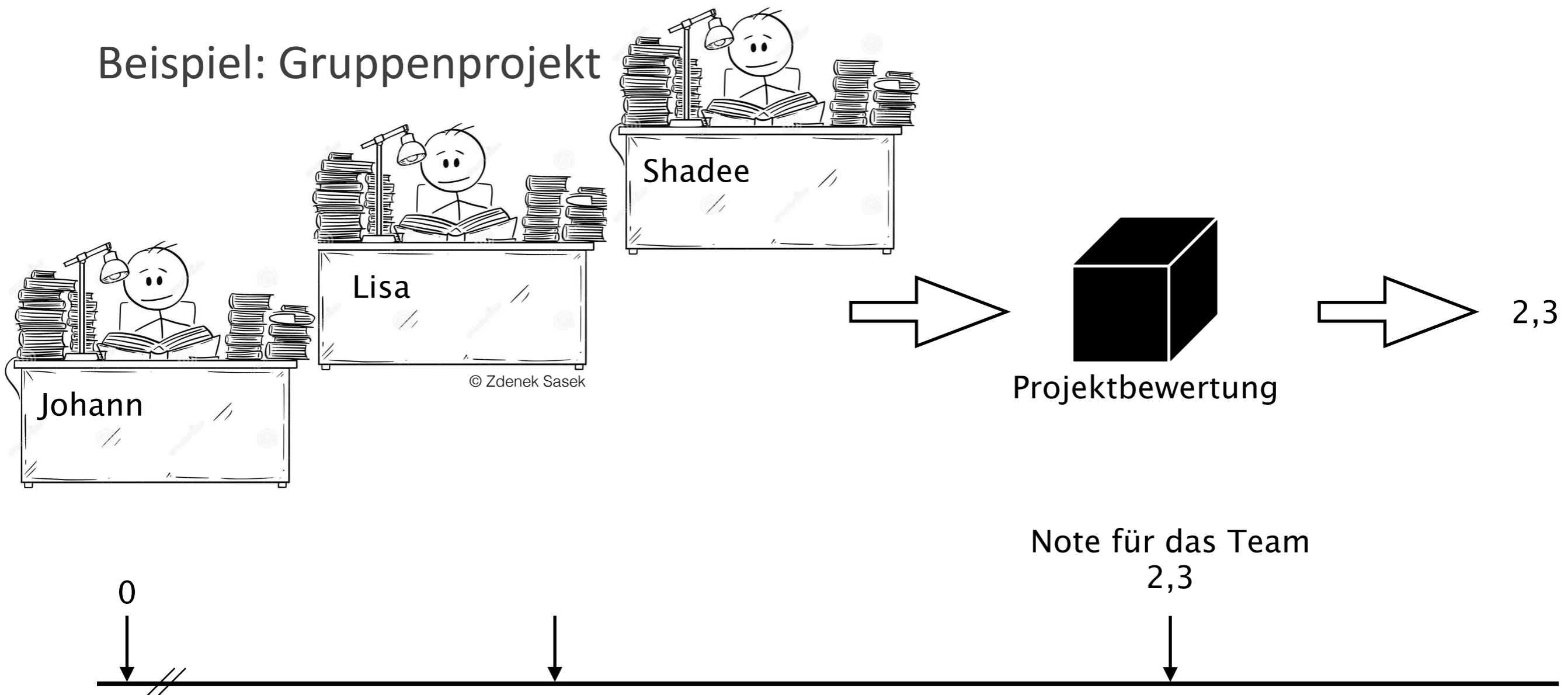
dazu: alle möglichen Zusammensetzungen betrachten

### Beispiel: Gruppenprojekt



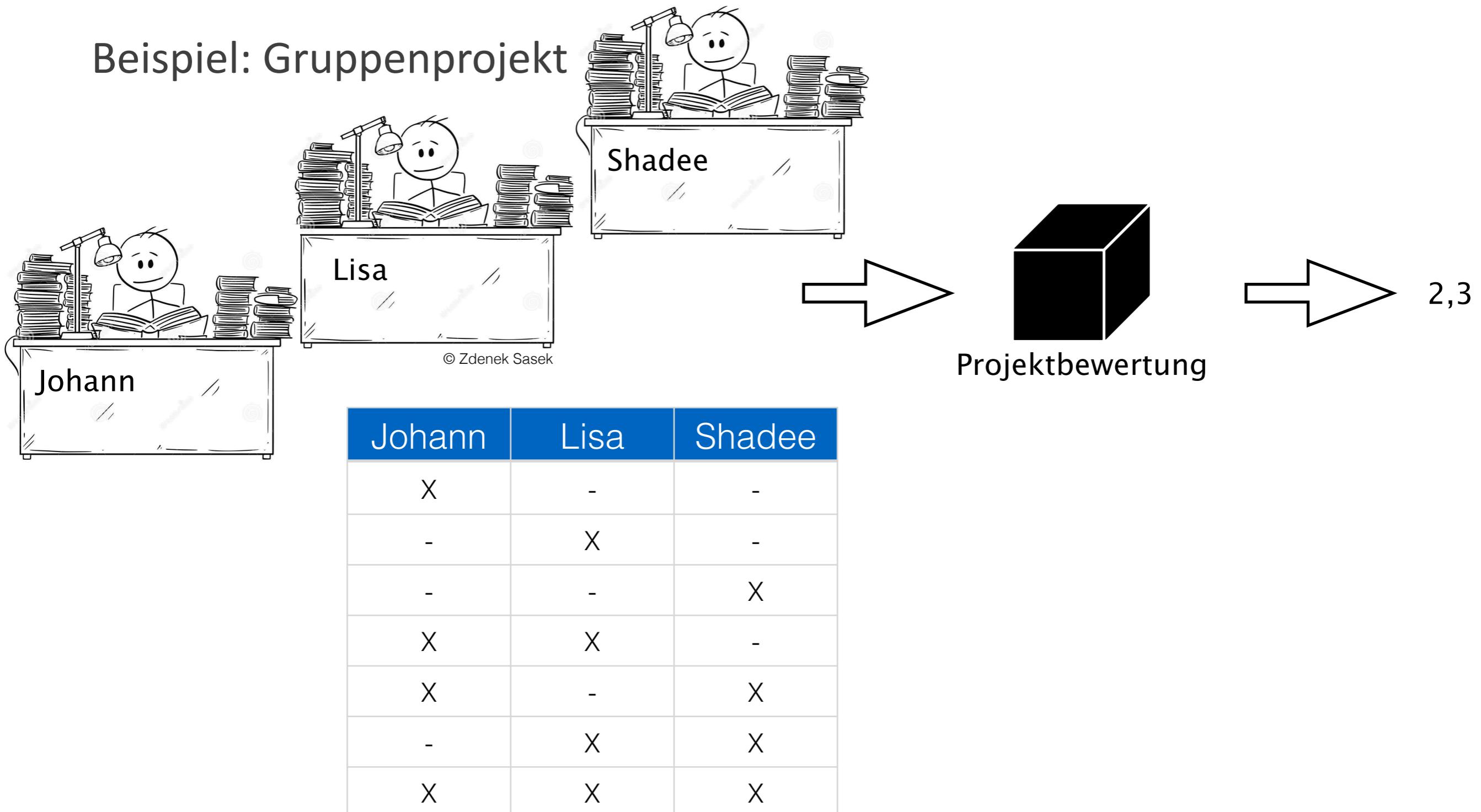
## SHapley Additive exPlanations

Beispiel: Gruppenprojekt



## SHapley Additive exPlanations

Beispiel: Gruppenprojekt



## SHapley Additive exPlanations

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} (f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S))$$

Johann	Lisa	Shadee
X	-	-
-	X	-
-	-	X
X	X	-
X	-	X
-	X	X
X	X	X

## SHapley Additive exPlanations

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} (f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S))$$

$$\pi_x(z') = \frac{M-1}{\binom{M}{|z'|} |z'| (M - |z'|)}$$

$$\mathcal{L}(f, g, \pi_x) = \sum_{z'} \pi_x(z') (f(z) - g(z'))^2$$

$$g(z') = \phi_0 + \sum_i \phi_i z'_i$$

# SHAP

## SHapley Additive exPlanations

