Duke

Great success !!!

sunset

Duke

sunset

Duke

post hoc analysis

sunset

concept based - human reason in concepts

Duke

# Post hoc analysis – concept based

- Single neuron (Zhou et al, 2014; 2018)



neuron 1

impure!

# Post hoc analysis – concept based

- Single neuron (Zhou et al, 2014; 2018)



neuron 1

distributed!

neuron 2

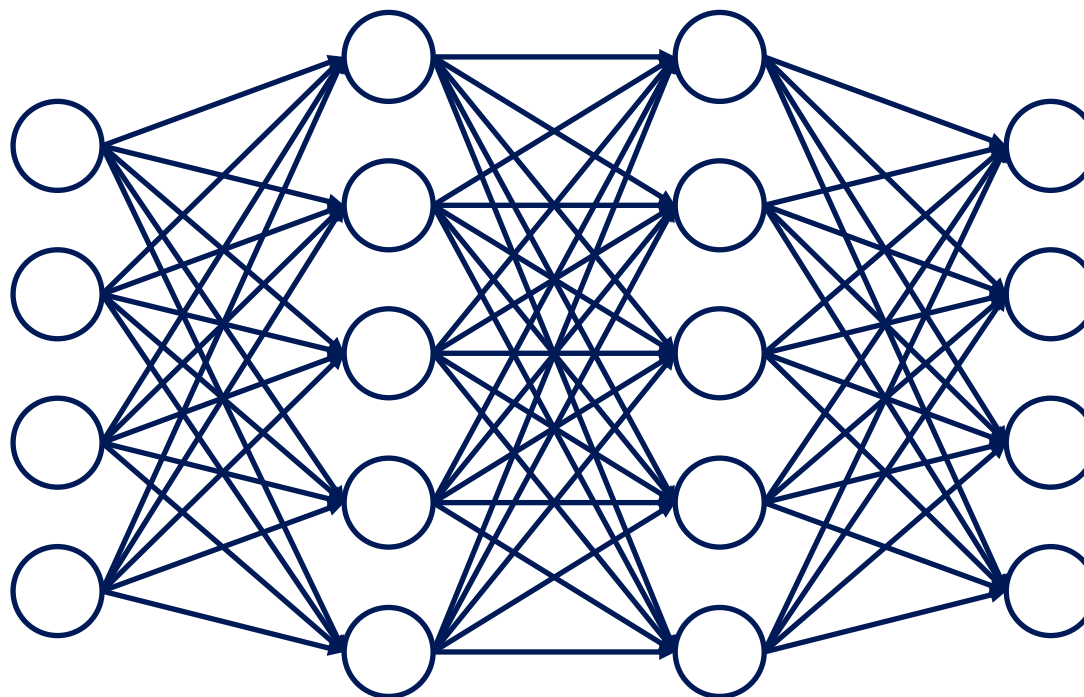single neuron of standard NN --✗ single concept

Duke

# Post hoc analysis – concept based

- Linear combination of neurons (Kim et al, 2017; Zhou et al, 2018)
  - better than single neuron



TCAV, Kim et al, 2017

# Post hoc analysis – concept based

- Linear combination of neurons (Kim et al, 2017; Zhou et al, 2018)
  - reality: concept vectors may point to the same direction

# The idea

- Why not do it by ourselves?
- Create a disentangled latent space that its axes represent known concepts

# Concept Whitening

- Step 1: Whitening transformation
  - Decorrelate the latent space
  - Separate the concepts
- Step 2: Rotation transformation
  - Align the concepts to corresponding axes
  - Maintain the decorrelation property

Duke

# Concept Whitening



concept 1
concept 2
concept 3
concept 4

neuron 2

neuron 1

Mean-centered latent features

$$\widetilde{\mathbf{Z}}_{d \times n} = \mathbf{Z}_{d \times n} - \boldsymbol{\mu} \, \mathbf{1}_{n \times 1}^{T}$$

Duke

# Concept Whitening



concept 1
concept 2
concept 3
concept 4

Whitening matrix $\boldsymbol{W_{d \times d}}$

neuron 2

neuron 1

neuron 2

neuron 1

Mean-centered latent features

$$\widetilde{\boldsymbol{Z}}_{d \times n} = \boldsymbol{Z}_{d \times n} - \boldsymbol{\mu} \, \boldsymbol{1}_{n \times 1}^{T}$$

Duke

After whitening transformation

$$\boldsymbol{W} \, \widetilde{\boldsymbol{z}}$$

$\boldsymbol{W}$ should obey $\boldsymbol{W}^{T}\boldsymbol{W} = \left( \dfrac{\widetilde{\boldsymbol{Z}} \, \widetilde{\boldsymbol{Z}}^{\mathrm{T}}}{n} \right)^{-1}$

# Concept Whitening



concept 1
concept 2
concept 3
concept 4

Whitening matrix $\boldsymbol{W}_{d \times d}$

Orthogonal matrix $\boldsymbol{Q}_{d \times d}$

Mean-centered latent features

$$\widetilde{\boldsymbol{Z}}_{d \times n} = \boldsymbol{Z}_{d \times n} - \boldsymbol{\mu} \, \mathbf{1}_{n \times 1}^{T}$$

After whitening transformation

$$\boldsymbol{W} \, \widetilde{\boldsymbol{Z}}$$

$\boldsymbol{W}$ should obey $\boldsymbol{W}^{T} \boldsymbol{W} = \left( \dfrac{\widetilde{\boldsymbol{Z}} \, \widetilde{\boldsymbol{Z}}^{\mathrm{T}}}{n} \right)^{-1}$

After rotation transformation

$$\boldsymbol{Q}^{T} \boldsymbol{W} \, \widetilde{\boldsymbol{Z}}$$

$\boldsymbol{W}' = \boldsymbol{Q}^{T} \boldsymbol{W}$ is also a whitening matrix

Duke

# Learning the parameters

- Sample mean $\boldsymbol{\mu}$ and whitening matrix $\boldsymbol{W}$
  - Training phase: compute on the fly, support back-propagation (Huang et al)
  - Testing phase: exponential moving average of mini-batches (Ioffe & Szegedy)

- Orthogonal matrix $\boldsymbol{Q}$
  - maximizing concept activation under orthogonality constraint

$$\max_{\boldsymbol{q_1},\boldsymbol{q_{2,\ldots,}}\boldsymbol{q_k}} \sum_{j=1}^{k} \frac{1}{n_j} \boldsymbol{q}_j^T \boldsymbol{W} \boldsymbol{Z}_{c_j} \mathbf{1}_{n_j \times 1}$$
$$s.t.\, \boldsymbol{Q}^T \boldsymbol{Q} = \boldsymbol{I}_d$$

$\boldsymbol{Z}_{c_j}$: samples of concept j

$\boldsymbol{W}\,\boldsymbol{Z}_{c_j}$: after whitening

$\boldsymbol{q}_j^T\,\boldsymbol{W}\,\boldsymbol{Z}_{c_j}$: projection on axis j

$\frac{1}{n_j}\boldsymbol{q}_j^T\,\boldsymbol{W}\,\boldsymbol{Z}_{c_j}\mathbf{1}_{n_j \times 1}$: average activation

$\boldsymbol{Q}$ can be trained by gradient descent on Stiefel manifold (Wen & Yin, 2013)

Duke

What's the cost of interpretability?
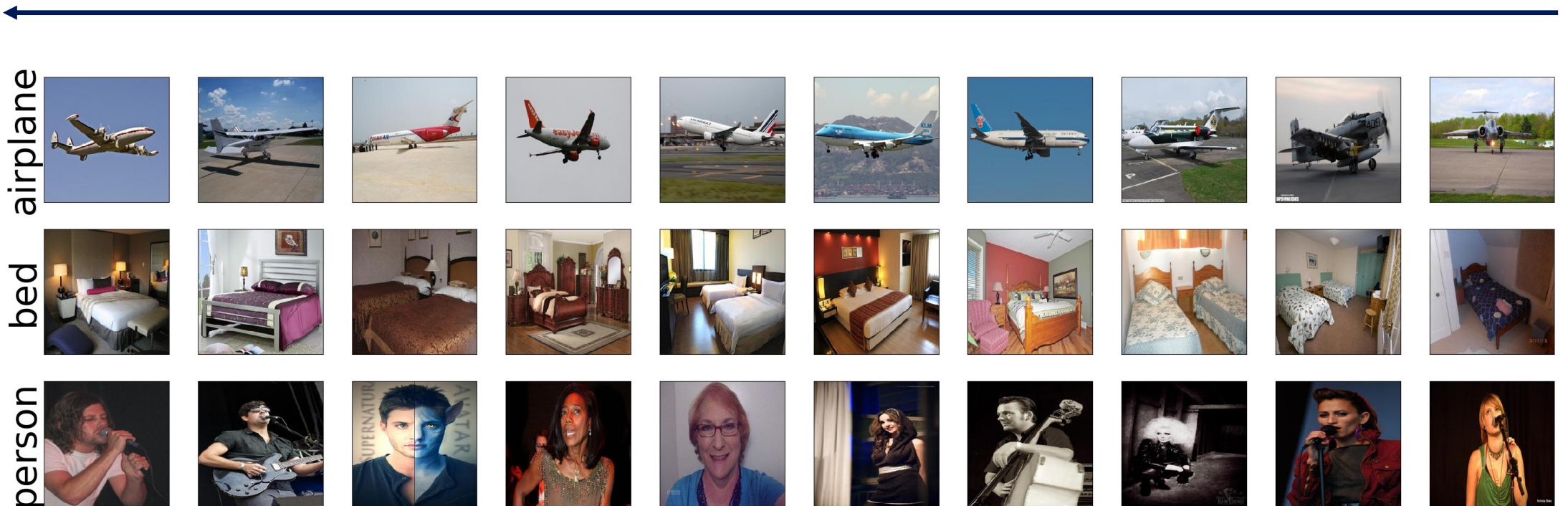
Duke

# Main task performance

- accuracy is on par with standard CNNs
  - different datasets, backbone architectures, layers, #concepts
- warm-start from pretrained model
  - replace BN with CW
  - *only one* additional epoch of further training

Duke

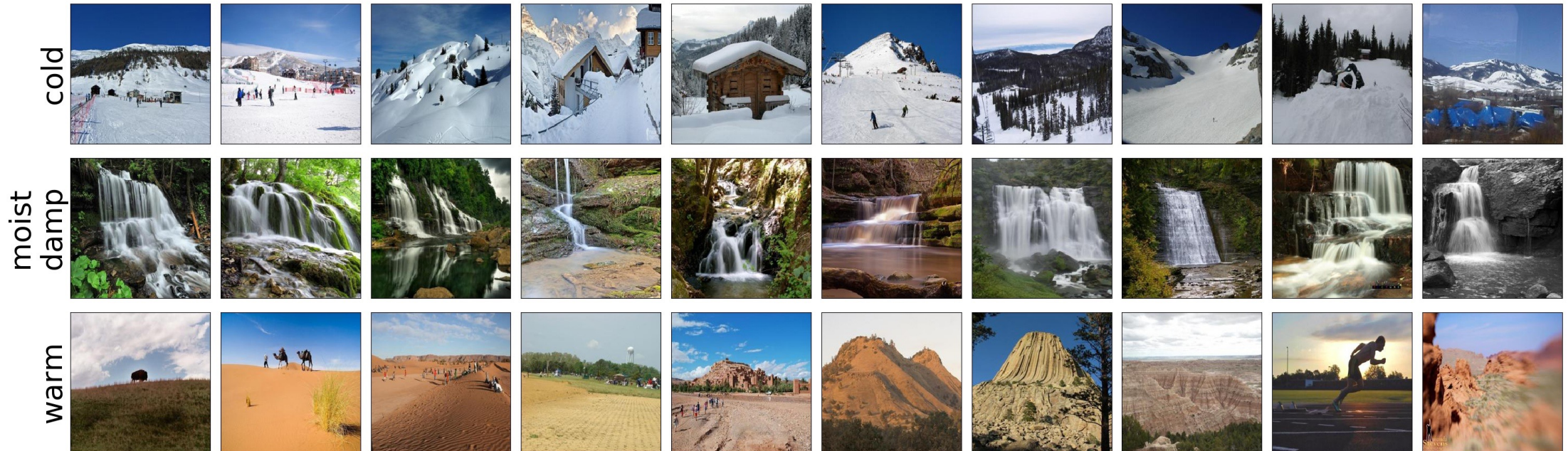What do the learned concepts look like?

# Visualize the concept axes

Most activated                                    16th layer
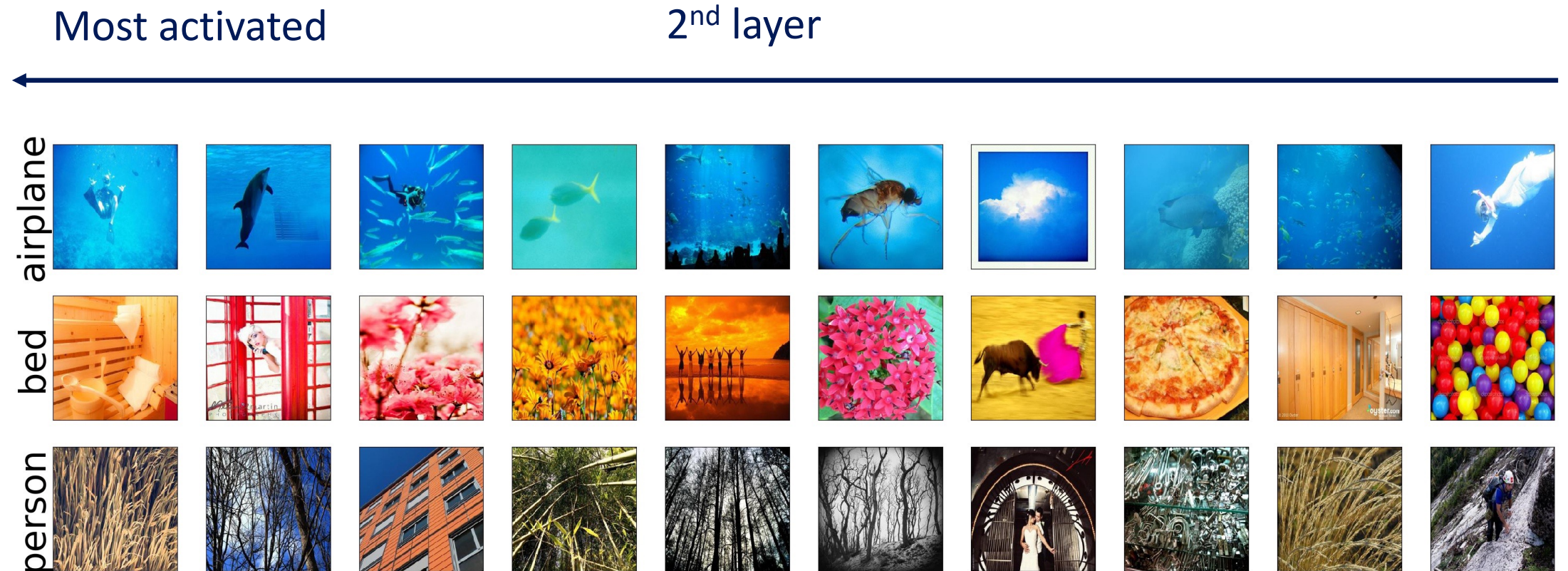
# Visualize the concept axes

- Not only objects - weather

# Visualize the concept axes

- Not only objects - material

# Visualize the concept axes

Most activated

2<sup>nd</sup> layer

How to quantitatively measure the quality
of the learned concepts?

# Concept separation

$$\text{inter- intra-concept ratio} = \frac{\text{avg cos similarity between concept } i \text{ and concept } j}{\sqrt{\text{avg cos similarity concept } i}\sqrt{\text{avg cos similarity concept } j}}$$

# Concept separation



directly build a concept classifier in the latent space

**a** BatchNorm (avg inter-sim = 0.74)

**b** auxiliary concept classification loss (avg inter-sim = 0.74)

**c** CW (avg inter-sim = 0.05)

# Concept purity



AUC of the activation measures concept purity

# Concept purity
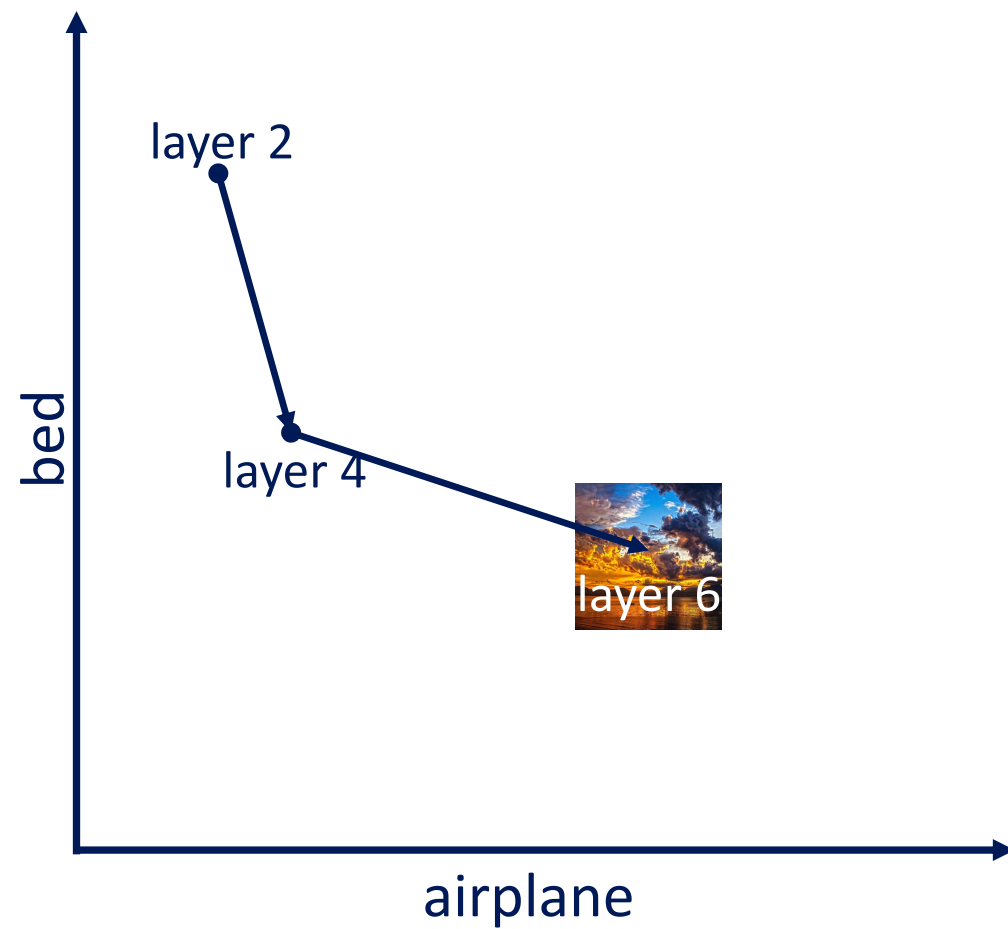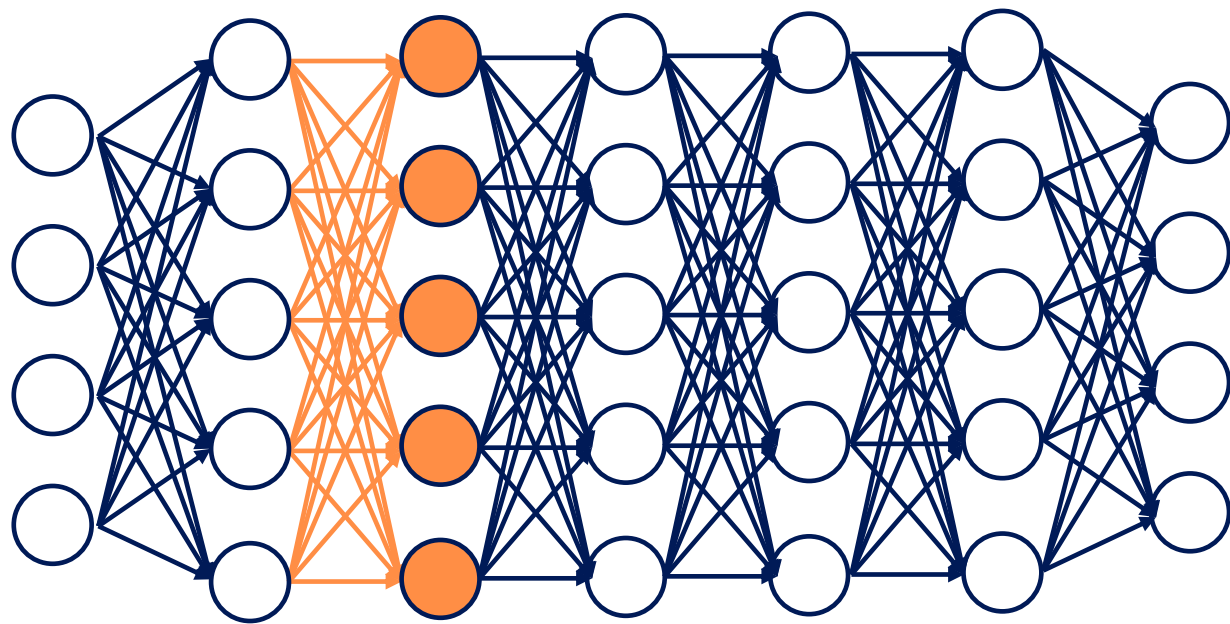

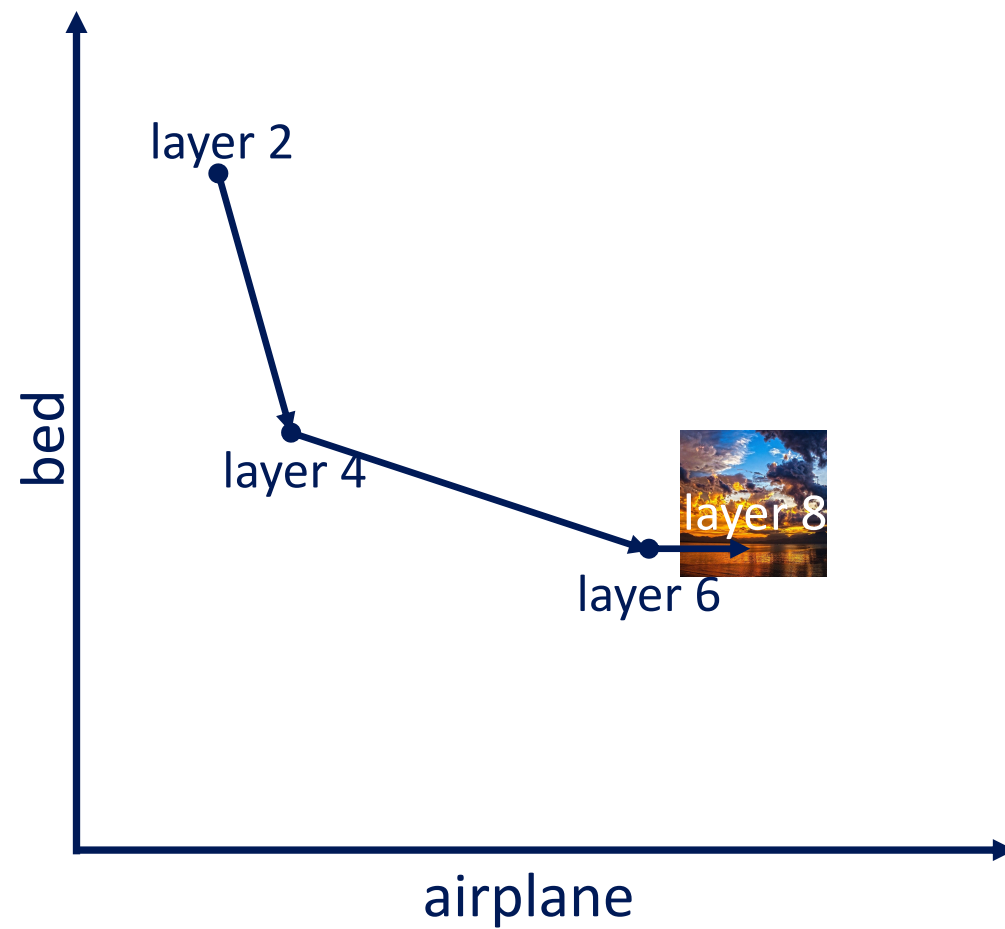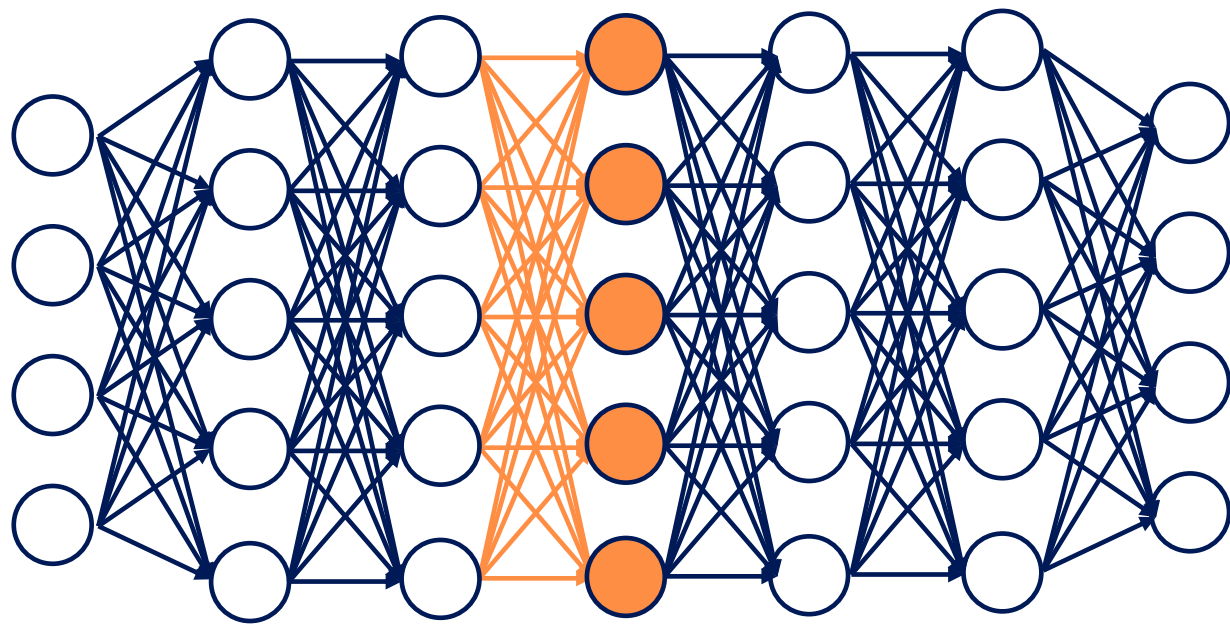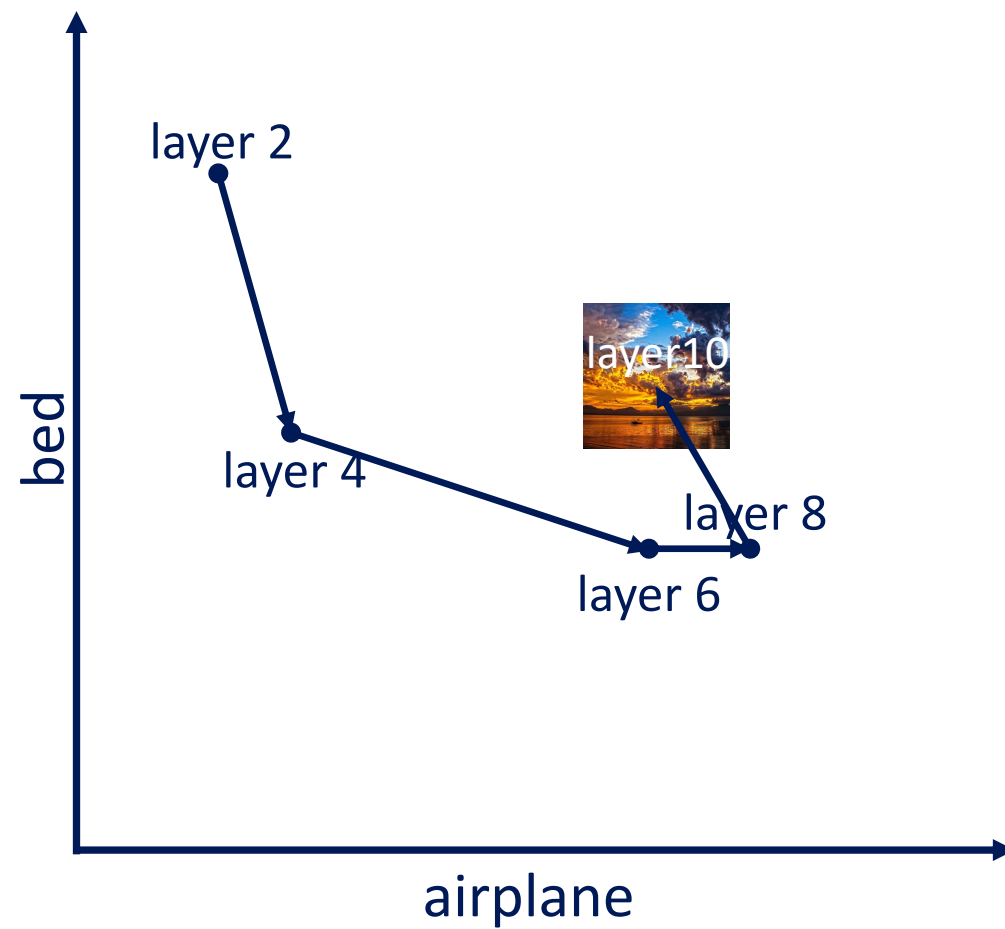
**a** airplane

**b** person

What can we use this model for?

Duke

# Reasoning process

# Reasoning process



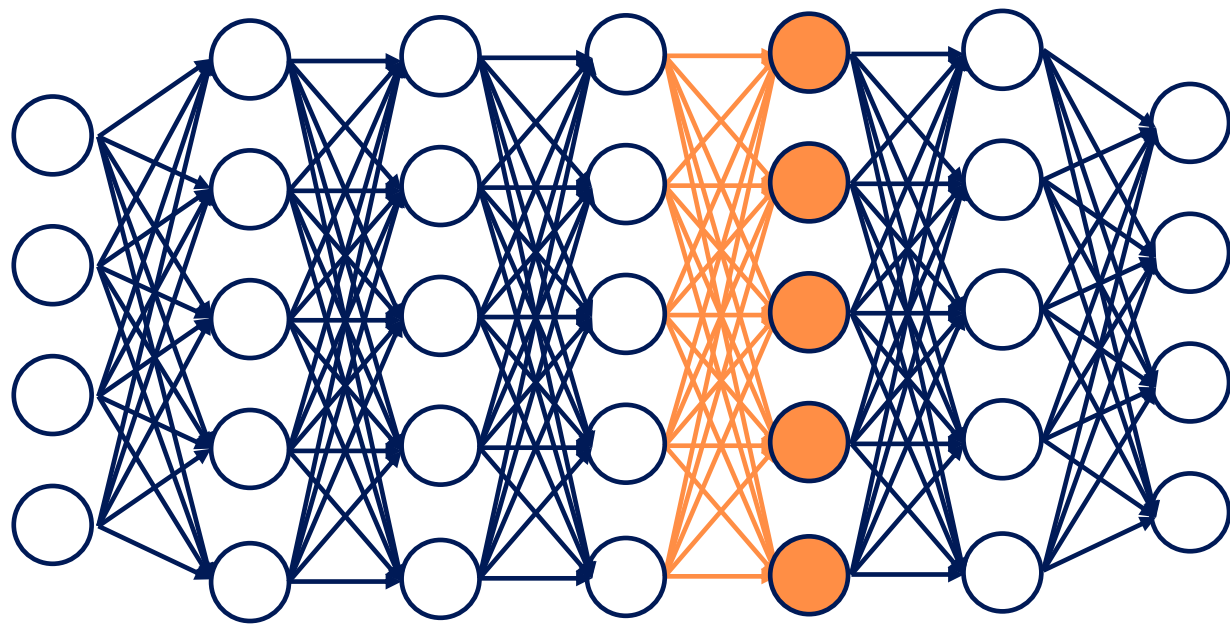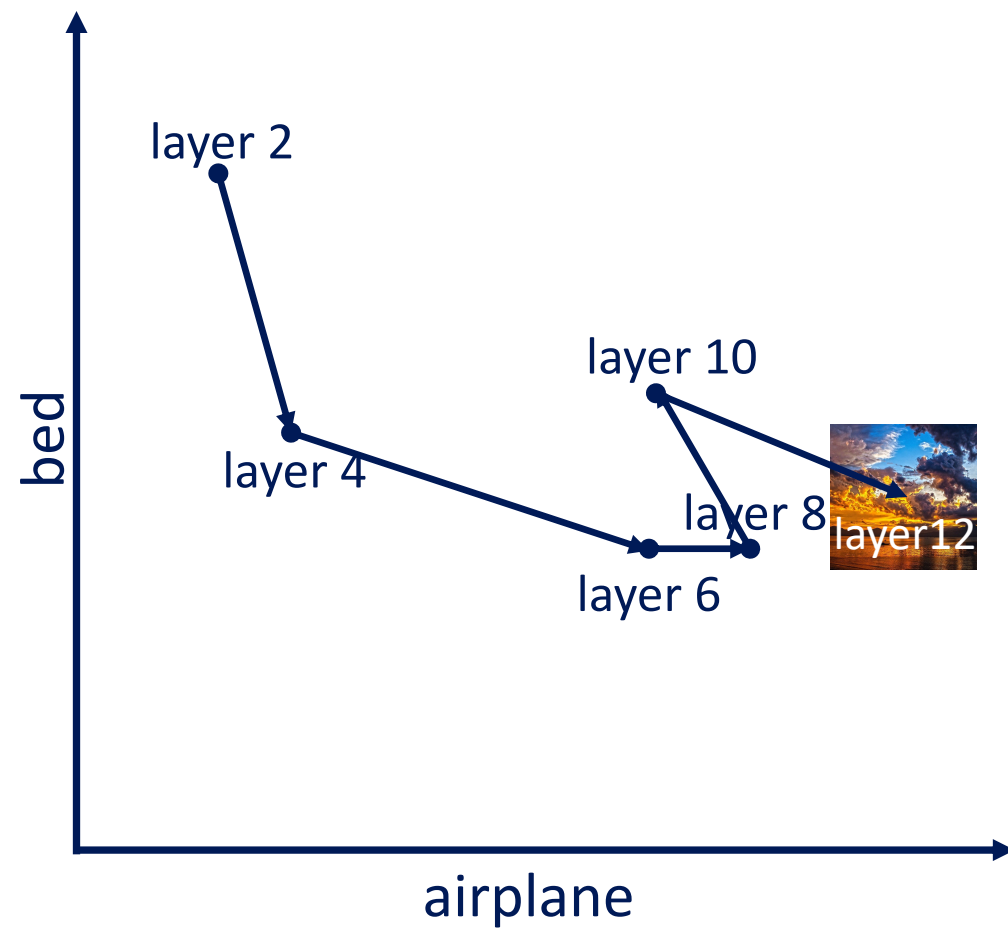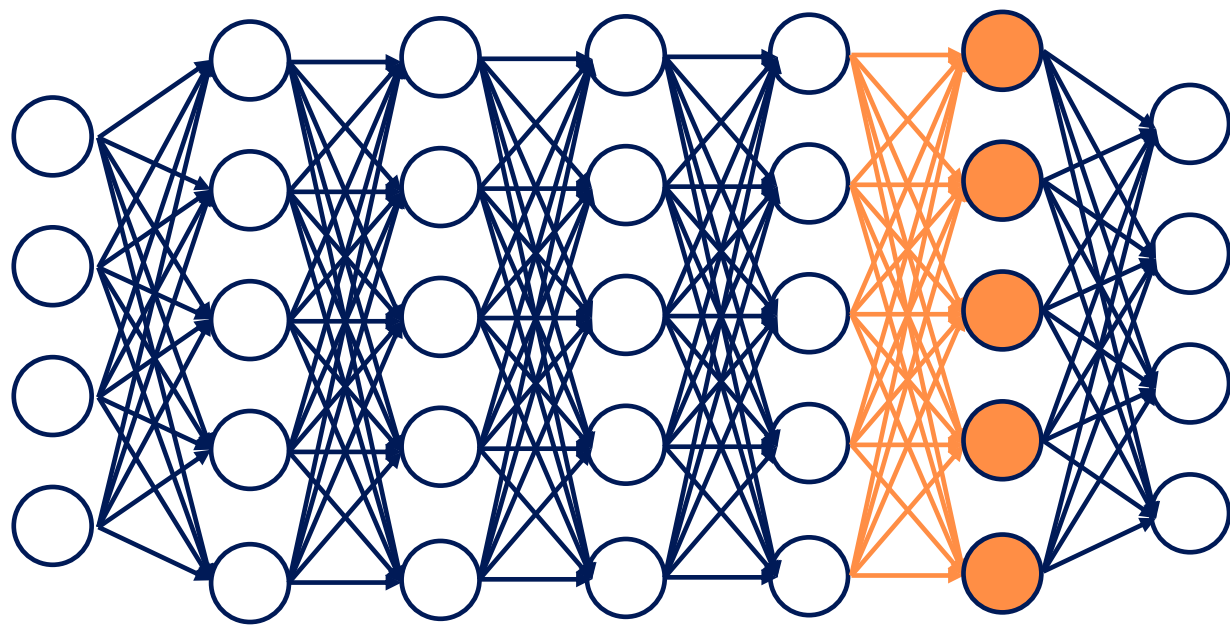Duke

# Reasoning process

# Reasoning process

# Reasoning process

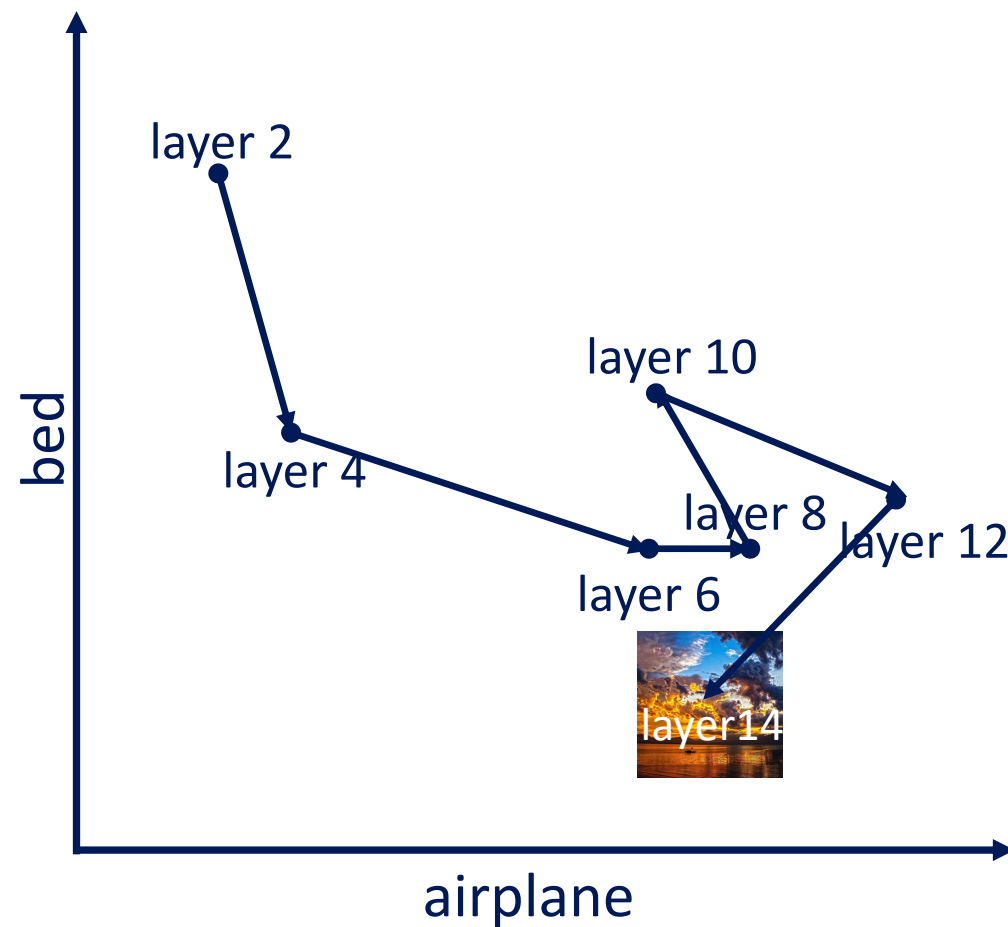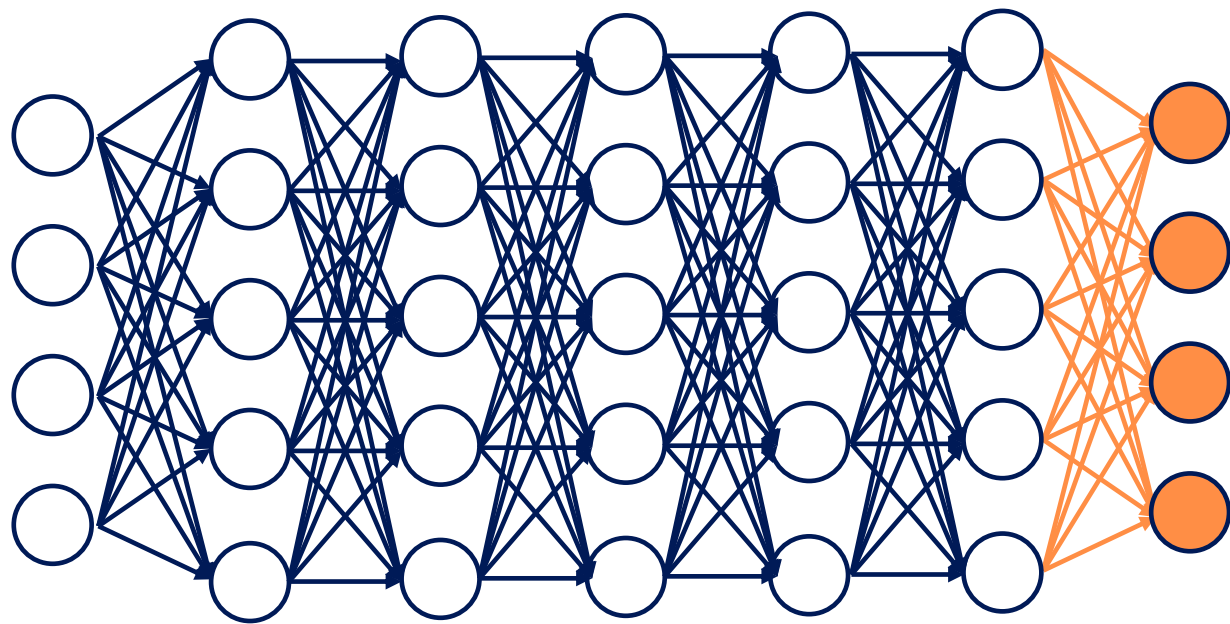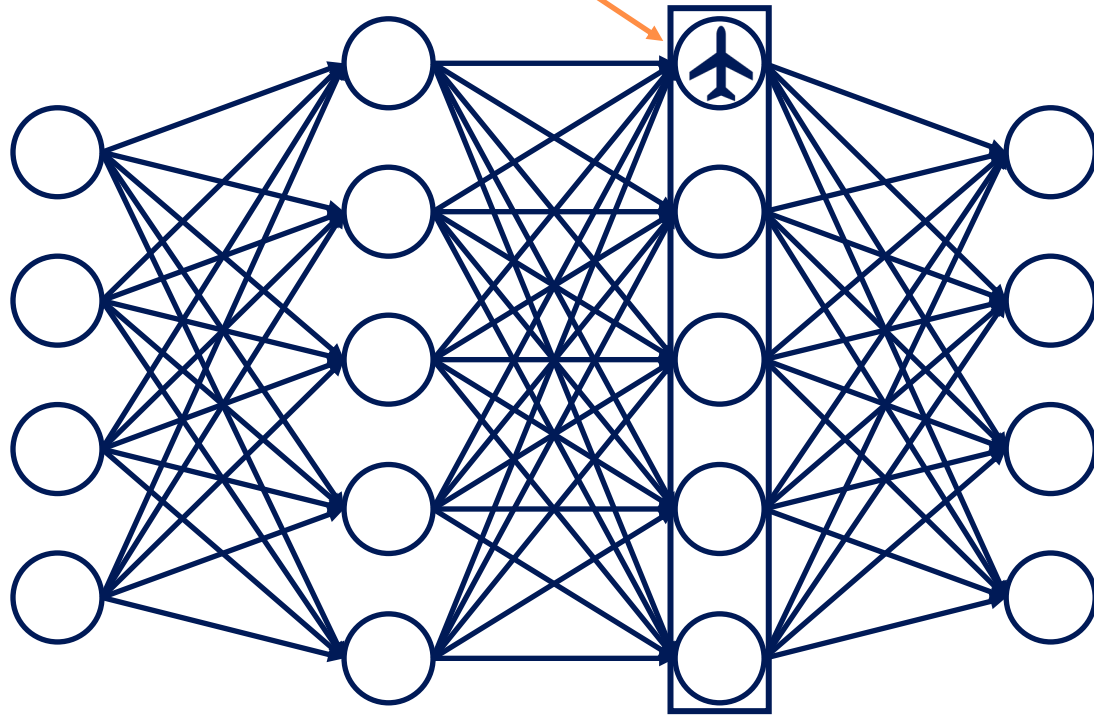# Reasoning process

# Reasoning process



Duke

# Concept importance
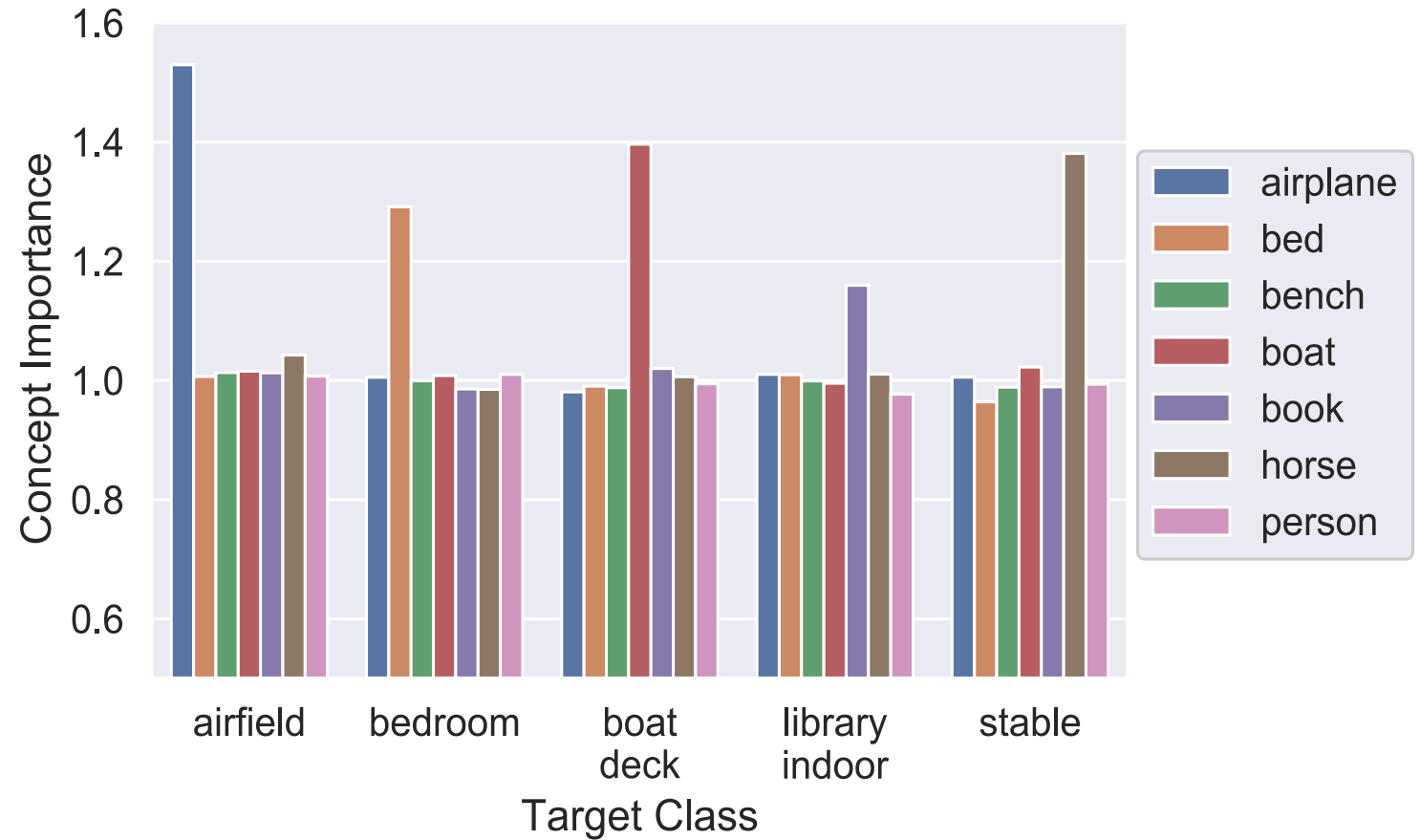
permute the output



- Variable importance of axis j

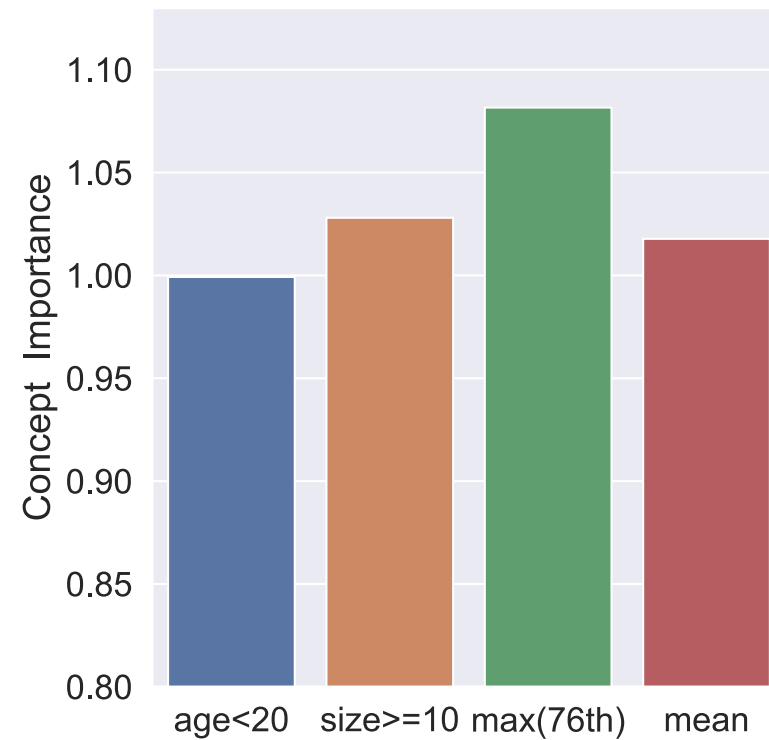$$CI_j = \frac{e_{\text{switch}}^{(j)}}{e_{\text{original}}}$$

Duke

# Concept importance
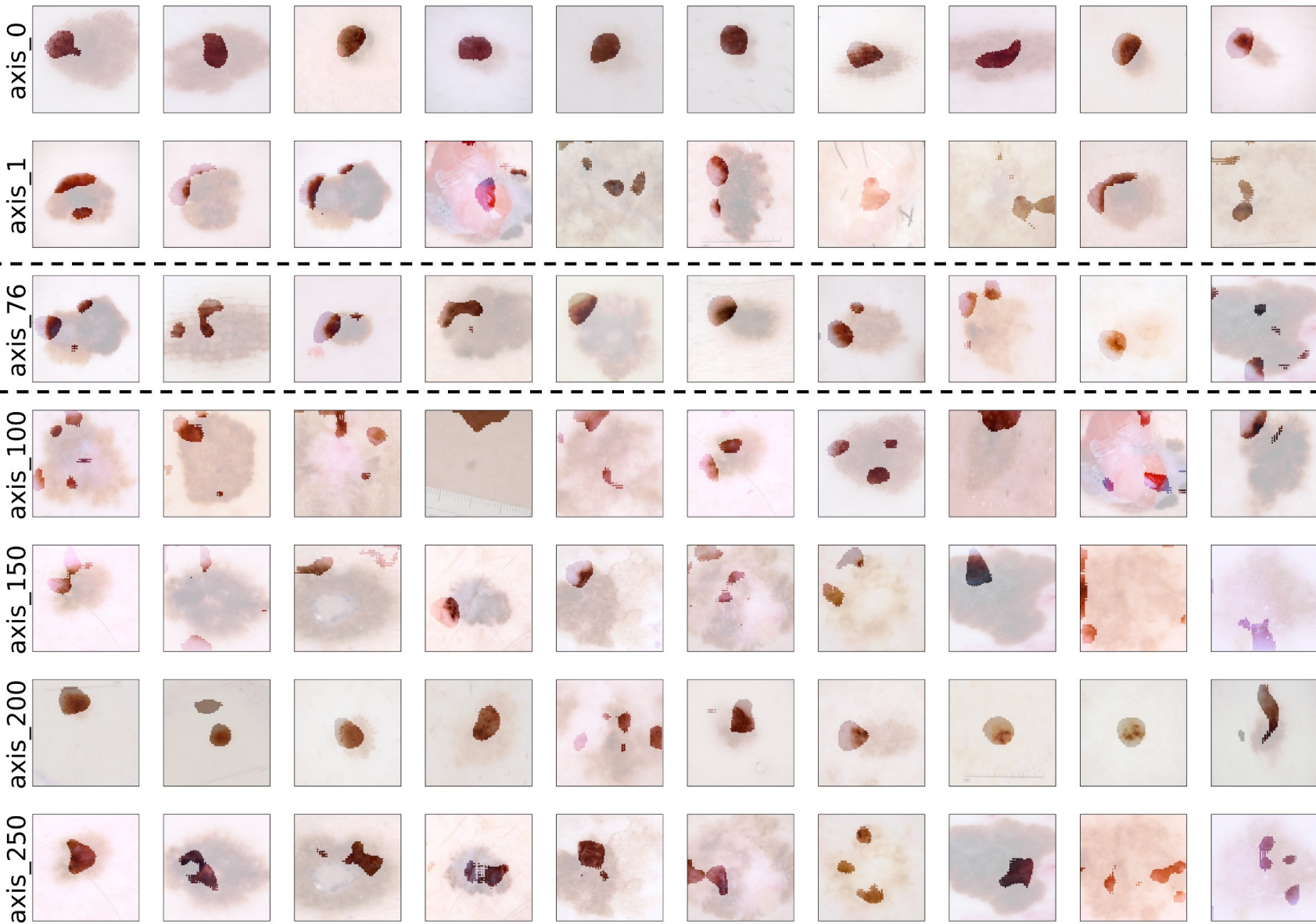
- Scene classification
  - Places365

# Concept importance

- Skin lesion malignancy
  - ISIC dataset
  - axis 1:  age < 20
  - axis 2:  size >= 10 mm
  - not most important



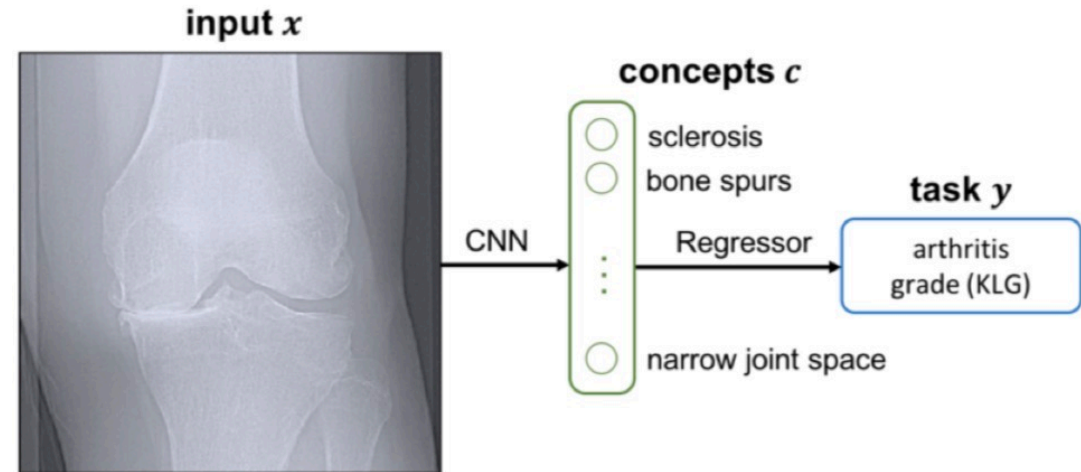Duke

Most Activated

axis_0
axis_1

focused on boundaries

axis_76

axis_100

doctors also agree
boundary > size > age

axis_150

axis_200

axis_250
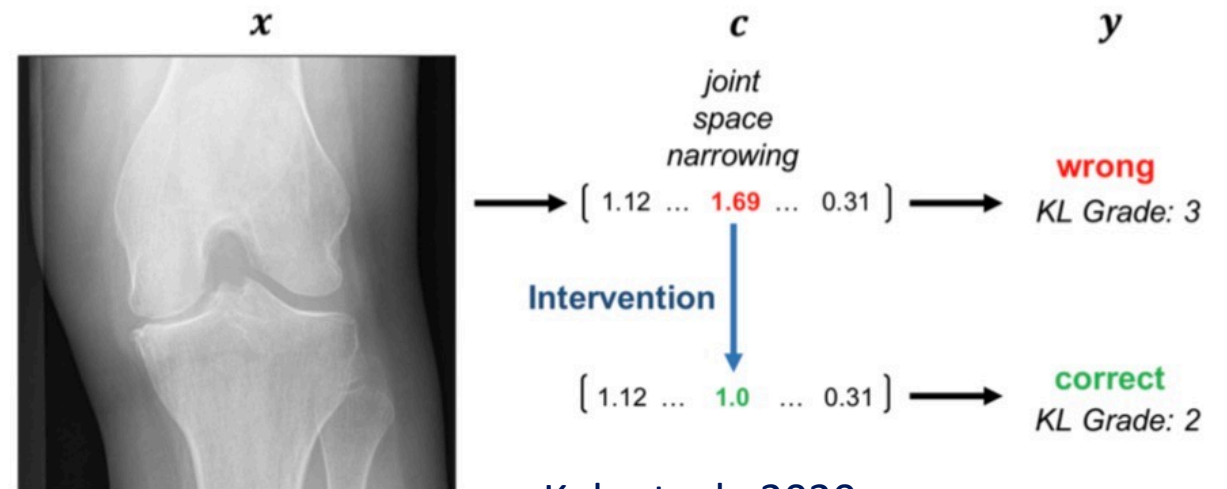
Duke

# Model intervention and editing

- Concept Bottleneck Model (Koh et. al , 2020)
  - they didn't disentangle
  - concept-based models can do test-time intervention



**doctors can change the model when it is wrong**
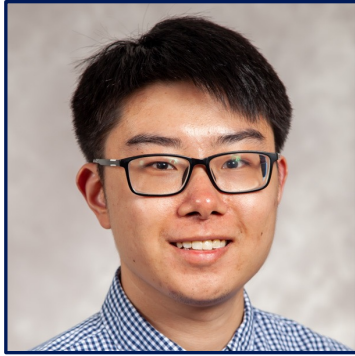


Koh et. al , 2020

Duke

# Summary: Concept Whitening

➢ Better interpretability

   - concepts are disentangled in the latent space

➢ No sacrifice in accuracy

   - accuracy is on par with standard CNNs

➢ Easy to use

   - warm-start from pretrained model requires only one
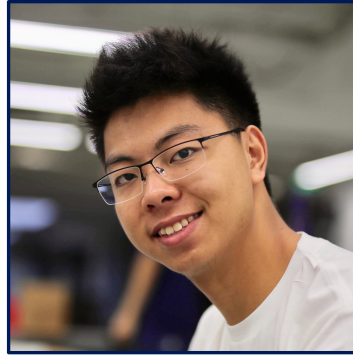     additional epoch of further training

Duke

# Links

➢ Nature Machine Intelligence paper

  - https://rdcu.be/cbOKj

➢ Code

  - https://github.com/zhiCHEN96/ConceptWhitening

Duke

# Thank you



Zhi Chen

Yijie Bei

Cynthia Rudin

Duke