# Introduction to Machine Learning

## Robert Haase

CENTER FOR SCALABLE DATA ANALYTICS AND
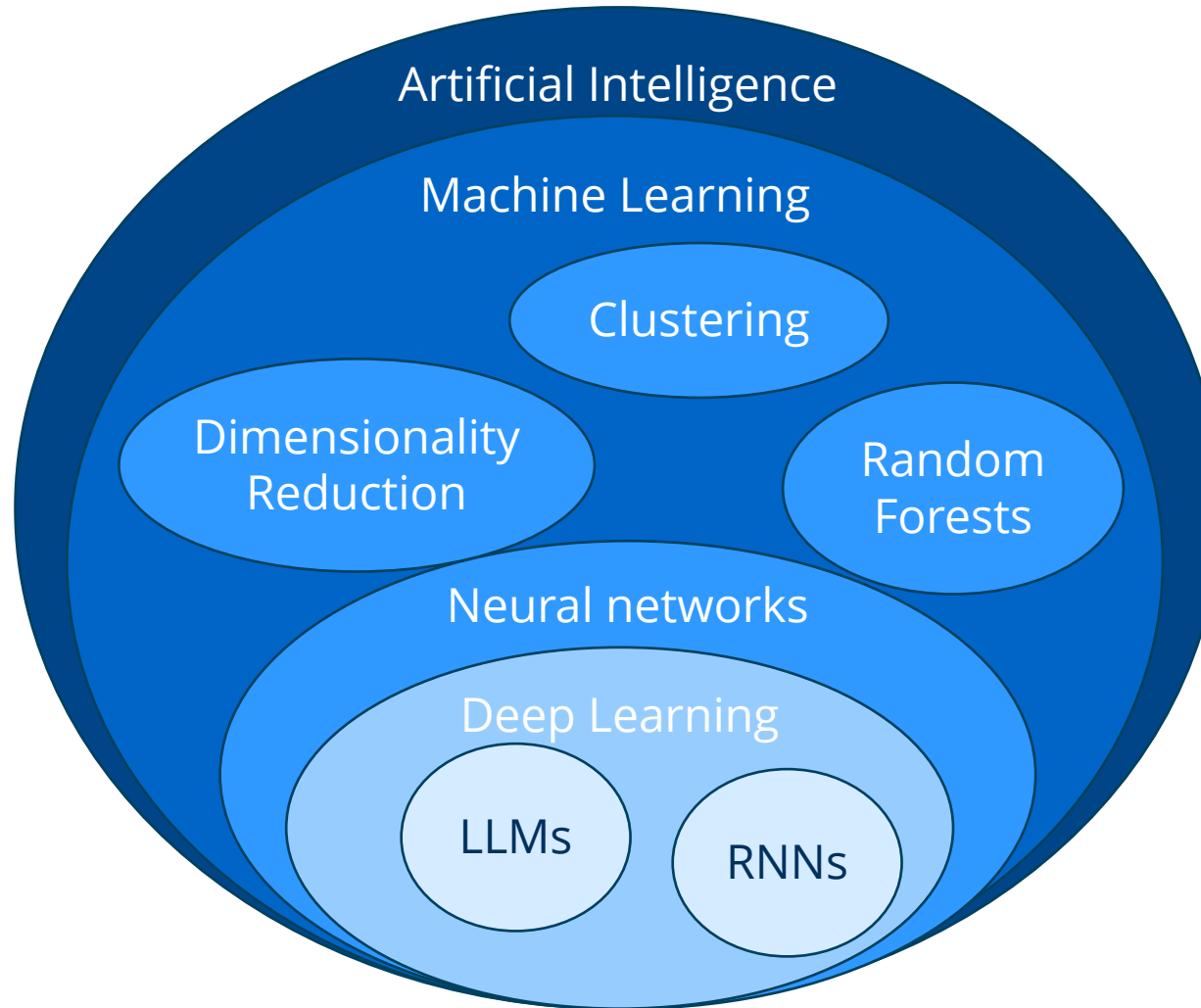ARTIFICIAL INTELLIGENCE

Funded by

Bundesministerium
für Bildung
und Forschung

SACHSEN

Diese Maßnahme wird gefördert durch die Bundesregierung
aufgrund eines Beschlusses des Deutschen Bundestages.
Diese Maßnahme wird mitfinanziert durch Steuermittel auf
der Grundlage des von den Abgeordneten des Sächsischen
Landtags beschlossenen Haushaltes.

Reusing materials from Johannes Soltwedel, Till Korten, Johannes Müller, Laura Žiguty (TU Dresden), Ryan Savill (MPI-CBG), Matthias Täschner (ScaDS.AI/Uni Leipzig) and the Scikit-learn community.

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Artificial intelligence

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Artificial intelligence

**Narrow AI**
- Application specific
- Trained on labelled data
- Reflexive tasks
- Cannot extrapolate
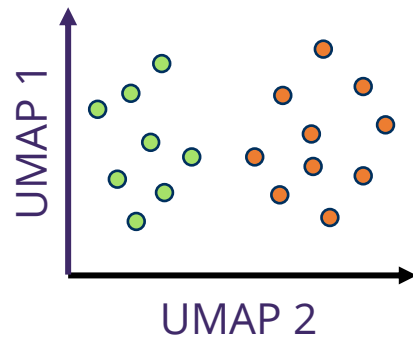
Great for data analysis tasks

**General AI**
- Human capabilities
- Access to knowledge of humanity, beyond individuals
- Can create *new* solutions by working creatively

Robert Haase
@haesleinhuepf
AI4Seismology
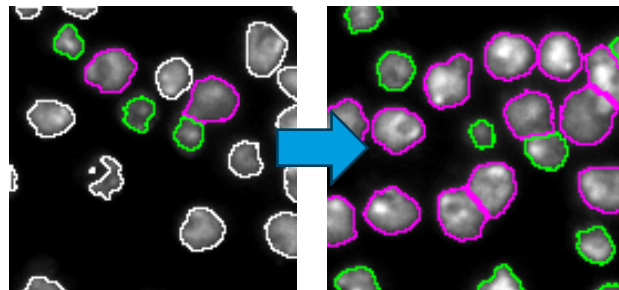May 5th 2025

# Artificial intelligence

## Unsupervised ML

- Dimensionality reduction
- Clustering
- Detecting patterns in unlabeled data
- Hypothesis generation



## Supervised ML

- Learning tasks otherwise only humans could do
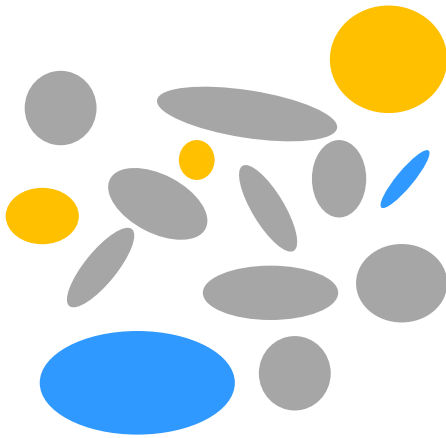- Train a model, predict a classification



## Generative AI

- Produces new data provided a context, often with human language prompts
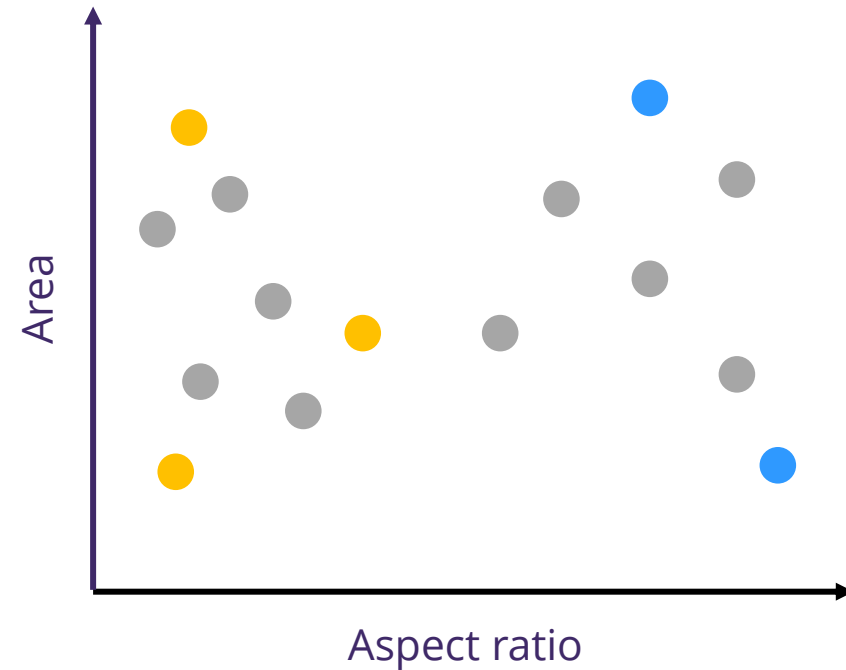- Hyped since 2022, with yet unclear limitations

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Labelled data

- E.g. for shape differentiation of objects
- Partially labelled data — Bias?

Elongated
Round
Unlabelled

Area

Aspect ratio

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025
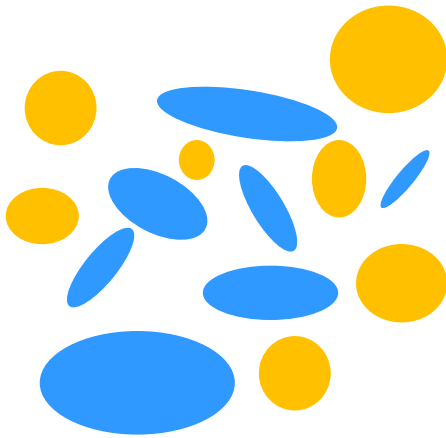
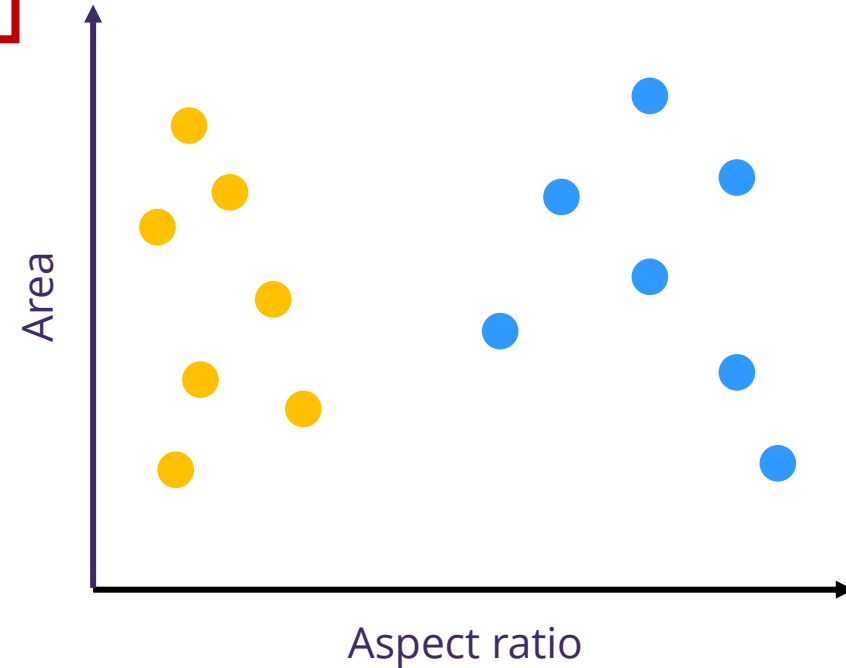# Labelled data

- E.g. for shape differentiation of objects
- Fully labelled data

Typically expensive



Elongated
Round
Unlabelled

Area

Aspect ratio

CENTER FOR SCALABLE DATA ANALYTICS AND
ARTIFICIAL INTELLIGENCE

# <u>Unsupervised</u> Machine Learning

## Robert Haase

Reusing materials from Johannes Soltwedel, Till Korten, Johannes Müller, Laura Žiguty (TU Dresden), Ryan Savill (MPI-CBG), Matthias Täschner (ScaDS.AI/Uni Leipzig) and the Scikit-learn community.

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Hypothesis-driven quantitative science

Hypothesis: The amplitude of a given signal is an indicator for upcoming earthquakes.

Null-Hypothesis: There is no relationship between the amplitude and future earthquakes.

Data download

> Shall we use a different dataset / sensor?

Data preprocessing

> Shall we use a different denoising algorithm?

> Shall we modify our measurement + hypothesis?

Amplitude measurement

Statistics

> Shall we use a different statistical test?

Reject / accept null-hypothesis

> Be careful going down this rabbit hole, you may be leaving good scientific practice behind.

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Data-driven quantitative science

~~Hypothesis: The amplitude of a given signal is an indicator for upcoming earthquakes.~~

Question: Which measurement is a good predictor for upcoming earthquakes?

**Which sensor / data is the most reliable?**
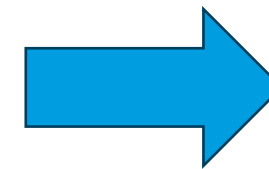
Data download (multiple sources, sensors, ...)

Data preprocessing using Method A, B, C

**Why?**

Amplitude, frequency, wavelength, ... measurement

**Which parameter shows any relationship with upcoming earthquakes?**

Statistics

➡️ **Hypothesis generation**

# Feature selection

- Which measurement / parameter / feature is related to the effect I'm investigating?

- Example goals:



Signal classification

Pixel classification

Object classification

- Amplitude
- Energy
- Duration
- …

- Silence
- Tourists jumping on a sensor
- Earthquake approaching

original    top_hat(10)

gaussian_sobel(1)    random

- Area
- Perimeter
- Aspect ratio
- …

- Round
- Elongated

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Feature selection

Question: Which features shall I analyse?

Challenges:

- Physical properties versus measurable features

- Correlation versus causation

- Too many features

If you have no idea -> unsupervised machine learning

- Dimensionality reduction

- Clustering

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Dimensionality reduction: Principal Component Analysis (PCA)

Linear transformation of high-dimensional data to concentrate information in a lower dimensional *embedding*
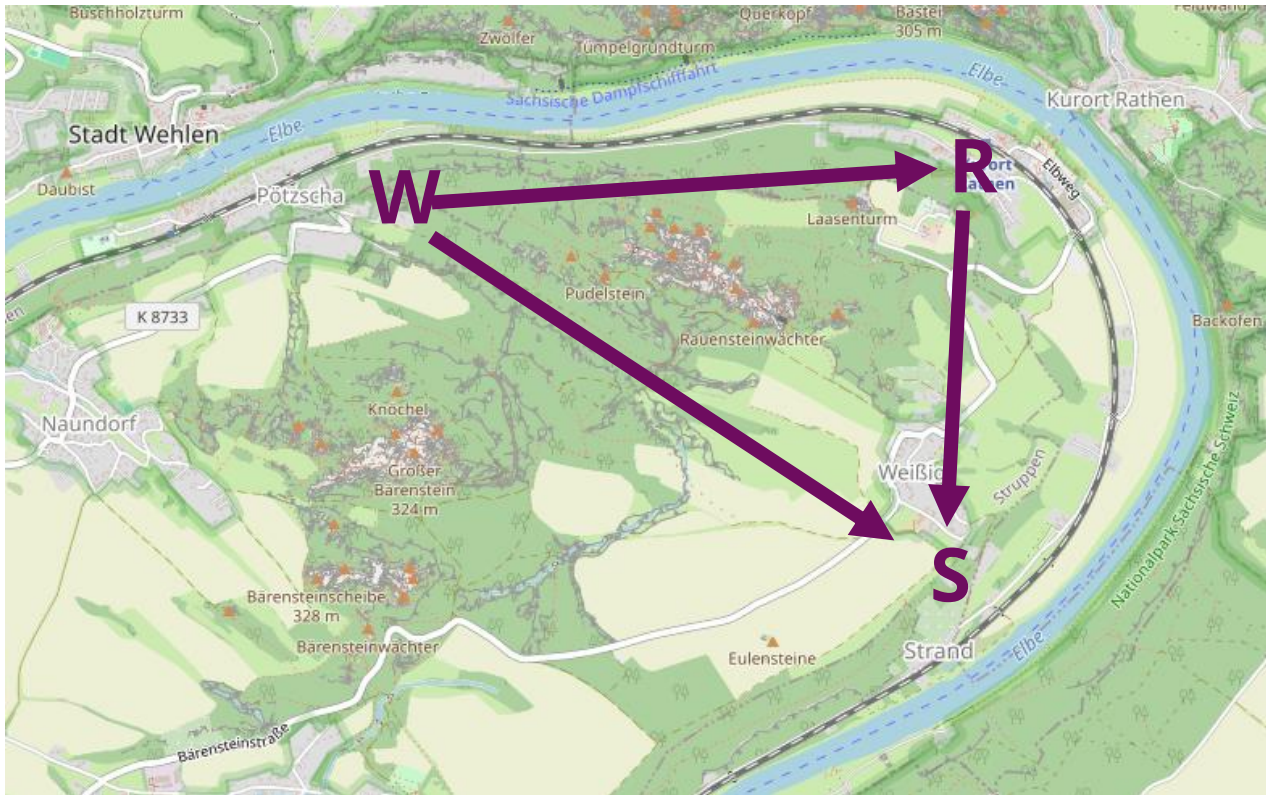
# Embeddings

- N-dimensional latent space

- Axes typically have no meaningful/physical name (PCA1, UMAP1, ...) and no physical unit

- Allow representing complex measurements, things, relationships in numeric space.

- Example:
  - You measure amplitude, frequency, wavelength, etc.,
  - derive a 2D-embedding from it,
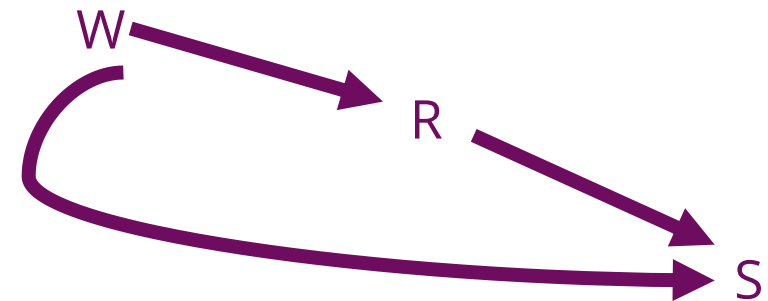  - to visualize the data or
  - to better process data

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Non-Euclidian spaces

## Not all features might be distances



Use travel time between W and S as metric for distance

→ Travelling from **W**ehlen to **S**trand by bike is probably faster if you make a detour through **R**athen
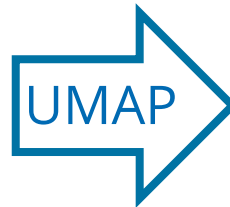
# Uniform Manifold Approximation Projection (UMAP)
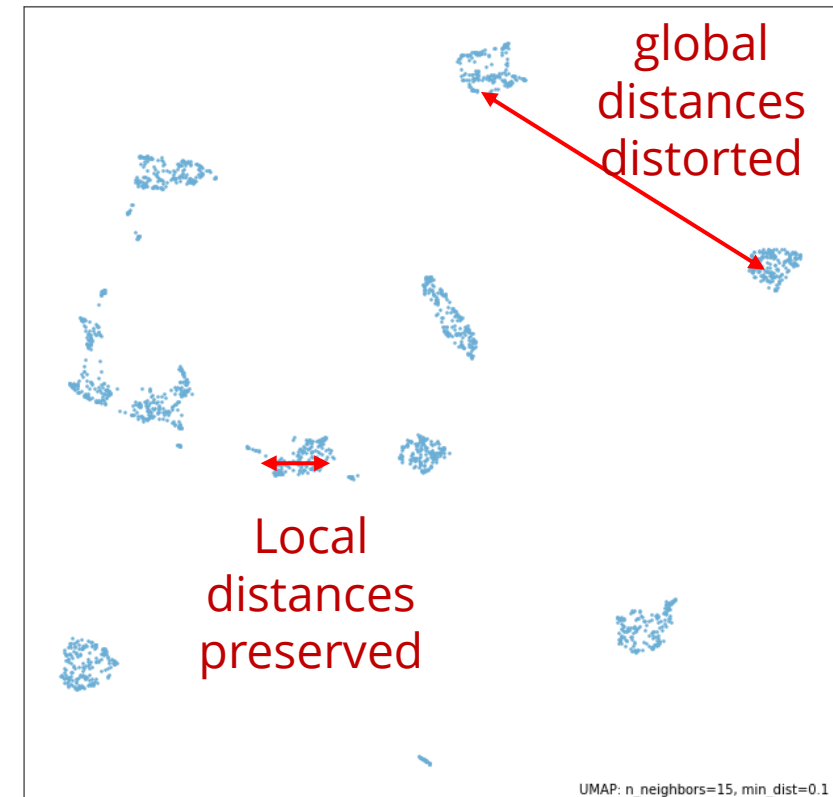
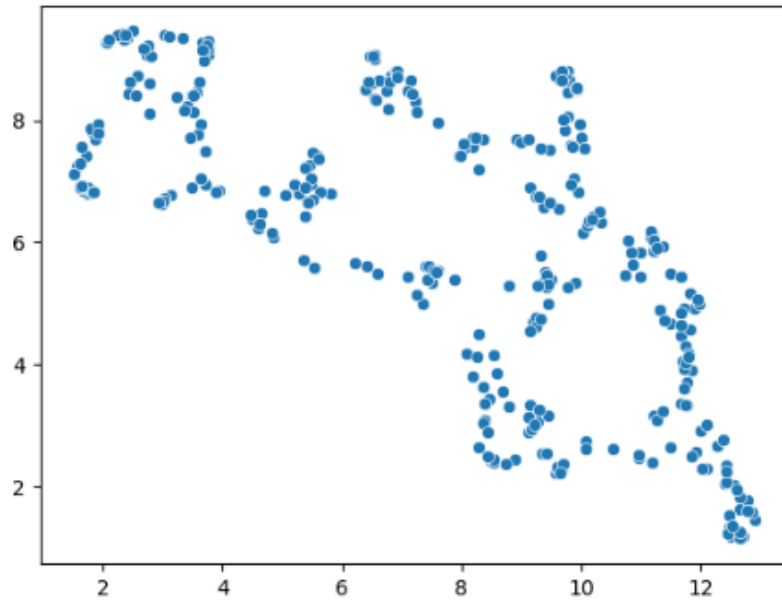Structural, hierarchical, non-linear transformation

Modifies density of data points.

# Uniform Manifold Approximation Projection (UMAP)

Non-deterministic algorithm: You execute it twice, you get different results.

```
[11]: reducer = umap.UMAP()
      embedding2 = reducer.fit_transform(scaled_statistics)

      seaborn.scatterplot(x=embedding2[:, 0],
                          y=embedding2[:, 1])
```

`[11]: <AxesSubplot: >`

```
[12]: reducer = umap.UMAP()
      embedding2 = reducer.fit_transform(scaled_statistics)

      seaborn.scatterplot(x=embedding2[:, 0],
                          y=embedding2[:, 1])
```

`[12]: <AxesSubplot: >`



Locally rotated

Globally rotated

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/47_clustering/umap.html?highlight=umap#a-note-on-repeatability

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

ScaDS.AI
DRESDEN LEIPZIG

# Dimensionality reduction

Uniform manifold approximation and projection (UMAP)
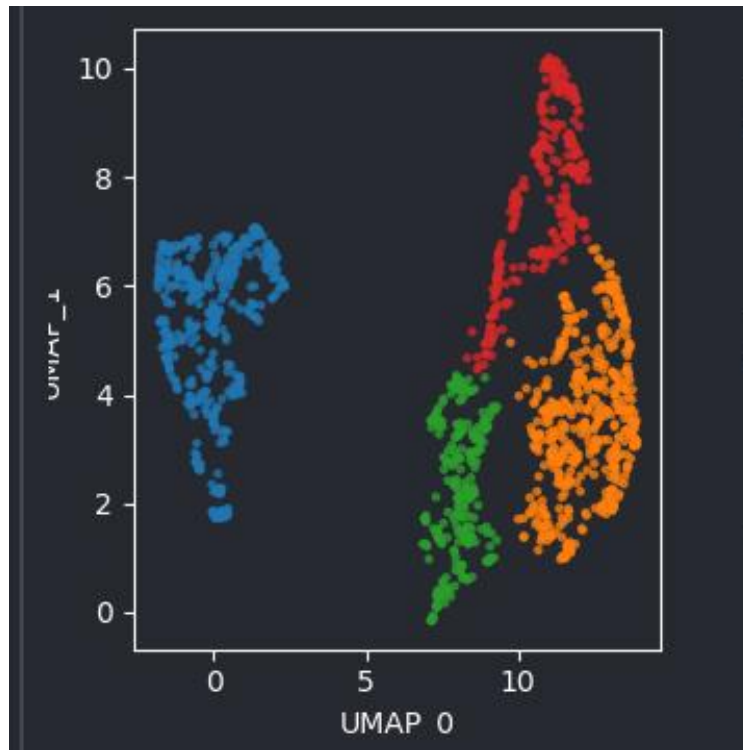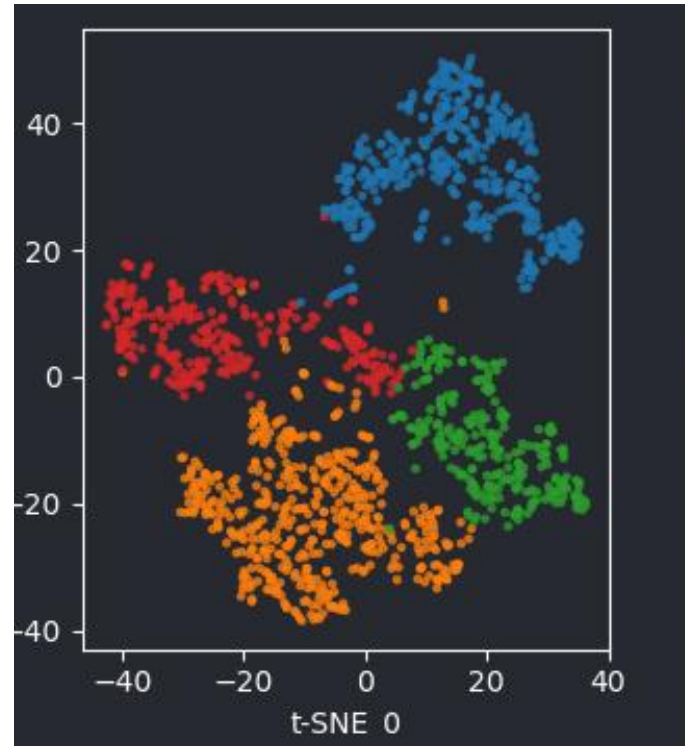
t-distributed stochastic neighbor embedding (t-SNE)

Principal component analysis (PCA)

# Clustering

Unsupervised machine learning may include grouping objects without given ground truth

# Clustering

Unsupervised machine learning may include grouping objects without given ground truth



Round  Elongated

Names given by human observer *after* grouping / clustering

UMAP 1

frequency

UMAP 2

frequency

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# K-Means Clustering

Clustering algorithm, where you *only* need to specify the number of clusters.

**Step1: Random initialization of cluster centers**

**Step2: Tessellation of space into cluster regions**

**Step3: Replace cluster center with centroids**

**Step4: Repeat 2&3 until convergence**

Slide adapted from Johannes Soltwedel, TU Dresden

# K-Means Clustering

Clustering algorithm, where you *only* need to specify the number of clusters.



Iteration #1

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

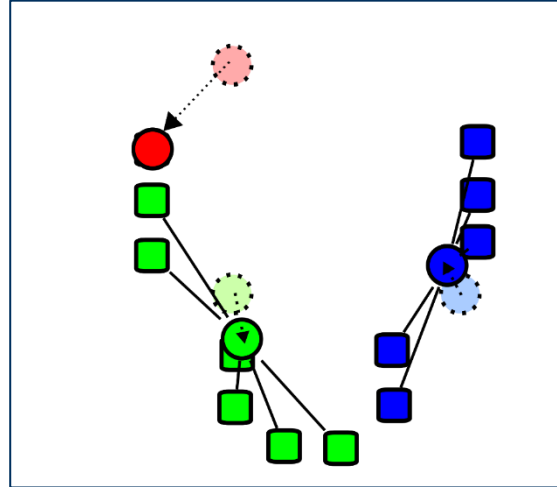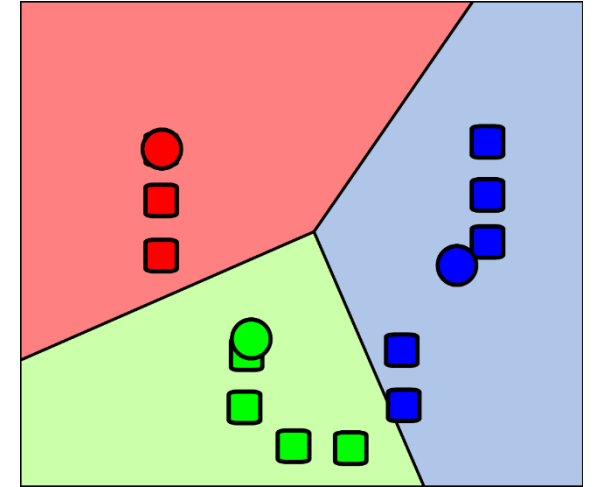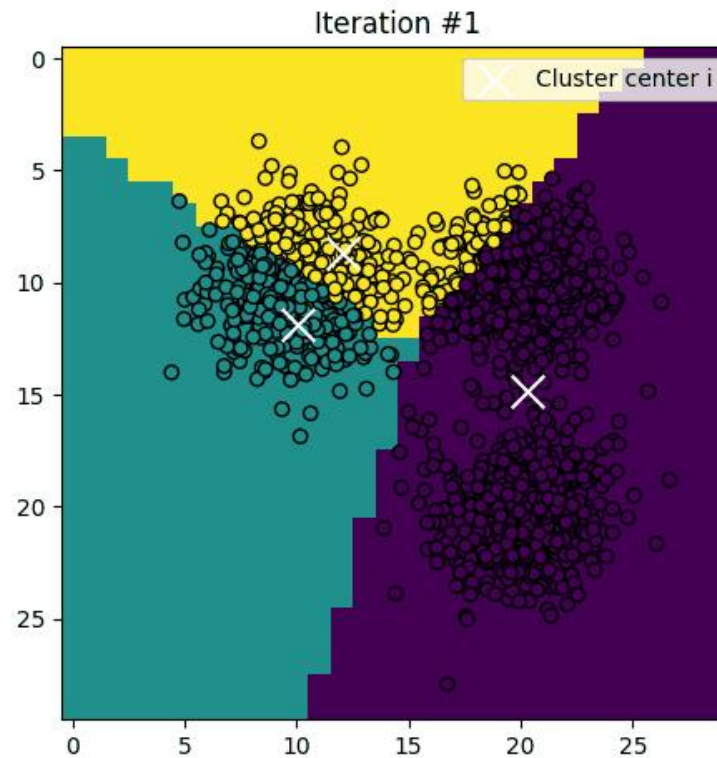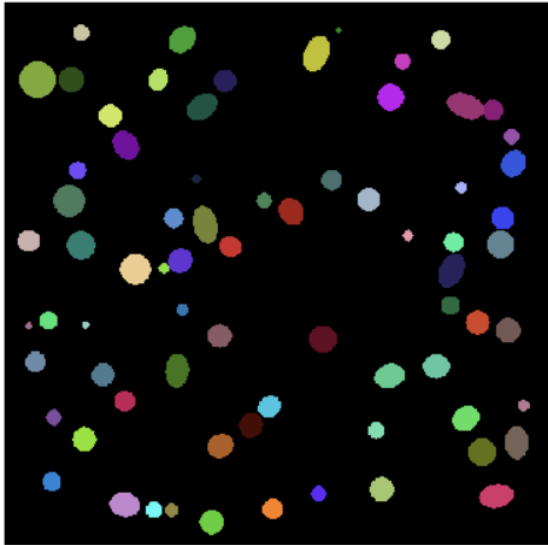Goal: Understand shape measurements

Data: Shape measurements from *randomly* shaped blobs.



| | label | area | perimeter | minor_axis_length | major_axis_length | circularity | solidity | aspect_ratio | elongation |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 97.0 | 32.970563 | 11.092860 | 11.092860 | 1.121318 | 0.788288 | 1.000000 | 0.000000 |
| **1** | 2 | 285.0 | 60.284271 | 19.052651 | 19.052651 | 0.985477 | 0.785116 | 1.000000 | 0.000000 |
| **2** | 3 | 473.0 | 79.597980 | 21.823280 | 27.594586 | 0.938138 | 0.785448 | 1.264456 | 0.209146 |
| **3** | 4 | 321.0 | 63.112698 | 19.033334 | 21.456036 | 1.012701 | 0.786033 | 1.127287 | 0.112915 |
| **4** | 5 | 407.0 | 72.769553 | 22.155138 | 23.384406 | 0.965839 | 0.785586 | 1.055485 | 0.052568 |

...

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 1: Dimensionality reduction (UMAP)

Observation: There appear to be *2 distinct groups*

Pinning the random seed is no solution to this general problem.

Beware: UMAPs are non-deterministic. Different runs lead to different results.

Run 1:

Run 2:

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 2: Clustering data into 2 clusters

Using K-Means clustering

Pinning the random seed is no solution to this general problem.

Beware: Clustering-algorithms are non-deterministic. Different runs lead to different results.

Run 1:

Run 2:

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 3: Feature selection

Based on correlation
with distance to cluster-centers

*Side note: beware of feature correlation.*

Hypothesis:
"Circularity and minor_axis_length
allow to predict classification."

*Hypothesis generation*

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 4: Train a classifier (supervised ML)

Goal: Eliminate non-determinism



Clustering result (non-deterministic)



Classification result (deterministic, repeatable)

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Supervised Machine Learning

## Robert Haase

CENTER FOR SCALABLE DATA ANALYTICS AND ARTIFICIAL INTELLIGENCE

Funded by

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Supervised Machine learning

Automatic construction of predictive models from given data

Pixels,     Objects,     Images, Audio, Sensor data, Text, Measurements, ...

Annotated raw data, often generated by humans

Classification (categorical)

Regression (continuous numerical)

n = 11

green_magenta_ratio=0.3

$P_{Cat}$ = 0.5
$P_{Microscope}$ = 0.4

Height = 80 cm

Cat ~~Dog~~

Earth-quake ~~Wind~~

Raw data

Training

Ground truth

Prediction

Model

Classification / regression

Quality

Accuracy, Precision, Recall, ...

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Goal

Guess classification (color) from position of a sample in parameter space.



Input data     Decision Tree     Random Forest     Neural Net

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

Adapted from https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Machine learning for image segmentation

*Supervised* machine learning: We give the computer some ground truth to learn from

The computer derives a *model* or a *classifier* which can judge if a pixel should be foreground (white) or background (black)

Example: Binary classifier



Training

Model / classifier

Raw image

Binary image

@haesleinhuepf
AI4Seismology
May 5th 2025

# Random forest based image segmentation

Decision trees are classifiers, they decide if a pixel should be white or black

Random decision trees are randomly initialized, afterwards evaluated and selected

Random forests consist of many random decision trees

Example: Random forest of binary decision trees

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Deriving random decision trees

For efficient processing, we randomly *sample* our data set
- Individual pixels, their intensity and their classification



Note: You cannot use a single threshold to make the decision correctly

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Deriving random decision trees

Decision trees combine several thresholds on several parameters

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Deriving random decision trees

Depending on sampling, the decision trees are different

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Random Forest Pixel Classifiers

By training many decision trees, errors are equilibrated

Sampling

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Random Forest Pixel Classifiers

Combination of individual tree decisions by voting or max / mean

Prediction

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Random Forest Pixel Classifiers

Typical numbers for pixel classifiers in microscopy

Available features:



- Gaussian blur image
- DoG image
- LoG image
- Hessian
- ….

Depth: 4

Number of trees: > 100

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Model validation

In order to assess model quality, we split the ground truth into two set
- Training set (50%-90% of the available data)
- Test set (10%-50% of the available data)



Training set

Ground truth → Training → Prediction → Classifier

Test set

Raw data → Prediction → Prediction → Ability to abstract → Ground truth

Typically done with hundreds or thousands of cells / images / objects / ...

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Model validation

Based on the theory of sets



Legend:
- (A) Prediction
- (B) Reference / ground truth
- [ROI] Region of interest
- TP True-positive
- FN False-negative
- FP False-positive
- TN True-negative

$$\text{Accuracy} = \frac{\text{correct classifications}}{\text{all classifications}}$$

This means: $= \dfrac{TP + TN}{FP + FN + TP + TN}$

$$\text{Precision} = \frac{\text{Relevant retrieved instances}}{\text{All } \mathbf{retrieved} \text{ instances}}$$

This *may* mean: $= \dfrac{TP}{FP + TP}$

# Model validation: Accuracy versus precision



Accurate and precise

Accurate and but not precise

Not accurate and but precise

Neither accurate nor precise

Lesson learned:
A single quality metric cannot describe the whole situation

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Model validation: Accuracy versus Jaccard Index

## Side-effect of number of true negatives

Nuclei

Reference

Segmentation result

$$A = \frac{TP + TN}{FN + FP + TP + TN}$$

$$J = \frac{TP}{FN + FP + TP}$$

Accuracy: 0.97
Jaccard Index: 0.73

Accuracy decreases because there are less correct black pixels (TN)

Accuracy: 0.95
Jaccard Index: 0.73

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/29_algorithm_validation/jaccard_index_versus_accuracy.html

Slide 41

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

41

CENTER FOR SCALABLE DATA ANALYTICS AND
ARTIFICIAL INTELLIGENCE

# Deep Learning

## Robert Haase

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Deep learning for image analysis

In deep learning, this selection becomes part of the black box



Neural networks

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Neural networks

- How biologists see neurons
- How computer scientists see neurons

"perceptron"

# Neural Networks

- Early form: "Multilayer Perceptron"
- fully connected class of feedforward artificial neural network

If there are *many* hidden layers, we speak of a *deep* neural network

Input layer          n hidden layer(s)          Output layer

# Convolutional neural networks

- Layer types

### Fully connected layer



### Convolutional layer

Field of View (FoV)



### Pooling layer ("Max pool", "Average pool")



| 3 | 15 | 1 | 13 |
|---|----|---|----|
| 9 | 7 | 0 | 10 |
| 11 | 5 | 5 | 3 |
| 1 | 8 | 9 | 6 |

Max pooling →

| 15 | 13 |
|----|----|
| 11 | 9 |

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Convolutional neural networks

- Assuming we had no activation functions in the network layers can be reduced by eliminating brackets!

input

$$y = w_5(w_1 x_1 + w_2 x_2) + w_6(w_3 x_2 + w_4 x_3)$$

$x_1$  $x_2$  $x_3$

$w_1$  $w_2$  $w_4$

$w_3$

$$y = w_5 w_1 x_1 + w_5 w_2 x_2 + w_6 w_3 x_2 + w_6 w_4 x_3$$

$$y = w_5 w_1 x_1 +$$

$w_6$

$w_5$

$$v_1 = w_5 w_1$$
$$v_2 = w_5 w_2 + w_6 w_3$$
$$v_3 = w_6 w_4$$

output  y

$$y = v_1 x_1 + v_2 x_2 + v_3 x_3$$

$x_1$  $x_2$  $x_3$

$v_1$  $v_2$  $v_3$

y

ScaDS.AI DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Activation functions

- Introduction of *non-linearity* and *activation functions* enabled what we call *deep-learning* today.



$$y = f(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Activation functions

- Introduction of *non-linearity* and *activation functions* enabled what we call *deep-learning* today.



Bias b

Activation function

Output

Inputs: $x_1 \rightarrow w_1$, $x_2 \rightarrow w_2$, $x_3 \rightarrow w_3$, $\Sigma$ $\rightarrow$ f $\rightarrow$ y

Weights

**Linear**  **Non-linear**

| Identity | | $x$ |
|---|---|---|
| Binary step | | $\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| Logistic, sigmoid, or soft step | | $\sigma(x) \doteq \dfrac{1}{1 + e^{-x}}$ |
| Rectified linear unit (ReLU)[8] | | $(x)^+ \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max(0, x) = x\mathbf{1}_{x>0}$ |

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Learning: Back propagation

- Step 0: Initialize the network randomly (weights, bias)

- Step 1: Forward pass the input through the network, get an initial prediction

- Step 2: Compare the output with the ground truth, compute the error (loss function)

  - The loss function can be freely defined.

  - Example: mean squared error

- Step 3: Update weights



- Silence
- Tourists jumping on a sensor
- Earthquake approaching

Prediction  0.3  0.4  0.4

Ground truth  0  0  1

Loss  0.18

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

ScaDS.AI DRESDEN LEIPZIG

# Learning: Back propagation

- Updating weights:
  - Set output to the error (per-parameter gradient)
  - Backward-pass: add/subtract gradients from weights, to push the network towards giving the right answer.
- Execute the same procedure for next sample
- Execute the same for multiple *epochs*

Set: 0.3   0.3   -0.6

- Silence
- Tourists jumping on a sensor
- Earthquake approaching

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Train- [validation]- Test-split

Training dataset (e.g. 80% of the data)

Used for training directly

Validation dataset (10% of the data)

After every iteration see if the model overfits

Test dataset (10% of the data)

Final evaluation after training is finished (once)



Underfitting

Overfitting

Loss (lower is better)

Training duration (epochs)

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

https://towardsdatascience.com/how-to-split-data-into-three-sets-train-validation-and-test-and-why-e50d22d3e54c

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Training NNs: Batch size & epochs

Problem:

- Assume you have $10^{10}$ samples and attempt to train for 1000 epochs

-> $10^{13}$ backprop steps required.

Data

epoch 1

# Training NNs: Batch size & epochs

Problem:

- Assume you have $10^{10}$ samples and attempt to train for 1000 epochs

-> $10^{13}$ backprop steps required.

Solution:

- Draw n=1000 random samples from the training data to train for one epoch.

- Next epoch: different n samples.

-> $10^{6}$ backprop steps required.



Data

Batch size

epoch 1

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Training NNs: Batch size & epochs

Problem:

- Assume you have $10^{10}$ samples and attempt to train for 1000 epochs
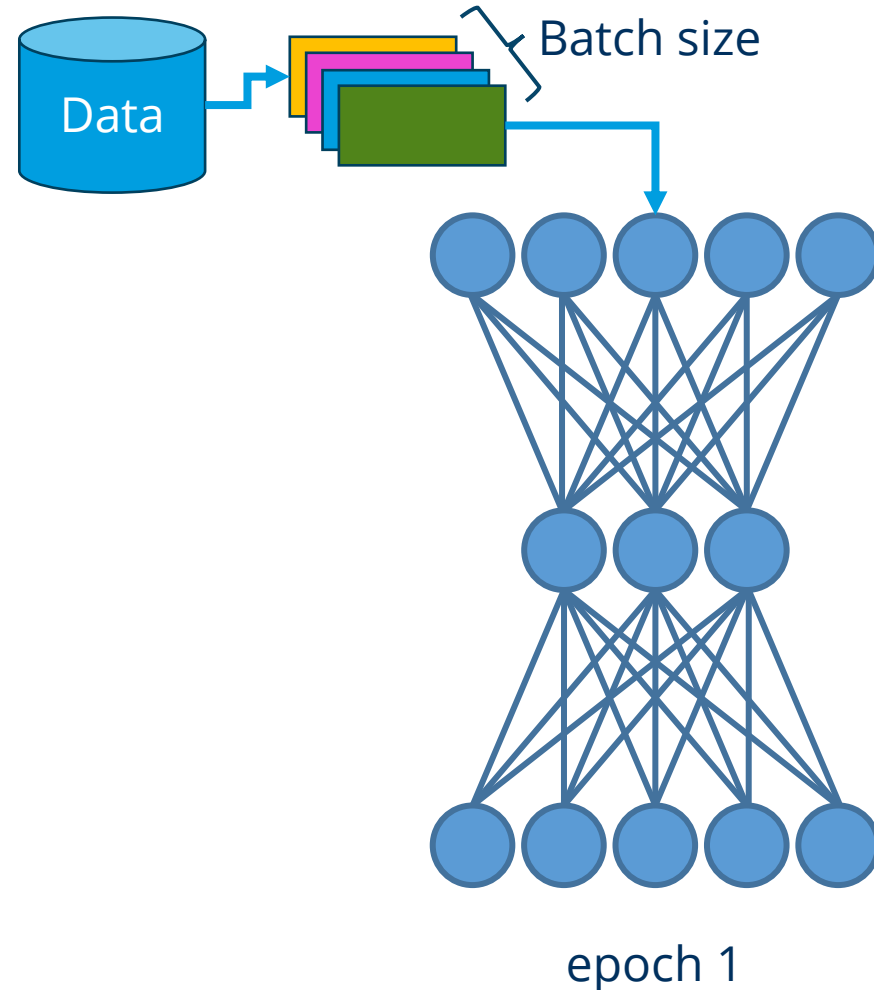
-> $10^{13}$ backprop steps required.

Solution:

- Draw n=1000 random samples from the training data to train for one epoch.

- Next epoch: different n samples.

-> $10^6$ backprop steps required.



Batch size

Data

epoch 2

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Training NNs: Drop-out

- Drop-out: deactivating individual neurons during training

- Helps with over-fitting, because the network cannot rely on individual neurons by chance being well trained, while others remain randomly initialized

- Example: drop-out-rate: 30%

epoch 1

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Training NNs: Drop-out

- Drop-out: deactivating individual neurons during training

- Helps with over-fitting, because the network cannot rely on individual neurons by chance being well trained, while others remain randomly initialized
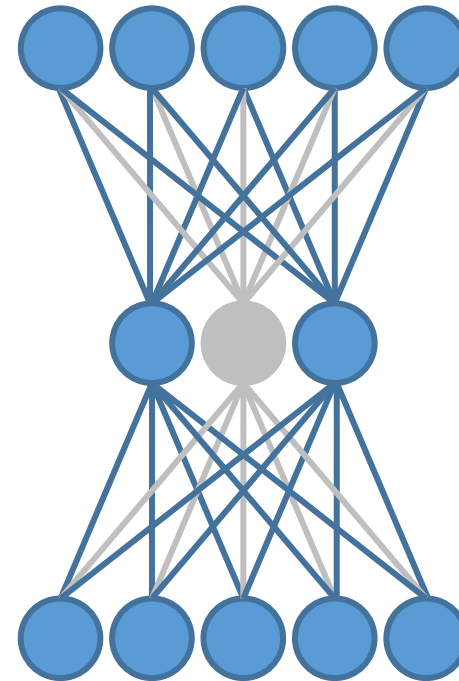
- Example: drop-out-rate: 30%



epoch 2

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# NN Architectures: Recurrent Neural Networks

Introducing some form of memory through additional connections and nodes.



Input    Output

Hidden layer with
self-feedback

Input    Output

Hidden layer with
context nodes
(Elman network, 1990)

Input    Output

Fully
Recurrent NN

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

Figure sources: Staudemeyer & Rothstein Morris
https://arxiv.org/pdf/1909.09586

# Training Recurrent Neural Networks

- Backpropagation through time

- Computationally expensive

- Unfolding through time

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

https://en.wikipedia.org/wiki/Backpropagation_through_time

# NN Architectures: Long Short-Term Memory (LSTM)

Differentiation between updating short-term memory (all the time) and updating long-term memory ([not] forgetting) thanks to separate input- and forget-gates.

# Traditional architecture: Encoder-Decoder Networks

Related: „Auto-encoder", „Variational Auto-Encoder", „U-Net"



Encoder

Decoder

Skip-connections (optional)

Input

Input: noisy image

10 µm

"Bottleneck", "Embedding", "Latent space"

Output

Input: denoised image

Images demonstrating noise2void, cropped from
https://github.com/HenriquesLab/ZeroCostDL4Mic/blob/master/Wiki_files/FigureS2.png license: MIT

# Traditional architecture: Encoder-Decoder Networks



"Embedding"

conv 3x3, ReLU
copy and crop
max pool 2x2
up-conv 2x2
conv 1x1

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

Figure source: Ronneberger et al 2015
https://arxiv.org/pdf/1505.04597

# NN Architectures: Transformers



Source: Vaswani et al (2017)
https://arxiv.org/abs/1706.03762

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Scaled dot-product attention

Attention score: How much related are two words?

Query: For which word are we calculating attention?

Key: To which word are we calculating attention

Value: Relevance of the query-key relationship



The cat is black and white.

attention score

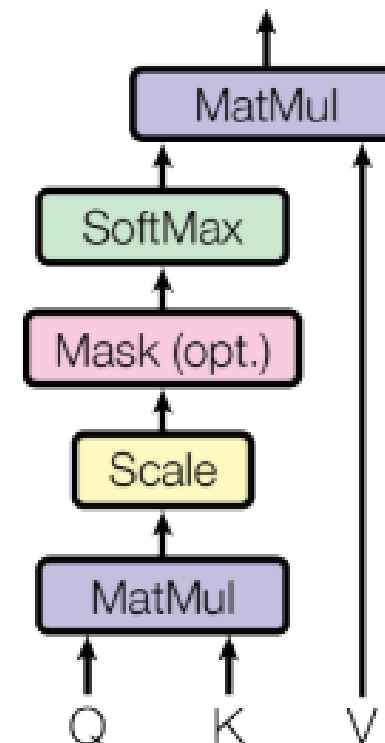Relevance value: 0.1

The cat is meowing.

attention score

Relevance value: 0.9



Scaled Dot-Product Attention

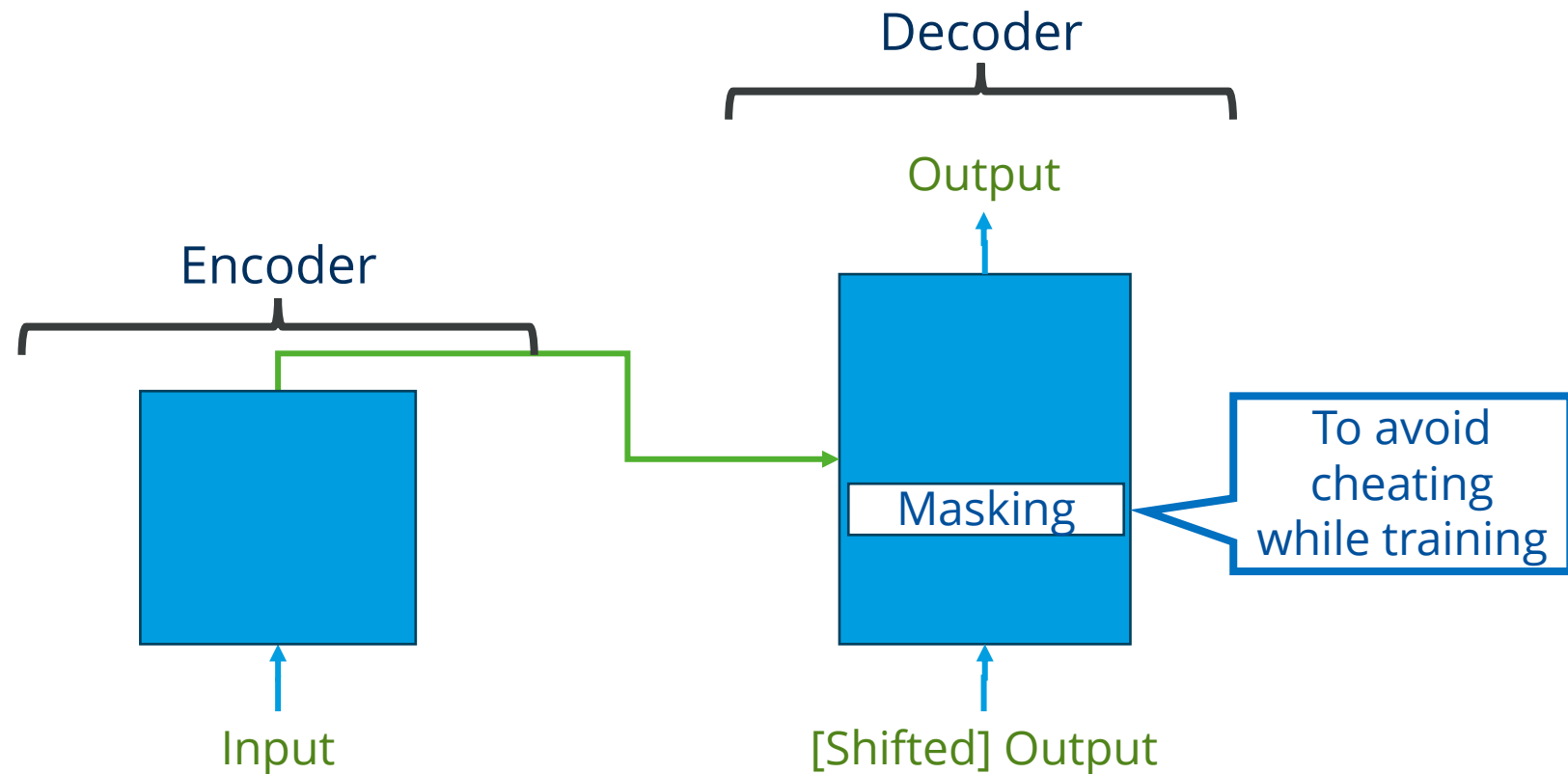Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Multi-head attention

## Multiple aspects represented by multiple attention heads

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

Slide 77

# NN Architectures: Transformers
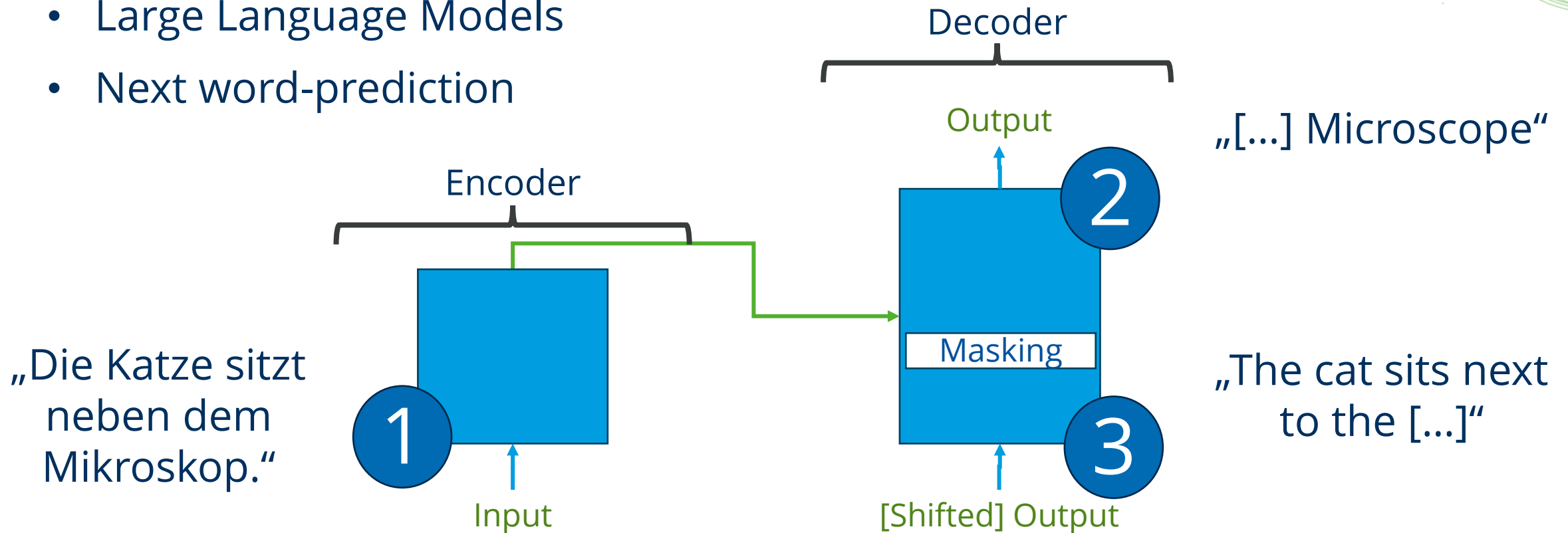
# NN Architectures: Transformers

Related terms:

- Generative Pretrained Transformer (GPT)

- Large Language Models

- Next word-prediction

Decoder

Encoder

Output

„[...] Microscope"

**2**

„Die Katze sitzt neben dem Mikroskop."

**1**

Masking

„The cat sits next to the [...]"

**3**

Input

[Shifted] Output

ScaDS.AI DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# NN Architectures: Vision Language Models

VLMs use combinations of traditional neural network architectures and transformers.

# NN Architectures: DNA Language Models

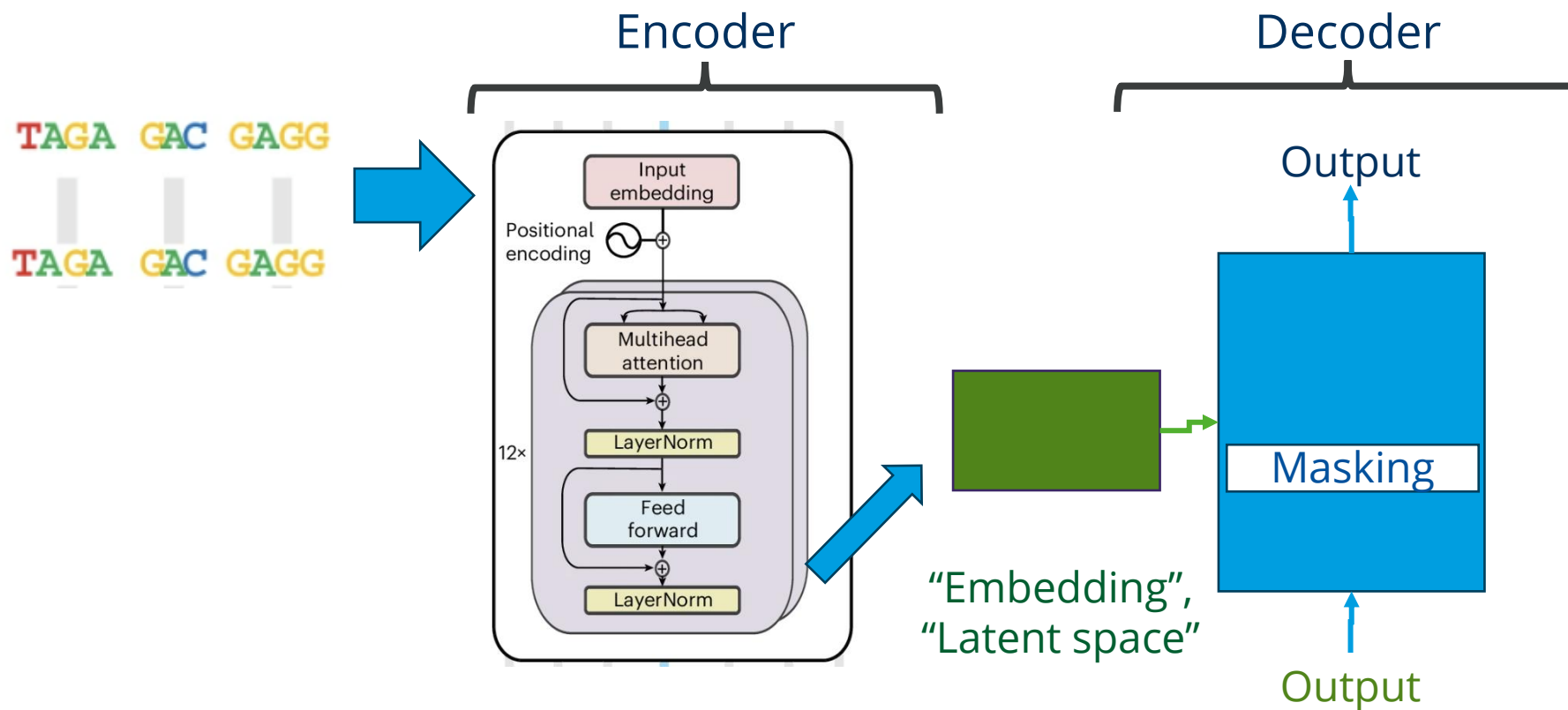DNA-LMs use a variation of the transformer architecture.



Encoder

Decoder

Input embedding

Positional encoding

Multihead attention

LayerNorm

12×

Feed forward

LayerNorm

"Embedding", "Latent space"

Output

Masking

Output

# Multi-modal Language Models

MMLMs use combinations and/or variations of traditional neural network architectures and  transformers.
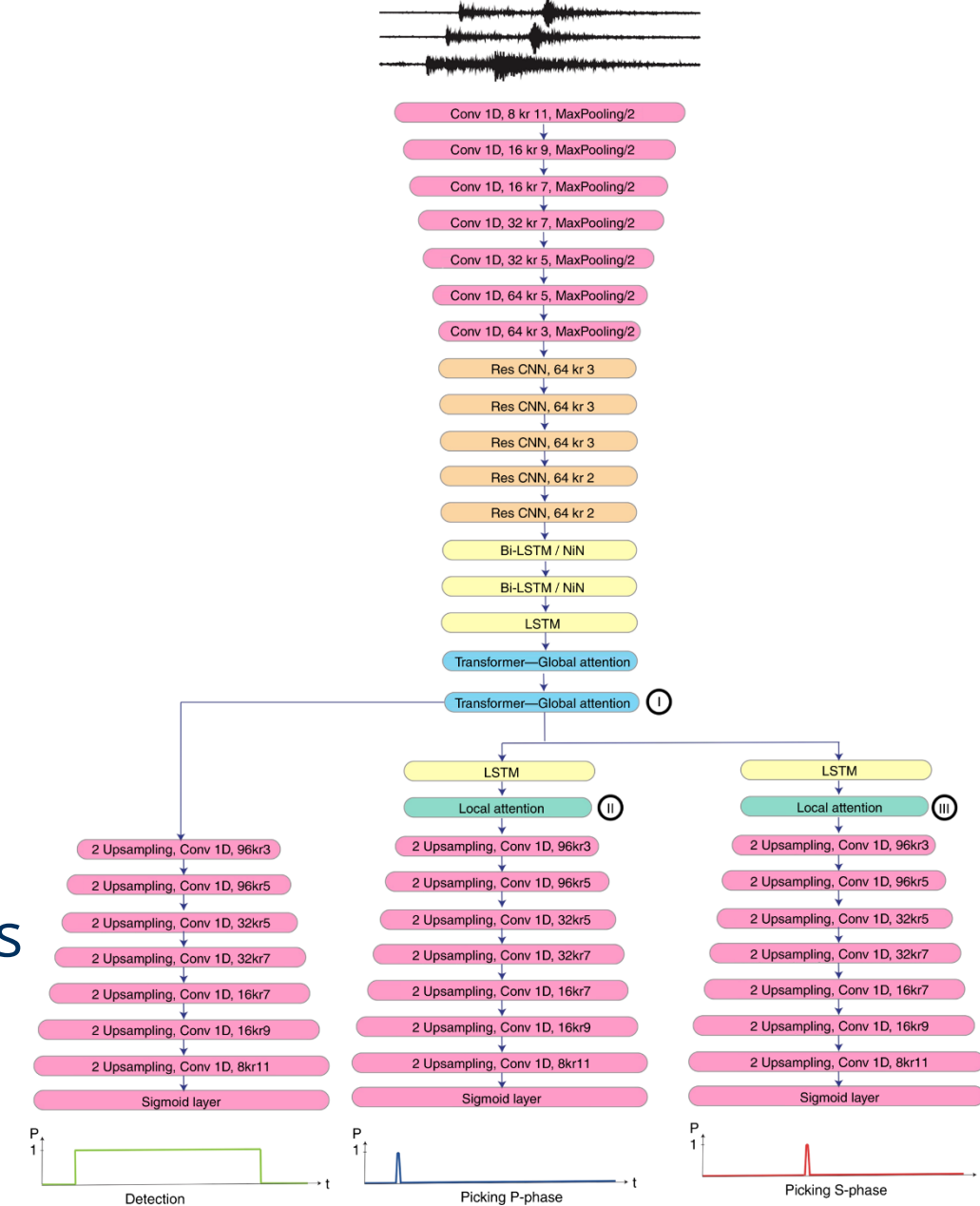
# NN Architectures

Modern NN architectures combine techniques quite freely. Example, for large earthquake detection:

- **LSTMs**

- **Transformers**

- **Convolutional**

- **Attention**

Combining architectures sometimes appears *more art than science*. Computer scientists world-wide struggle comparing different architectures.

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025

# Summary

Unsupervised ML: Explorative data science , Embeddings

Supervised ML / DL: Preduction: classification / regression , Embeddings

Explainability: SHapleys Additive exPlanations (SHAP-Analysis)

Neural networks

- Many hidden layers -> *deep* learning, Embeddings

- Training: Drop-out, batch-size, epochs

- RNNs / LSTMs -> Memory

- Transformers -> Attention, Embeddings

Good scientific practice

- Train-test-split

- Overfitting / underfitting

Robert Haase
@haesleinhuepf
AI4Seismology
May 5th 2025