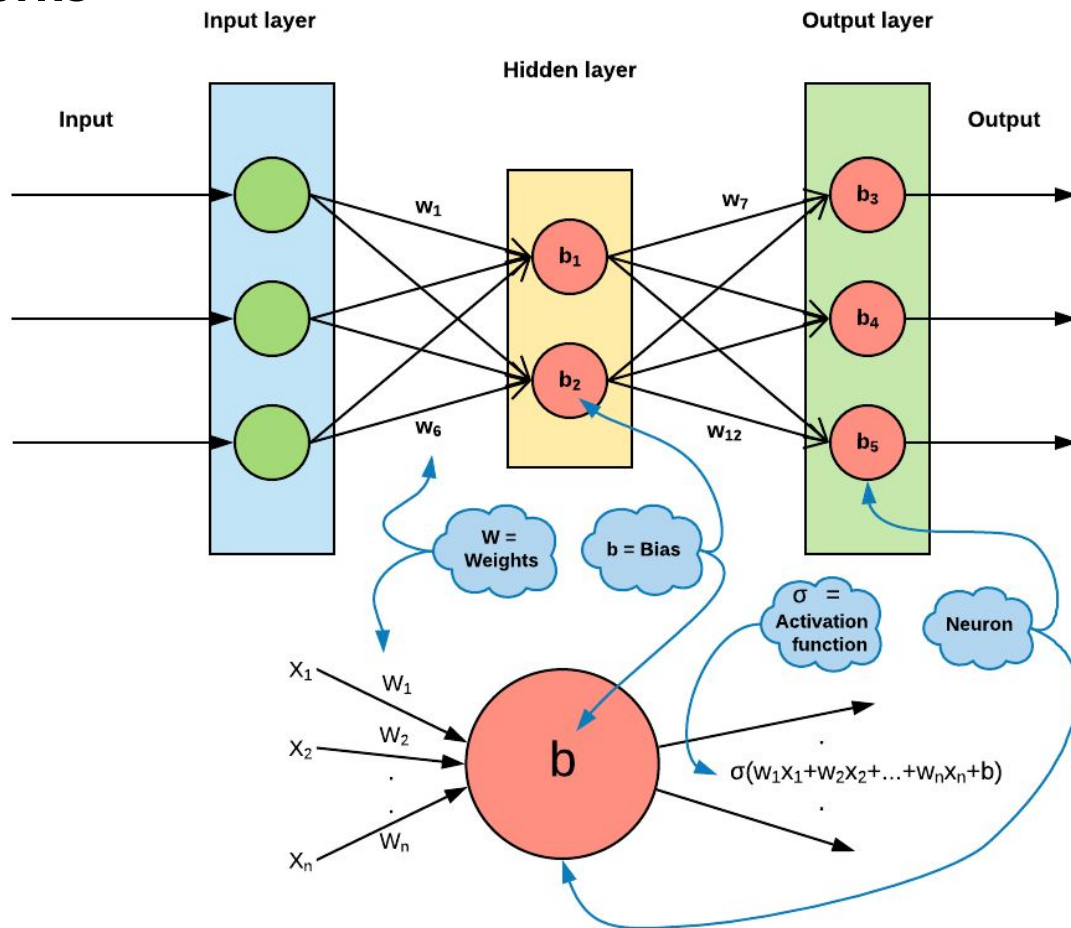# Attractors of autoencoders
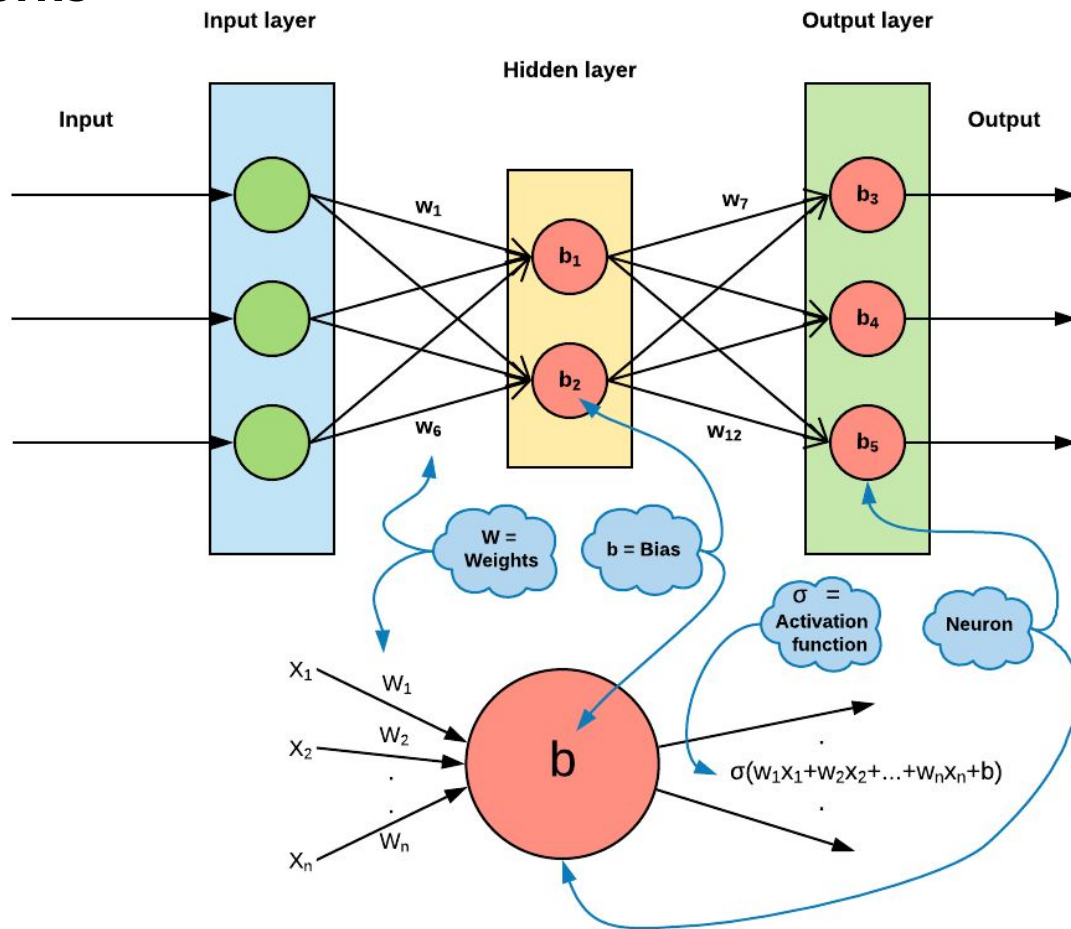### - Memorization in neural networks

Jonas Strandqvist

# Autoencoder neural networks

- Input & output layers: same number of neurons

# Autoencoder neural networks

- Input & output layers: same number of neurons

- **Depth** = number of hidden layers in neural network
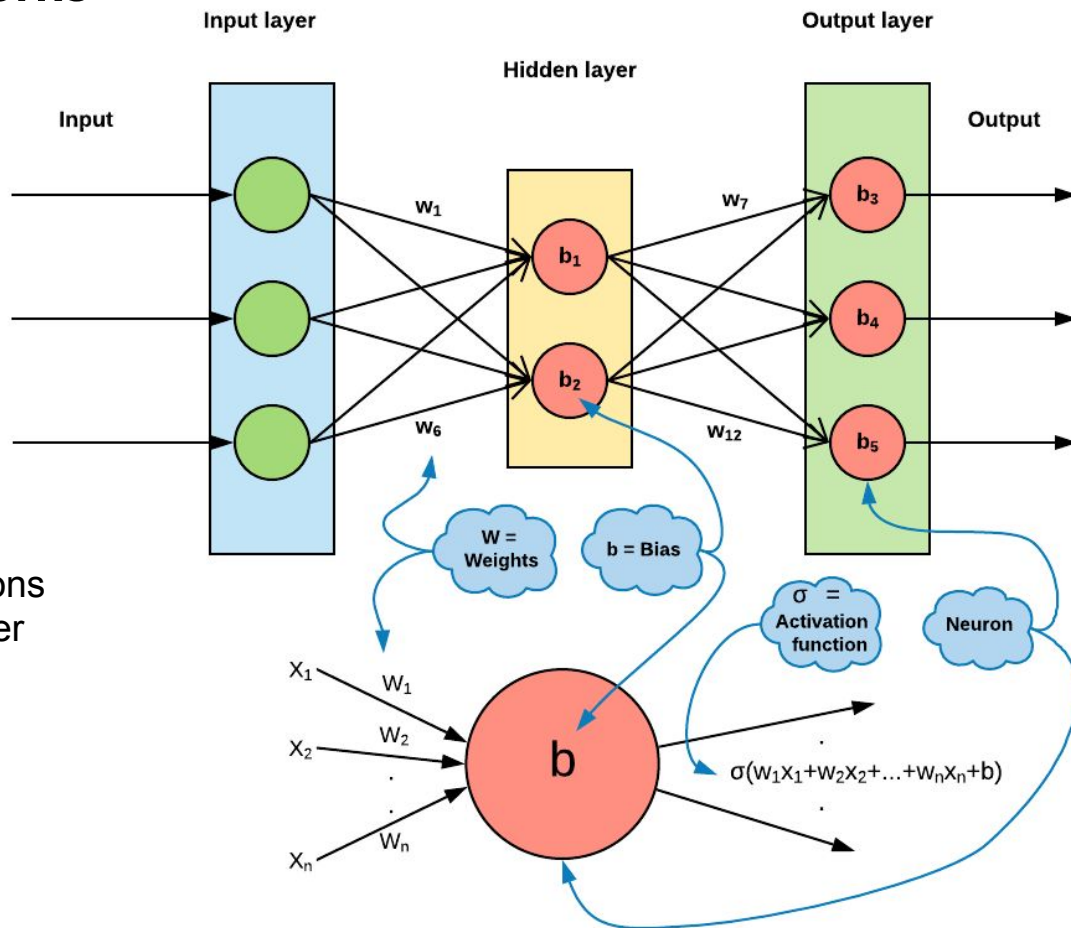  - Example: depth = 1

# Autoencoder neural networks

- Input & output layers: same number of neurons

- **Depth** = number of hidden layers in neural network
  - Example: depth = 1

- In this talk:
  - Every hidden layer has the same number of neurons
  - **Width** refers to this number
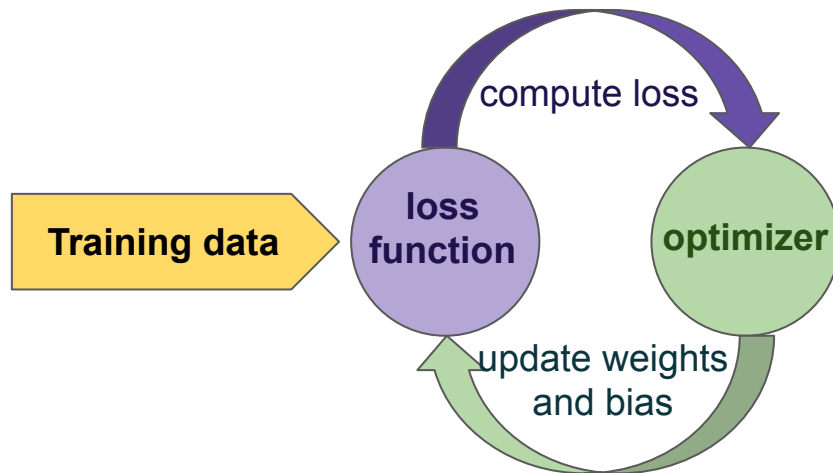  - Example: width = 2

# Machine learning with autoencoders

**Training** is the process that improves weights and biases such that our neural networks produces better results.

**Goal of an autoencoder:**



Input            Output

Encoding    Decoding

# Machine learning with autoencoders

**Training** is the process that improves weights and biases such that our neural networks produces better results.
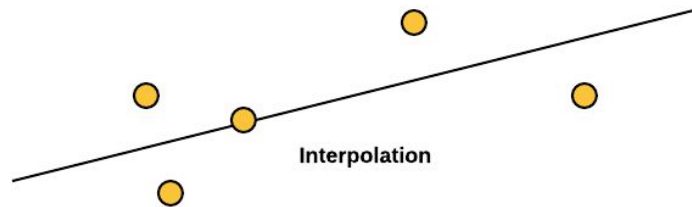
Requires three ingredients:

- a set of training data

- a loss function: measures how well the autoencoder achieves its task

- an optimizer: changes weights and biases to improve loss

**Goal of an autoencoder:**



Input → Encoding → Decoding → Output



Training data → loss function → compute loss → optimizer → update weights and bias →

# Memorization in autoencoders

Few parameters lead to interpolation.

Interpolation

Too many parameters allow the network to learn the dataset.
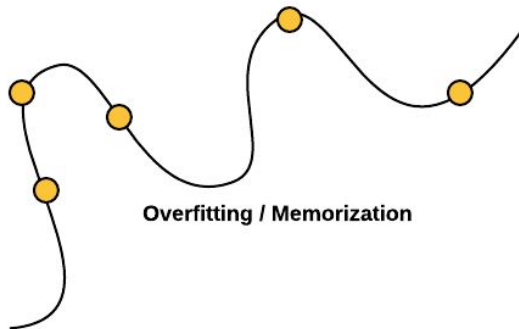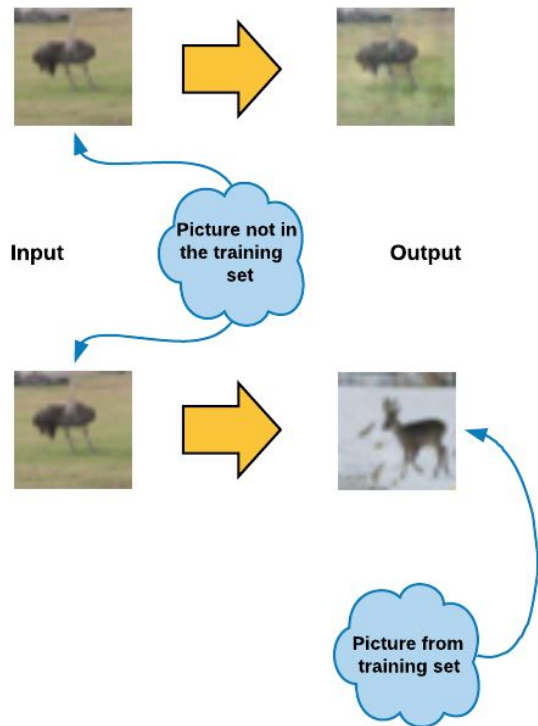
Overfitting / Memorization
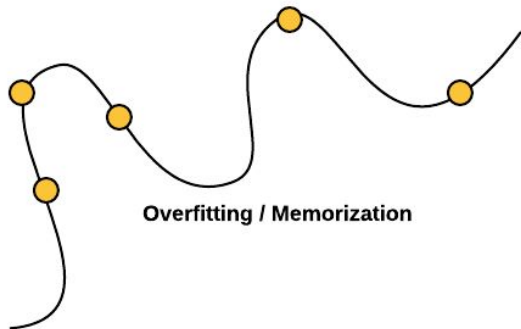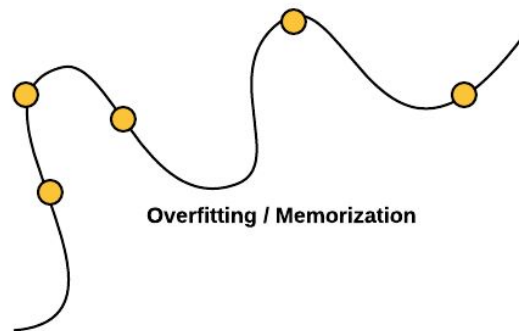
# Memorization in autoencoders

Few parameters lead to interpolation.

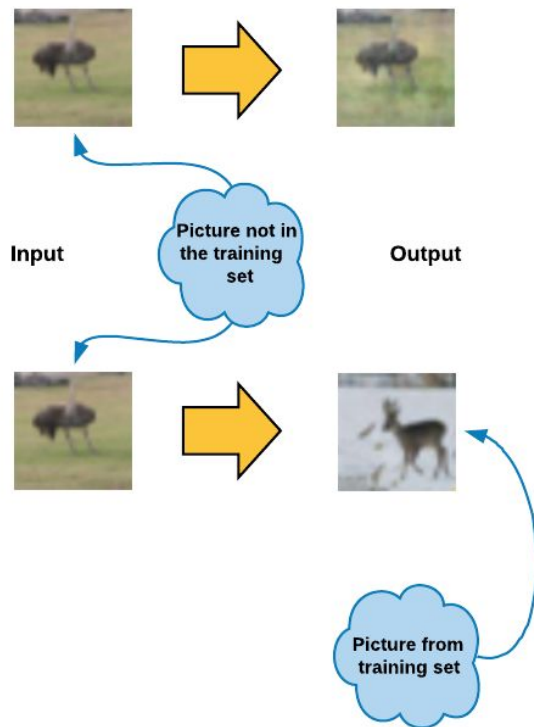Too many parameters allow the network to learn the dataset.



Interpolation



Overfitting / Memorization



Input    Picture not in the training set    Output



Picture from training set

# Memorization in autoencoders

Few parameters lead to interpolation.



Interpolation

Too many parameters allow the network to learn the dataset.

Overfitting / Memorization



Input

Picture not in the training set

Output

Picture from training set

It is still an open question how neural networks memorize data.
For this, articles [1] and [2] suggest to study **attractors** of autoencoders.

[1] "Memorization in overparameterized autoencoders" - A. Radhakrishnan, K.D. Yang, M. Belkin and C. Uhler
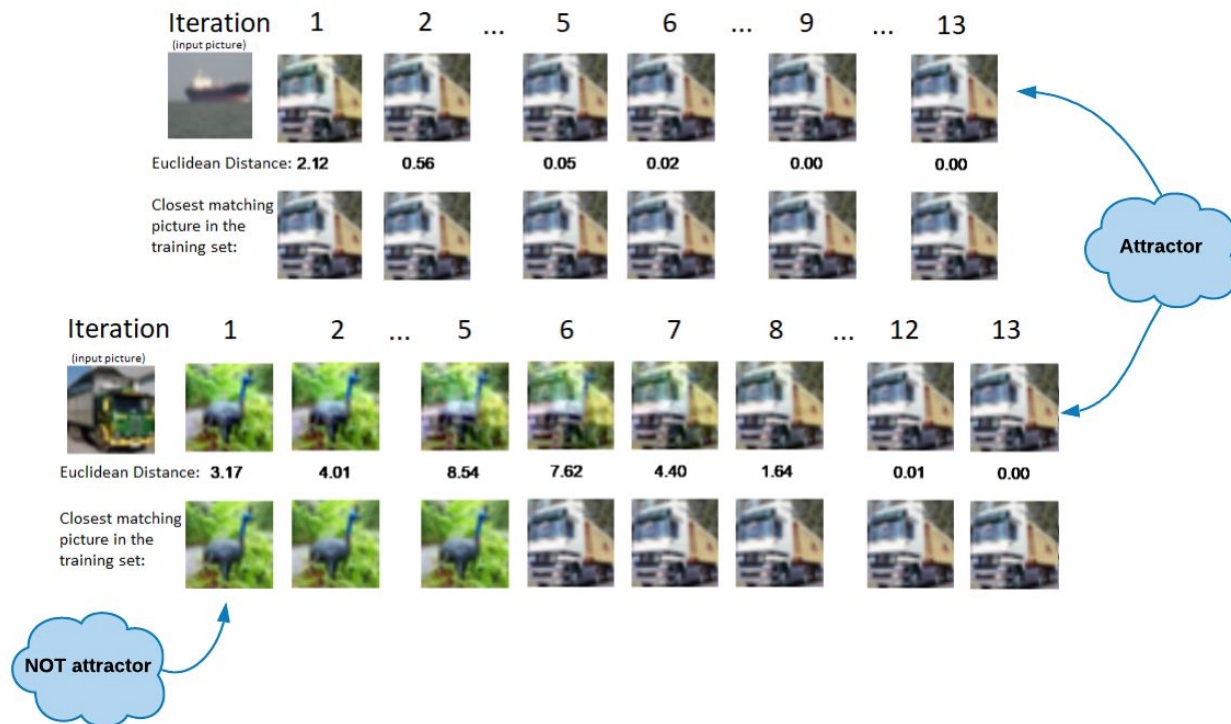[2] "Overparameterized Neural Networks Can Implement Associative Memory" - A. Radhakrishnan, M. Belkin, C. Uhler

# Attractors of autoencoders

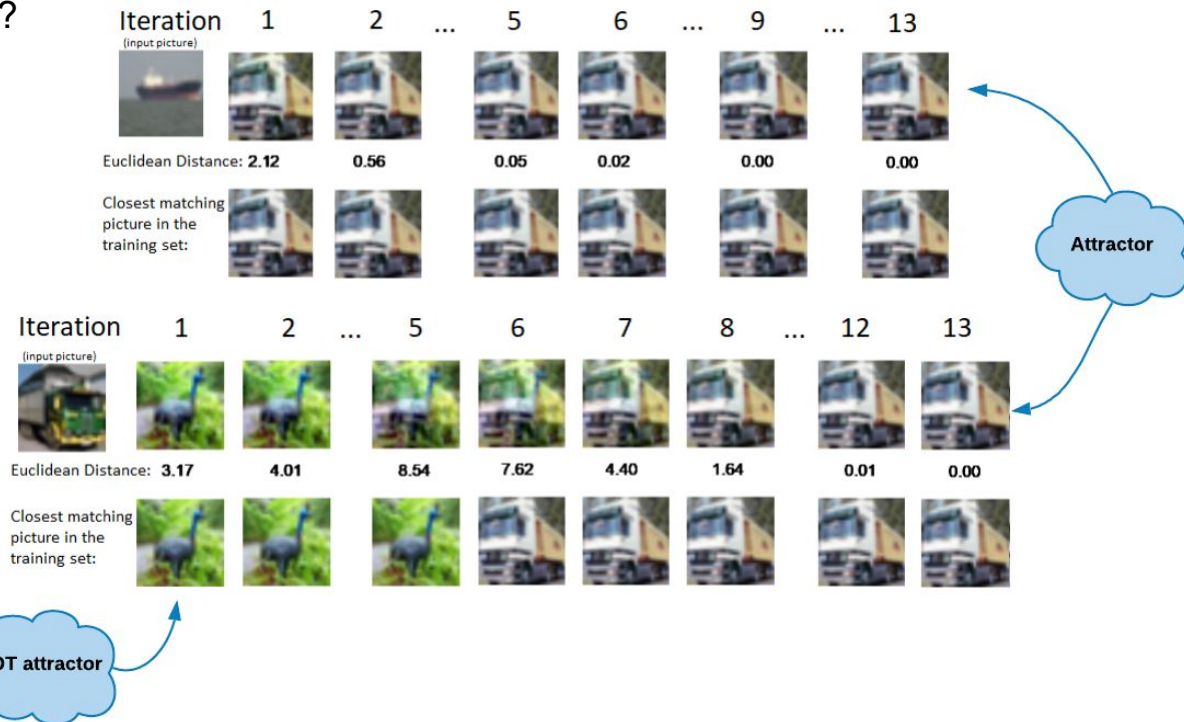# Attractors of autoencoders

# Attractors of autoencoders

# Attractors of autoencoders

When is a training image 🟨 an attractor?

**Two** things can go wrong:

- 🟨 becomes another image when iterating autoencoder

  → we say that 🟨 is not an **iterative fixed point**

# Attractors of autoencoders

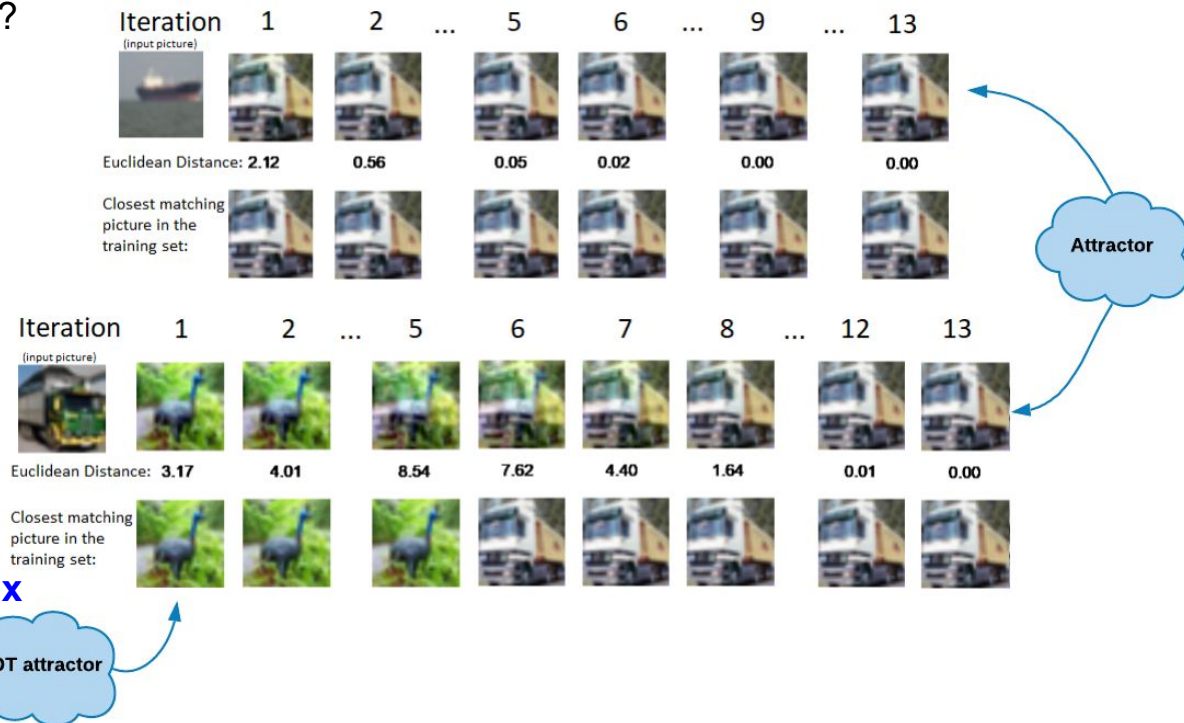When is a training image 🟨 an attractor?

**Two** things can go wrong:

- 🟨 becomes another image when iterating autoencoder

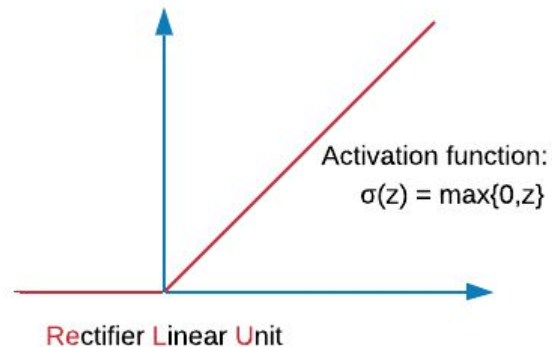  → we say that 🟨 is not an **iterative fixed point**

- **other** images never become 🟨

  (can be checked with the **eigenvalues of the Jacobian matrix** of the autoencoder at 🟨)

# Problem formulation

We want to extend the experiments in [2] to **ReLU** autoencoders.

Activation function:
$\sigma(z) = \max\{0, z\}$

Rectifier Linear Unit

[2] "Overparameterized Neural Networks Can Implement Associative Memory" - A. Radhakrishnan, M. Belkin, C. Uhler
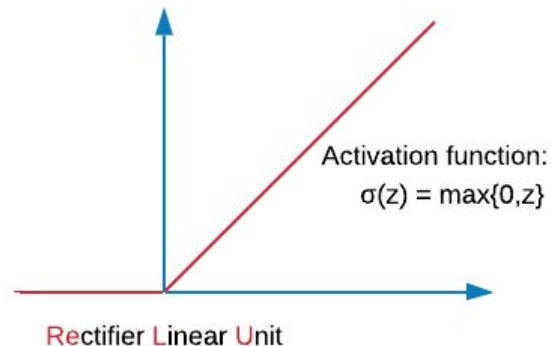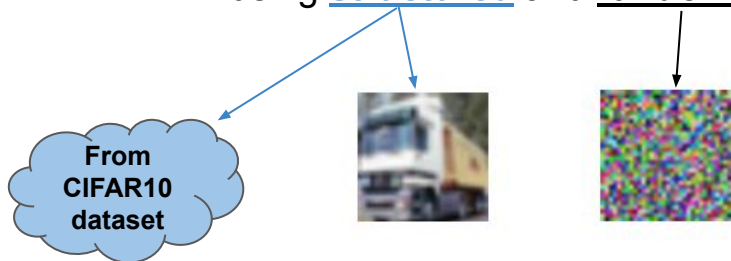
# Problem formulation

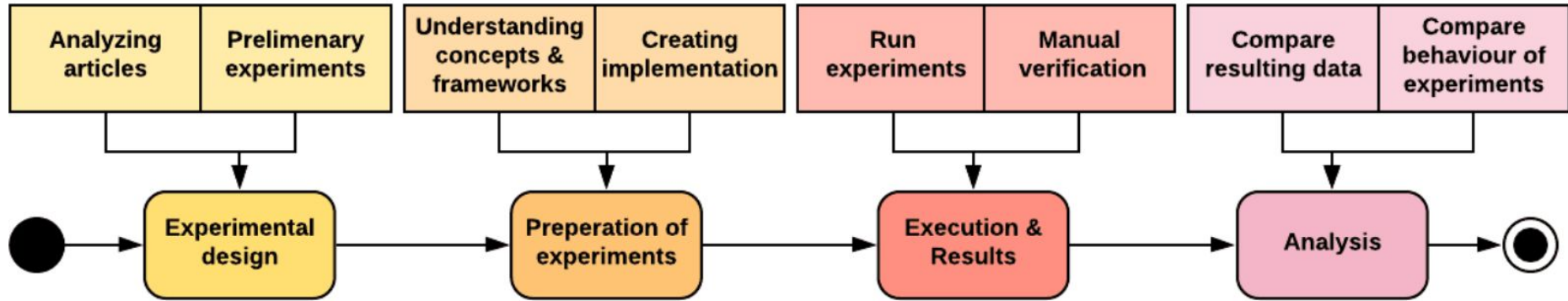We want to extend the experiments in [2] to **ReLU** autoencoders.

We investigate the impact on the amount of attractors:

- changing **depth** and **width**

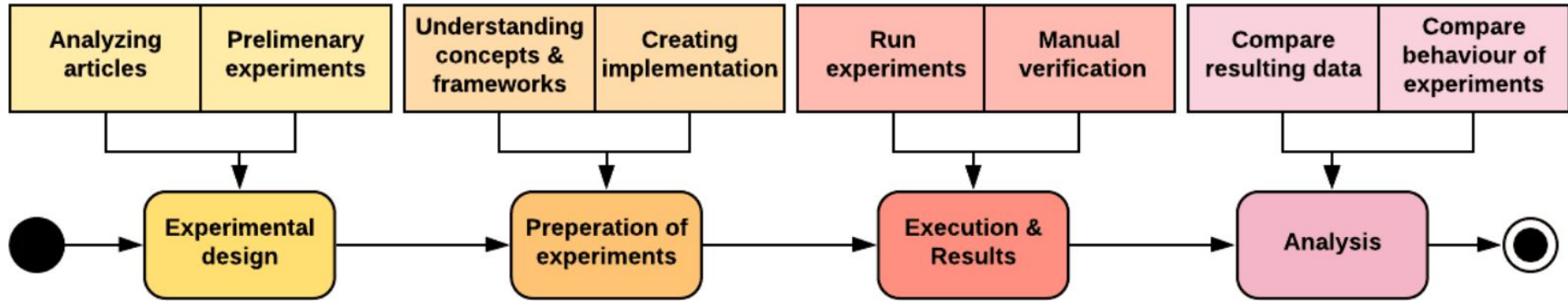- with and without **bias**

- using **structured** and **random pictures**

From CIFAR10 dataset

Activation function:
$\sigma(z) = \max\{0, z\}$

Rectifier Linear Unit

[2] "Overparameterized Neural Networks Can Implement Associative Memory" - A. Radhakrishnan, M. Belkin, C. Uhler
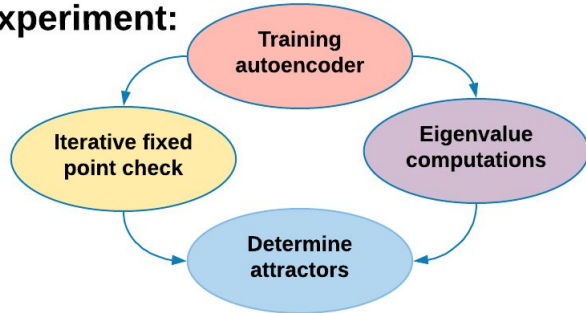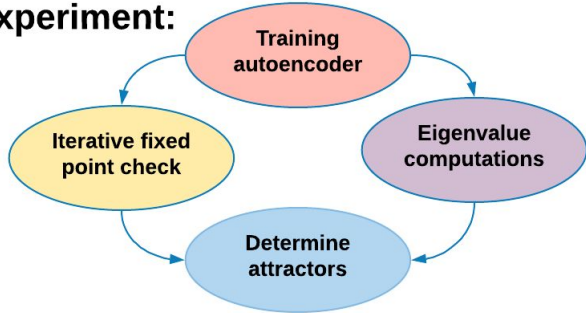
# Method : Controlled experiment(s)

# Method : Controlled experiment(s)



Single experiment:

# Method : Controlled experiment(s)



**Single experiment:**

**Each experiment repeats 4x, totalling 240 experiments**

**Scope:**

| width depth | 128 | 64 | 32 |
|---|---|---|---|
| 11 | | | |
| 6 | | | |
| 3 | | | |
| 2 | | | |
| 1 | | | |

with and without bias

training pictures: random and CIFAR10

# Implementation
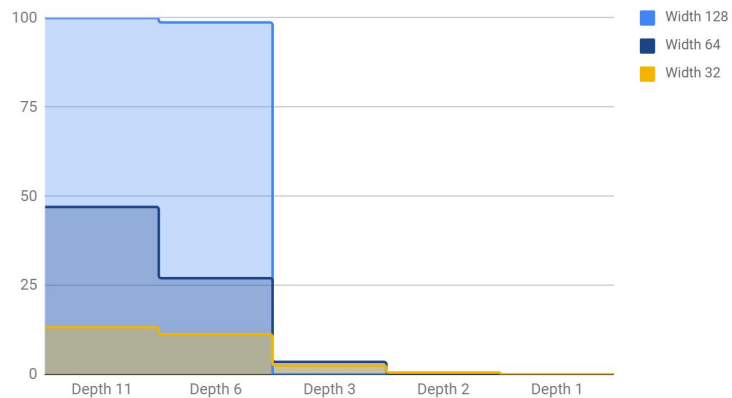


**Experiment environment**

**Manual verification**

# Results & Analysis
# (trained models with bias)



Attractors with CIFAR10 pictures

Attractors with Random pictures

# Results & Analysis
# (trained models with bias)



Attractors with CIFAR10 pictures

Width 128
Width 64
Width 32

Attractors with Random pictures

Width 128
Width 64
Width 32

# Results & Analysis
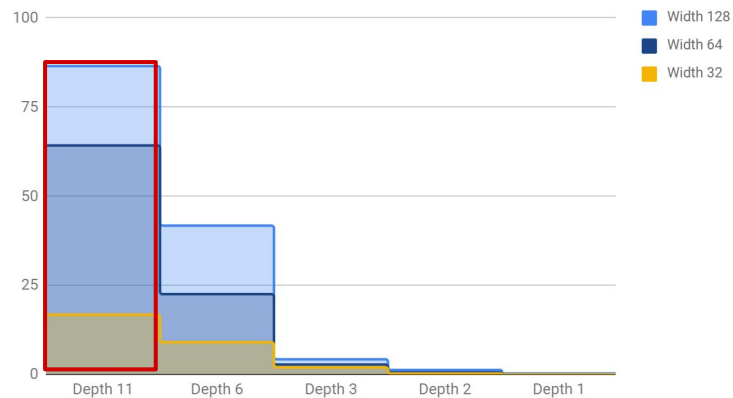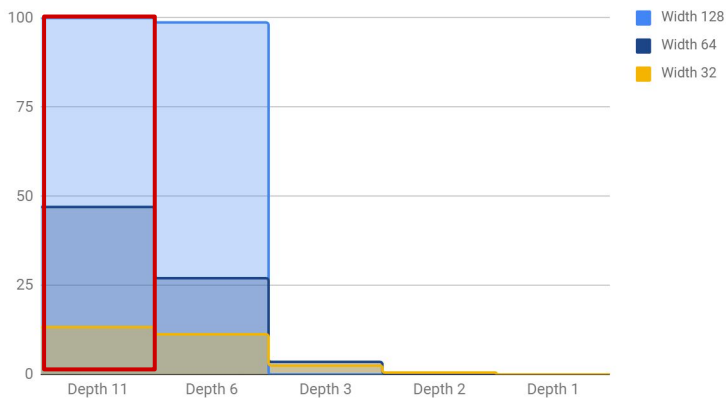# (trained models with bias)



Attractors with CIFAR10 pictures

Attractors with Random pictures

# Results & Analysis
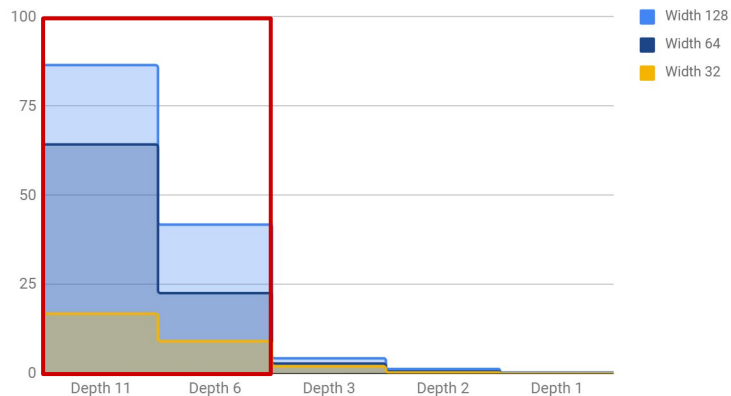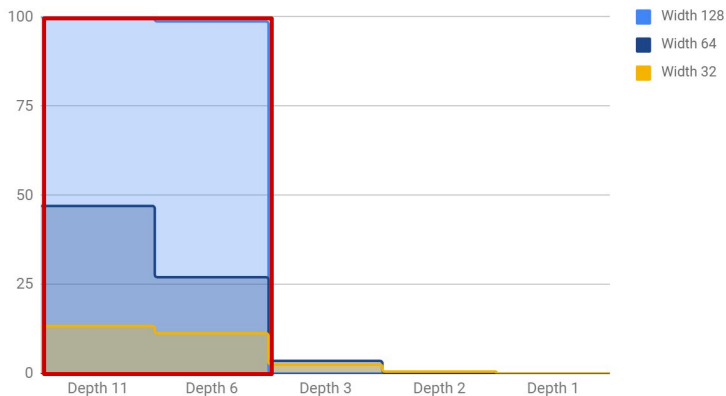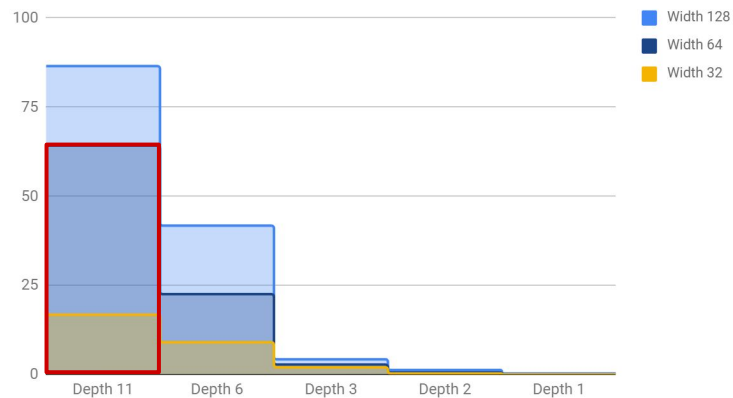# (trained models with bias)



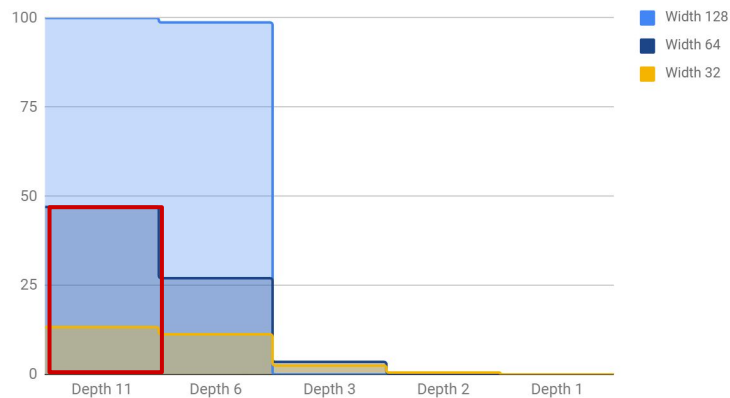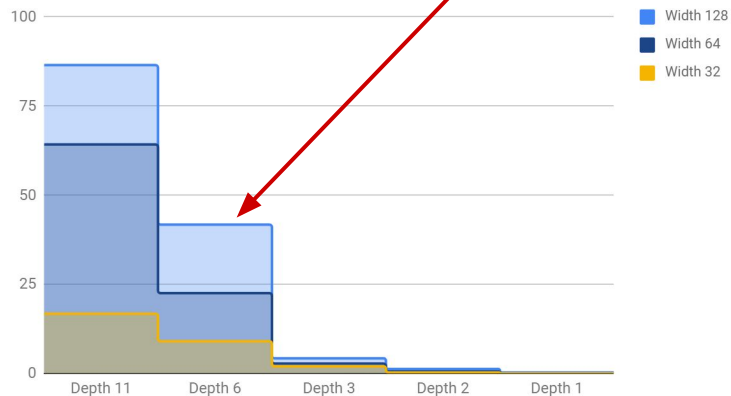Attractors with CIFAR10 pictures

Attractors with Random pictures

# Results & Analysis
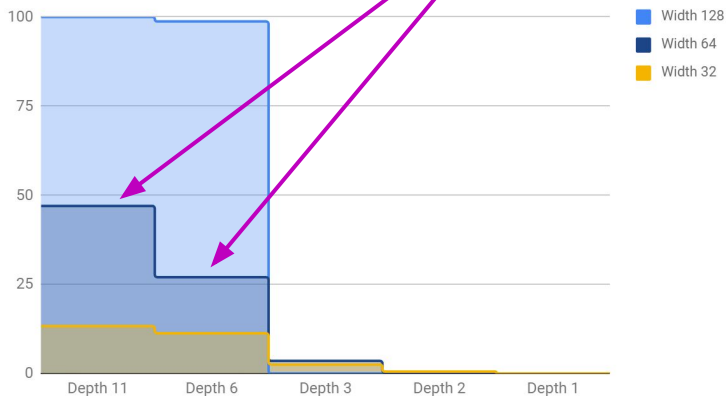# (trained models with bias)

- CIFAR10 pictures:
  greater impact when changing depth

- Random pictures:
  greater impact when changing width



Attractors with CIFAR10 pictures



Attractors with Random pictures
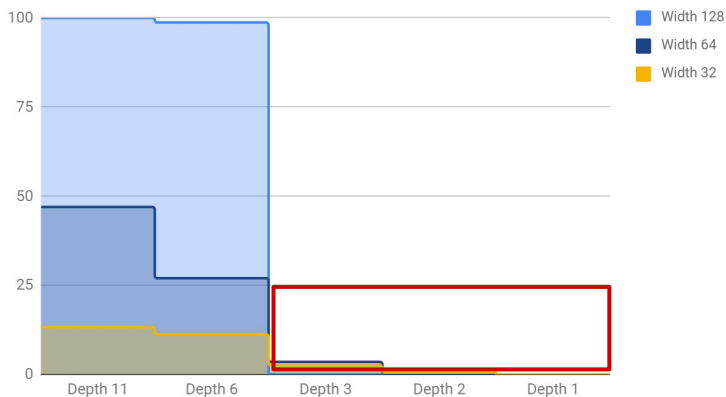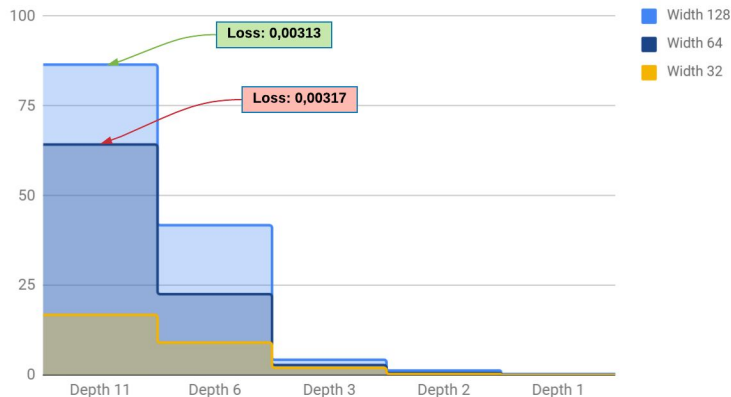
# Results & Analysis
# (trained models with bias)

- CIFAR10 pictures:
  greater impact when changing depth

- Random pictures:
  greater impact when changing width

- Sufficient depth is required for creating attractors

Attractors with CIFAR10 pictures

Width 128
Width 64
Width 32

Attractors with Random pictures

Width 128
Width 64
Width 32

# Results & Analysis
# (trained models with bias)

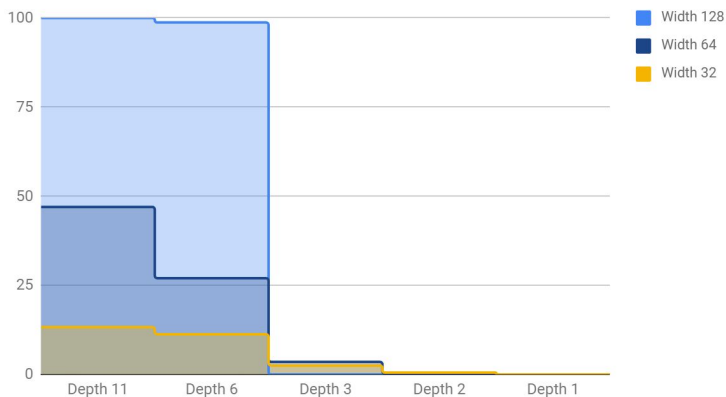- CIFAR10 pictures:
  greater impact when changing depth

- <span style="color:red">Good loss and number of attractors
  are not necessarily related</span>

- Random pictures:
  greater impact when changing width
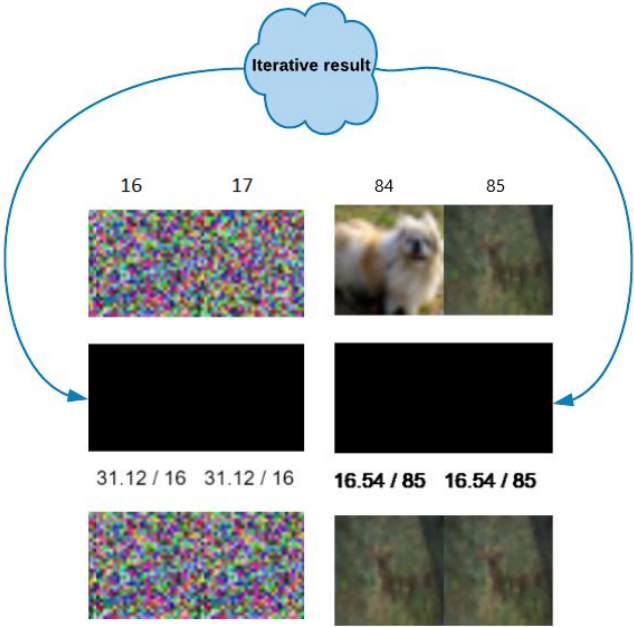
- Sufficient depth is required for
  creating attractors

Attractors with CIFAR10 pictures
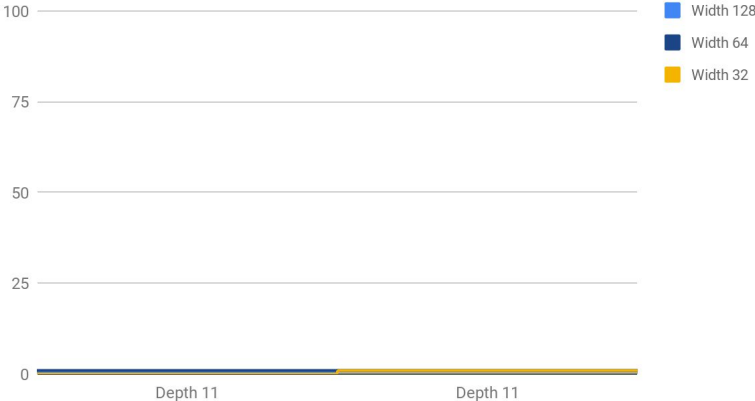
Loss: 0,00313

Loss: 0,00317

- Width 128
- Width 64
- Width 32

Depth 11   Depth 6   Depth 3   Depth 2   Depth 1

Attractors with Random pictures

- Width 128
- Width 64
- Width 32

Depth 11   Depth 6   Depth 3   Depth 2   Depth 1

# Results & Analysis
# (trained models without bias)

Not using bias, CIFAR10 & Random pictures





Iterative result

16      17          84      85

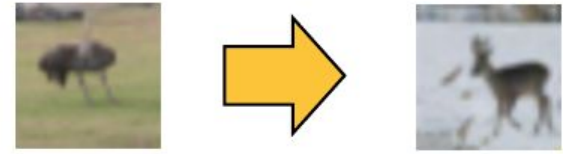31.12 / 16    31.12 / 16        16.54 / 85    16.54 / 85

**Not using bias, produces no attractors**

# Future work

- increase the scope (width and depth) of autoencoders

- the effect of epochs on trained autoencoders was not covered in the thesis (only trained 50k).

- training pictures was always 100, the impact of changing this amount would be interesting.

- training with pictures with a more specific structure could give interesting results.

- optimizing the thesis' notion of iterative fixed points, excluding the "false positives".

- a mathematical proof showing that every data point becomes an attractor for a sufficiently generic training set for a sufficiently large ReLU autoencoder

An autoencoder is a mathematical function; in this thesis:



$$\alpha : \mathbb{R}^{3072} \to \mathbb{R}^{3072}$$

At every picture, we can compute the Jacobian matrix $J_\alpha(\text{🦌})$ → this is a 3072 x 3072 matrix.

[1] and [2] show for a **perfectly trained autoencoder with loss 0**:

- 🦌 is an attractor if the highest absolute value of the eigenvalues of $J_\alpha(\text{🦌})$ is smaller than 1.
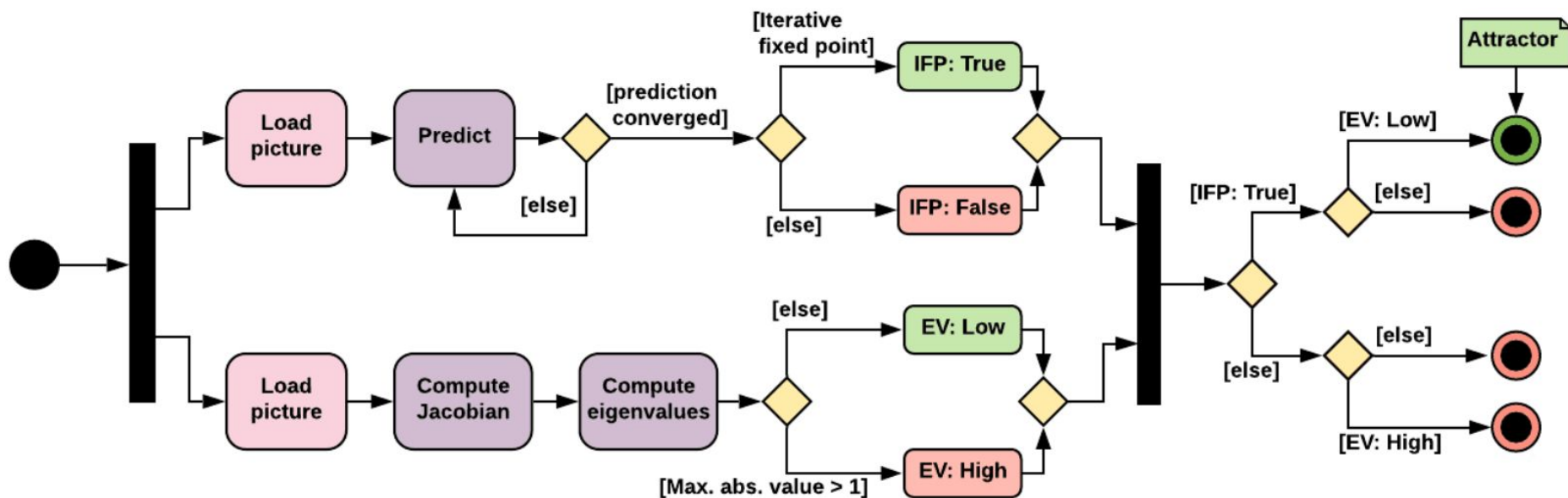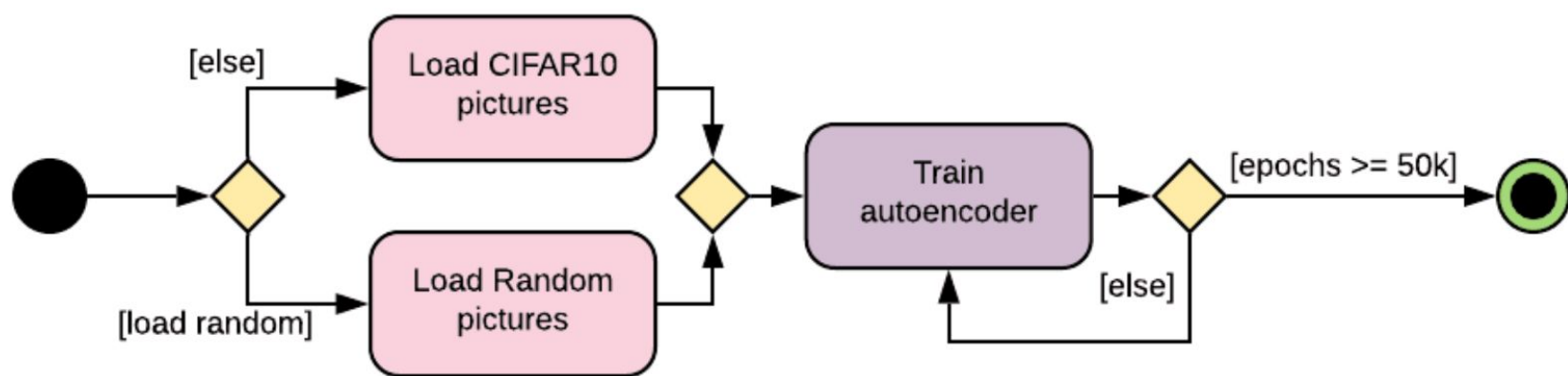
- 🦌 is not an attractor if this highest absolute value is larger than 1.
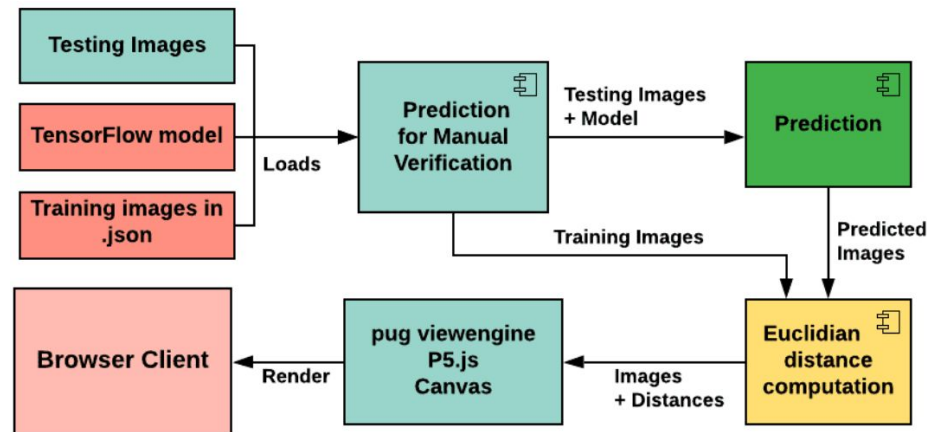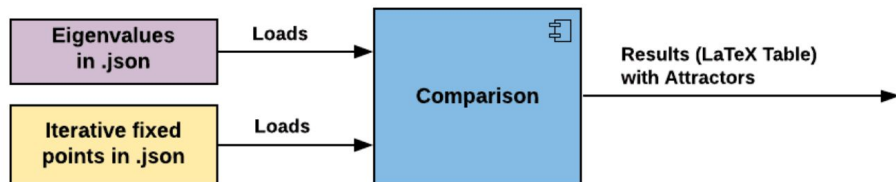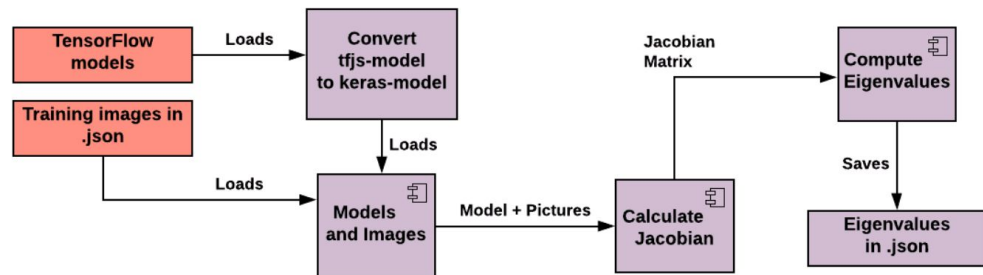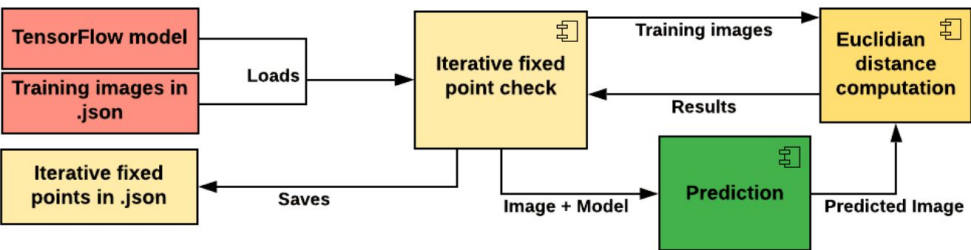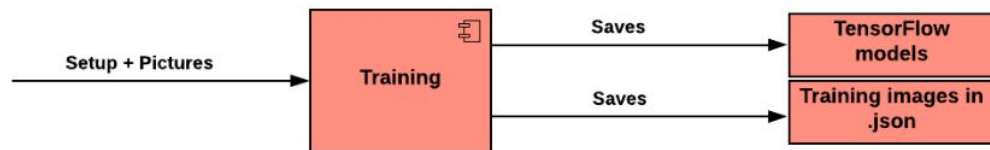
For an **autoencoder with loss > 0**, 🦌 is an attractor if this highest absolute value is smaller than 1 and 🦌 is an iterative fixed point.

| | | largest absolute value of all eigenvalues: | |
| --- | --- | --- | --- |
| | | < 1 | > 1 |
| iterative fixed point? | yes | attractor | not attractor |
| | no | not attractor | not attractor |

[1] "Memorization in overparameterized autoencoders" - A. Radhakrishnan, K.D. Yang, M. Belkin and C. Uhler
[2] "Overparameterized Neural Networks Can Implement Associative Memory" - A. Radhakrishnan, M. Belkin, C. Uhler
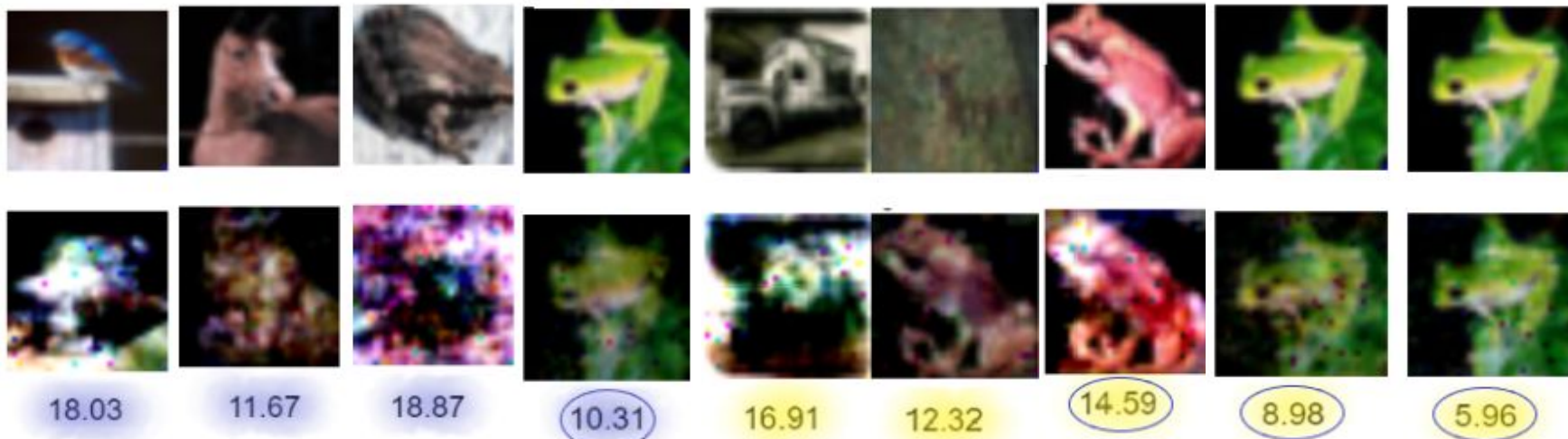
# Results & Analysis
## (iterative fixed points at lower depths
## & manual verification)

'false predictions':
- 82% at depth 1
- 36% at depth 2
- 19% at depth 3



18.03    11.67    18.87    10.31    16.91    12.32    14.59    8.98    5.96

Examples depth 2

Examples depth 1