

# Intelligente Informationssysteme

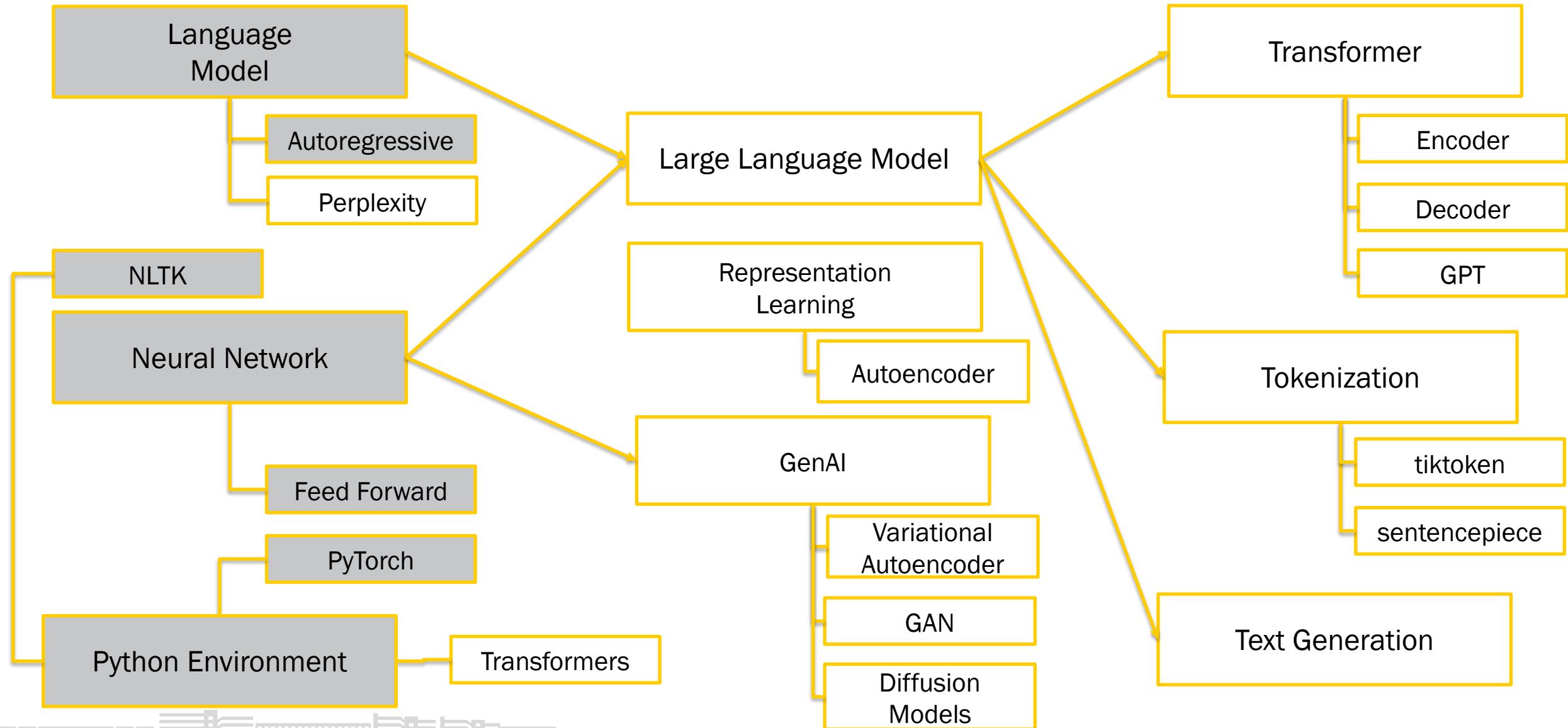
## Block 1 – Language Models

Dominik Neumann

## Language Model

1. Why is **language** so important in Artificial Intelligence?
2. What is a model?
3. What is a Language Model?
4. How to build a Language Model?

# Overview



# 01

## Language is the jewel in the crown of cognition

Steven Pinker



1. Language is a system of communication that uses symbols in a regular way to create meaning.
2. Language gives us the ability to communicate to others by talking, reading, and writing.
3. Human language is powerful and has been transformative to our species because it gives groups of people a way to network human brains together.



## 02

# Concept of a Model in Machine Learning



What does it mean to “model something”?

Imagine that we have a model of a physical world.

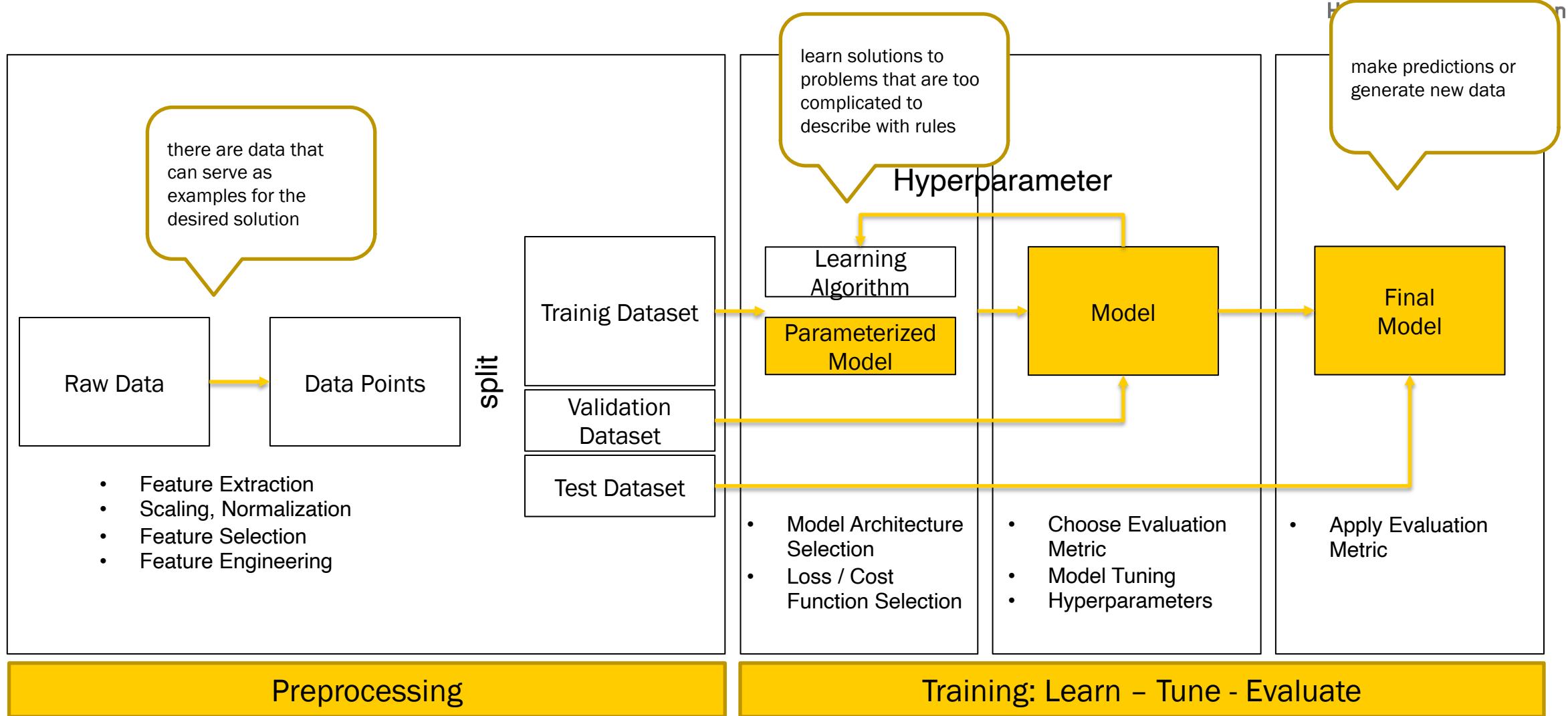
What do you expect it to be able to do?

- A good model would predict what happens next, given some description of “context”.
- A good model would simulate the behavior of the real world.
- It would “understand” which events are in better agreement with the world. Which of them are more likely.



# Concept of a Model in Machine Learning

- learned in some machine learning process



# Concept of a Model in Machine Learning

As a mathematical parameterized model that explain some data

Machine learning (ML) develops, researches and uses statistical algorithms, also known as learning algorithms.

Learning algorithms can learn solutions to problems that are too complicated to describe with rules, but there are data that can serve as examples for the desired solution.

A learning algorithm maps given example data onto a mathematical parameterized model. The learning algorithm adapts the parameters of the model so that it can generalize from the example data to new cases. This process is called training. After training, the solution found is stored in the model. It is not explicitly programmed.

The trained model can make predictions or generate new data.

There are always three ingredients for Machine Learning:

1. some data
2. a parameterized model that can explain the data, given an appropriate choice of parameters.
3. at least an evaluation procedure, the loss function, that allows us to adjust the parameters of our model.



# Concept of a Model in Machine Learning

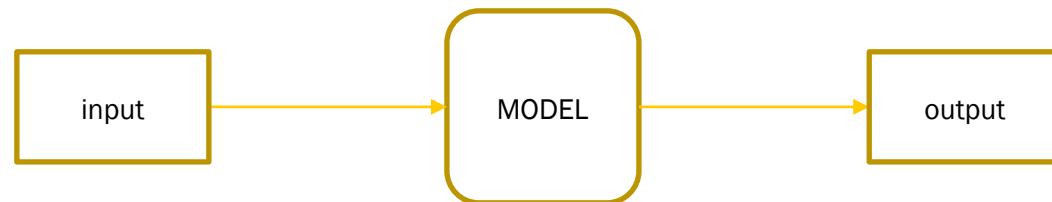
- as a „black box“ function with input and output performing some task

The universal Function Approximation Theorem tells us, that every function can be approximated by a deep neural network:

A ML **task** is a type of prediction or inference being made by the model. It is based on the problem or question that is being asked, and the available data.

For example, the classification task assigns data to categories, and the clustering task groups data according to similarity.

A Model is like a function with input and output performing some task.

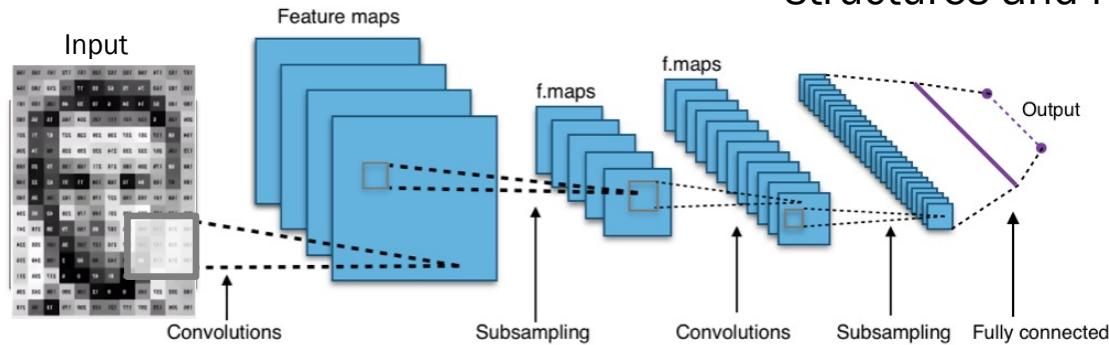
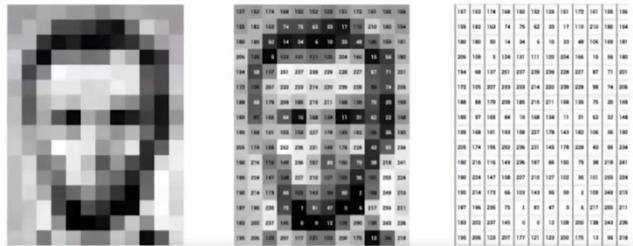


# Concept of a Model in Machine Learning

Models and their Deep Learning Network Architectures

While designing a model, we make assumptions about our data depending on the geometry of the data: called **Inductive Bias**

## 1. Images are represented as Matrix or Tensor



### Inductive Bias

- **Nearby pixels have similar colors**
- Use convolutional filters to extract structures and hidden features

# Concept of a Model in Machine Learning

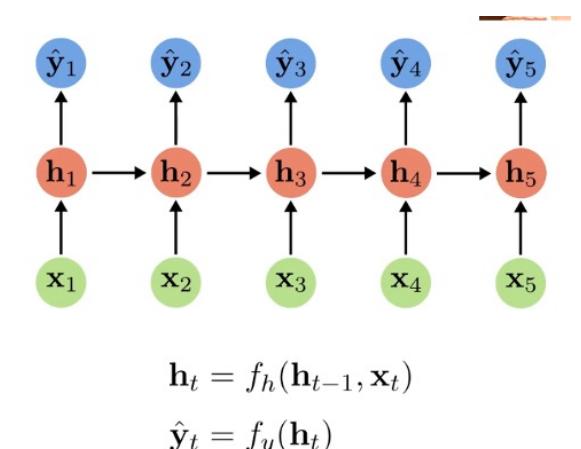
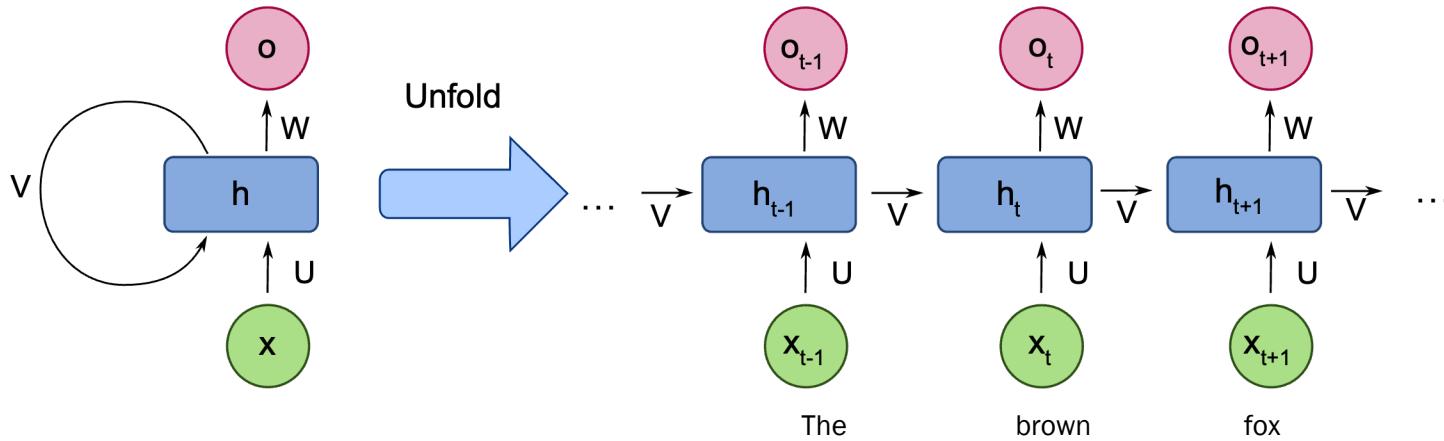
Models and their Deep Learning Network Architectures

While designing a model, we make assumptions about our data depending on the geometry of the data: called Inductive Bias

2. Language is represented as a Sequence: Inductive Bias

[“The”, “brown”, “fox”, ....]

- Nearby words belong to the same semantic meaning
- sequence-2-sequence models like recurrent neural networks could be used to model context with hidden states



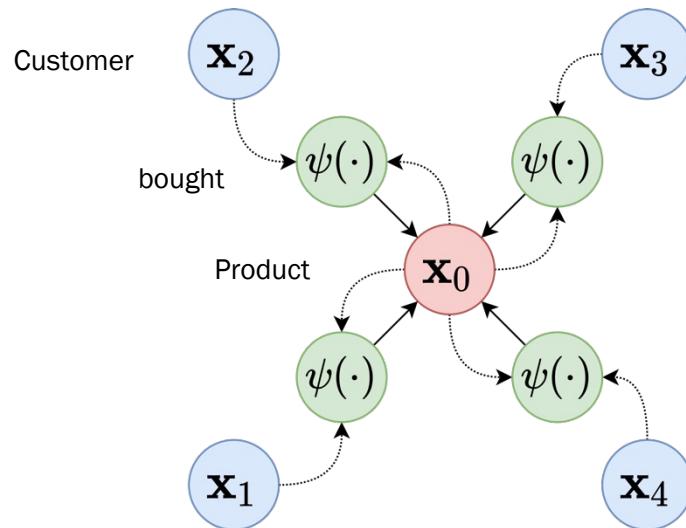
Source: Wikipedia - <https://commons.wikimedia.org/w/index.php?curid=60109157>

# Concept of a Model in Machine Learning

Models and their Deep Learning Network Architectures

While designing a model, we make assumptions about our data depending on the **geometry of the data**: called **Inductive Bias**

3. Data like objects with relationships to each other can be represented as a graph



## Inductive Bias

- Use the graph structure of the data to build a graph neural network.
- In a Message Passing GNN the network nodes update their representations by aggregating the messages received from their neighbors.
- Node **representation** could be used to downstream task, like classification, clustering or link prediction.

# 03

# Language Models



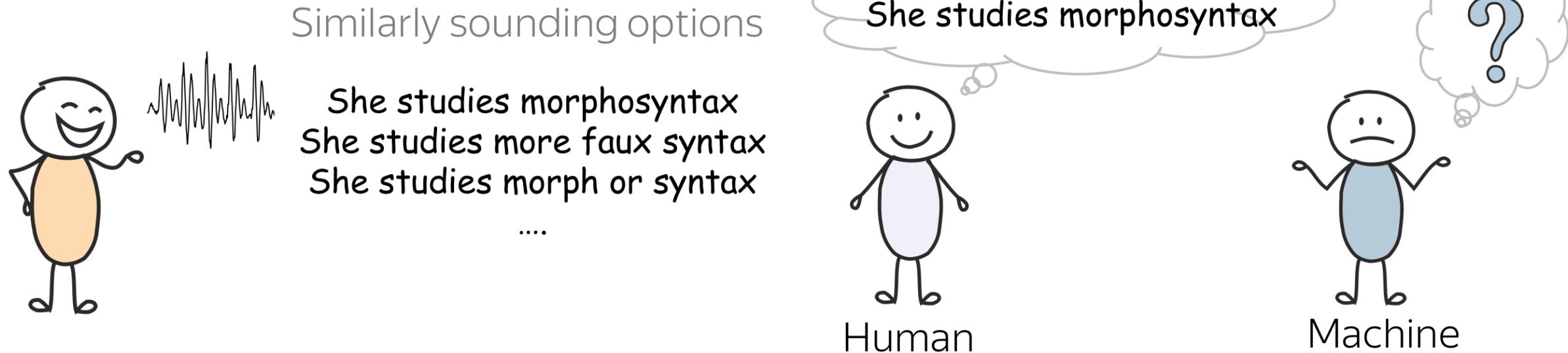
In language, an event is a linguistic unit: text, sentences, word, token, symbol

The goal of a language model is to estimate the probabilities of these events.

Definition: A Language model is a probability distribution over linguistic units. It estimate the probability of different linguistic units: symbols, tokens, token sequences.



What is easy for humans, can be very hard for machines



Source: the morphosyntax example is from the slides by Alex Lascarides and Sharon Goldwater, Foundations of Natural Language Processing course at the University of Edinburgh.

But how a machine is supposed to understand?

- A machine needs a language model, which estimates the probabilities of sentences.
- If a language is good, it will assign a larger probability to a correct sentence.

$P(\text{``this is a correct sentence.'}') > P(\text{``htis is a coreckt sentence''}.)$



## History of (Large) Language Models

ELIZA  
(Joseph Weizenbaum)

```
Welcome to
EEEEE LL IIII ZZZZZZ AAAA
EE EE LL II ZZ AA AA
EEEEE LL IIII ZZZZZZ AA AA
EEEEE LLLLLL IIII ZZZZZZ AA AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landesleiter 2005.

ELIZA: Is something troubling you ?
YOU: Men are all alike.
ELIZA: I'm glad you mention it, do you suppose ?
YOU: They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU: Well, they're always pestering us.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU: He says I'm depressed much of the time.
ELIZA: I'm sorry to hear that. You are depressed.
YOU: It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
```

N-gram:

N-gram models emphasized the importance of context in language

$P(\text{word} \mid \text{given some history})$

(Sufficient digitally available text)

1991 - The World Wide Web provides a Massive Source of Data

1966

1990

1995

2013

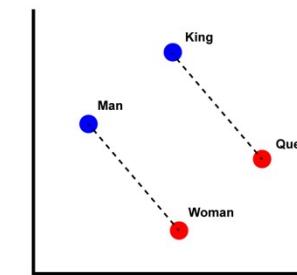
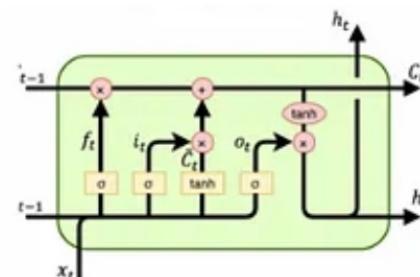
2016

2017

Rule-based  
models

Statistical language  
models

Neural  
Networks



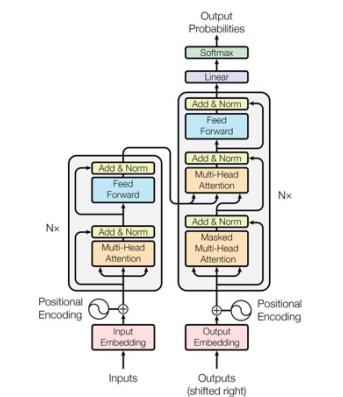
LSTM:

Long Short-Term Memory  
Sepp Hochreiter  
Jürgen Schmidhuber

Google's Neural Machine  
Translation System

Word2Vec:  
Word Embeddings  
Thomas Mikolov  
(Google)

Transformer  
Attention is all you need  
(Google)



# 04

## How to build a Language Model



A language model is a probability distribution over sequences of a vocabulary.

## General Framework:

- Our goal is to estimate probabilities of text fragments.
- Without loss of generality (w.l.o.g.) we deal with sentences.
- We want these probabilities to reflect knowledge of a language.
- We want sentences that are “more likely” to appear in a language to have larger probability according to our language model.

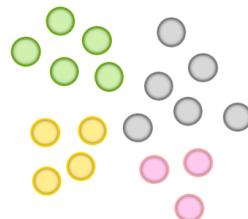
Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



# How to build a Language Model - General Framework

## Could simple probability theory help?

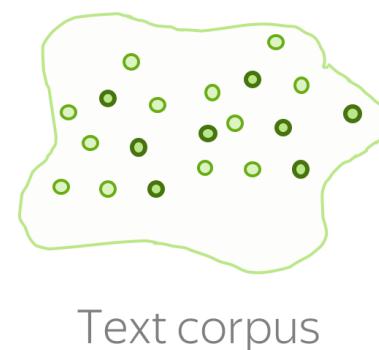
What is the probability to pick a green ball?



$$\frac{5}{5 + 6 + 4 + 3} = \frac{5}{18}$$

The probability to pick a ball of a certain color (let's say green) from this basket is the frequency with which green balls occur in the basket.

Can we do the same for sentences?



$$P(\text{the mut is tinning the tebn}) = \frac{0}{|\text{corpus}|} = 0$$

$$P(\text{mut the tinning tebn is the}) = \frac{0}{|\text{corpus}|} = 0$$

With this approach, sentences that never occurred in the corpus will receive zero probability

But the first sentence is “more likely” than the second!  
This method is not good!



Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



## Solution: Decompose a sentence into smaller parts!

We can not reliably estimate sentence probabilities if we treat them as atomic units. Instead, let's decompose the probability of a sentence into probabilities of smaller parts.

For example, let's take the sentence **I saw a cat on a mat** and imagine that we read it word by word. At each step, we estimate the probability of all seen so far tokens.

$$P(\mathbf{I}) =$$

$$\underline{P(\mathbf{I})}$$

Probability of **I**

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



# How to build a Language Model - General Framework

Formally, let  $(x_1, \dots, x_n)$  be tokens in a sentence, and  $P(x_1, \dots, x_n)$  the probability to see all these tokens (in this order). Using the product rule of probability (aka the chain rule), we get:

$$\begin{aligned} P(x_1, x_2, \dots, x_n) &= P(x_1) * P(x_2 | x_1) * P(x_3 | x_1, x_2) * \dots * P(x_n | x_1, x_2, \dots, x_{n-1}) \\ &= \prod_{t=1}^n P(x_t | x_{<t}). \end{aligned}$$

- We decomposed the probability of a text into conditional probabilities of each token given the previous context.
- What we got is the **autoregressive left-to-right language modeling framework**.

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



# How to build a Language Model - General Framework

Need to define, how to compute  $P(x_n|x_1, x_2, \dots, x_{n-1})$ .

- N-Gram
- Neural Models (neural networks)

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



## Generate text using a Language Model?

I \_\_\_\_\_

Once we have a language model, we can use it to generate text.

We do it one token at a time: predict the probability distribution of the next token given previous context, and sample from this distribution.

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



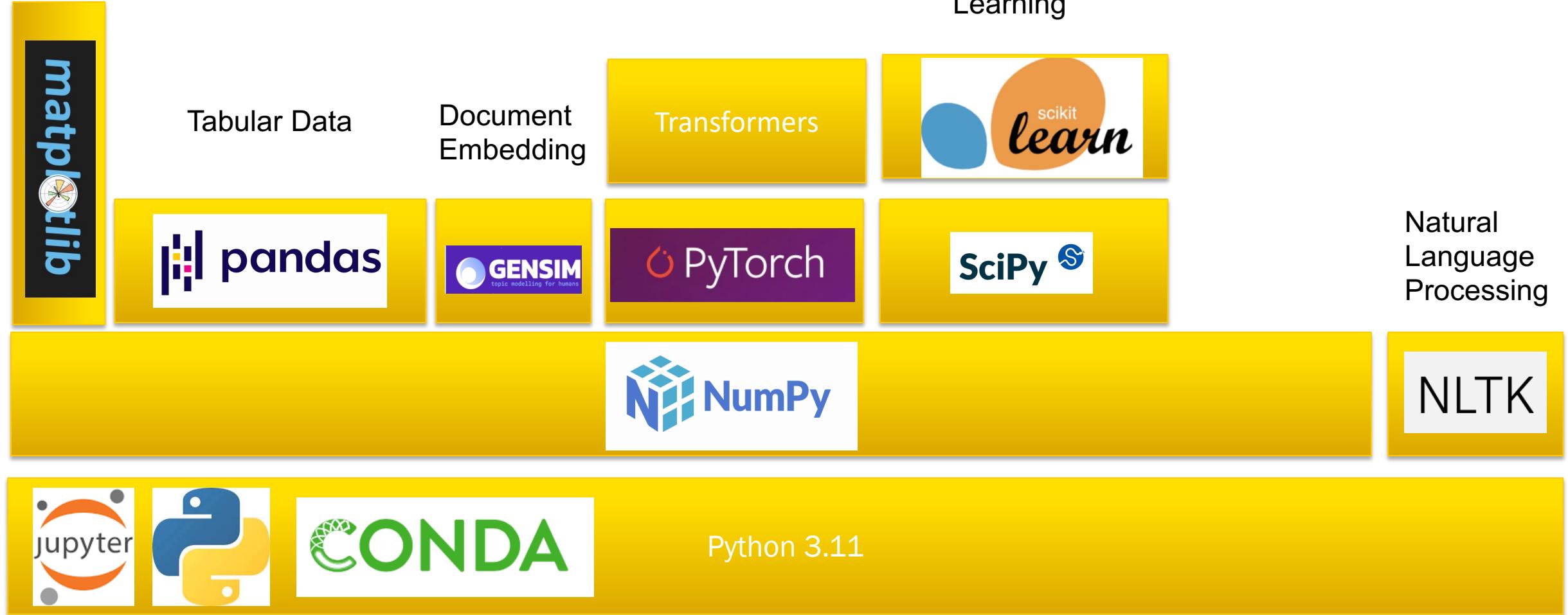
# 05

## Build a

# N-Gram Language Model with Python



# Install Python Environment



# Install Python

- Python: <https://www.python.org>
- Python on Mac with homebrew: <https://docs.python-guide.org/starting/install3/osx/>
- Python on Windows: <https://docs.python-guide.org/starting/install3/win/>



Create for each project its own Python environment to omit dependency conflicts:

\$ python -m venv *envpath* creates a virtual environment (in *envpath* directory)

Beispiel:

```
(base) done@Dominiks-MBP ~ % python -m venv Documents/tmpenv/test1
(base) done@Dominiks-MBP ~ % cd Documents/tmpenv/test1
(base) done@Dominiks-MBP test1 % ls -l
total 8
drwxr-xr-x 12 done  staff  384  2 Okt 10:46 bin
drwxr-xr-x   2 done  staff    64  2 Okt 10:46 include
drwxr-xr-x   3 done  staff    96  2 Okt 10:46 lib
-rw-r--r--   1 done  staff    89  2 Okt 10:46 pyvenv.cfg
```



# Python Environment Management

```
(base) done@Dominiks-MBP ~ % tree -dL 4 Documents/tmpenv/test1
```

```
Documents/tmpenv/test1
```

```
└── bin   ← Link to system python
    ├── include
    └── lib
        └── python3.10
            └── site-packages
                ├── _distutils_hack
                ├── pip
                ├── pip-22.3.1.dist-info
                ├── pkg_resources
                ├── setuptools
                └── setuptools-65.5.0.dist-info
```



## Activate and deactivate the virtual environment

On Unix/MacOS:

\$ source envpath/bin/activate	activates the environment
\$ source envpath/bin/deactivate	deactivates the environment

On Windows:

C:\> envpath\Scripts\activate.bat

Activation does many things:

1. Add virtualenv's *bin* directory at the beginning of shell PATH
2. Defines a deactivate command to return the environment to its former state
3. Modifies shell prompt: puts the environment name in front of it
4. Defines a VIRTUAL\_ENV environment variable as the path to virtual environments root directory.

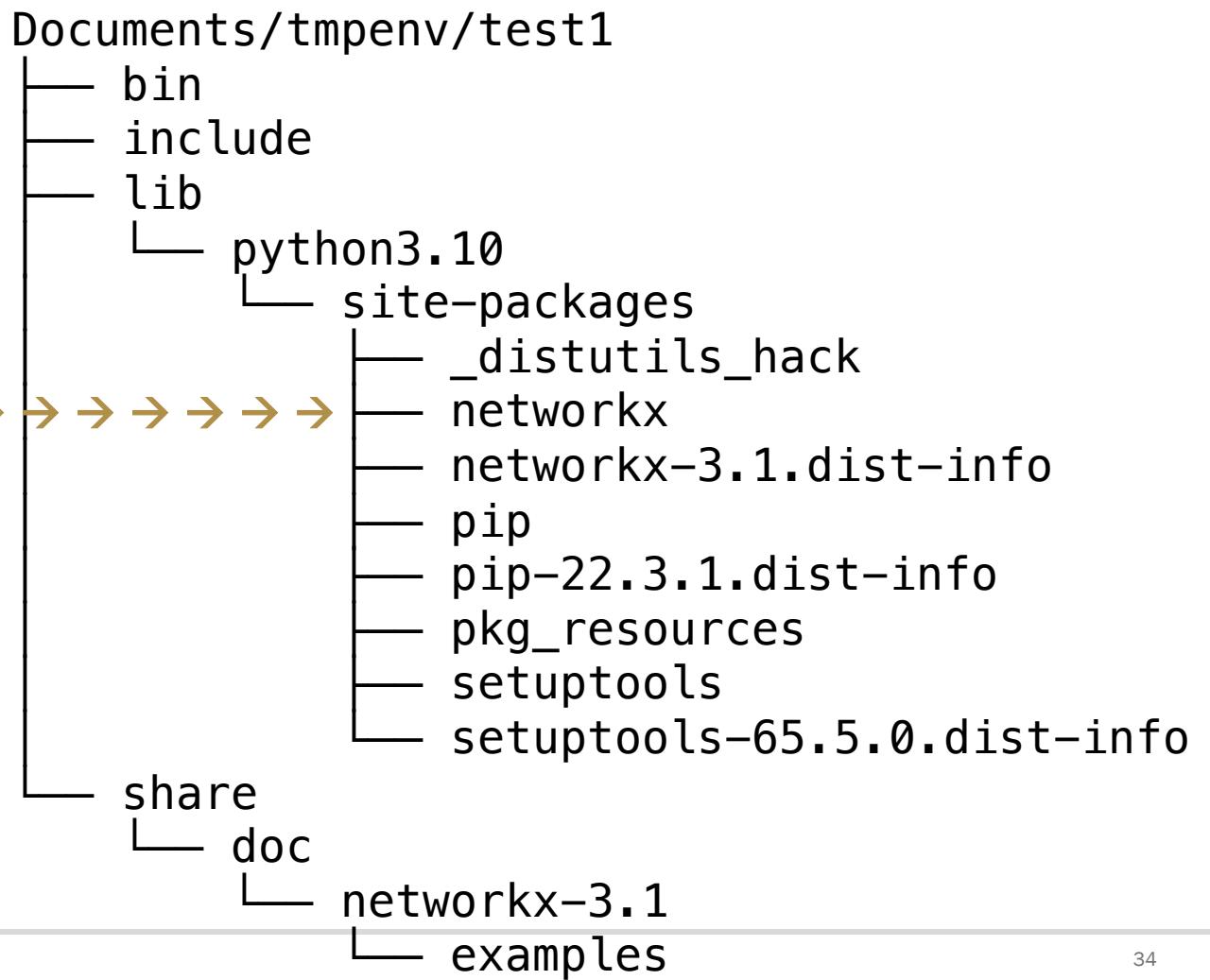


Install a package from <https://pypi.org> into the virtual environment:

`pip install packegname`

Example:

`pip install networkx`



A smart way to manage environments is conda

<https://docs.conda.io/en/latest/>

```
conda create --name >envname< python=3.11
```

```
conda activate >envname<
```

```
conda deactivate >envname<
```

**CONDA** <https://docs.conda.io/projects/conda/en/latest/index.html>

# Python Environments

- <https://pypi.org> pip install jupyter
- <https://jupyter.org> pip install nltk
- <https://www.nltk.org> pip install torch
- <https://pytorch.org> pip install transformers
- <https://huggingface.co>



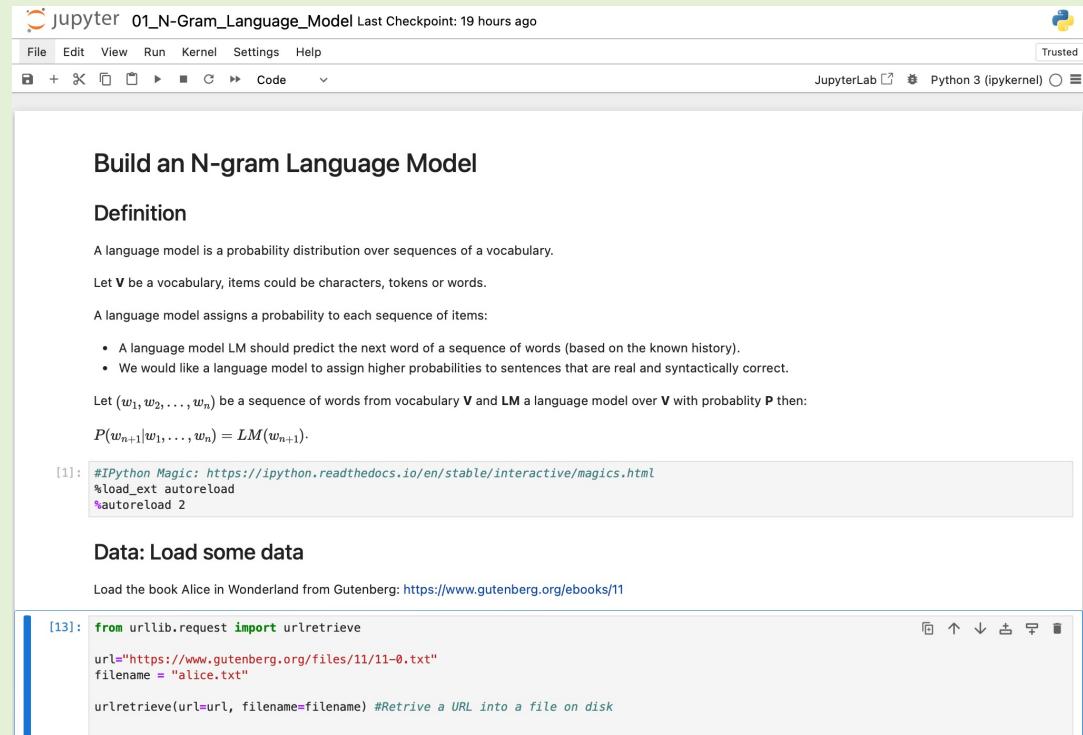
# Github Repository

Our github Repo: <https://github.com/dneumnn/ws24>



# N-Gram Notebook

## 01\_N-Gram\_Language\_Model.ipynb



jupyter 01\_N-Gram\_Language\_Model Last Checkpoint: 19 hours ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

### Build an N-gram Language Model

#### Definition

A language model is a probability distribution over sequences of a vocabulary.

Let  $V$  be a vocabulary, items could be characters, tokens or words.

A language model assigns a probability to each sequence of items:

- A language model LM should predict the next word of a sequence of words (based on the known history).
- We would like a language model to assign higher probabilities to sentences that are real and syntactically correct.

Let  $(w_1, w_2, \dots, w_n)$  be a sequence of words from vocabulary  $V$  and LM a language model over  $V$  with probability P then:

$$P(w_{n+1}|w_1, \dots, w_n) = LM(w_{n+1})$$

```
[1]: #IPython Magic: https://ipython.readthedocs.io/en/stable/interactive/magics.html
%load_ext autoreload
%autoreload 2
```

#### Data: Load some data

Load the book Alice in Wonderland from Gutenberg: <https://www.gutenberg.org/ebooks/11>

```
[13]: from urllib.request import urlretrieve
url="https://www.gutenberg.org/files/11/11-0.txt"
filename = "alice.txt"
urlretrieve(url=url, filename=filename) #Retrive a URL into a file on disk
```

# 06

## Neural Language Models



Recap: Need to define, how to compute  $P(x_n|x_1, x_2, \dots, x_{n-1})$ .

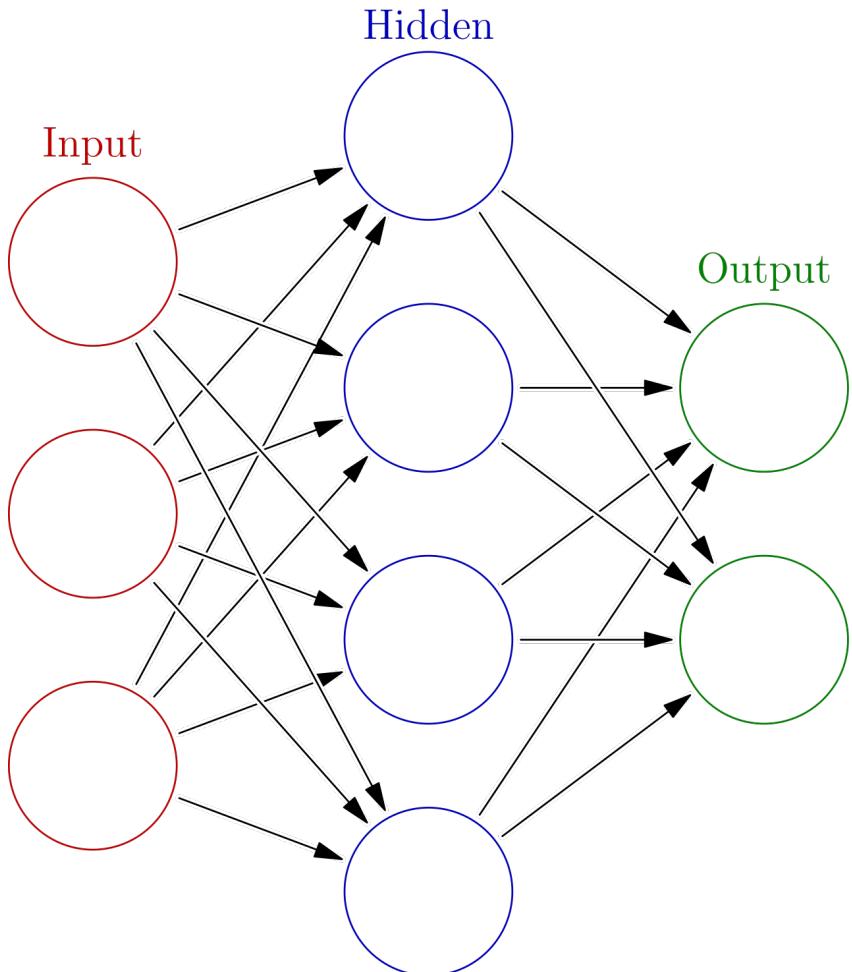
 history or context

- N-Gram
- Neural Models (neural networks)

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)



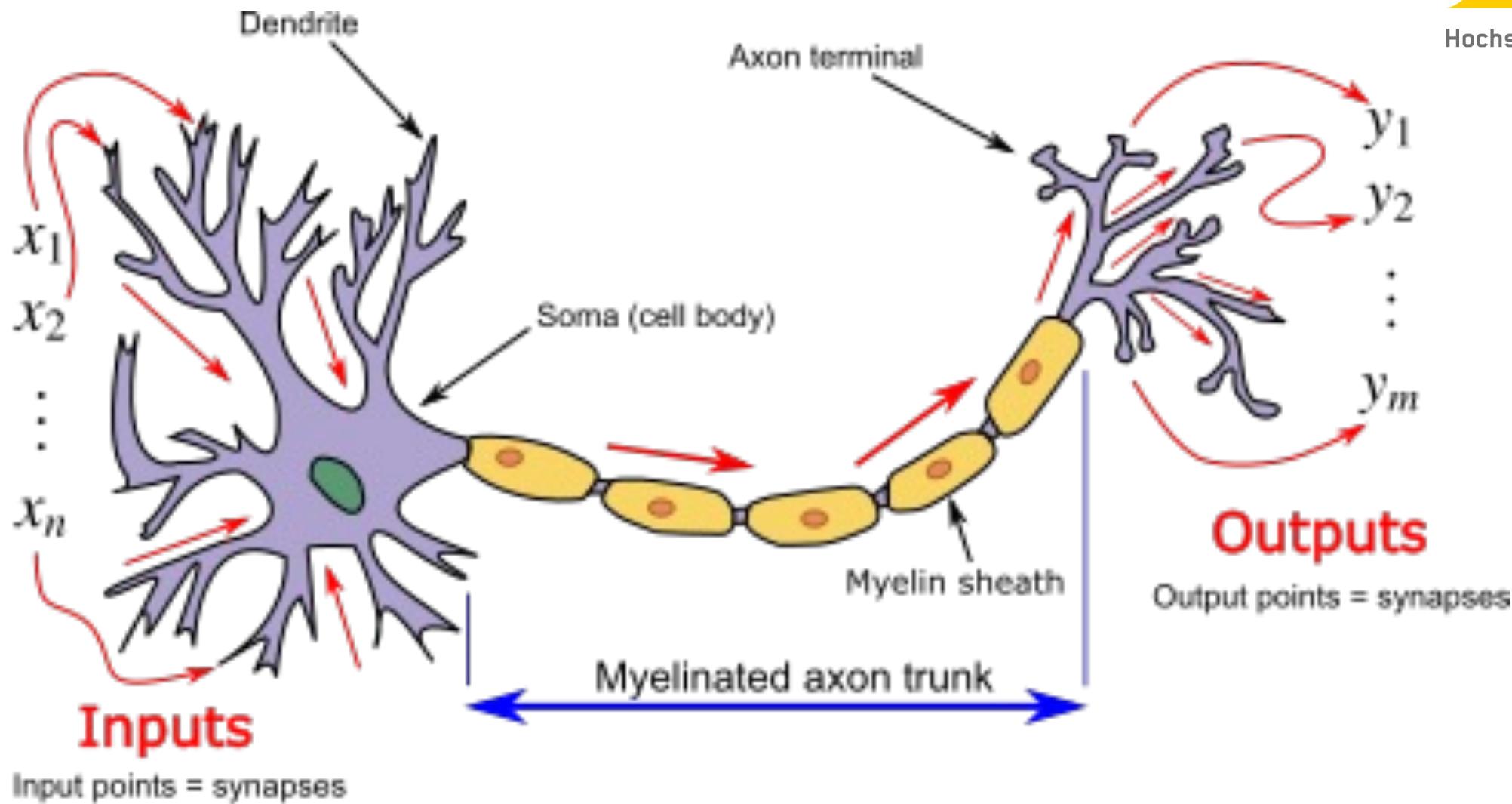
# A very short Introduction into Neural Networks



Source: [https://en.wikipedia.org/wiki/Neural\\_network\\_\(machine\\_learning\)](https://en.wikipedia.org/wiki/Neural_network_(machine_learning))

- In machine learning, a **neural network (NN)** is a model inspired by the structure and function of biological neural networks in animal brains.
- An NN consists of connected units or **nodes** called artificial neurons, which loosely model the neurons in the brain.
- These are connected by **edges**, which model the synapses in the brain.
- Each artificial neuron receives signals from connected neurons, then processes them and sends a signal to other connected neurons. The "signal" is a real number, and the output of each neuron is computed by some non-linear function of the sum of its inputs, called the **activation function**.
- The strength of the signal at each connection is determined by a **weight**, which adjusts during the learning process.

# A very short Introduction into Neural Networks

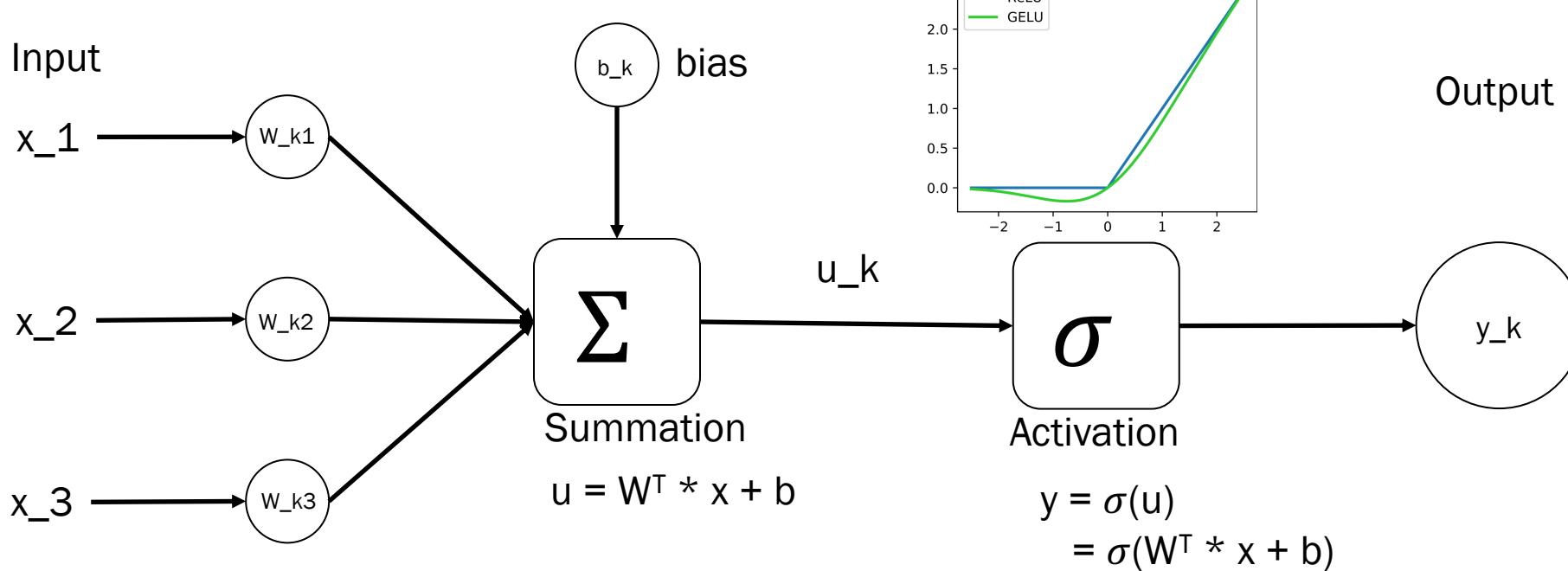


Input points = synapses

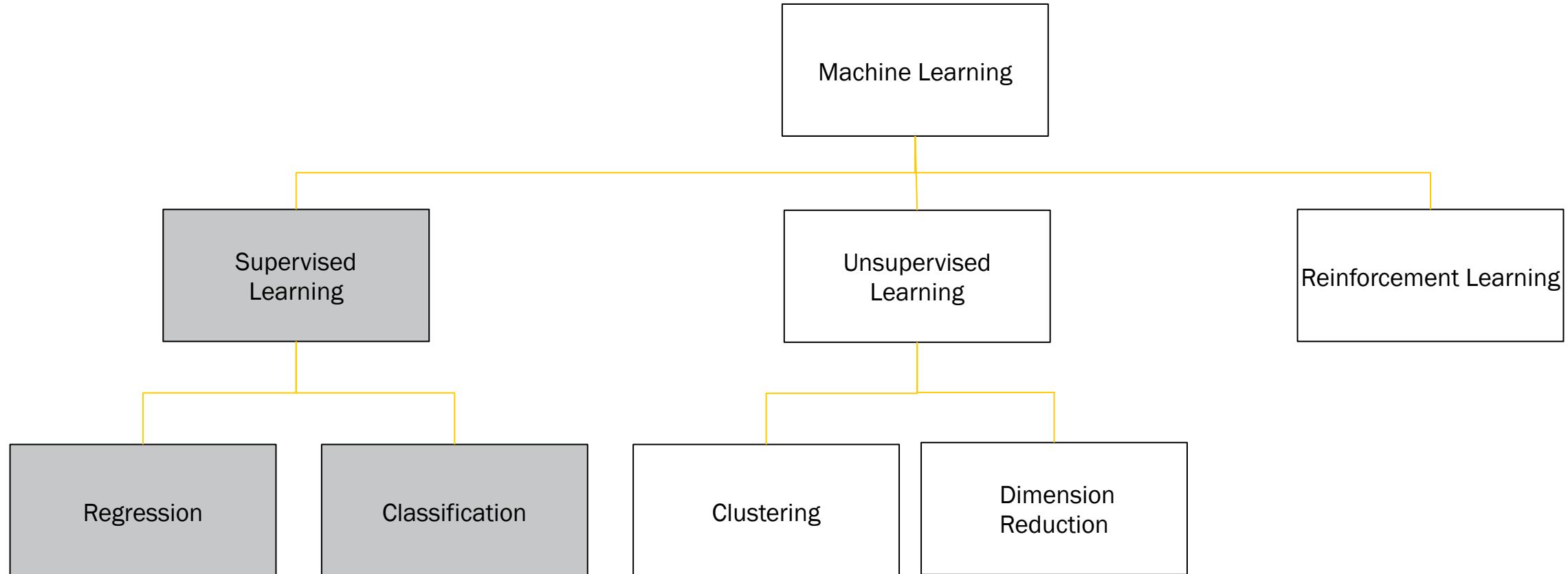
Source: [https://en.wikipedia.org/wiki/Neural\\_network\\_\(machine\\_learning\)](https://en.wikipedia.org/wiki/Neural_network_(machine_learning))

# A very short Introduction into Neural Networks

Each neuron of the hidden layer does the following computation



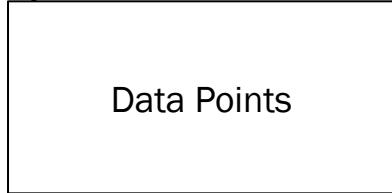
# A very short Introduction into Neural Networks



## Supervised Learning

### (2) Collect (labeled) Data

Training Data



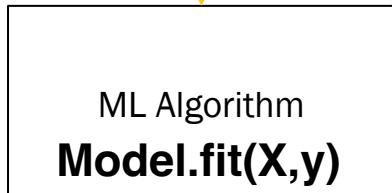
FEATURES  
LABELS

### (1) Define the problem to be solved

Unknown Pattern

$$f: X \rightarrow Y: f(x) = y$$

### (3) Choose Algorithm and Optimization Metric



Hyperparameter:  
- Learning Rate  
- epochs  
...

Model  
(Hypothesis)

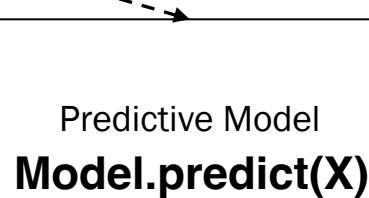
Cost / Loss  
-function

$$h: X \rightarrow Y: h(x; \theta) = \hat{y}$$

$$J(\theta) = \sum L(y, h(x; \theta))$$

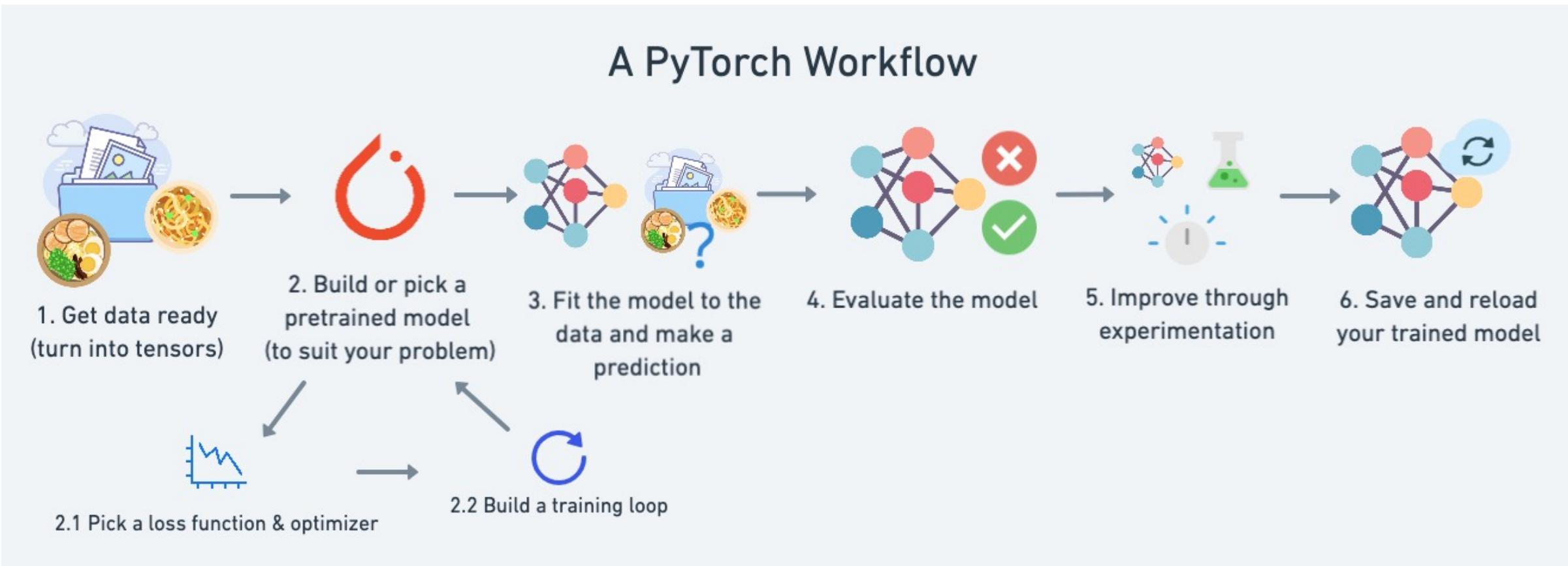
### (4) Choose a Metric to evaluate the final model

New Data →



→ Prediction

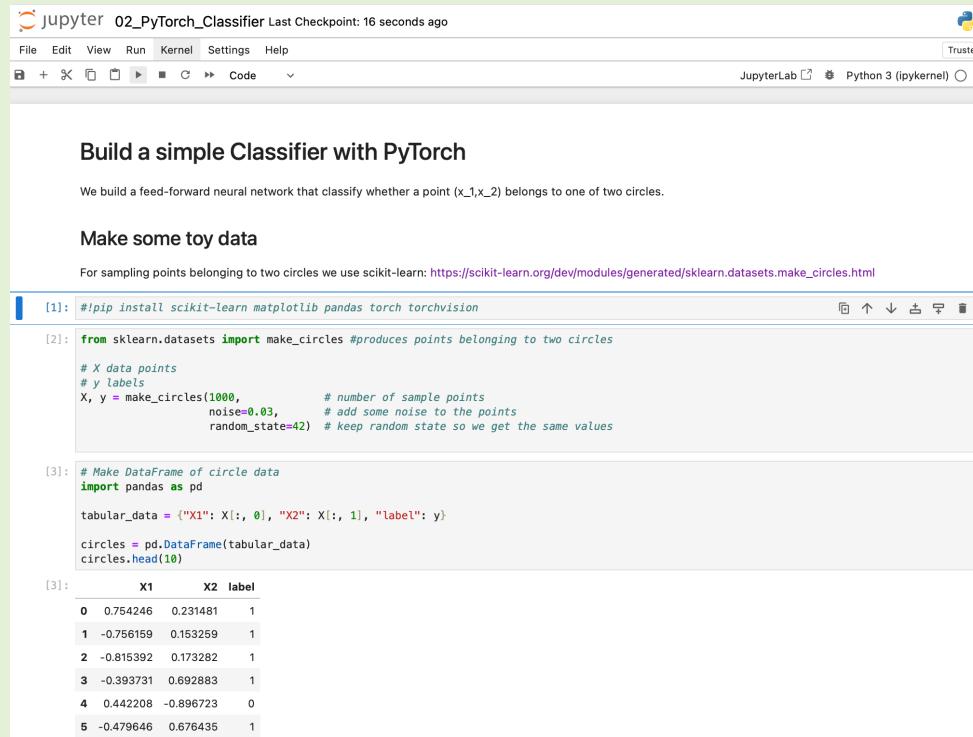
## Neural Networks with PyTorch



Source: [https://www.learnpytorch.io/01\\_pytorch\\_workflow/](https://www.learnpytorch.io/01_pytorch_workflow/)

# A very short Introduction into Neural Networks

## Build a Neural Classifier with PyTorch



The screenshot shows a Jupyter Notebook interface with the title "Build a simple Classifier with PyTorch". The notebook contains the following content:

We build a feed-forward neural network that classifies whether a point ( $x_1, x_2$ ) belongs to one of two circles.

**Make some toy data**

For sampling points belonging to two circles we use scikit-learn: [https://scikit-learn.org/dev/modules/generated/sklearn.datasets.make\\_circles.html](https://scikit-learn.org/dev/modules/generated/sklearn.datasets.make_circles.html)

```
[1]: #!pip install scikit-learn matplotlib pandas torch torchvision
[2]: from sklearn.datasets import make_circles # produces points belonging to two circles
# X data points
# y labels
X, y = make_circles(1000,           # number of sample points
                    noise=0.03,      # add some noise to the points
                    random_state=42) # keep random state so we get the same values
[3]: # Make DataFrame of circle data
import pandas as pd
tabular_data = {"X1": X[:, 0], "X2": X[:, 1], "label": y}
circles = pd.DataFrame(tabular_data)
circles.head(10)
```

	X1	X2	label
0	0.754246	0.231481	1
1	-0.756159	0.153259	1
2	-0.815392	0.173282	1
3	-0.393731	0.692883	1
4	0.442208	-0.896723	0
5	-0.479646	0.676435	1

Source: [https://www.learnpytorch.io/01\\_pytorch\\_workflow/](https://www.learnpytorch.io/01_pytorch_workflow/)

Intuitively, neural Language Models do two things:

1. process context  $(x_1, x_2, \dots, x_{n-1})$  → model-specific

The main idea here is to get a vector representation for the previous context. Using this representation, a model predicts a probability distribution for the next token. This part could be different depending on model architecture (e.g., RNN, CNN, whatever you want), but the main point is the same - to encode context.

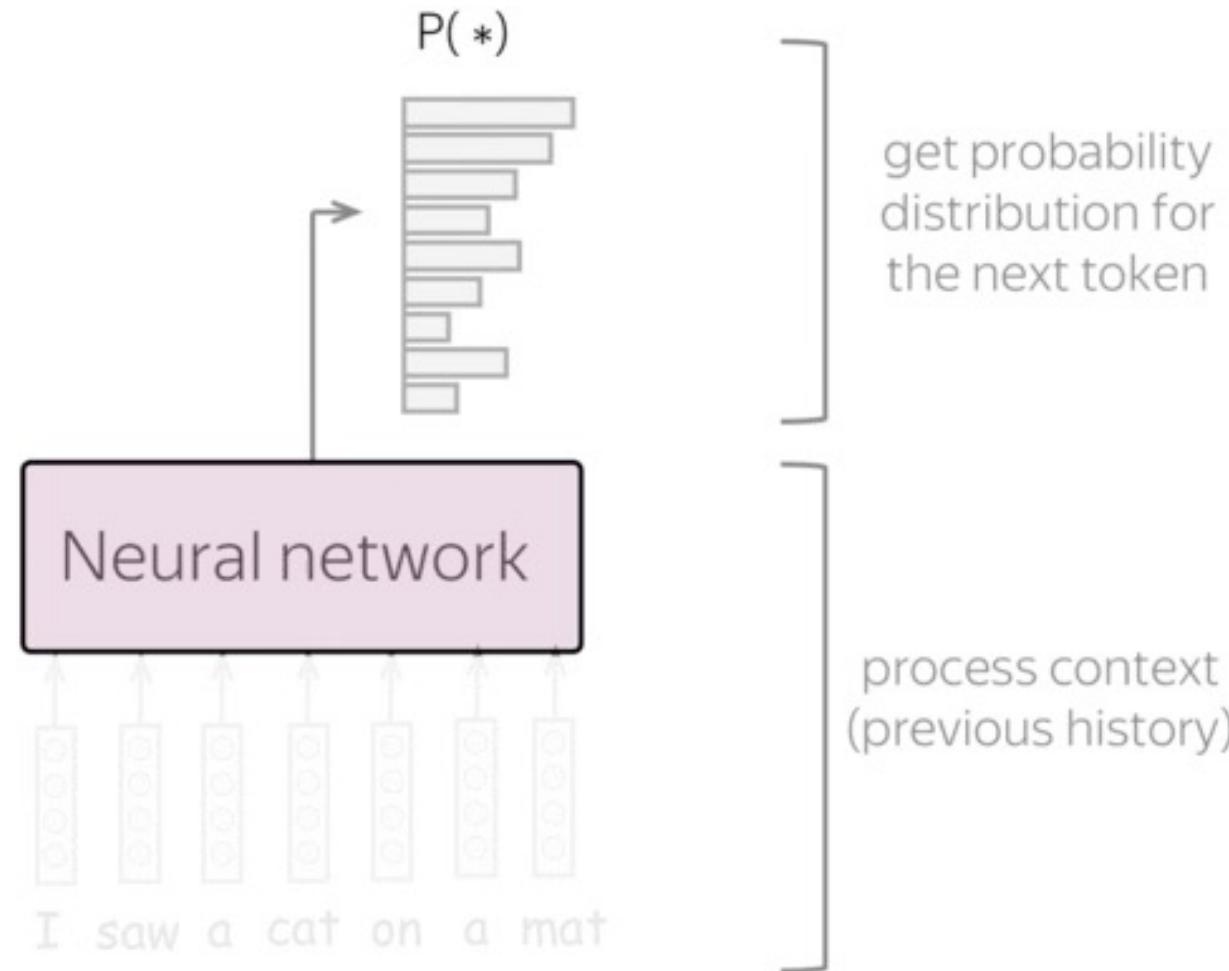
2. generate a probability distribution for the next token  $(x_n)$  → model-agnostic

Once a context has been encoded, usually the probability distribution is generated in the same way - see below.

Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)

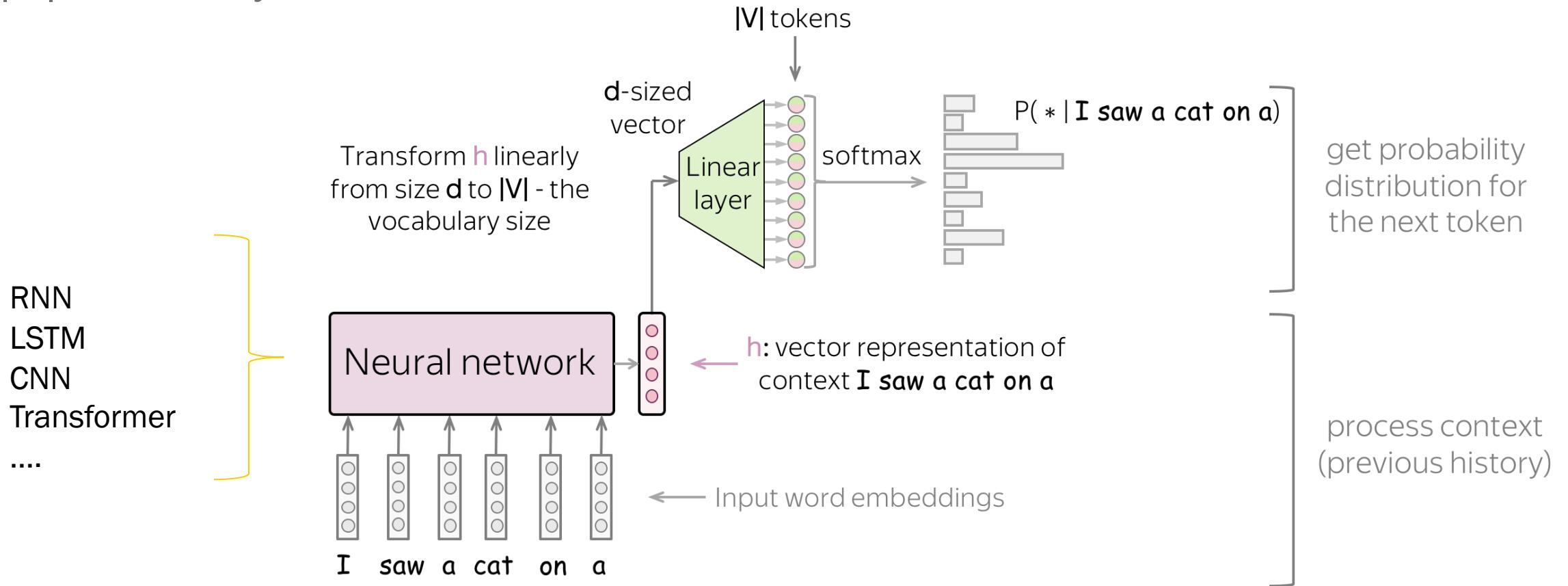


# Neutral Language Models



Source: [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)

We can think of neural language models as **neural classifiers**, where the classes are  $|V|$  vocabulary tokens.



## Pipeline

- feed word embedding for previous (context) words into the network
- get vector representation of context back from the network
- from this vector representation, predict a probability distribution for the next token.

We have:

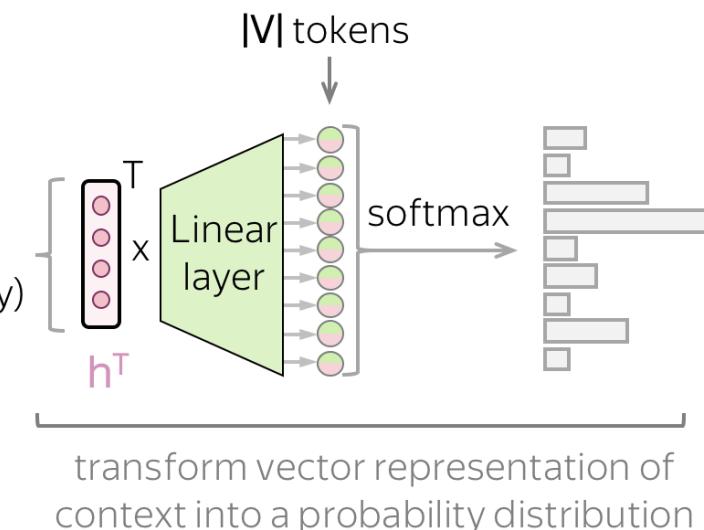
- $h$  - vector of size  $d$

We need:

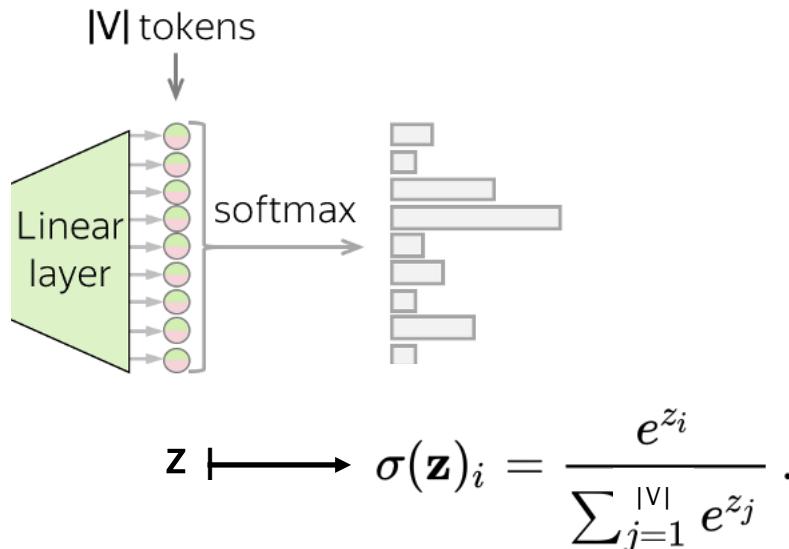
- vector of size  $|V|$  – probabilities for all tokens in the vocabulary

Transform linearly  
from size  $d$  to size  $|V|$

$d$  features  
(vector  
dimensionality)



Using Softmax to convert logits into a probabilities:



```
import numpy as np

def softmax(z, beta=1.0):
    return np.exp(beta * z) / np.sum(np.exp(beta * z))

z = np.array([1.0, 2.0, 3.0, 4.0, 1.0, 2.0, 3.0])
p = softmax(z)
p_max = np.argmax(p)
```

the softmax applies the standard exponential function to each element  $z_i$  of the input vector  $z$  (consisting of  $|V|$  real numbers), and normalizes these values by dividing by the sum of all these exponentials.

See: <https://pytorch.org/docs/stable/generated/torch.nn.Softmax.html>

## Training and the Cross-Entropy Loss

See: [https://lena-voita.github.io/nlp\\_course/text\\_classification.html](https://lena-voita.github.io/nlp_course/text_classification.html)

→ Equivalence to minimizing cross-entropy



# Further Reading



# Further Reading

- Raymond Lee, “Natural Language Processing”, 2024, Springer
- Manning & Schütze, Foundations of Statistical Natural Language Processing, 1999, MIT Press
- Lena Voita, “Language Modeling” [https://lena-voita.github.io/nlp\\_course/language\\_modeling.html](https://lena-voita.github.io/nlp_course/language_modeling.html)
- Zero to Mastery Learn PyTorch for Deep Learning: <https://www.learnpytorch.io>
- Dive into Deep Learning: <https://d2l.ai/index.html>



Dominik Neumann  
Hochschule Reutlingen, Alteburgstraße 150, 72762 Reutlingen  
[www.reutlingen-university.de](http://www.reutlingen-university.de)  
T. +49 172 9861157  
[dominik.neumann@reutlingen-university.de](mailto:dominik.neumann@reutlingen-university.de)  
[dominik.neumann@exxeta.com](mailto:dominik.neumann@exxeta.com)