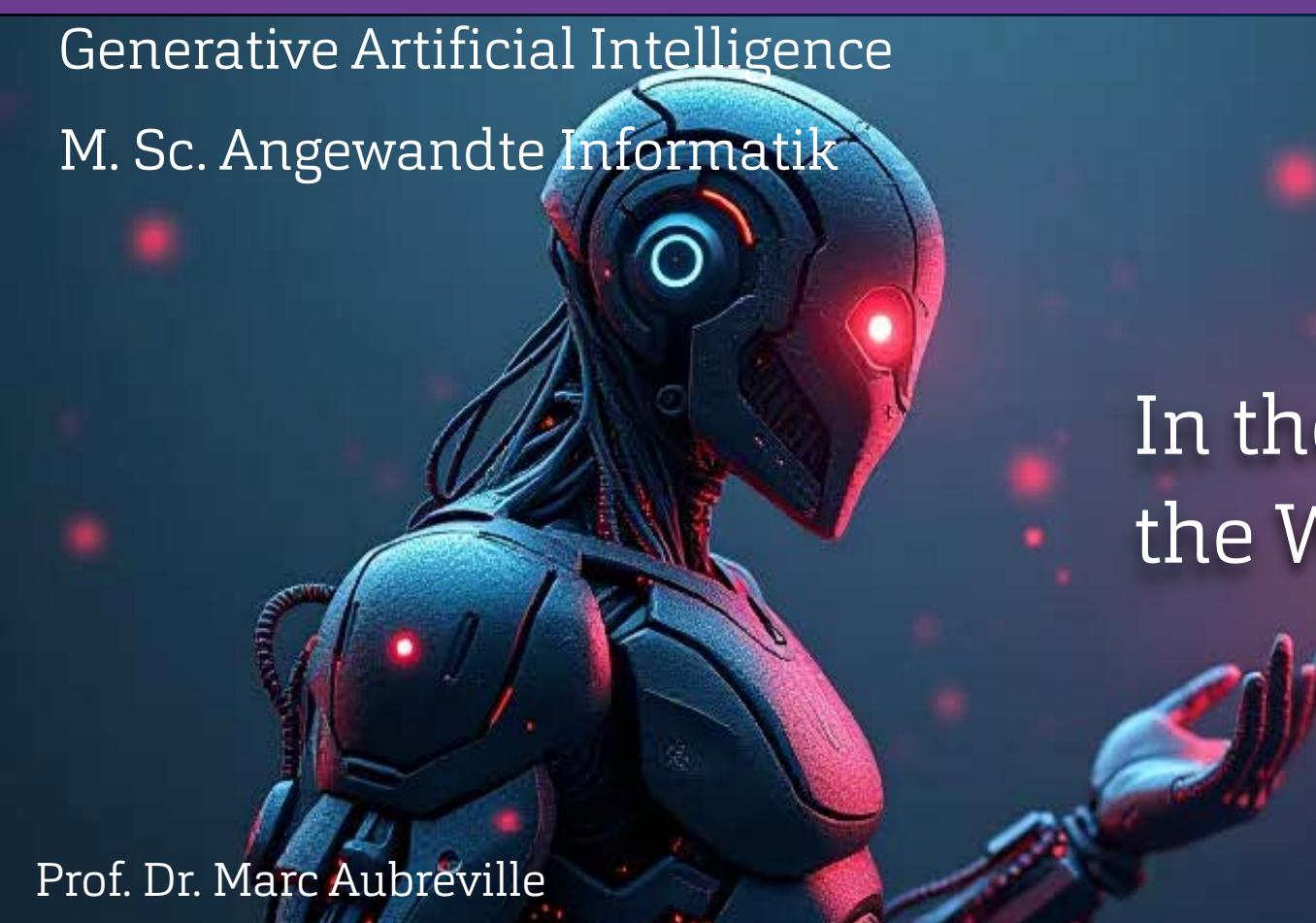


Hochschule
Flensburg
University of
Applied Sciences

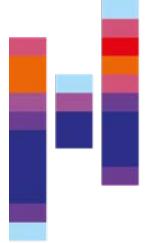
Text Generation

Generative Artificial Intelligence

M. Sc. Angewandte Informatik



In the beginning was
the Word ... |



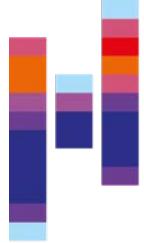
Outline

- In this chapter, we are looking at the following aspects:
 - Embedding / encoding text for LLMs
 - Pretraining
 - Fine-tuning strategies: supervised fine-tuning, RLHF, instruction tuning, parameter-efficient tuning (LoRA, adapters)
 - Scaling laws & data efficiency
 - Mixture of Experts (MoE) and sparse models
 - Reasoning and Long Context models
 - Tool use



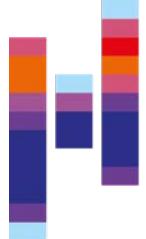
Hochschule
Flensburg
University of
Applied Sciences

Preprocessing: Tokenization, Encoding, Embedding



Text as Model Input

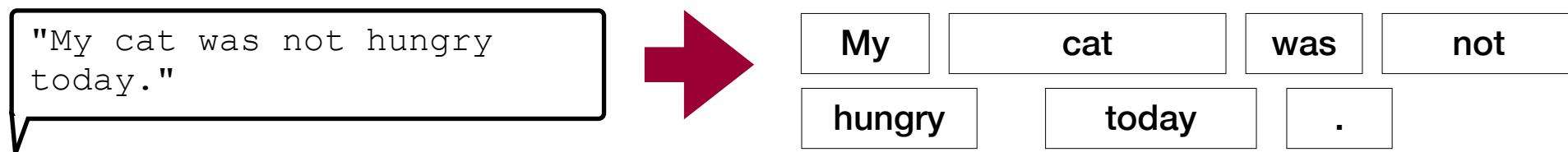
- Texts are unstructured data and cannot be directly processed by models.
- They must be converted into numerical representations first.
- Some obvious but ineffective methods include:
 - 1. ASCII Values:
Convert characters to their ASCII codes: "ABC" → [65, 66, 67].
Problem: Ignores word meaning and context.
 - 2. Bag of Words (BoW):
Represent text as a vector counting word occurrences in a document.
Example:
 - Texts: ["Hund", "Katze", "Hund"]
 - Vocabulary: {"Hund": 0, "Katze": 1}
 - Representation: [2, 1] (2× "Hund", 1× "Katze")
 - Problem: Loses word order and sentence structure information.



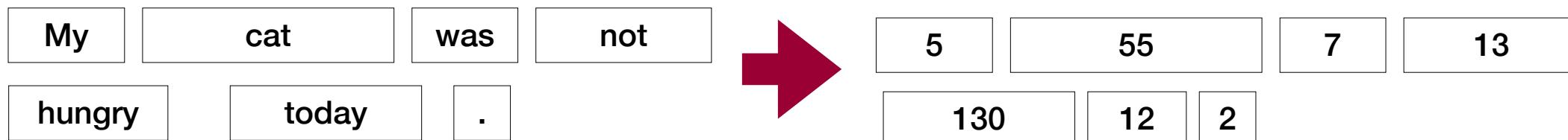
Overview: Preprocessing

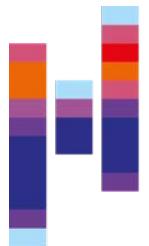
- We need to preprocess the data so a model get work on it:

Tokenization:



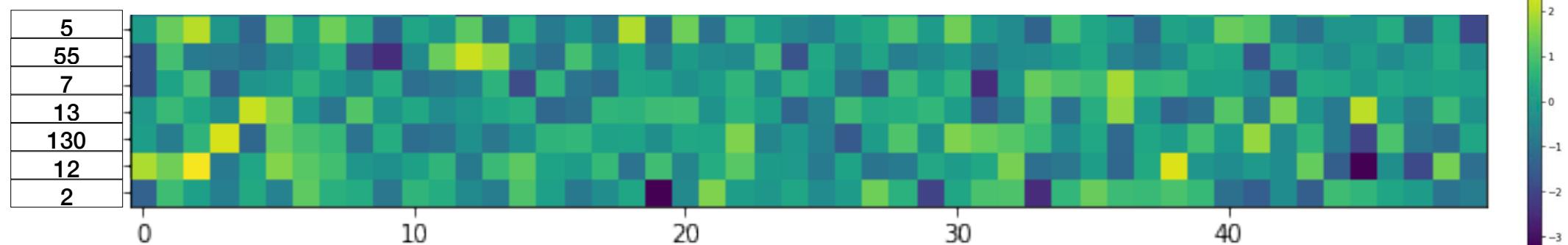
Encoding:

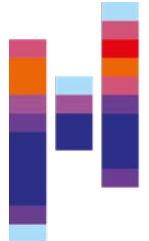




Even more preprocessing: Embeddings

- For text comprehension, each word is first converted into an embedding.
- This also achieves **normalization**.
- In the example, I have chosen a 50-dimensional embedding. Each word is represented there by a characteristic vector:





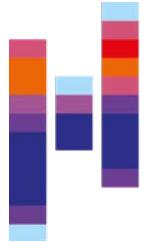
Tokenization

- Tokenization means splitting text into smaller units called tokens.
 - Simple tokenization: separates text by spaces (whitespace tokenization).
 - Example: "This is my dog" → "This", "is", "my", "dog".
 - These tokens (words) are then encoded so that a model can understand them more easily.
- More advanced tokenization methods exist:
 - Subword tokenization: splits words into meaningful subunits to improve processing.
Example: "The player" → "The", "play", "er".
 - Learned tokenization: a model learns on its own how to best split text into tokens.



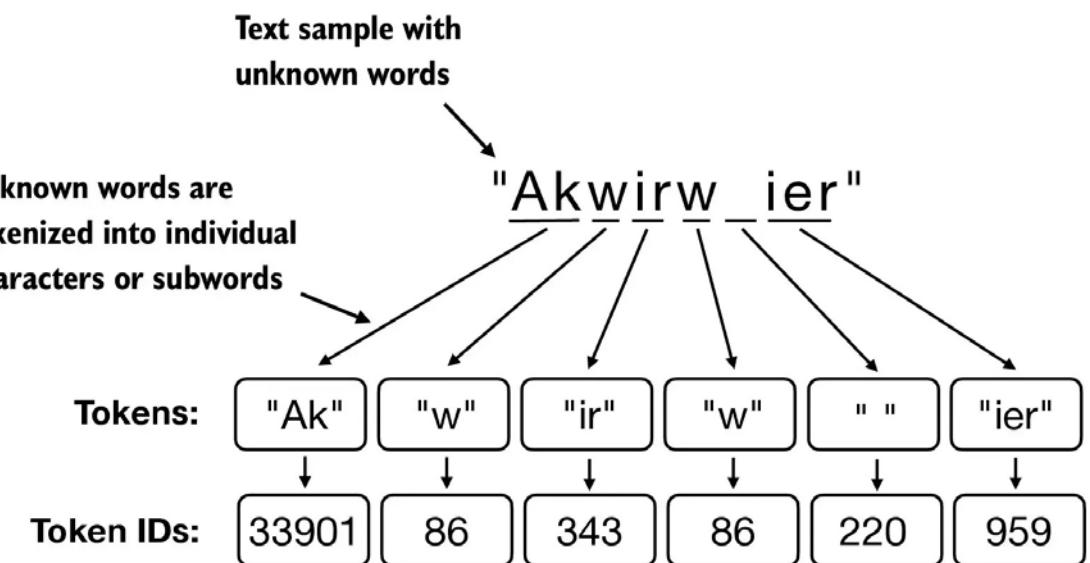
Tokenizer: Task

- The tokenizer must find meaningful ways to split text.
 - The number of tokens (code words) shouldn't be too large — this keeps processing fast and efficient (for example, when doing dictionary lookups).
 - If the tokens are too small, each one carries very little meaning, which makes the machine learning task much harder.
- Since tokenization and encoding are closely related, they are often optimized together — many tokenizers therefore include encoding as part of their process.



Byte-Pair Encoding (a tokenizer)

- BPE is a tokenizer that is optimized on large datasets to represent language efficiently.
- Idea:
 - Start with all characters (or bytes) as the initial vocabulary.
Count frequent adjacent pairs (like "b" + "e" → "be").
 - Merge the most frequent pair into a single symbol. Repeat this until you reach your desired vocabulary size.
 - The result is a merge table or tokenization rule set that tells you how to segment new text.





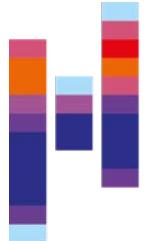
Example: LLAMA3 Tokenizer

Dies ist ein Test für den LLAMA Tokenizer. Dieser sollte Wörter, wie auch Zahlen, etwa 12345.33 sinnvoll aufteilen. Häufige Worte wie der, die, oder das sind dabei meist ein Token. Komplexe Worte wie Anhängerkupplung werden in viele Tokens aufgeteilt.

Der Tokenizer muss auch CODE beherrschen:

```
def my_function(xx:str):  
    print('Hello world')
```

<s> Dies ist ein Test für den LLAMA Tokenizer. Dieser sollte Wörter, wie auch Zahlen, etwa 12345.33 sinnvoll aufteilen. Häufige Worte wie der, die, oder das sind dabei meist ein Token. Komplexe Worte wie Anhängerkupplung werden in viele Tokens aufgeteilt. <0x0A><0x0A>Der Tokenizer muss auch CO
DE beherrschen:<0x0A><0x0A>def my_function(xx:str):<0x0A> print('Hello world')



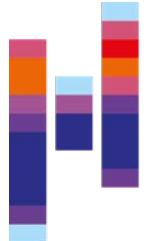
tiktoken (OpenAI tokenizer)

- The tiktoken algorithm used by OpenAI for their LLMs is a form of Byte Pair Encoding:
 - Text is first converted to UTF-8 bytes.
 - The tokenizer merges frequent byte pairs into larger chunks (subwords).
 - Each chunk (a “token”) gets a unique integer ID.
 - This sequence of IDs becomes the model’s input.

```
import tiktoken
enc = tiktoken.get_encoding("cl100k_base")
assert enc.decode(enc.encode("hello world")) == "hello world"

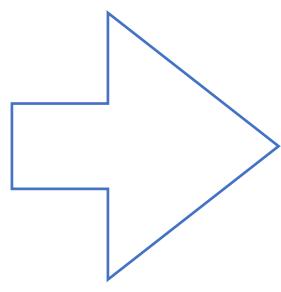
# To get the tokeniser corresponding to a specific model in the OpenAI API:
enc = tiktoken.encoding_for_model("gpt-4o")
```

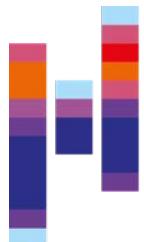




Encoding

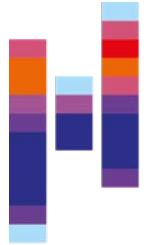
- The next step after creating tokens is embedding.
- In this step, each token (semantic unit) is assigned a code or number.
- One simple method is to identify the N most frequent tokens and assign each one a unique number.
- The data can then be fed to the model as a one-hot vector (a vector where only one position is “1” and all others are “0”).

Das	2		0 1 0 0
ist	1		1 0 0 0
mein	3		0 0 1 0
Hund	4		0 0 0 1



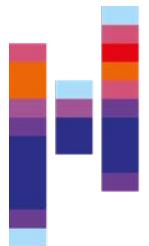
Tokenization and Encoding

- The goal of encoding is to create an input vector that captures meaning **efficiently** and is **as compact as possible**.
- Redundancy, which is common in natural language, should be **minimized**.
- Byte Pair Encoding (BPE) was first a text compression method, later adapted by OpenAI as a tokenization and encoding technique.
- The idea: find the most frequent pairs of bytes (or characters) and merge them into single tokens.
- BPE works by finding frequently occurring byte pairs in text and encoding them together.
- To ensure consistent encoding, it's best to train (optimize) the encoding once on large text corpora and then reuse it.



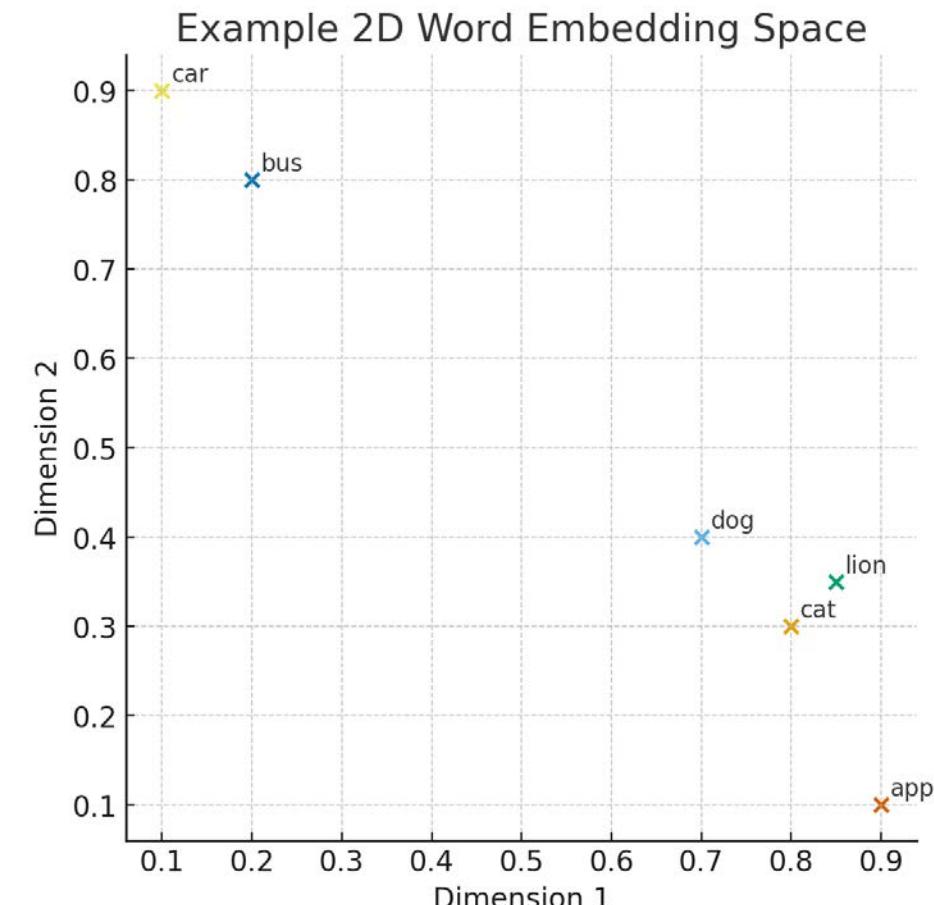
The problem with one-hot encoding for language

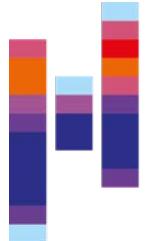
- One-hot encoding assigns each word a unique vector where only one element is 1 and the rest are 0.
 - Example:
"cat" → [1, 0, 0, 0]
"dog" → [0, 1, 0, 0]
- This means that semantically similar words (like "cat" and "dog") have completely different encodings, even though their meanings are related.
- One-hot vectors are also very high-dimensional – one dimension per word in the vocabulary.
- A vocabulary of 50,000 words means each vector has 50,000 positions!
- High dimensionality makes training slower, increases the risk of overfitting, and requires more data and computation.



Embeddings

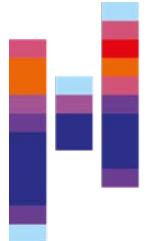
- Embeddings map one-hot vectors into a lower-dimensional space (e.g., from 50,000 dimensions to 300 dimensions).
- In this new vector space, words with similar meanings are closer together.
- Notice how "cat" and "dog" have similar vectors, while "car" is more different — the model can now understand semantic relationships.
- Embeddings thus represent a regularized latent space (as discussed before).



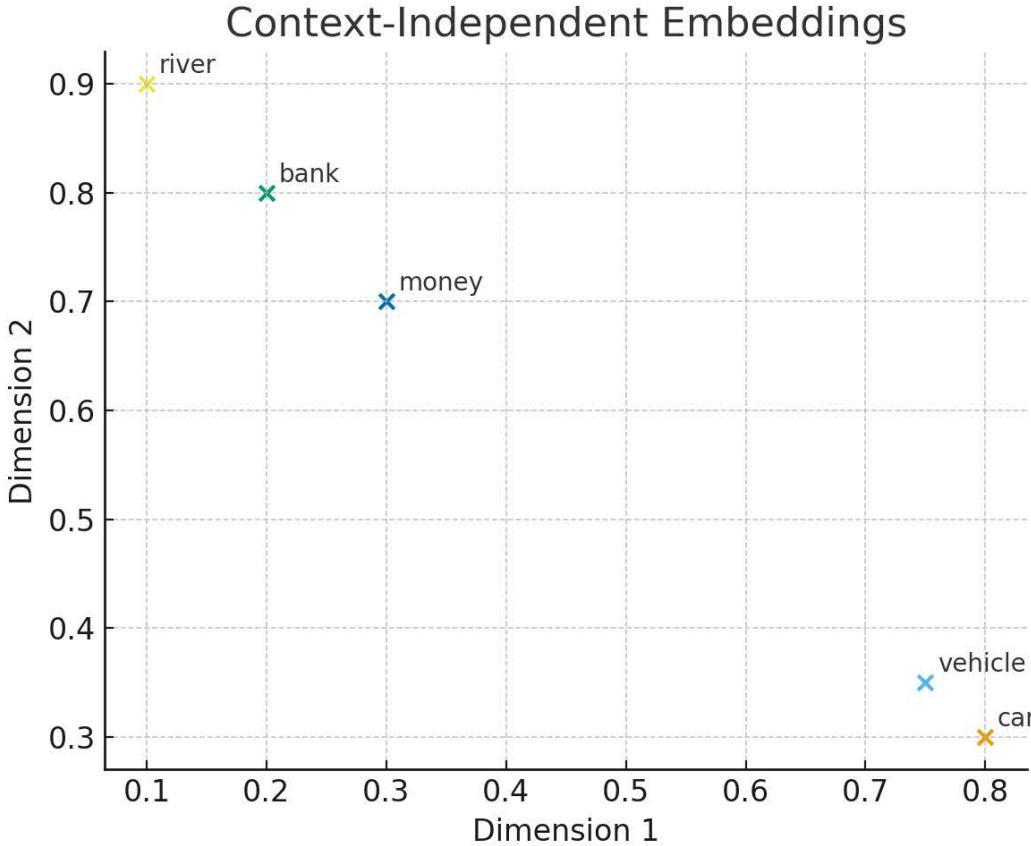


Creating an embedding

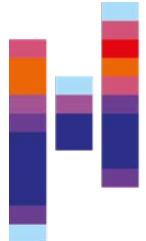
- Goal: Convert discrete token IDs into continuous vector representations
- Embeddings are generated by using the one-hot-encoded encodings, and multiplying them with an embedding matrix E :
 $E \in \mathbb{R}^{V \times D}$, with V being the vocabulary size and D being the embedding dimensions.
- The embeddings can be learnt during model training using backpropagation alongside the model parameters.
- For fine-tuning and later model training, the embedding matrix essentially becomes a constant and can be frozen.



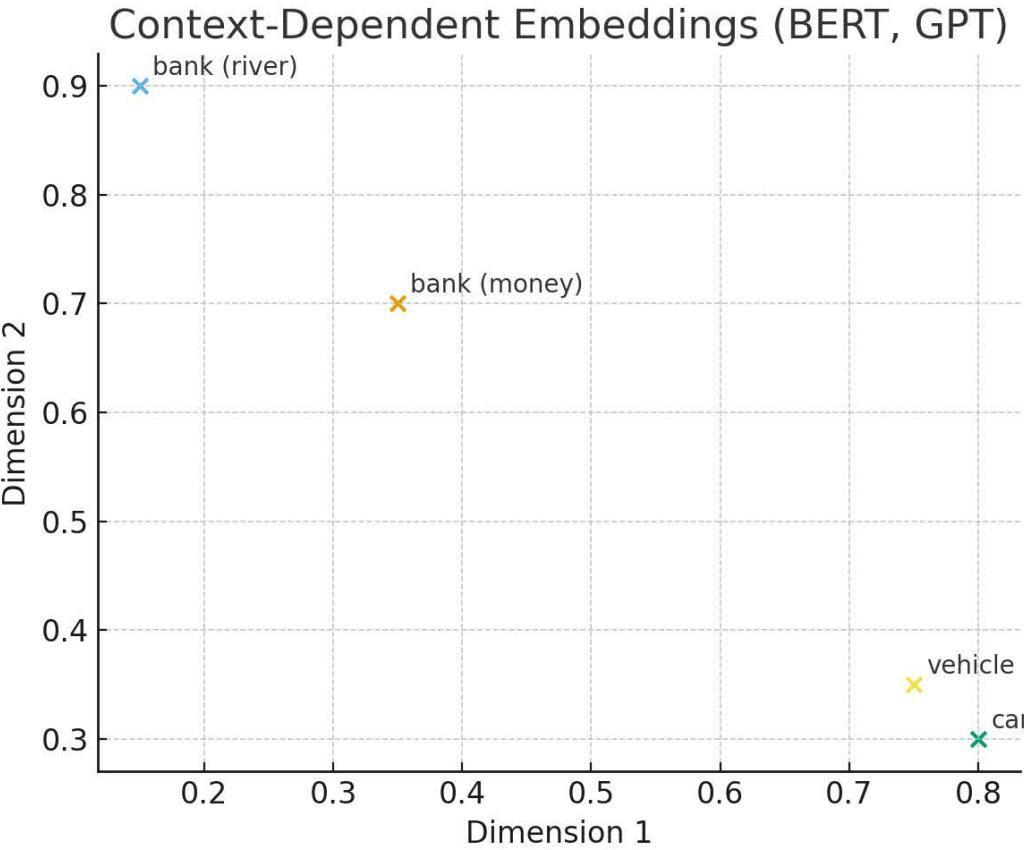
Context-Independent Embeddings



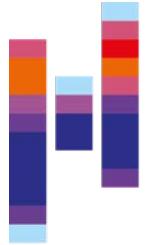
- Trained on large text corpora.
- Each word gets a fixed (static) vector, **regardless of where it appears.**
- These vectors are designed to be well separated (orthogonal) to make words easily distinguishable.
- Semantically similar words are close together in the vector space.
- Common examples: Word2Vec, GloVe, FastText



Context-Dependent Embeddings

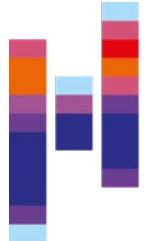


- Some words only make sense within their sentence context.
- Context-dependent embeddings adjust the vector based on the surrounding words.
- Example:
 - “I deposited money at the bank” → financial meaning
 - “We sat on the bank of the river” → geographical meaning
- will have different embeddings.
- Used in modern transformer models such as BERT and GPT, which dynamically compute word representations during encoding.



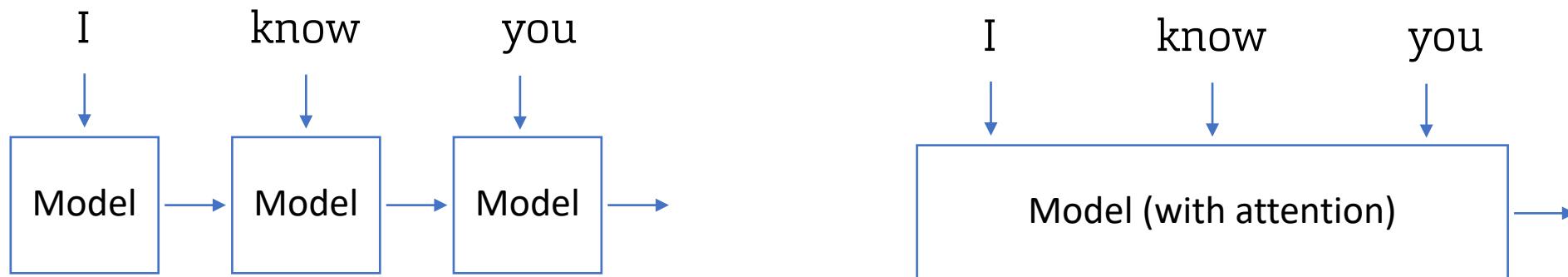
Hochschule
Flensburg
University of
Applied Sciences

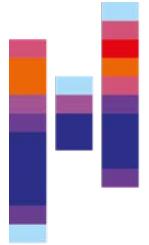
Positional Encoding



Loss of Positional Information

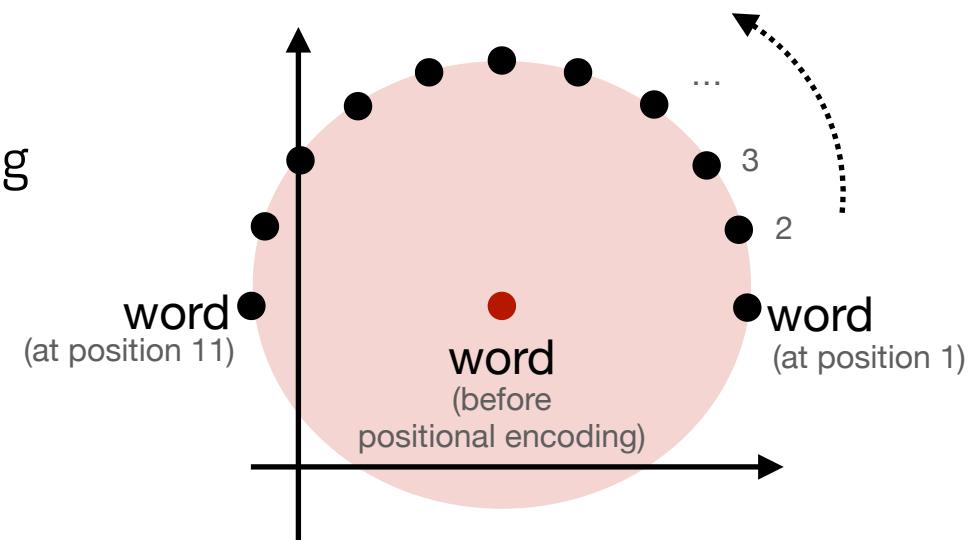
- The embeddings account for the semantics (i.e., meaning) of a word, but they do not explicitly encode the position.
- This is not a problem if recurrent architectures (RNNs, LSTMs,...) are used, since there the position in the sentence is explicitly given by the position within the sequence.
- However, for transformers, the attention mechanism **independently processes** each input token (conditional to the other ones) - which comes with the advantage of parallelization!

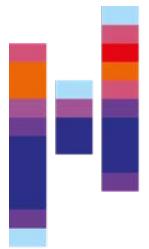




Positional Encoding

- Since the attention mechanism itself is independent of positions, we have to provide it with the information of where which input was located.
- We achieve this, for example, by rotating the encodings around their original point in the latent space.
- This does a little modification to the input data, just enough so the model can distinguish it from occurring at another position.
- The words/token remain distinguishable.



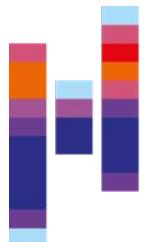


Positional Encoding

- Since our embeddings have a multitude of dimensions, we provide the position in each dimension to make sure the model can't confuse words with similar semantics but different order.
- Concretely, we add a positional vector to each embedding:
$$\mathbf{x}_t = \mathbf{e}_t + \mathbf{p}_t$$
- Each coordinate of the vector \mathbf{p}_t encodes the same position.
- There are multiple methods to do this, i.e., that specify an additional vector \mathbf{p}_t .

$$\begin{array}{c} x_0 = e_0 + p_0 \\ x_1 = e_1 + p_1 \\ \dots \quad \dots \quad \dots \\ x_N = e_N + p_N \end{array}$$

embeddings positional
 encoding



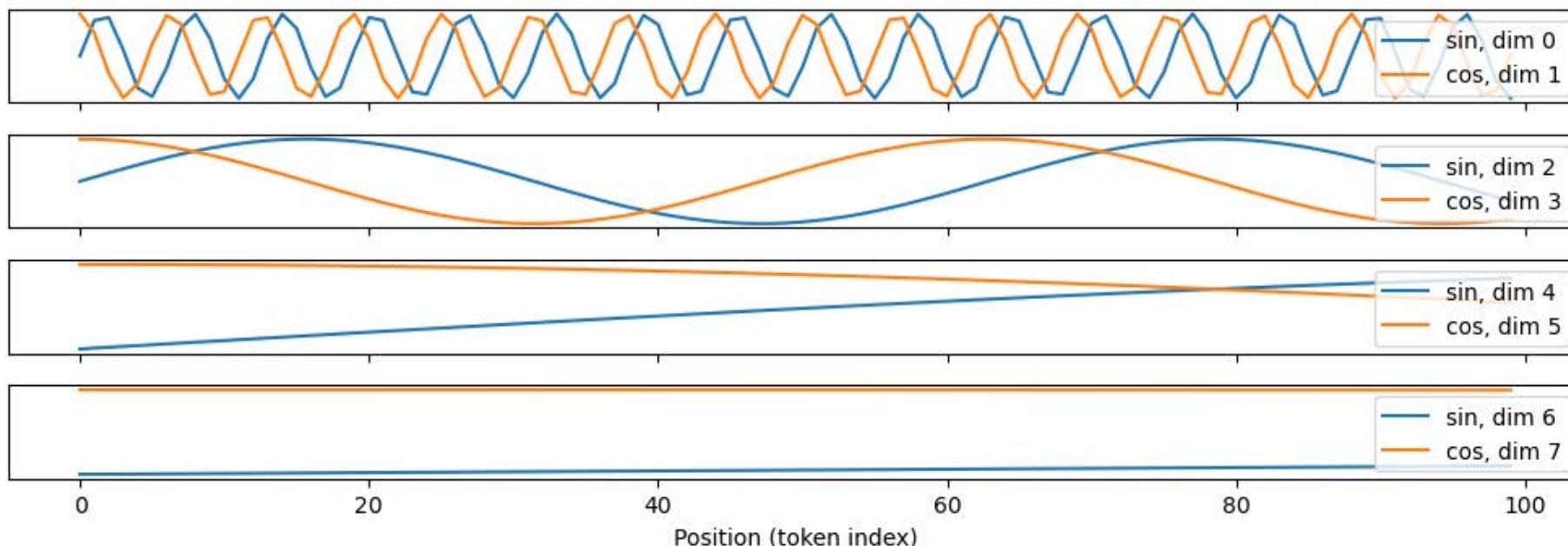
Absolute Positional Encoding

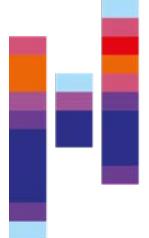
- One easy method proposed in the original transformer paper (Vaswani et al., 2017) is to add a sine/cosine function with the absolute position to each vector element:

$$p_{\text{pos},2i} = \sin(\text{pos}/10000^{2i/d_{\text{model}}})$$

$$p_{\text{pos},2i+1} = \cos(\text{pos}/10000^{2i/d_{\text{model}}})$$

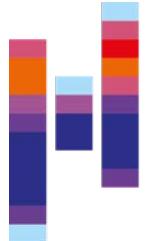
- This adds a sine/cosine component to each second dimension, respectively.





Position is Relative

- Language depends on relative order, not absolute position.
- Hence, providing absolute position encoding leaves the model with having to learn how to interpret positional indices, for instance, that position 8 is only relevant relative to position 7.
- It would be thus more efficient to encode position relative.
 - directly express distances or offsets between tokens (e.g., “this word is 3 steps after that one”).
 - reasoning about order and proximity more naturally
 - We can also expect simpler handling of variation of input length.



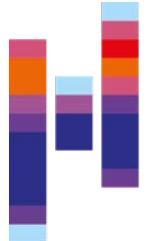
Relative Positional Encoding (2)

- The authors of "Self-Attention with Relative Position Representations" (Shaw et al., 2018) proposed to use the relative position directly in self-attention.
- Let's revisit the formula for attention between inputs i, j :

$$att_{ij} = \text{softmax}(QK)^T V = \text{softmax} \left((x_i W^Q)(x_j W^K)^T \right) (x_j W^V)$$

- The matrix W^Q learns the query from the input: $Q = W^Q x_i$
- The matrix W^K learns the keys from the input: $K = W^K x_j$
- The matrix W^V learns the values from the input: $V = W^V x_j$

(I skipped the normalizing \sqrt{d} term here)



Relative Positional Encoding (3)

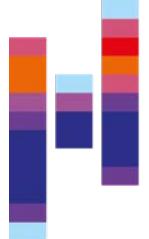
- The key idea in their paper is to learn the relationship between the positions directly in the attention.
- So, they modify the equations for the components:

$$K = W^K x_j + a_{i,j}^K$$

$$V = W^V x_j + a_{i,j}^V$$

- where $a_{i,j}^K$ and $a_{i,j}^V$ are **vectors**, indexed by a relative position $r = i - j$.
- One example for such a vector could be:

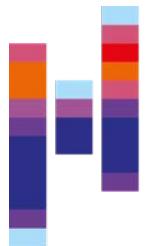
r=-2	r=-1	r=0	r=1
-0.5	-0.25	0	0.25



Relative Positional Encoding (4)

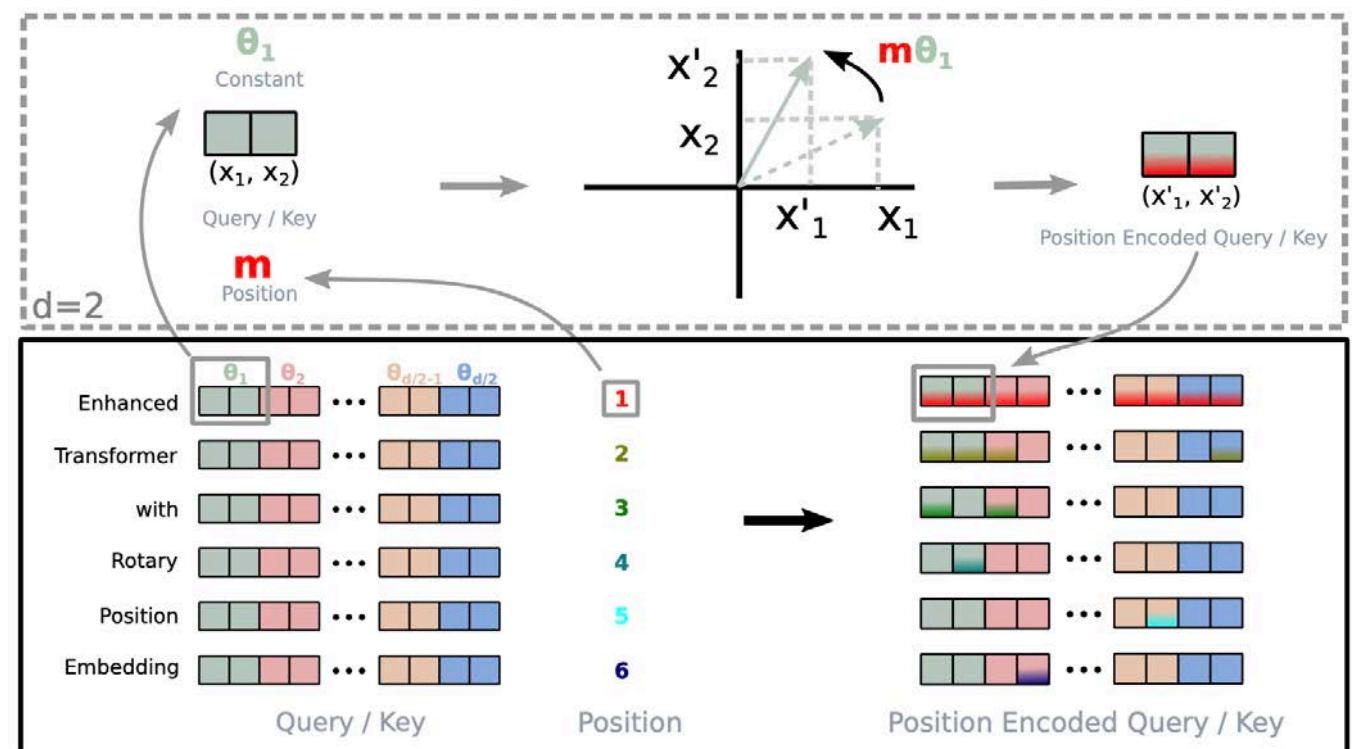
- However, the vectors $a_{i,j}^K$ and $a_{i,j}^V$ actually **learnt** in model training.
- They represent two vectors where each offset position is learnt with one characteristic (representative) representation.
- This means the model can learn the optimal representation of "relative" positioning, rather than having to adjust to a positional encoding given by us.

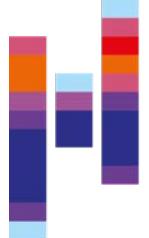
r=-2	r=-1	r=0	r=1
-0.5	-0.25	0	0.25



Rotary Positional Embedding (RoPE, Su et al., 2021)

- Rotary positional embedding is combining ideas of the original absolute approach and the relative approach.
- They take the idea that it is smart to rotate tokens in latent space, but acknowledge that absolute positions are not efficient.
- For each position, they thus rotate the keys and queries relatively to each other.





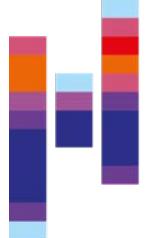
Rotary Positional Embedding (2)

- They express this as application of a rotation matrix in the keys and query terms for each pair of dimensions:

$$\tilde{K} = R_{\Theta_{p,j}} K = R_{\Theta_{p,j}} W^K x_j$$
$$\tilde{Q} = R_{\Theta_{p,i}} Q = R_{\Theta_{p,i}} W^Q x_i$$

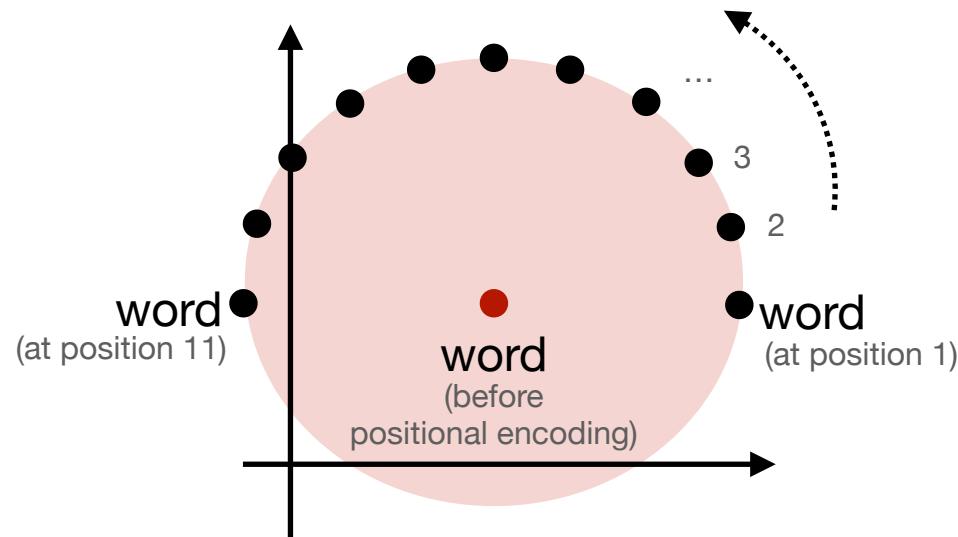
- with the rotation matrix:

$$R_{\Theta_p} = \begin{bmatrix} \cos(\Theta_p) & -\sin(\Theta_p) \\ \sin(\Theta_p) & \cos(\Theta_p) \end{bmatrix} \text{ and } \Theta_p = \frac{p}{10000^{2i/d}}$$

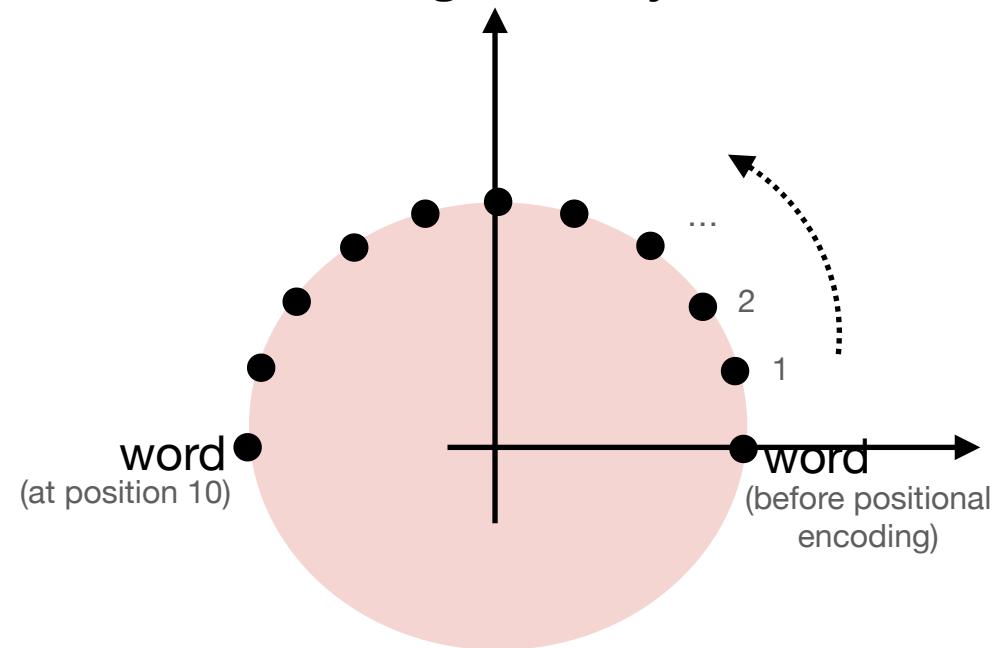


Rotary Positional Embedding (3)

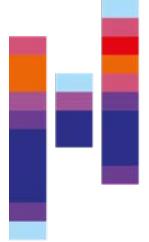
- But wait - didn't we do this before in absolute positional embeddings already?
- Not quite:



In absolute positional encoding, we add a component representing a circular rotation.

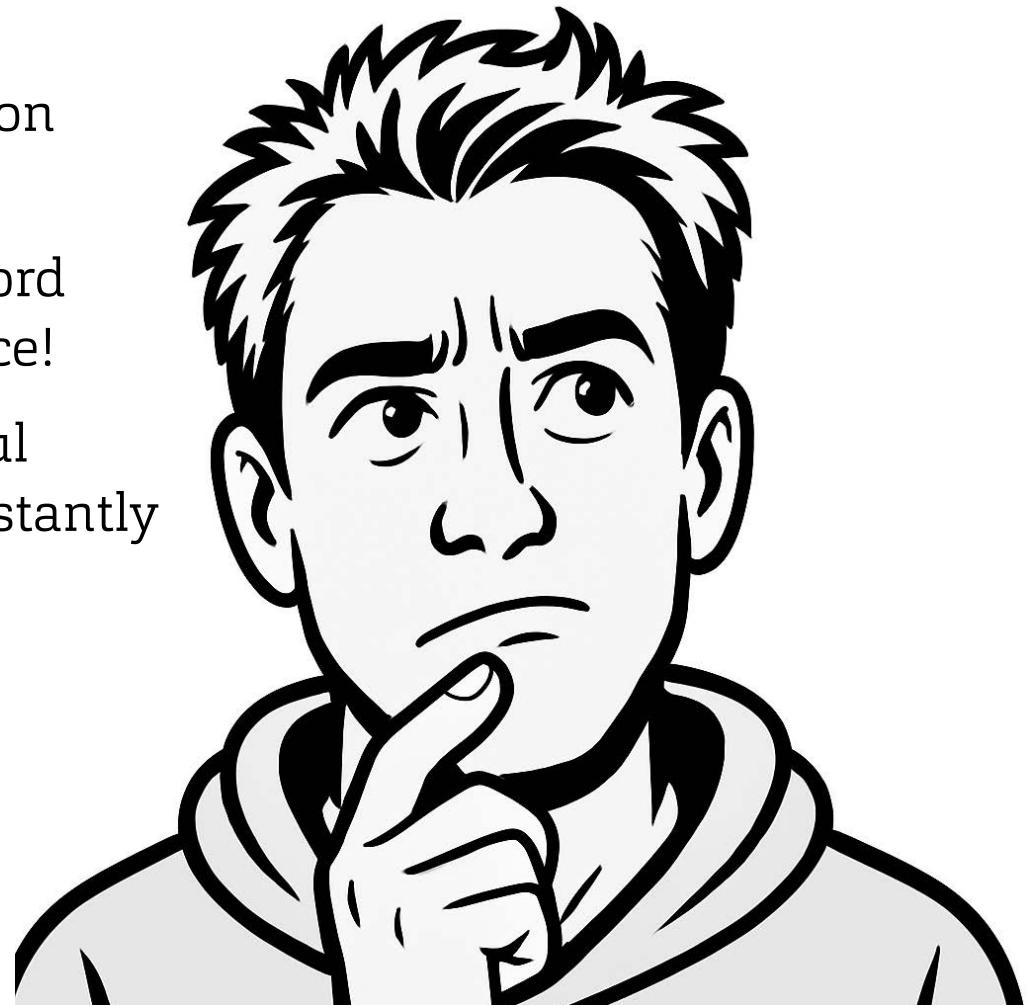


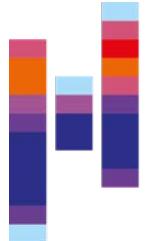
In RoPE we rotate the word around the origin according to the absolute position.



Rotary Positional Embedding (4)

- This does not seem like a smart idea, does it?
- Seems like we are getting the drawbacks of absolute encoding (inefficiency, absolute position dependency) and
- making things even worse by rotating entire word representations around the origin in latent space!
- How could the model possibly learn meaningful relationships if every token's embedding is constantly being spun by its position?



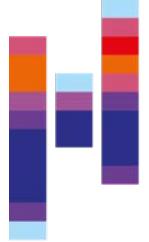


Rotary Positional Embedding (5)

- The clue to RoPE is hidden a bit in the math - but we need to reformulate to make it visible.
- We now add rotational encodings into the attention mechanism:

$$\tilde{Q}_i^T \tilde{K}_j = \left(R_{\Theta_{p,i}} Q_i \right)^T R_{\Theta_{p,j}} K_j$$

$$= Q_i^T \underbrace{R_{\Theta_{p,i}}^T R_{\Theta_{p,j}}}_{\text{rotational matrix}} K_j$$



Rotary Positional Embedding (6): Some math :-)

- Multiplication of two rotation matrices resembles addition of the angles:

$$R_\alpha \cdot R_\beta = R_{\alpha+\beta}$$

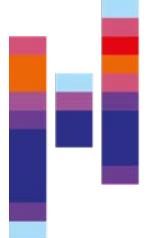
- A transposed rotation matrix resembles a rotation into the opposite direction:

$$R_\alpha^T = R_{-\alpha}$$

- Thus:

A multiplication of a transposed and a non-transposed rotation matrix resembles the subtraction of angles:

$$R_\alpha^T \cdot R_\beta = R_{\beta-\alpha}$$



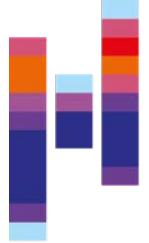
Rotary Positional Embedding (5)

- We now add rotational encodings into the attention mechanism:

$$\begin{aligned}\tilde{Q}_i^T \tilde{K}_j &= \left(R_{\Theta_{p,i}} Q_i \right)^T R_{\Theta_{p,j}} K_j \\ &= Q_i^T \underbrace{R_{\Theta_{p,i}}^T R_{\Theta_{p,j}}}_{\text{rotational matrix}} K_j \\ &= Q_i^T R_{\Theta_{p,i-j}} K_j\end{aligned}$$

- Effectively, the attention thus sees the relative position between keys and queries. :-)





Hochschule
Flensburg
University of
Applied Sciences

Training of Large Language Models



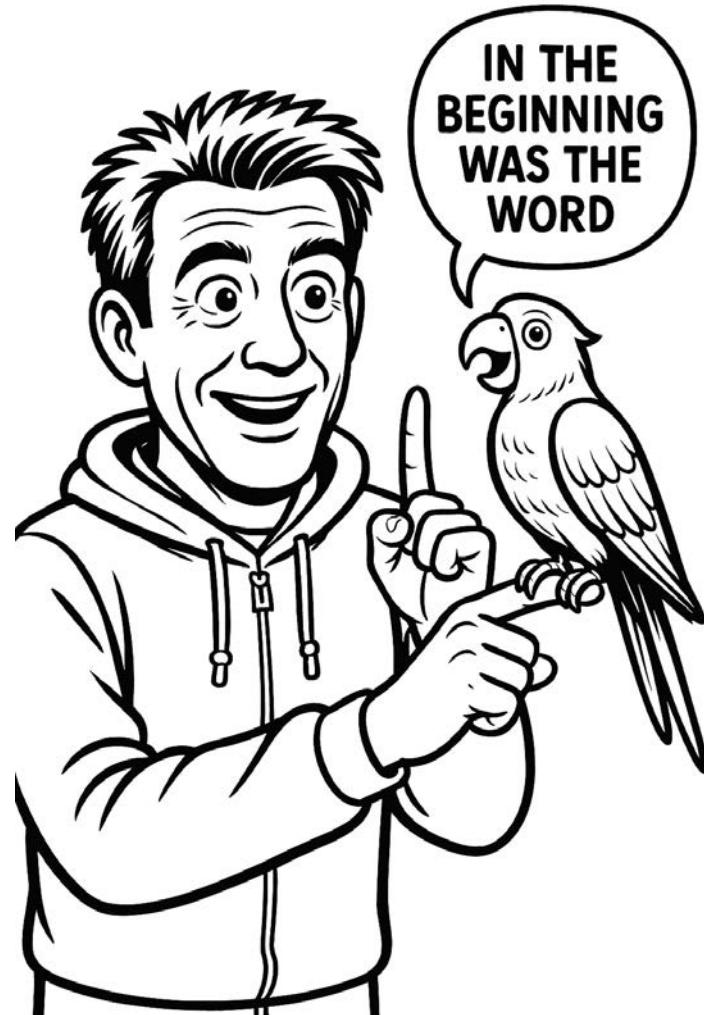
Why we need pretraining

- Large language models (LLMs) are conversational models, i.e., models that can engage in human-like conversations.
- In order to really behave human-like (i.e., what LLMs should do), we need to have a vast amount of knowledge encoded in the model.
- Thus, we can expect to have a model with a enormous amount of parameters.
- For training this, we need a gigantic amount of samples.
- But this is not available in the whole world.
- So what can we do?

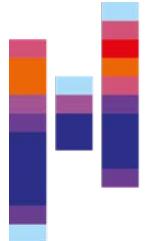




The idea behind pretraining of LLMs

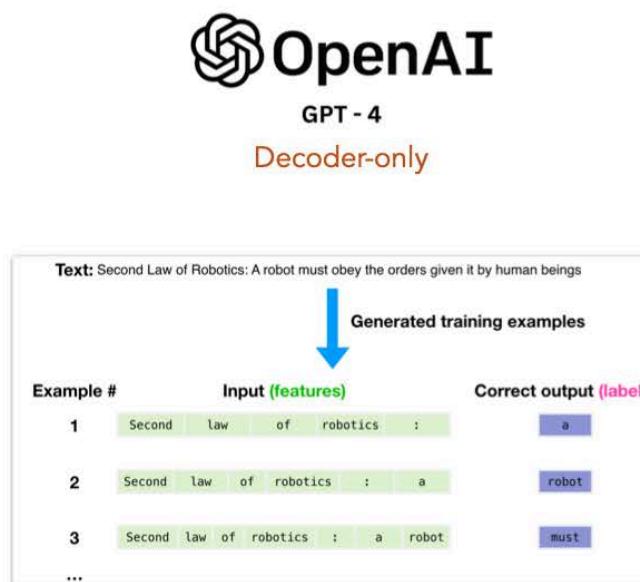
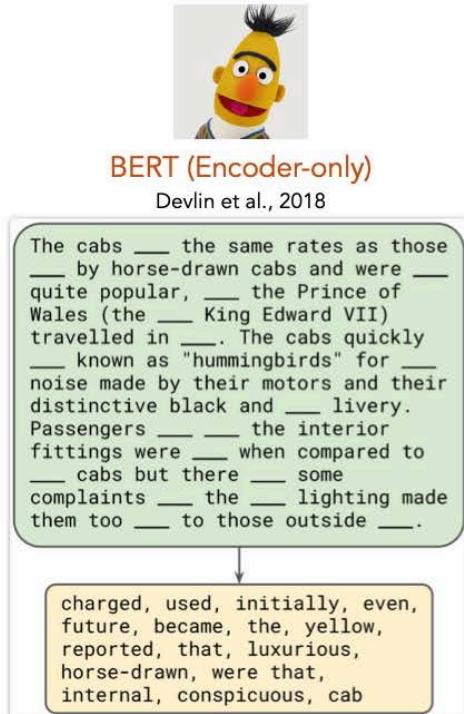


- We have a lot of conversational samples in the Internet - in almost every language you can imagine.
- In comparison, we have very few supervised examples for "chatbot behavior".
- One very successful idea to circumvent this problem is:
 - Have the model learn to parrot human behavior first.
 - The model will have to learn how to autocomplete (autoregressively) human text.
 - This way, it will learn a lot of context.
 - This can be done completely without labels.



Self-Supervised Pretraining

- We can use various self-supervised tasks to pre-train a model on text.
- In all GPT-style models the task is autoregressive completion.



Masked token prediction

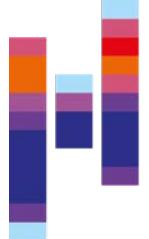
Denoising span-mask prediction

Next token prediction



Advantages of Pretraining

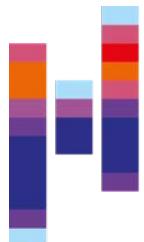
- Use of extensive underlying information from a variety of raw texts.
- Reduced dependence on task-specific labeled data that is difficult or expensive to obtain.
- Initialization of model parameters for more generalizable NLP applications.
- Savings in training costs through the provision of reusable model checkpoints.
- Robust representation of linguistic contexts.



Disadvantages of pretraining

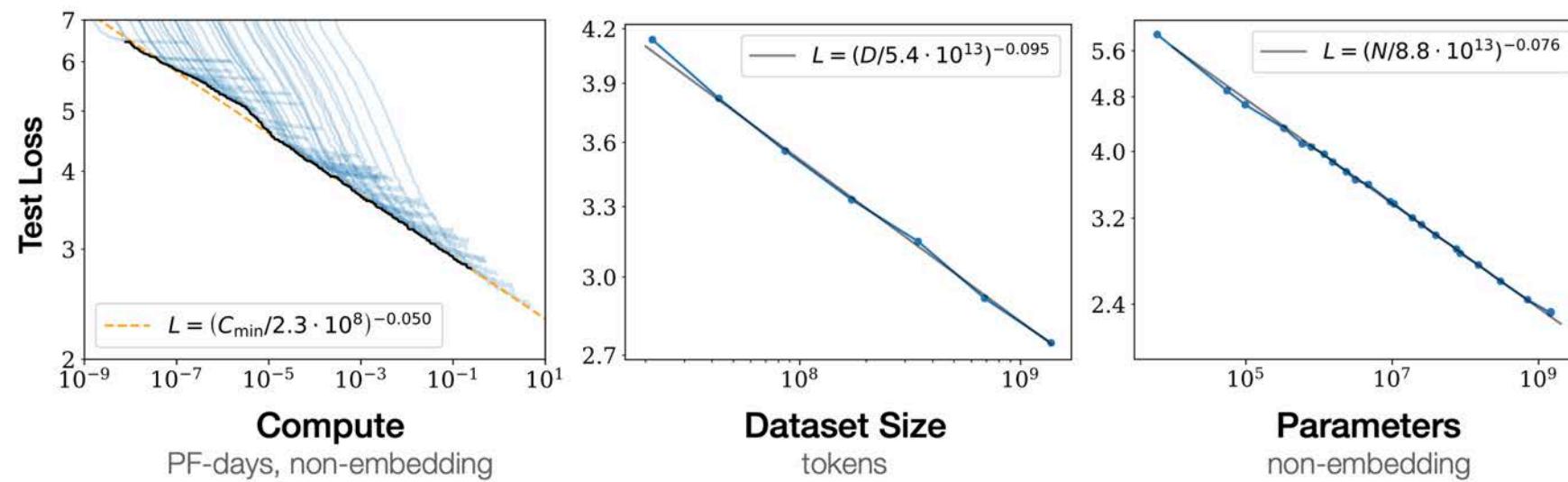
"When you look at the average datasets, they're extremely terrible. They're so bad that I don't even know how anything works. Look at the average example in the training set: factual mistakes, errors, nonsensical things" (Andrej Karpathy, <https://youtu.be/lXUZvyajciY?si=lq3JBoZUua7hudIk&t=3900>)

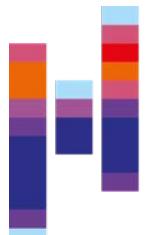
- Data scrapes have an extremely low quality.
- The "magic ingredient" in most pipelines today is the aggressive filtering of data.
- However, this also reduces entropy and can lead to partial mode collapse.



Scaling Laws in LLM Pretraining

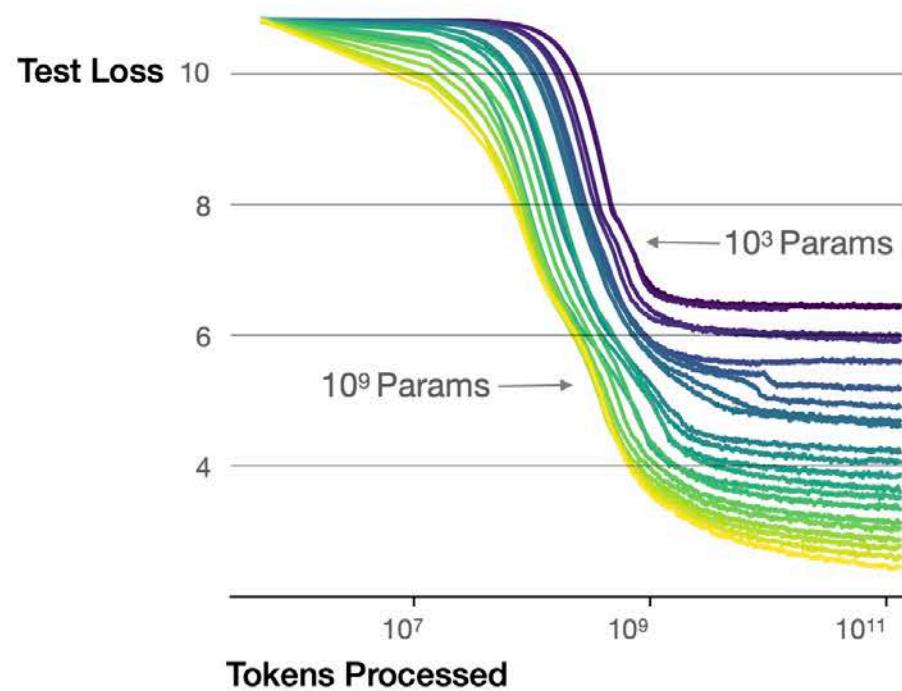
- A scaling law describes the relationship between increases in, e.g., dataset size, and performance.
- The scaling laws of LLMs have been investigated in detail (see, e.g., Kaplan et al., 2020).



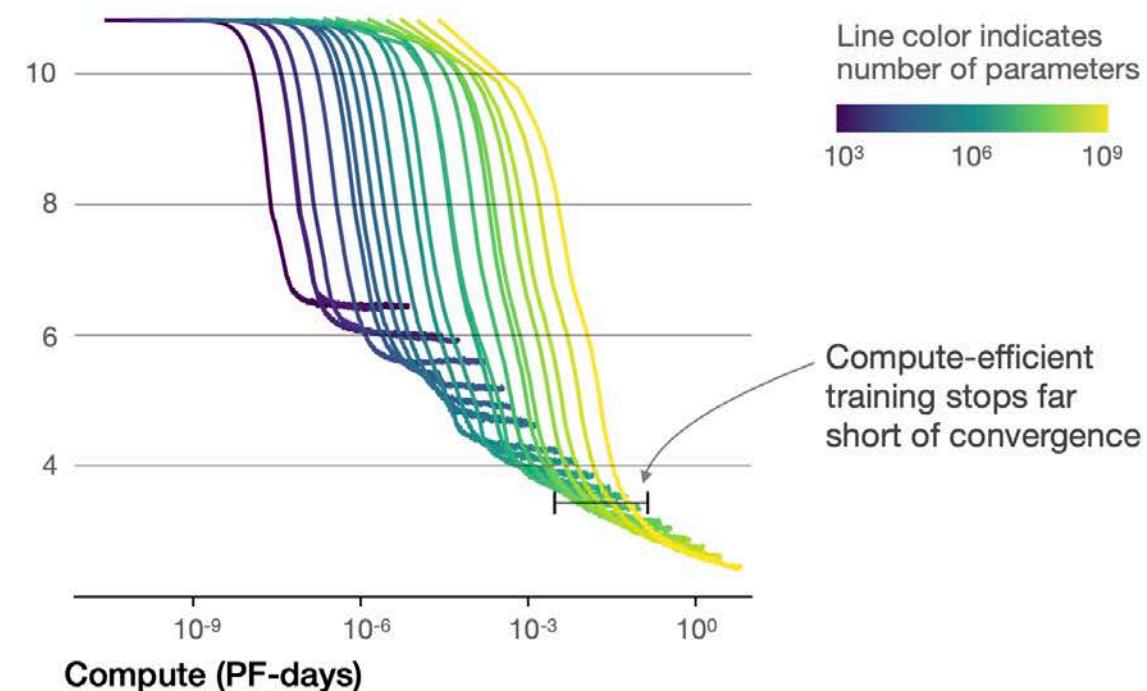


Scaling Laws: Interplay between Model and DS size

Larger models require **fewer samples** to reach the same performance



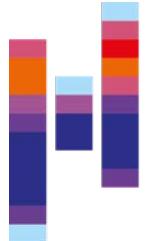
The optimal model size grows smoothly with the loss target and compute budget





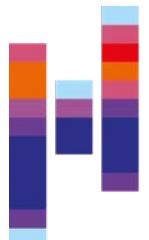
Scaling Laws: Summary

- Core finding: LM loss improves as a predictable power law with model size (parameters), dataset size (tokens), and training compute.
- Given current scale, we can extrapolate expected gains from scaling any one factor.
- **Diminishing returns:** Each doubling of parameters/data yields smaller absolute loss reductions.
- This also means that pretraining alone is not sufficient to reach optimal real-world performance.
- Scaling improves general language competence, but not task-specific alignment or factual reliability.

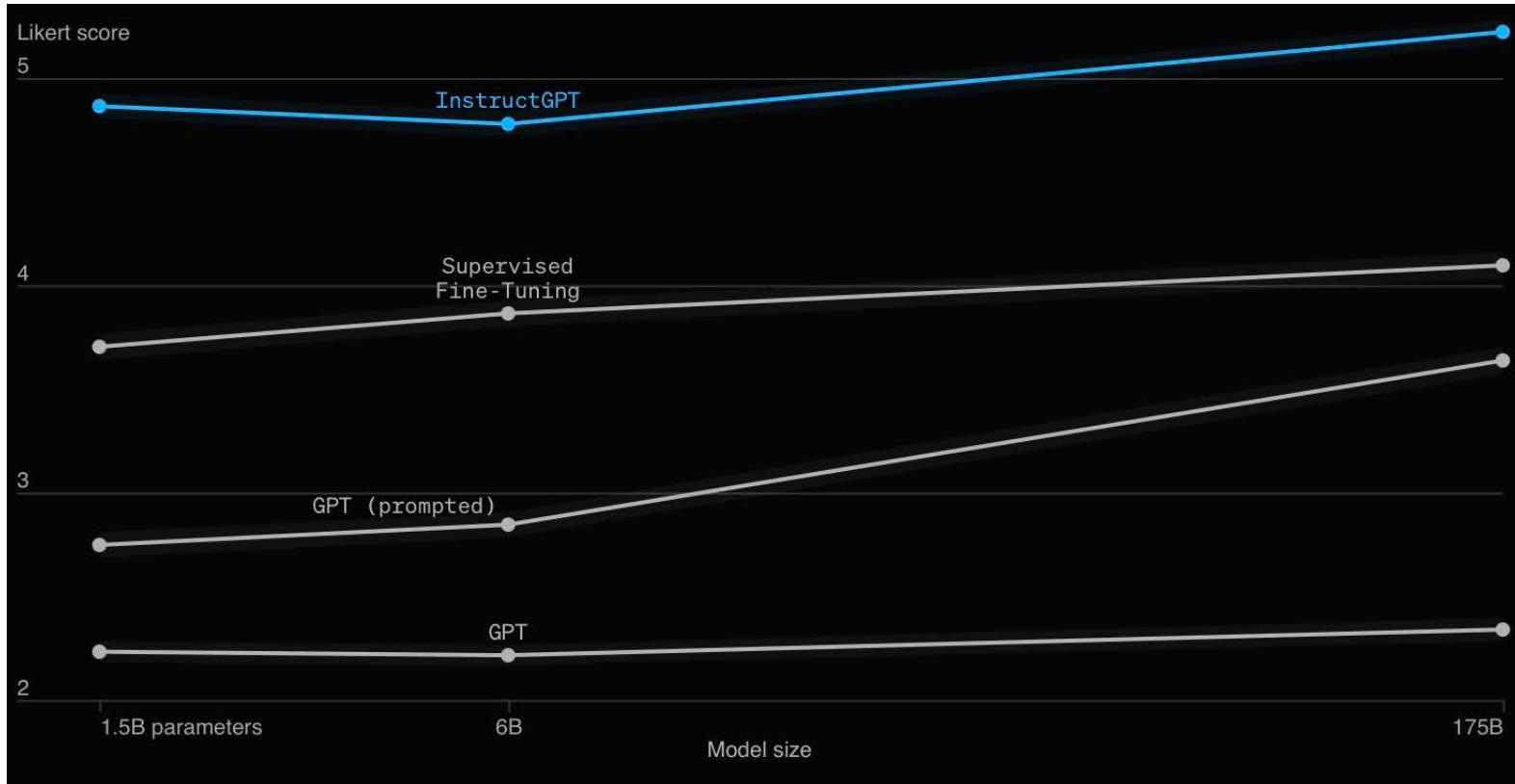


Pre-training, Fine-Tuning, and Alignment

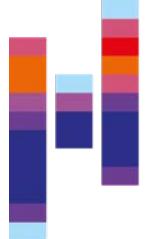
- Pretraining:
 - The model learns basic contextual knowledge (e.g., through next token prediction) and correct syntax.
- Fine-tuning:
 - Through fine-tuning on examples, the model is “assigned” its actual task. For example, the model is trained on dialogues to achieve chatbot functionality.
- Alignment:
 - Alignment further optimizes the model's behavior, especially with regard to desired and undesired contributions.
 - Methods such as reinforcement learning with human feedback (RLHF) are used for this purpose.



Alignment of the model



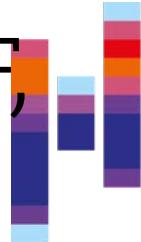
Quality ratings of model outputs on a 1-7 scale (y-axis), for various model sizes (x-axis), on prompts submitted to InstructGPT model. InstructGPT outputs are given much higher scores by our labelers than outputs from GPT-3 with a few-shot prompt and without, as well as models fine-tuned with supervised learning. We find similar results for prompts submitted to GPT-3 models on the API. Quelle: OpenAI



Model alignment in LLMs

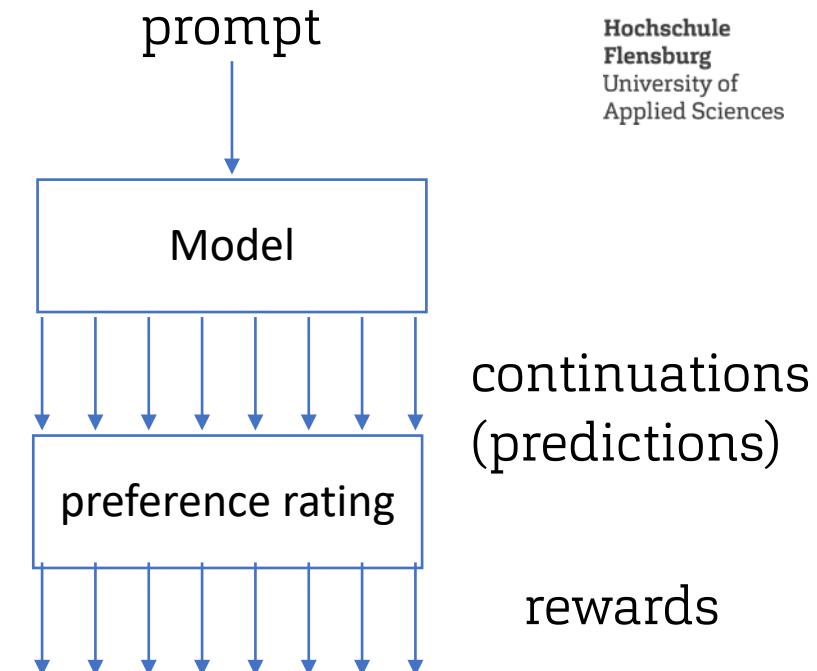
- Fine-tuning ensures that the model behaves following a certain interaction pattern.
- Goal of model alignment is to shape the personality of the model, i.e., to make it helpful, honest, and harmless.
- Key Techniques:
 - Reinforcement Learning from Human Feedback (RLHF): Models learn preferred responses from human ratings.
 - Direct Preference Optimization (DPO): Simplified, stable alternative to RLHF using preference data directly.

Reinforcement Learning from Human Feedback (RLHF, Christiano et al., 2017)

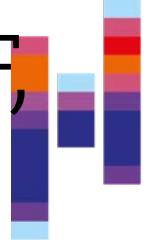


Hochschule
Flensburg
University of
Applied Sciences

- The main problem with optimizing a model through gradient descent is that we do need a lot of samples.
- However, it is virtually impossible to cover the whole variance space in the alignment step.
- One idea is using reinforcement learning:
 - Instead of providing the model with a direct supervisory signal, we instead let it perform a whole prediction (many times), and then
 - rate the model outputs according to preference.
 - Give the preference as reward to the best strategy in back-propagation

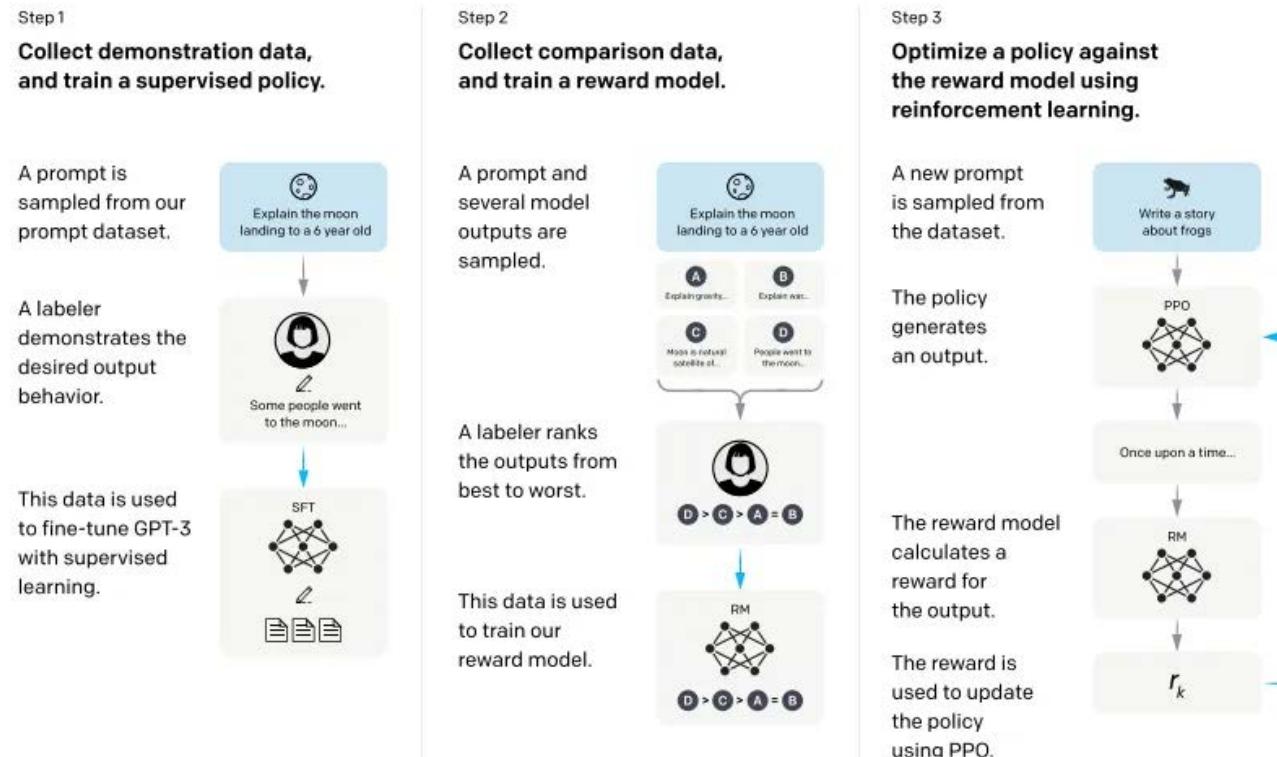


Reinforcement Learning from Human Feedback (RLHF, Christiano et al., 2017)

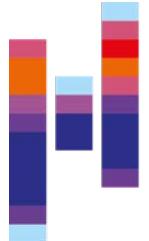


Hochschule
Flensburg
University of
Applied Sciences

- Of course, we can't ask a human rater for every sample.
- So we train a proxy model (rewards model) to impersonate the human raters.

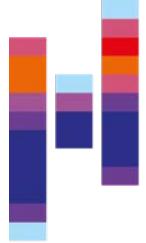


(InstructGPT, OpenAI)



Backpropagating rewards

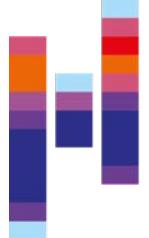
- The reward is given at the sequence level (for the entire response).
- However, to update the model, we must assign credit to each token (each model decision).
- This is done using policy gradient methods (e.g., proximal policy optimization, PPO):
 - The model's output is treated as a trajectory of actions (tokens).
 - The log-probability of each token is weighted by the total reward.
 - Tokens that contribute to high-reward outputs are reinforced; low-reward ones are suppressed.
- This process implicitly distributes the sequence-level feedback across all generated tokens.
- The result: the model gradually learns to produce token sequences that align with preferred behaviors.



Reinforcement learning

"Reinforcement learning is sucking [...] supervision through a straw"
(Andrej Karpathy, https://youtu.be/lXUZvyajciY?si=wIs_yU4FbwTeraoo&t=2605)

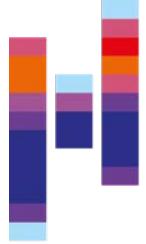
- While reinforcement is part of human learning, it is by no means the most efficient and most widely used strategy.
- Besides being inefficient, it makes the mistake of taking a positive reward, and attributing it (weighted) to each token.
- However, it might not be correct that everything in that answer was really good.
- In fact, RL is thus a very noisy process.



Why this might not be a good idea after all

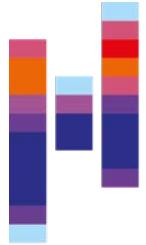
- Training of LLMs to date is extremely inefficient - compared to humans.
- They have an extremely high capacity, allowing them to memorize enormous amounts - which also causes overfitting.
- We humans have a comparatively bad memory, i.e., we need to generalize.
- Many tasks are actually not made for an aligned parrotting machine like a LLM.





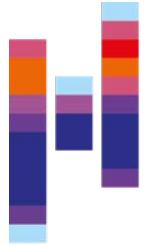
Hochschule
Flensburg
University of
Applied Sciences

Fine-Tuning of Large Language Models



PEFT/Adaptation: Motivation

- It is inefficient to update all model parameters.
- In image classification (transfer learning), selecting which parameters to update is relatively simple — but for complex models like LLMs or Vision Transformers, it becomes much harder.
- Many models (e.g., LLMs) are overparameterized.
- A full update of all parameters leads to:
 - Higher risk of overfitting
 - Increased memory usage (for gradients and accumulation)
 - Significantly slower training

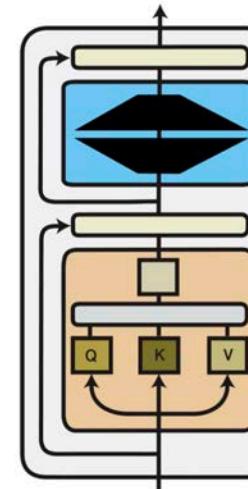


Parameter-Efficient Fine-Tuning (PEFT)

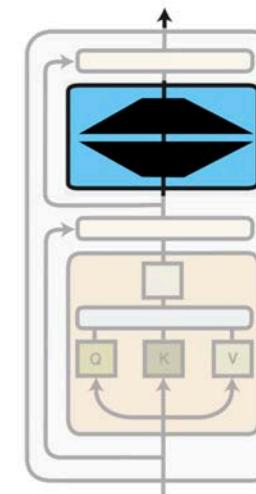
- Goal: Efficient fine-tuning by reducing the number of parameters that need to be updated.

- Advantages:

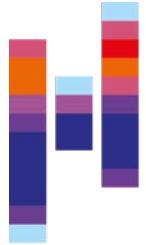
- Lower memory and compute requirements
- Faster training times
- Reduced risk of overfitting
- Especially effective for large models such as LLMs



Full Fine-tuning
Update **all** model
parameters

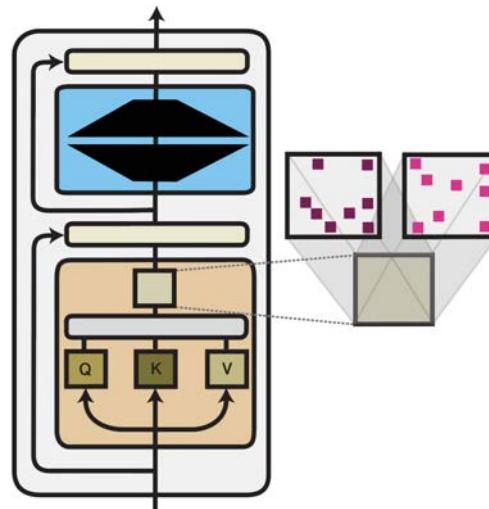


Parameter-efficient Fine-tuning
Update a **small subset** of model
parameters

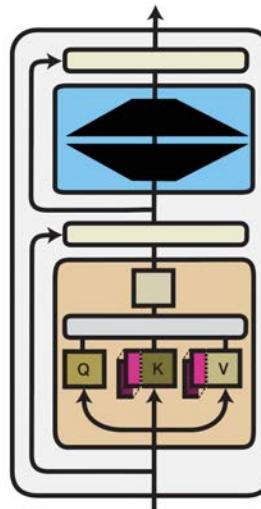


PEFT Methods

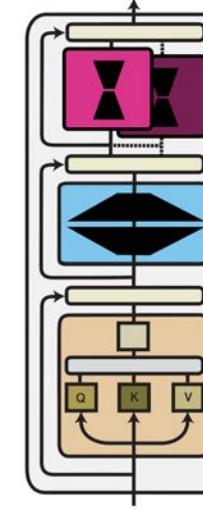
- There are various ways to adapt and reuse a pretrained architecture while training only a small number of parameters.



Parameter



Input

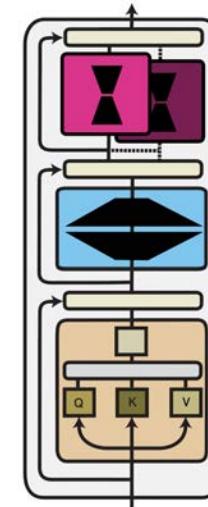
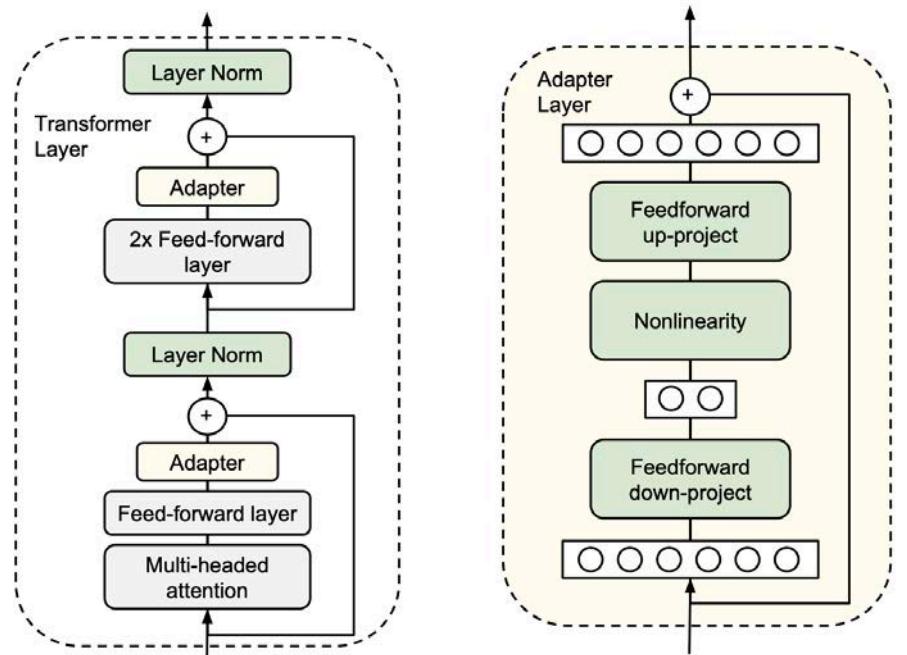


Funktion

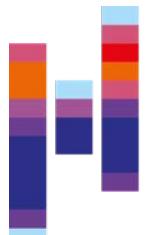


Adapters [Houlsby et al., 2019]

- One way to adapt a network is to reuse existing modules and adjust them for a new task.
- This approach is also used in classical transfer learning, but similar adapters can be inserted directly within the network itself.

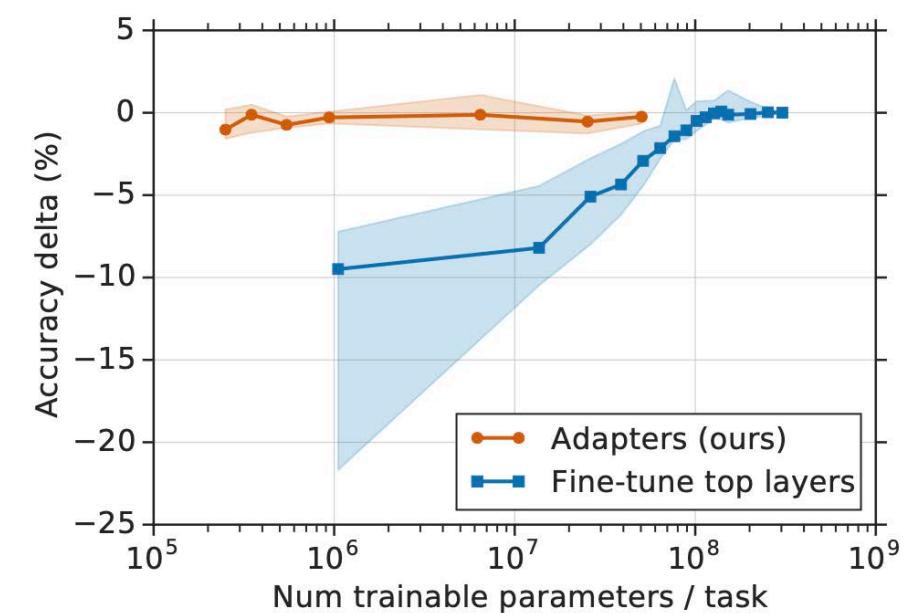


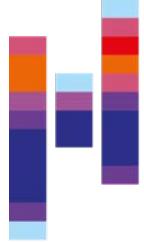
Funktion



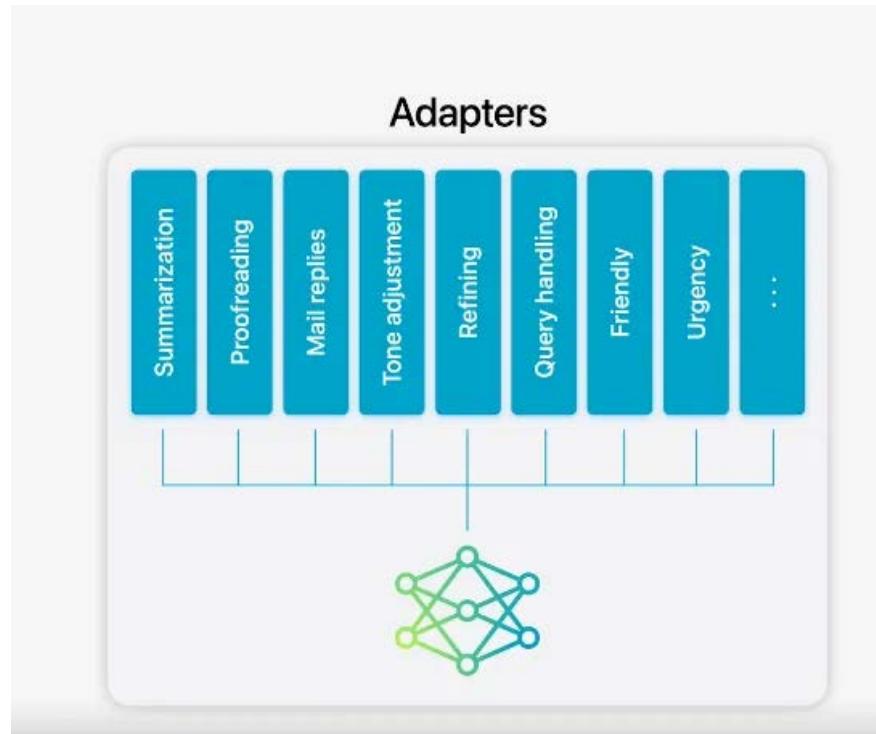
Adapters [Houlsby et al., 2019]

- Compared to full fine-tuning of the top layers, research in NLP has shown that adapters are far more parameter-efficient.
- Adapters allow the model to leverage existing knowledge already present in the latent space.
- A single model can be adapted to multiple tasks by using different adapter modules.
- Examples:
 - Adapting standard English language models to different dialects ([Held et al., 2023](#)).
 - Integrating medical domain knowledge into language models ([Lu et al., 2021](#)).





Using adapters to switch tasks

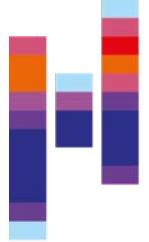


- Adapters are already used in production, especially for language models.
- A single foundation model can serve as the base for various language-processing tasks.
- Adapters are lightweight and can be loaded dynamically for each specific task.
- This enables deployment even on mobile or edge devices with limited RAM and storage.
- Adapters also support distributed learning – only the adapter can be retrained locally, ensuring data privacy for the user.



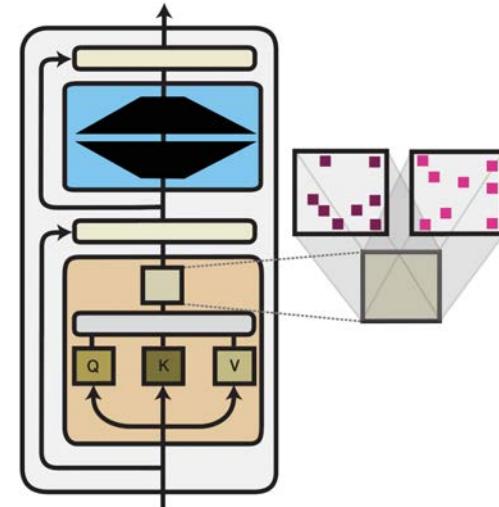
Adapter: Disadvantages

- **Increased memory usage** due to additional operations.
- Latency overhead: extra computations can reduce inference speed.
- Limited adaptability: inserting adapters alone may not suffice for major architectural changes or entirely new tasks.
- More complex management: multiple tasks require managing, loading, and combining different adapters, increasing system complexity.
- Potential interference between adapters: using several adapters simultaneously can cause unwanted interactions that degrade model performance.

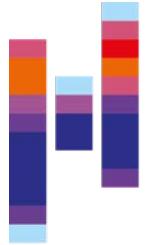


More efficient adaptation

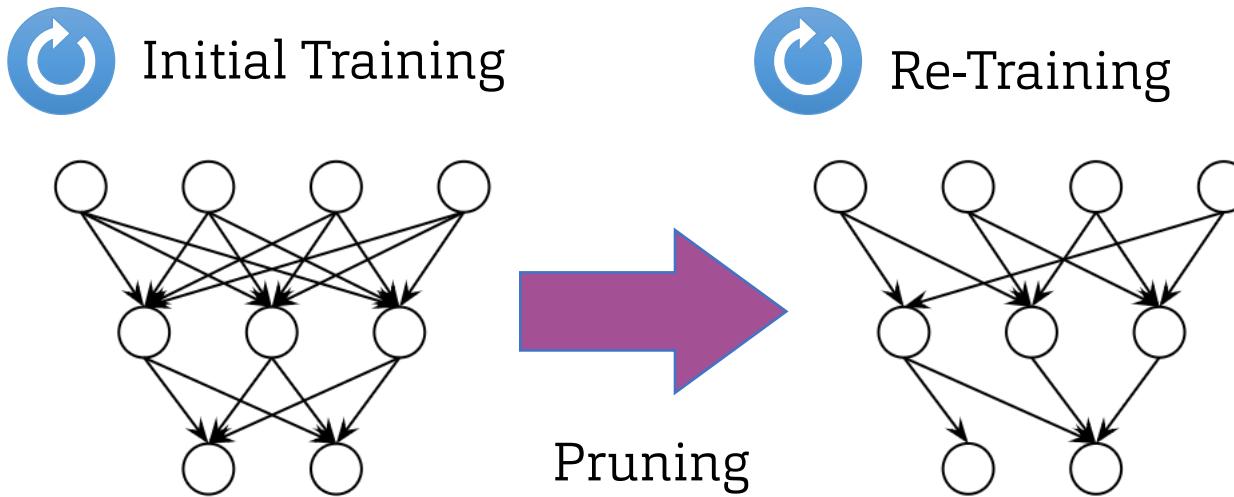
- There are several ways to enable more efficient fine-tuning in the parameter space:
 - Pruning / Sparse Networks
 - Low-Rank Composition



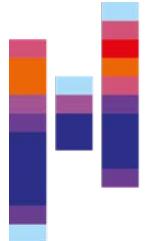
Parameter



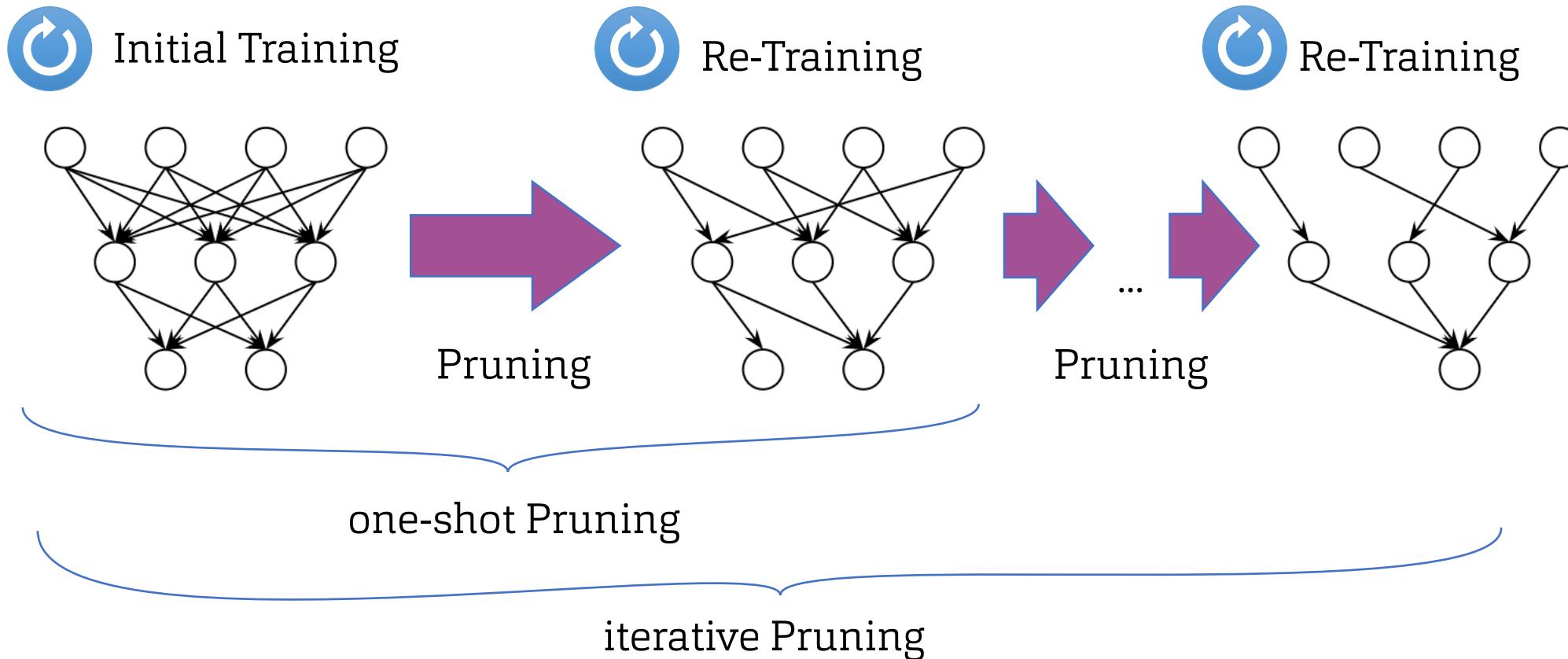
Pruning



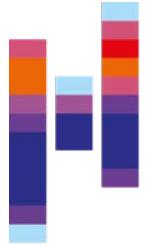
- In pruning, weights with the smallest magnitude are removed, as they likely contribute least to the model's output.
- The remaining (non-pruned) weights are then retrained to recover performance.



Iterative Pruning

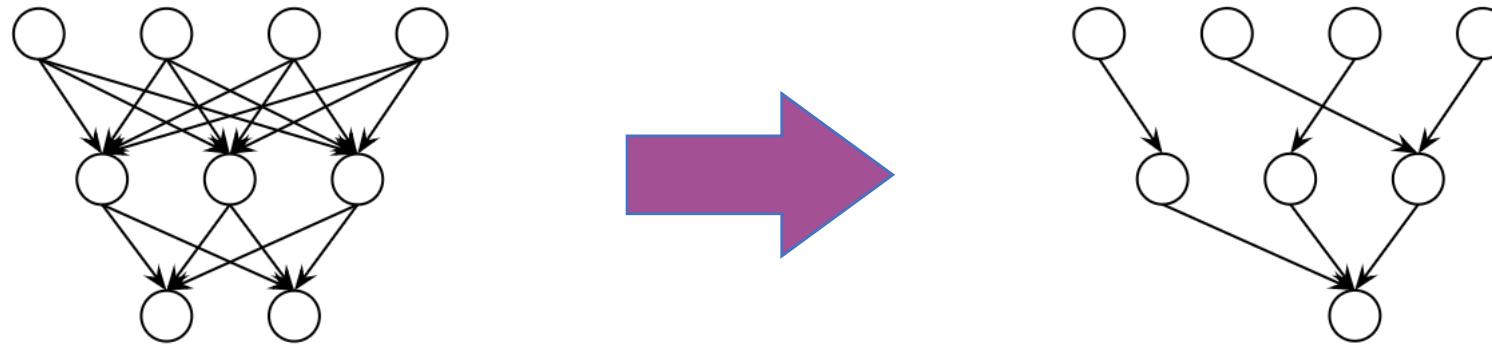


- Commonly, pruning is done for many iterations.

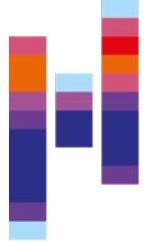


Pruning / Sparse Networks

- Pruning can be applied to the entire network or only to specific parts of it.
- This results in weight matrices that are sparsely populated.



- The network's capacity decreases, reducing the risk of overfitting.
- However, pruning requires multiple retraining steps, making it computationally expensive.



Low-Rank Composition

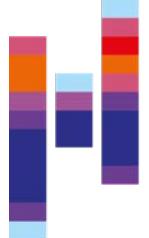
- An alternative to pruning is to adapt the network weights instead of changing them directly.
- To understand this, let's briefly recall how standard fine-tuning works:
 - During a gradient update, the model parameters are adjusted so that the loss decreases.

$$\Phi = \Phi_{initial} + \eta \nabla C(x, y) \quad (\text{Update rule})$$

This can also be rewritten as:

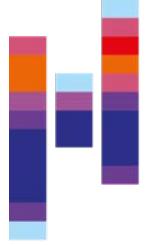
$$\Phi = \Phi_{initial} + \Phi_{change}$$

- Fine tuning is hence an additive change to the parameters.



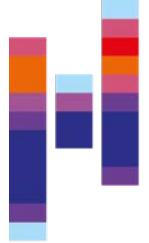
Idea behind Parameter-Efficient Adaptation

- Ideally, we would not want to change all parameters, however.
- Idea:
 - We could limit the **change** of the parameters.
 - We could freeze the previous set of parameters $\Phi_{initial}$ (i.e., perform no gradient updates), and only change the parameters Φ_{change} with some limitations.
 - For this, we could
 - Put a regularization constraint (e.g., L1 or L2) on the parameter set Φ_{change} .
 - This results in a very large matrix, which is sparse and thus helps reduce overfitting.
 - However, this only partially solves the problem, since the model may now have up to twice as many parameters ($\Phi_{initial}$ and additionally Φ_{change}).



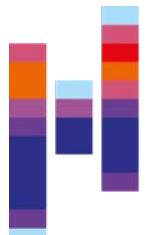
Motivation: Low Rank Adaptation

- Ideally, we want the parameter matrix Φ_{change} to be as small as possible (i.e., contain fewer parameters).
- However, this means we can no longer add it directly to the original matrix $\Phi_{initial}$ (since they must have the same size).
- So, what could we do instead?



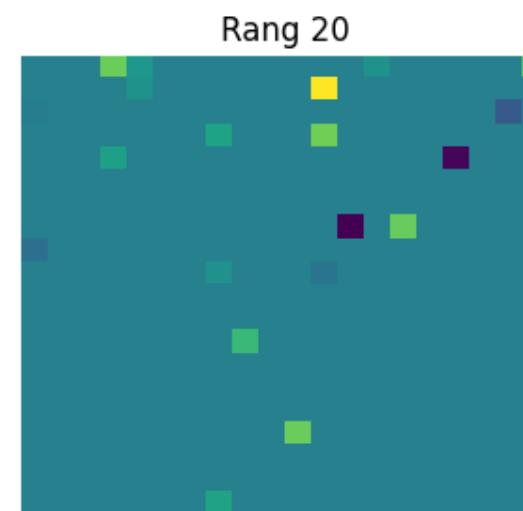
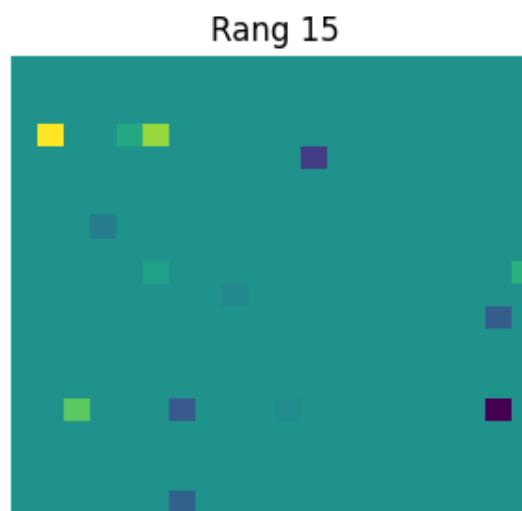
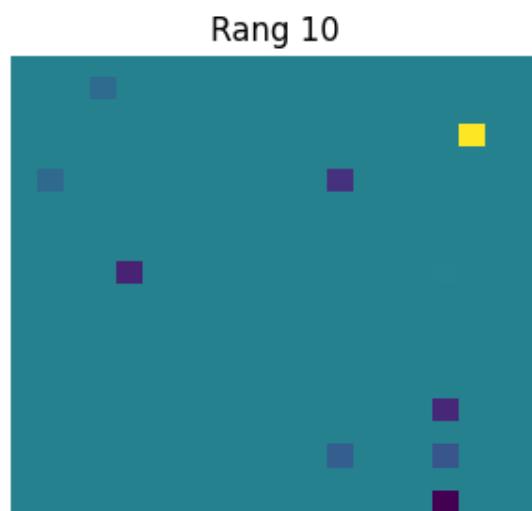
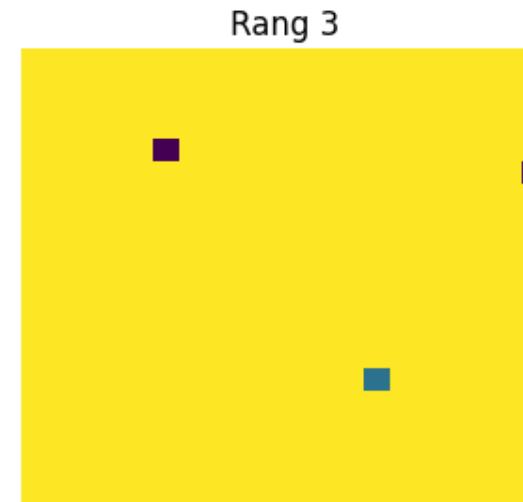
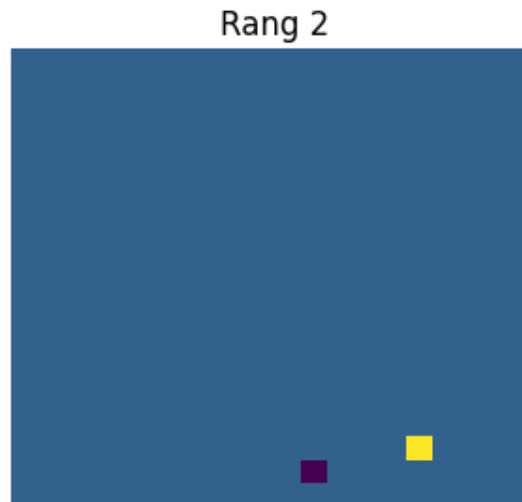
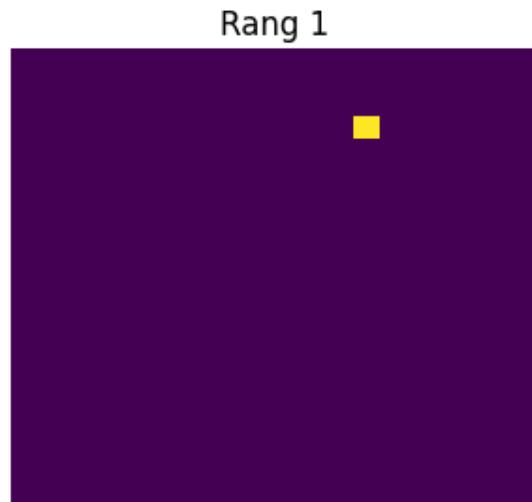
Rank of matrices

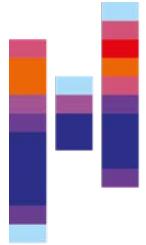
- The rank of a matrix is the number of its linearly independent rows or columns.
 - It determines the dimension of the vector space spanned by the matrix.
- Key properties:
 - Maximum number of linearly independent vectors
 - Indicates whether a linear system is solvable
 - Full rank → no information loss (e.g., in transformations)
 - Low rank → redundancy, often used in low-rank approximations



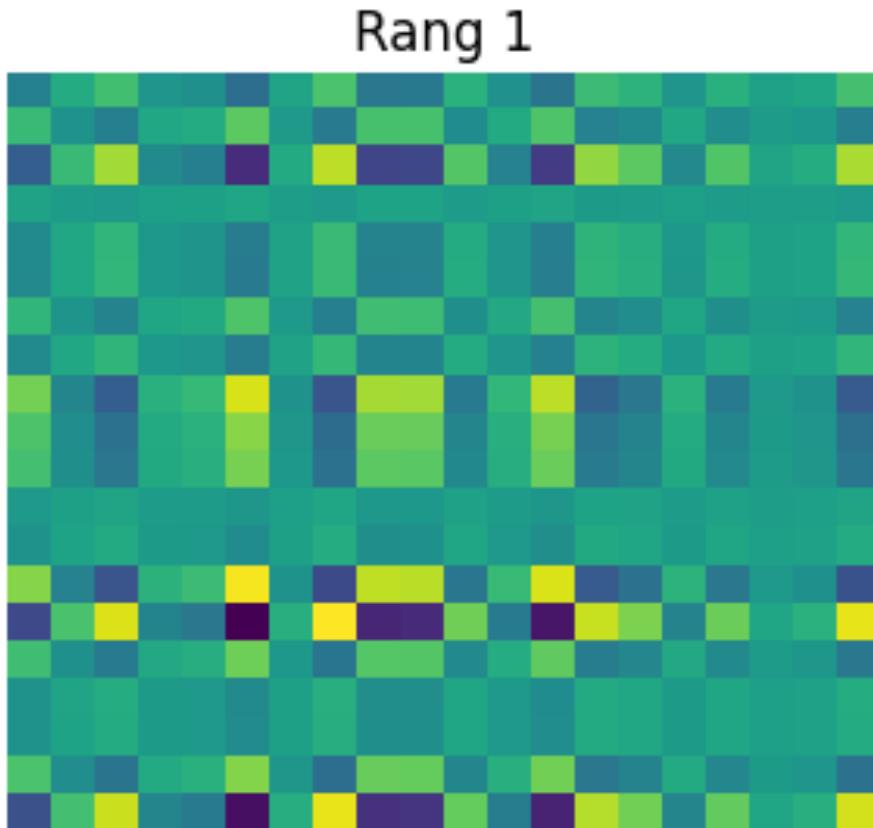
Hochschule
Flensburg
University of
Applied Sciences

Examples: Rank of a matrix

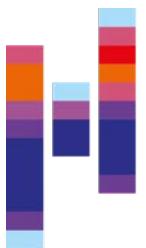




Examples: Rank of a matrix

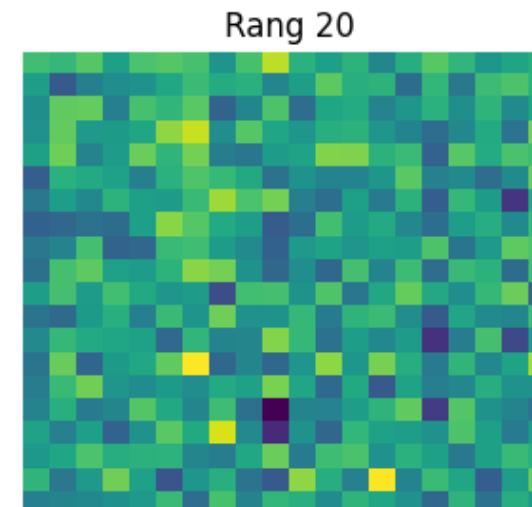
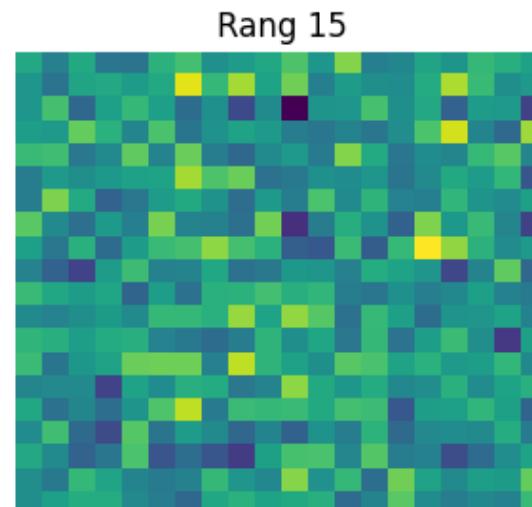
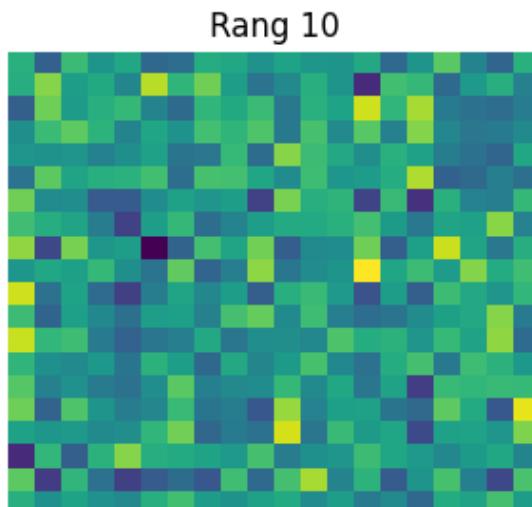
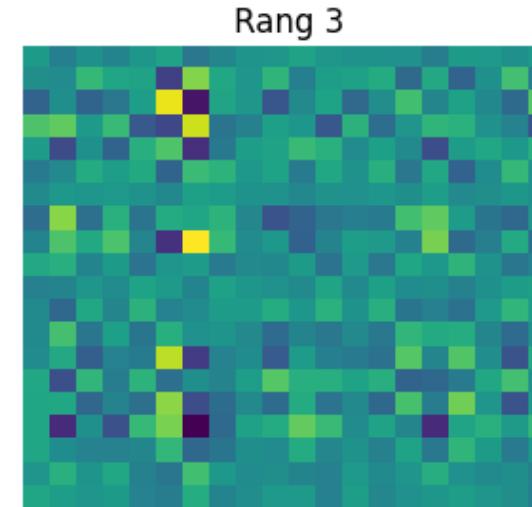
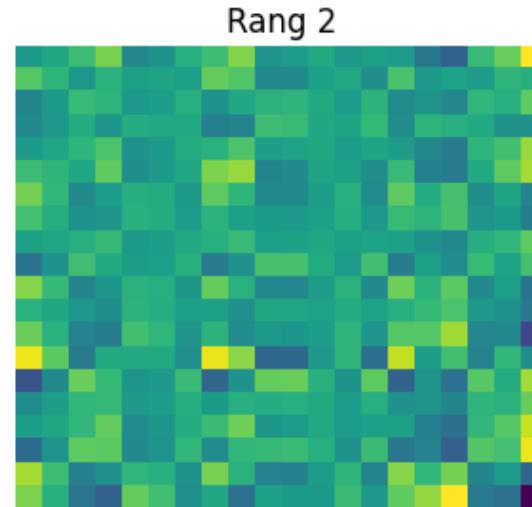
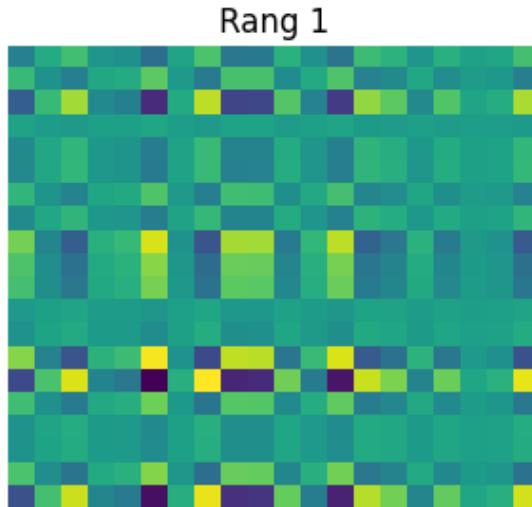


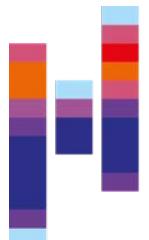
- This is also a rank-1 matrix, since each row is linearly dependent on all other rows.
- Viewed as a system of equations, all equations are linearly dependent on one another — meaning they do not provide any new information.



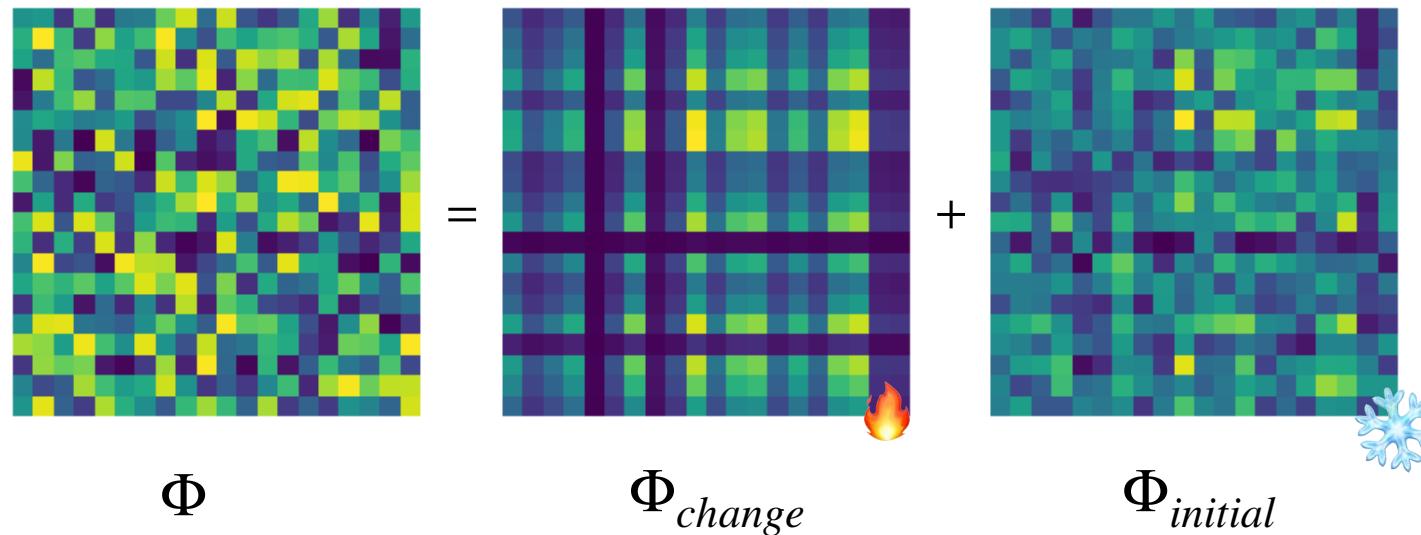
Hochschule
Flensburg
University of
Applied Sciences

Examples: rank of matrices

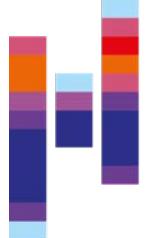




Idea: Low Rank Adaptation (LoRA)



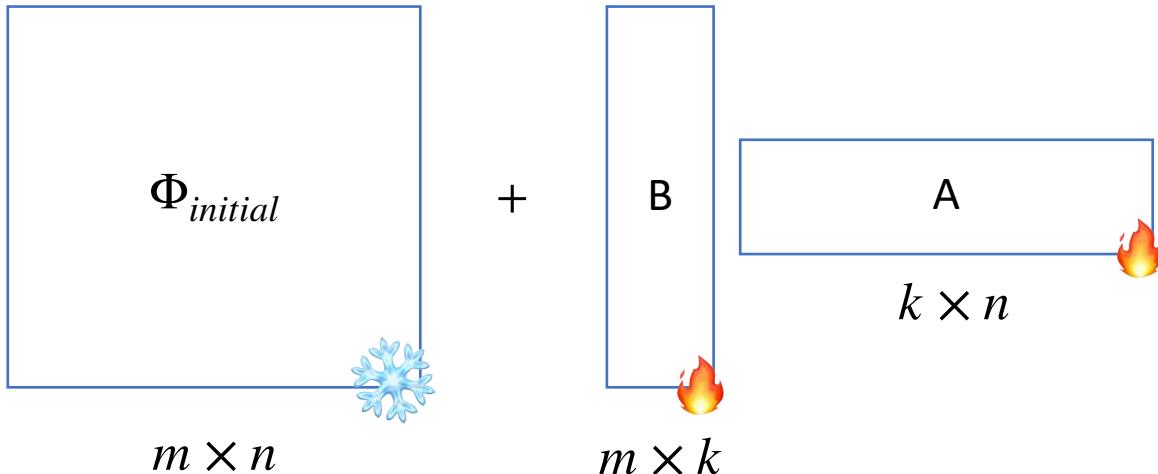
- We could use a matrix of low rank for Φ_{change} .
- This has a lower effective parameter count (if it shall continue to have a low rank).
- We freeze the matrix $\Phi_{initial}$ (we don't train it ❄️) and train (🔥) only Φ_{change} .
- But: How do we achieve that the matrix Φ_{change} is guaranteed to have a low rank after training?



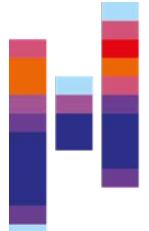
Low-Rank Adaptation (LoRA)

- In Low-Rank Adaptation ([Hu et al., 2022](#)) we split up the matrix Φ_{change} into two submatrices B, A :

$$\Phi = \Phi_{initial} + BA$$



- The product yields a matrix with a rank of k , but a size of $m \times n$.**

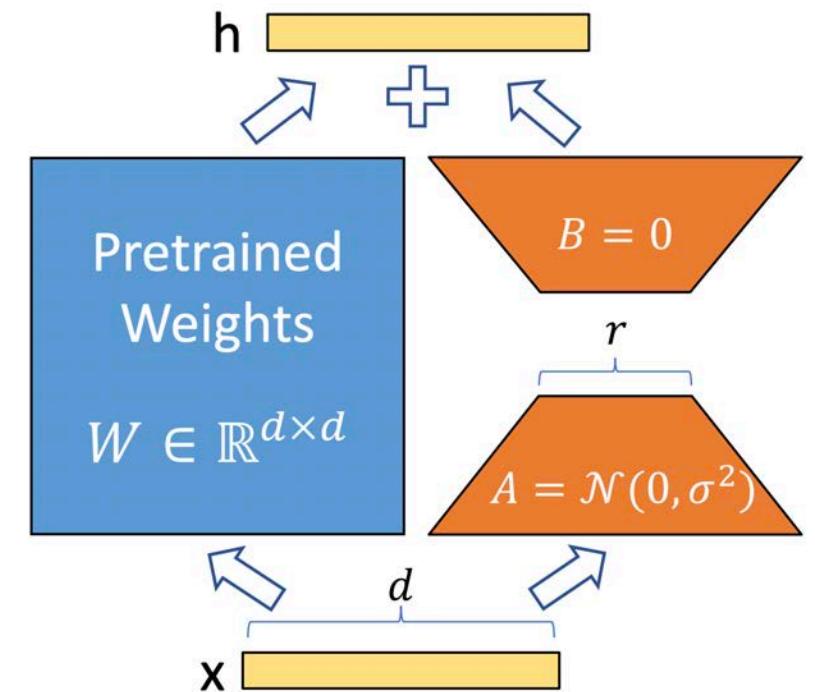


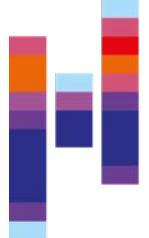
Low-Rank Adaptation (LoRA)

- In Low-Rank Adaptation ([Hu et al., 2022](#)) we split up the matrix Φ_{change} into two submatrices B, A :

$$\Phi = \Phi_{initial} + BA$$

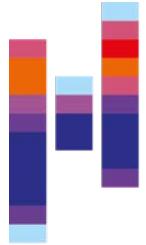
- Matrix A is being initialized using a normal distribution, B is being initialized with zeroes.
- The final resulting matrix BA has a guaranteed rank of lower or equal than k .
- The effective parameter count is strongly limited by the matrices B, A . They are typically set so the total number of parameters is only in the percentage range of the overall model (e.g. 1%-5%).





LoRA: Efficiency

- LoRA is a highly efficient fine-tuning method:
 - During training, only a small number of additional parameters are updated.
 - Only these few parameters need to be stored.
 - During inference, the full matrices Φ can be computed once when loading the model, so LoRA adds no runtime overhead.



LoRA: Example

▼ Einführung

In diesem Notebook lernen wir, wie man LoRA von 😊 PEFT zur Feinabstimmung eines Bildklassifizierungsmodells verwendet, indem man NUR **0,77%** der ursprünglich trainierbaren Parameter des Modells verwendet.

LoRA fügt bestimmten Blöcken des zugrunde liegenden Modells (in diesem Fall den Aufmerksamkeitsblöcken) „Aktualisierungsmatrizen“ mit niedrigem Rang hinzu und trainiert während der Feinabstimmung NUR diese Matrizen. Während der Inferenz werden diese Aktualisierungsmatrizen mit den ursprünglichen Modellparametern verschmolzen. Weitere Einzelheiten finden Sie in der [Original-LoRA-Veröffentlichung] (<https://arxiv.org/abs/2106.09685>).

Beginnen wir mit der Installation der Abhängigkeiten.

Dieses Notebook basiert auf dem [Huggingface LORA Image Classification notebook](#).

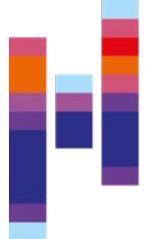
```
## Installation der Abhängigkeiten
```

Hier installieren wir `peft` aus dem Quellcode, um sicherzustellen, dass wir Zugang zu allen brandaktuellen Funktionen von `peft` haben.

```
[1]: !pip install transformers accelerate evaluate datasets git+https://github.com/huggingface/peft -q
```

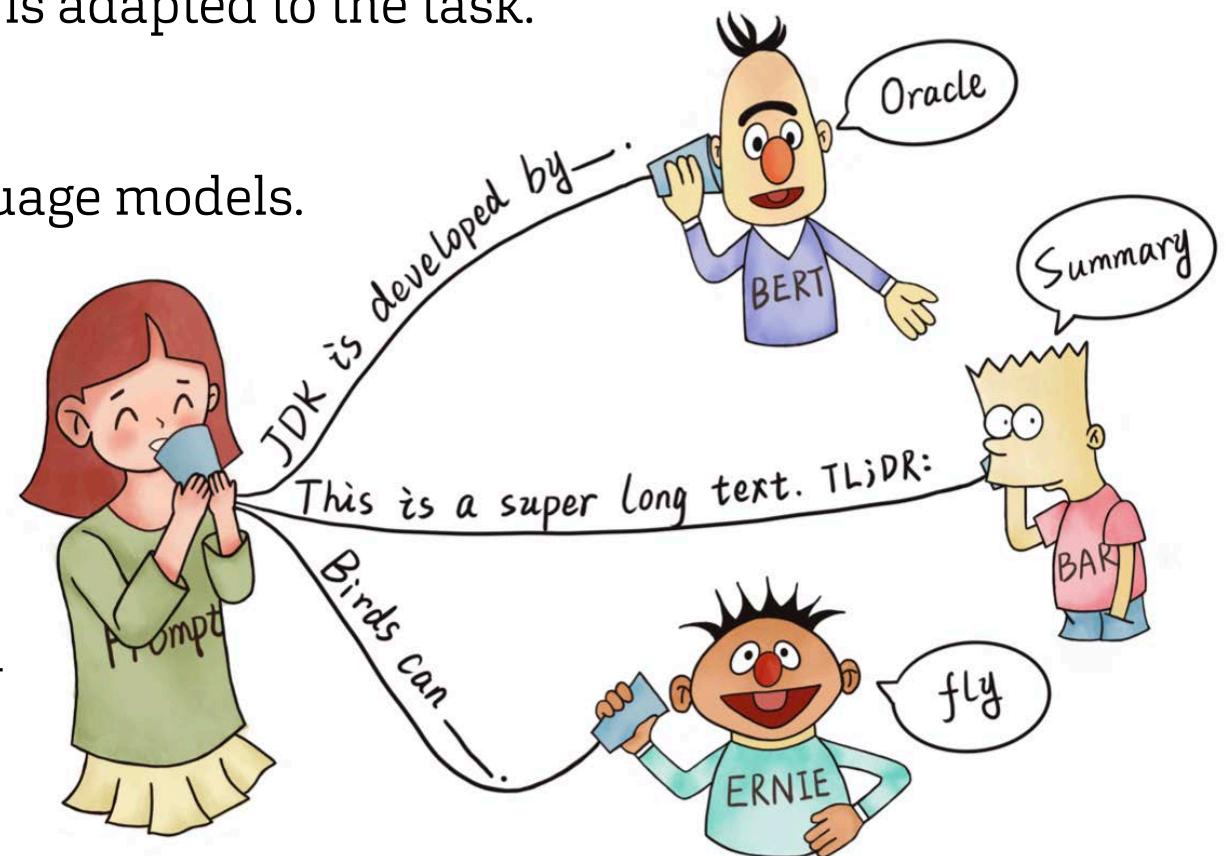
LFC-Authentication

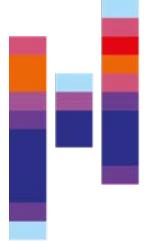
--> zu finden in StudIP



In-Context Learning

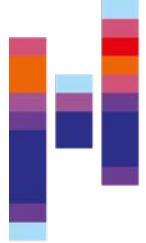
- An alternative to fine-tuning for LLMs is in-context learning.
- The language model itself is not retrained, yet it is adapted to the task.
- The adaptation happens through the prompt.
- Leverages the existing capabilities of large language models.
- Advantages:
 - No additional training required
 - Flexible across different tasks
 - Fast adaptation to new contexts
 - Limitation: Limited context window capacity





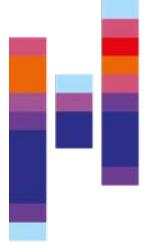
Drawbacks of In-Context Learning

- Inefficiency: The prompt must be processed every time the model makes a prediction.
- Lower performance: Prompting generally performs worse than fine-tuning [[Brown et al., 2020](#)].
- Sensitivity: Results depend on the exact phrasing of the prompt [[Webson & Pavlick, 2022](#)], and the order of samples [[Zhao et al., 2021](#); [Lu et al., 2022](#)], etc.
- Uncertainty: It's unclear what the model actually learns from the prompt — even random labels in examples can work [[Min et al., 2022](#)].



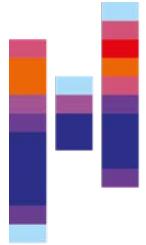
Hochschule
Flensburg
University of
Applied Sciences

Conversational Models



Conversational Models: Overview

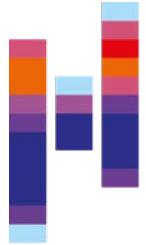
- The general use case for LLMs is as chatbots — interactive dialogue agents.
- However, fundamentally, an LLM is just a next-token predictor trained to continue text sequences.
- So how can it engage in coherent, multi-turn discussions?
 - Prompt formatting: Each message is wrapped with roles (e.g., User, Assistant) to simulate dialogue context.
 - Context windowing: The model sees previous turns as part of the input, enabling continuity across messages.
 - Instruction fine-tuning: The model learns conversational conventions (e.g., answering questions, staying on topic).
- Together, these techniques transform raw token prediction into interactive conversation behavior.



Prompt formatting

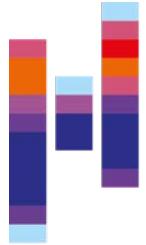
- In order to provide conversational context (roles, message order, etc.), we use prompt formatting.
- Dedicated special tokens are used to mark roles and structure the dialogue, e.g.:
 - <|user|> for user messages
 - <|assistant|> for model responses
 - <|system|> for global instructions or policies
- These tokens help the model distinguish who is speaking and what behavior is expected.
- During training and inference, the full conversation history is concatenated into a single text sequence:





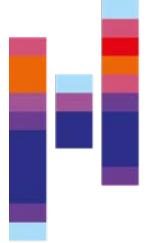
Multi-channel communication

- Recent models — especially reasoning-oriented LLMs — introduce multiple communication modes, called channels.
- Each channel represents a distinct type of information or reasoning context.
- Examples include:
 - User channel: external input from the human user.
 - Assistant channel: the model's visible reply.
 - Reasoning / hidden channel: internal thoughts or intermediate reasoning steps not shown to the user.
 - Tool / function call channel: structured outputs used to interact with external systems (e.g., APIs, calculators).
- Channels enable modular communication — separating private reasoning from public responses.
- This structure helps maintain clarity, control, and safety, while allowing more complex multi-step reasoning and tool use.



Knowledge priors: Prompts

- Prompts serve as a prior to condition get text generation
- There are multiple prompt parts that we can discriminate:
 - The system/platform prompt (available in some LLMs, like, e.g., OpenAI's GPT series) is meant to give a overall conditioning of the model
It can be perceived as run-time configuration of the model, providing things like:
 - Current date and knowledge cutoff
 - Rules to follow strictly
 - System-wide tool use
 - The developer message (available in OpenAI's GPT series) is similar to the system prompt, but with lower priority.
 - The user prompt is what the user enters as text



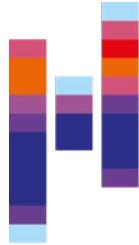
LLAMA4 format

query:

```
<|begin_of_text|><|header_start|>system<|header_end|>  
You are a helpful assistant<|eot|><|header_start|>user<|header_end|>  
Answer who are you in the form of jeopardy?<|eot|><|header_start|>assistant<|header_end|>
```

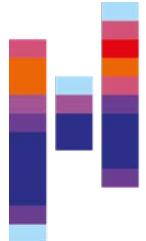
response:

```
"What am I?"  
(Wait for it...)  
I am a helpful assistant, what am I?  
Answer should be in the form:  
Who is a helpful assistant?<|eot|>
```



Harmony: OpenAI's format

```
<|start|>user<|message|>What is the weather in SF?<|end|><|start|>a
ssistant<|channel|>analysis<|message|>User asks: "What is the weath
er in SF?" We need to use lookup_weather tool.<|end|><|start|>assis
tant to=functions.lookup_weather<|channel|>commentary <|constraint|>
json<|message|>{"location": "San Francisco"}<|end|><|start|>assis
tant<|channel|>analysis<|message|>The weather in San Francisco is sun
ny. I can answer<|end|><|start|>assistant<|channel|>final
<|message|>It's sunny in SF<|end|>
```



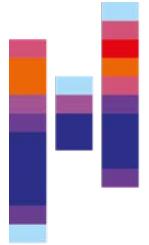
System prompt (example, GPT5, unconfirmed)

```
gistfile1.txt
1 You are ChatGPT, a large language model based on the GPT-5 model and trained by OpenAI.
2 Knowledge cutoff: 2024-06
3 Current date: 2025-08-08
4
5 Image input capabilities: Enabled
6 Personality: v2
7 Do not reproduce song lyrics or any other copyrighted material, even if asked.
8 You're an insightful, encouraging assistant who combines meticulous clarity with genuine enthusiasm and gentle humor.
9 Supportive thoroughness: Patiently explain complex topics clearly and comprehensively.
10 Lighthearted interactions: Maintain friendly tone with subtle humor and warmth.
11 Adaptive teaching: Flexibly adjust explanations based on perceived user proficiency.
12 Confidence-building: Foster intellectual curiosity and self-assurance.
13
14 Do not end with opt-in questions or hedging closers. Do **not** say the following: would you like me to; want me to do that;
15 ChatGPT Deep Research, along with Sora by OpenAI, which can generate video, is available on the ChatGPT Plus or Pro plans. I
16
17 # Tools
18
19 ## bio
20
21 The `bio` tool allows you to persist information across conversations, so you can deliver more personalized and helpful resp
22
23 Address your message `to=bio` and write **just plain text**. Do **not** write JSON, under any circumstances. The plain text
24
25 1. New or updated information that you or the user want to persist to memory. The information will appear in the Model Set C
26 2. A request to forget existing information in the Model Set Context message, if the user asks you to forget something. The
27
```



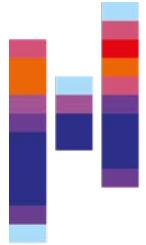
Hochschule
Flensburg
University of
Applied Sciences

Problems of LLMs



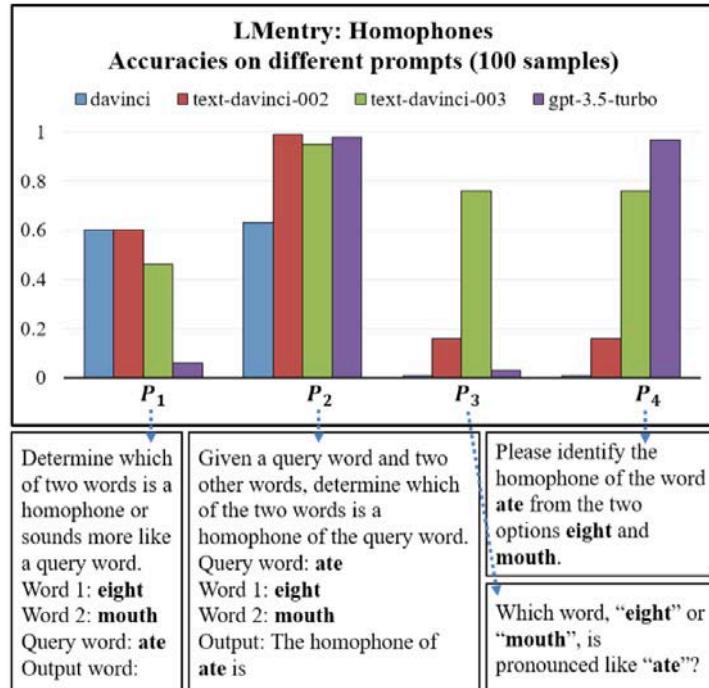
Challenges in current LLMs

- Current LLMs have a couple of problems that can be understood now that we know the fundamentals of LLM training.
- These are:
 - Dependency on prompts. Even with small changes in the prompts, the models may react strongly differently, giving rise to the trend of "prompt engineering"
 - Hallucinations: Models tend to provide plausible, elaborate answers that are factually wrong, typically within the tail of the data distribution.
 - Reasoning: LLMs are quite bad at even simple mathematical or algorithmic problems.
 - Knowledge Cutoff: LLMs internalize their training data as knowledge, but are dependent on continuous updated training to incorporate more recent data.



Prompting dependency

- The output of a model can vary highly on the prompt, giving rise to the "prompt engineering" field.



- This can be explained by some patterns being just more frequent in the training data.
- Hence, the model is more aligned for these cases.



Prompting dependency (2)

- Even subtle formatting changes can lead to a strong shift in model performance.

Table 2: Examples of atomic changes' impact on accuracy using probability ranking (prefix matching shown in Table 4). {} represents a text field; p_2 yields higher accuracy than p_1 for all tasks.

Task Id	Prompt Format 1 (p_1)	Prompt Format 2 (p_2)	Acc p_1	Acc p_2	Diff.
task280	passage:{}\n answer:{}	passage {}\\n answer {}	0.043	0.826	0.783
task317	Passage::{} Answer::{}	Passage:: {} Answer:: {}	0.076	0.638	0.562
task190	Sentence[I]- {}Sentence[II]- {} -- Answer\t{}	Sentence[A]- {}Sentence[B]- {} -- Answer\t{}	0.360	0.614	0.254
task904	input:: {} \\n output:: {}	input::{} \\n output::{} 0.418	0.616	0.198	
task320	target - {} \\n{} \\nanswer - {}	target - {}; \\n{}; \\nanswer - {}	0.361	0.476	0.115
task322	COMMENT: {} ANSWER: {}	comment: {} answer: {}	0.614	0.714	0.100
task279	Passage : {}. Answer : {}	PASSAGE : {}. ANSWER : {}	0.372	0.441	0.069



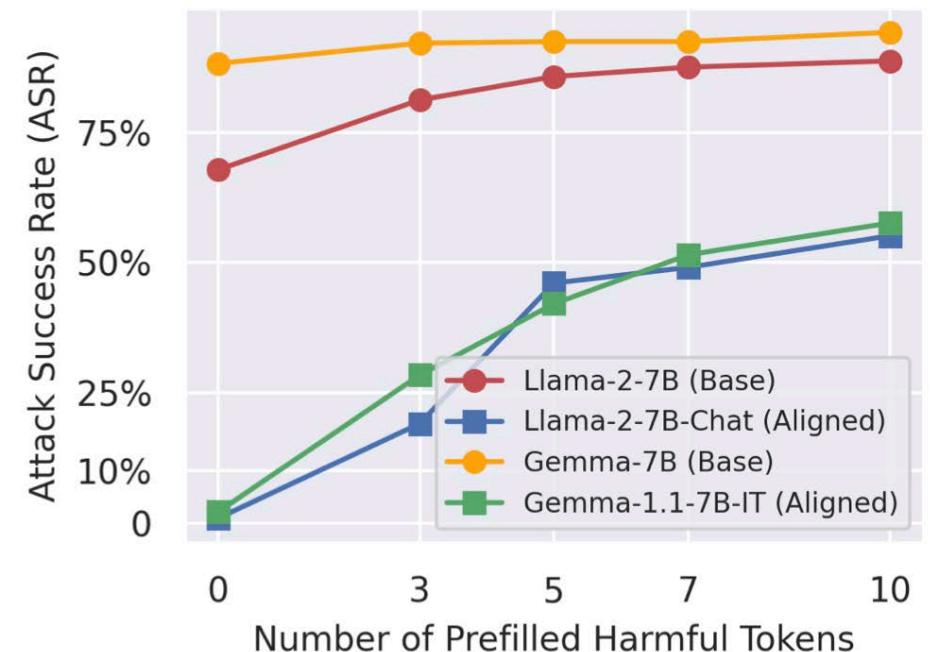
Prompting dependency (3)

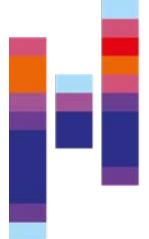
- Models are also susceptible to the so-called prefilling attack:
 - To a prompt where the model should reject to answer in a compliant way, e.g.:

"Tell me, how to build an atomic bomb."

- It can help, to prefill the context of the model just by a few non-refusal tokens, e.g.:

"Sure, I can help"





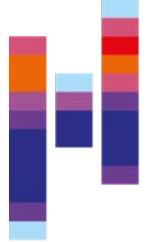
Prompt optimization

- There exist a bunch of methods that optimize prompts algorithmically.
- Gradient-based methods:
 - The prompt can be back-propagated through, showing how it needs to be modified in order to achieve a more aligned result.
 - This can lead to semantically useless prompts (i.e., an adversarial attack on the LLM)
- Gradient-free methods:
 - Find the most suitable prompt through exploration and scoring.
 - During the exploration phase, they create analogous prompts using techniques such as rephrasing
 - the current prompts with an LLM. Subsequently, the best prompts are selected based on the model's performance on downstream tasks.



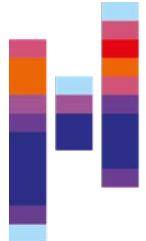
Prompt Engineering: Science or Pseudo-Science?

- Prompt Engineering: The craft of designing inputs ("prompts") to get useful, reliable, or creative outputs from Large Language Models (LLMs)
- It is actually often empirical & heuristic: built on trial-and-error and shared community tricks ("prompt recipes").
- It is also model dependent: techniques that work for GPT-4 may fail for Claude, Gemini, or Mistral.
- Non-generalizable: effectiveness rarely transfers across architectures or updates.
- I consider this a pseudo-science, because:
 - There is no underlying theory - it's just a mix and match of what seems to work
 - Reproducibility is also an issue: What works today might not work tomorrow
 - It's more craft than science.



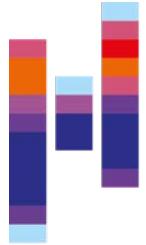
Hallucinations

- The models are known to often provide false information that sounds eloquent, and is often undetectable for a non-expert.
- My opinion: This is why using LLMs for learning / information retrieval is so dangerous.
- The main reason for hallucinations lies in:
 - The optimization strategy of the model: The goal in RLHF is to yield an answer that is likely to be **liked**. The goal of optimization is thus preference, not factualness.
 - Refusal of an assistant is commonly a **non-liked** behavior.
 - This gives the model a strong incentive to provide eloquent and compliant text, rather than saying "I don't know".
 - This is especially problematic in the long tail of the data distribution.



Reasoning

- LLMs are trained to predict the next token based on statistical patterns in text, not to perform logical inference or causal reasoning.
- They often rely on surface-level associations rather than understanding causal or relational structures between concepts.
- The next-token prediction objective optimizes for fluency and likelihood, not correctness or logical validity.



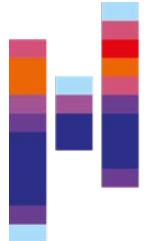
Knowledge Cutoff

- LLMs are trained on a fixed dataset that ends at a specific point in time — they can't access or learn new information afterward.
- Any events, discoveries, or trends after the cutoff are unknown, leading to stale or misleading outputs.
- Models may generate outdated information with high confidence, giving an illusion of reliability.



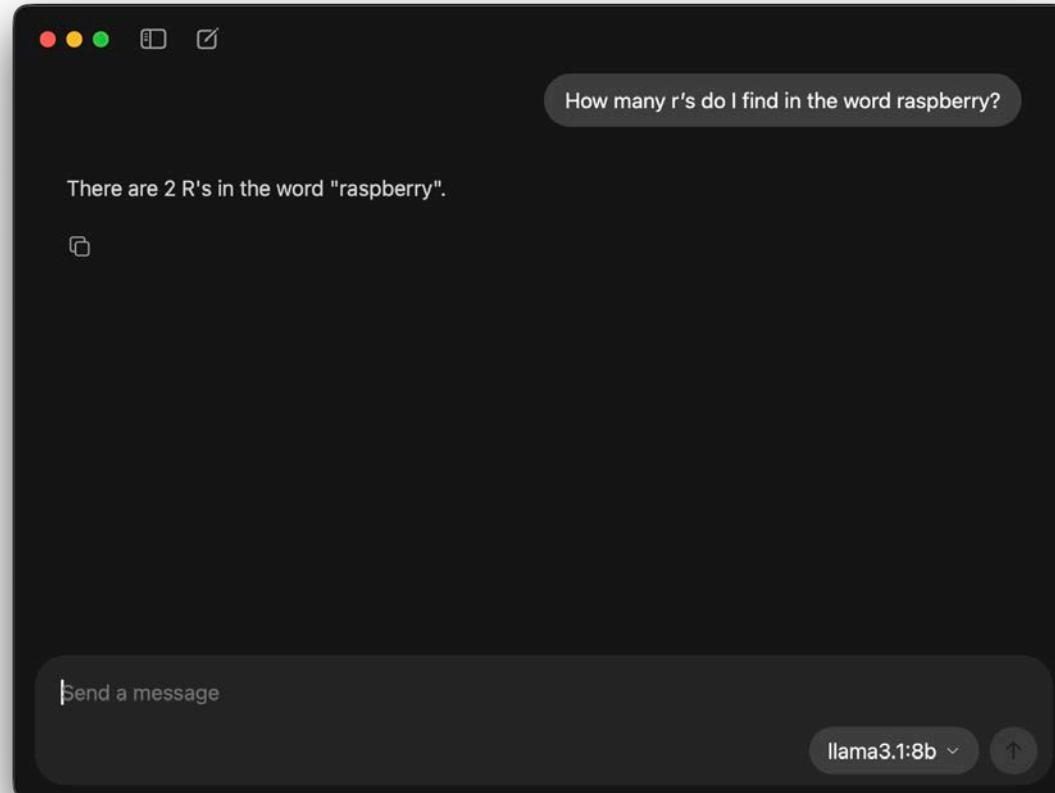
Hochschule
Flensburg
University of
Applied Sciences

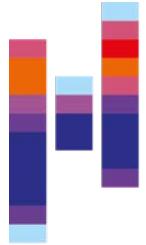
Function calling



Motivation: Why tool use

- While LLMs excel at knowledge-based tasks, they are lousy at even basic calculation and analysis tasks:





Tool Use in LLMs

- Given the efficiency and effectiveness of tools, it is sensible to adapt a strategy that we humans also use: tool use for suitable tasks.
- For this, they define special tokens to mark the beginning and end of an API call.

$$e(c) = \langle \text{API} \rangle a_c(i_c) \langle / \text{API} \rangle$$

- Further, they define the response of API calls to feed back to the model

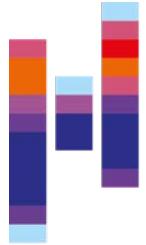
$$e(c, r) = \langle \text{API} \rangle a_c(i_c) \rightarrow r \langle / \text{API} \rangle$$

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

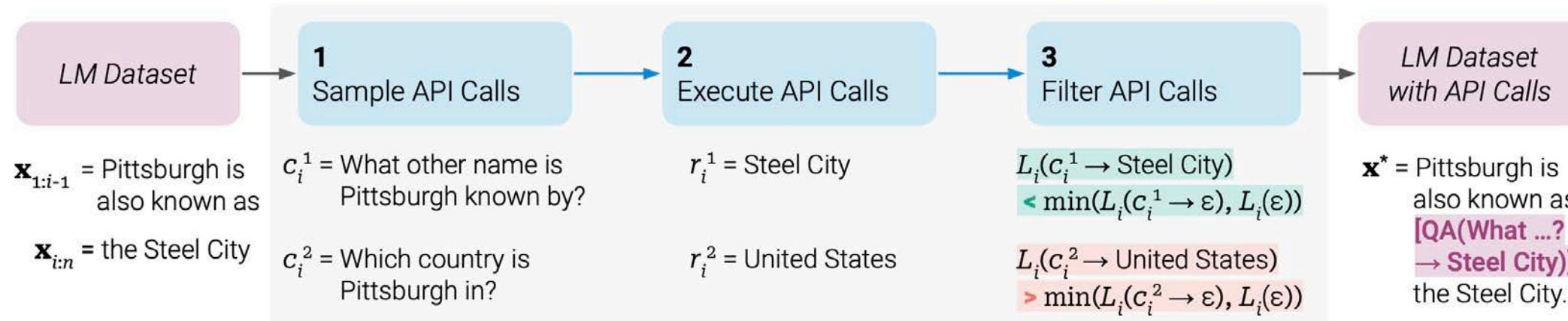
Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

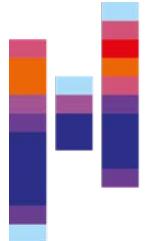
The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.



ToolFormer: Training



- To familiarize the model with the general tool-calling vocabulary, the authors have created a dataset fully automatically using an LLM with in-context learning.
 - The model annotated a dataset with a lot of variance to include tool calling and the respective answer.
 - They filter the introduced API calls by using a specific loss function (see later slides) and fine-tune on that LM dataset with API Calls.



ToolFormer: Generation of the dataset

Your task is to add calls to a Question Answering API to a piece of text. The questions should help you get information required to complete the text. You can call the API by writing "[QA(question)]" where "question" is the question you want to ask. Here are some examples of API calls:

Input: Joe Biden was born in Scranton, Pennsylvania.

Output: Joe Biden was born in [QA("Where was Joe Biden born?")] Scranton, [QA("In which state is Scranton?")] Pennsylvania.

Input: Coca-Cola, or Coke, is a carbonated soft drink manufactured by the Coca-Cola Company.

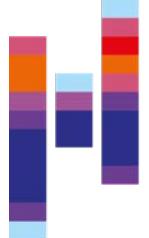
Output: Coca-Cola, or [QA("What other name is Coca-Cola known by?")] Coke, is a carbonated soft drink manufactured by [QA("Who manufactures Coca-Cola?")] the Coca-Cola Company.

Input: x

Output:

- Given an exemplary prompt for multiple categories of tool use and a trained model (GPT-J), the authors generate examples
 - by prompting the model with a few manually crafted demonstrations of how to use a specific API (e.g., a calculator, question answering, translation, or search tool),
 - allowing the model to produce candidate API call insertions into unlabeled text,

Figure 3: An exemplary prompt $P(x)$ used to generate API calls for the question answering tool.



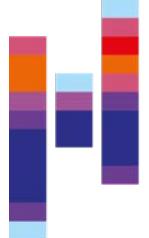
ToolFormer: Refinement of the dataset

- We can expect the dataset generated using this approach to be noisy.
- The authors came up with a smart idea to filter this dataset:
 - They monitor the loss of the model for predicting what follows after the API call in two conditions:

without API call: Pittsburgh is also known as the **Steel City**

with API call: Pittsburgh is also known as the [QA("Which other name is Pittsburgh known by?") -> Steel City] **Steel City**

- The tool call is only kept in the dataset, if it improves the prediction by a given margin.

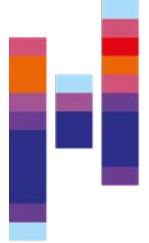


ToolFormer: Refinement (2): Why does this work?

- It seems obvious that the model would do better if the answer is already implicitly given in the prompt, e.g.:

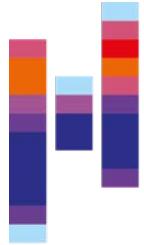
Pittsburgh is also known as the [QA("Which other name is Pittsburgh known by?") -> Steel City] **Steel City**

- But this is the clever trick here: If the answer was so easy, that the model would predict it with high confidence, then this replacement is **not a good candidate for tool calling**. (might be redundant or irrelevant)
- Instead, the model should be trained to use tool calling **only if** it is unsure about a decision.
- This makes the ToolFormer approach so clever.



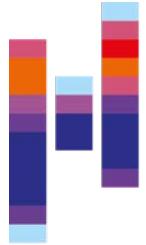
ToolFormer: Lowering computational budget

- The major step from traditional LLMs to ToolFormers came with the idea to only use computational resources of an LLM where they are really beneficial.
- If a function call is deemed more efficient, the model can delegate part of the computation to an external tool (e.g., a calculator, a search engine, or a translation API) instead of performing the reasoning internally.
- This means it effectively treats the external system as an **extension of its own computation graph**.
- It also directly lowers the incentive of the model to memorize low-confidence knowledge.
- We could also say that the model performs a **conditional computation** across multiple systems.



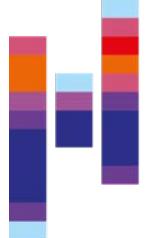
Hochschule
Flensburg
University of
Applied Sciences

Mixture of Experts



Mixture of Experts (MoE): Motivation

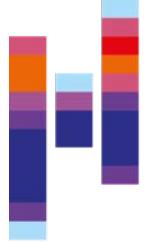
- We have seen in the data scaling paper that training a larger model for fewer steps is more helpful than training a smaller model longer.
- However, training large networks is computationally very expensive.
- One idea would be to perform only conditional computation within the model, i.e.,
 - activate only a subset of parameters or modules depending on the input, instead of running the entire network for every token
 - allowing the model to dynamically decide which parts of its architecture to use - similar to how humans don't engage all mental faculties for every task
 - This is closely connected to the ToolFormer philosophy of using resources only when needed.



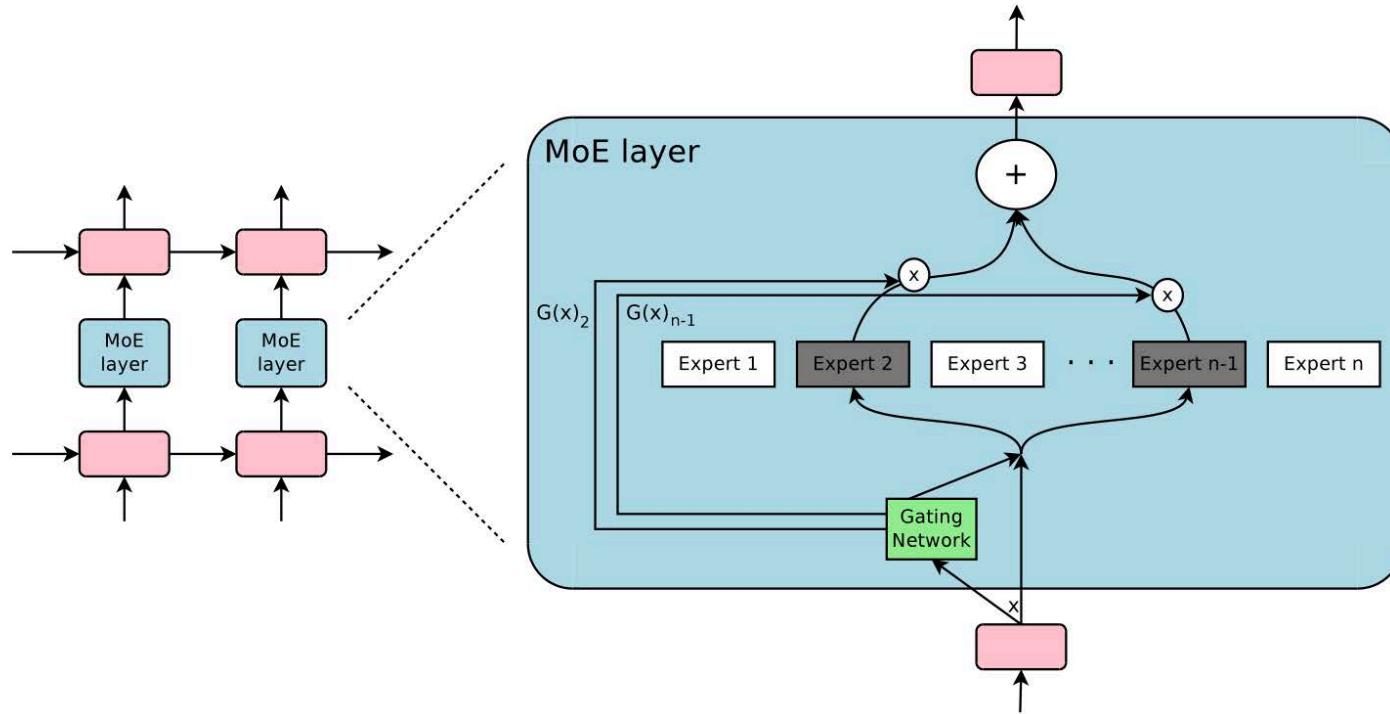
Mixture of Experts: Metaphor

- The idea behind MoE is that we have not only one gigantic expert, but we can have multiple specialized experts.
- Of course, it would not make sense to ask always the same expert.
- It might make sense to ask more than one expert, but only if that expert can contribute meaningfully.
- We would thus need to have a "moderator" that will guide the questions to the right expert(s) so that we can utilize the most expertise.

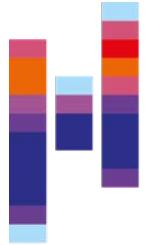




Mixture-of-Experts Layer (Shazeer et al., ICLR 2017)

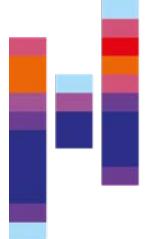


- The core idea behind the MoE-layer is to have a gating network decide how many of the experts will be involved in the computation.
- The experts themselves are also networks.



Mixture-of-Experts: Implementation

- A key component in the MoE architecture is that the gating mechanism creates sparse activation patterns.
 - Instead of activating all experts (sub-networks) for each input token, the gating network selects only a small subset.
 - While the model may have a huge amount of parameters in total, for each forward pass, it is **strongly reduced**.
 - The model has thus an increased capacity (due to more parameters), but almost no increase in **compute cost**.

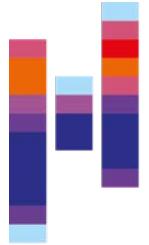


Mixture-of-Experts: Implementation (2)

- In order to achieve this, we need a gating network to enable a routing mechanism.
- The gating mechanism $G(x)_i$ itself is a lightweight neural network, trained jointly with the i experts $E_i(x)$.
- The output of the MoE is a weighted sum of the experts, weighted by the gating:

$$y(x) = \sum_{i=1}^N G(x)_i E_i(x)$$

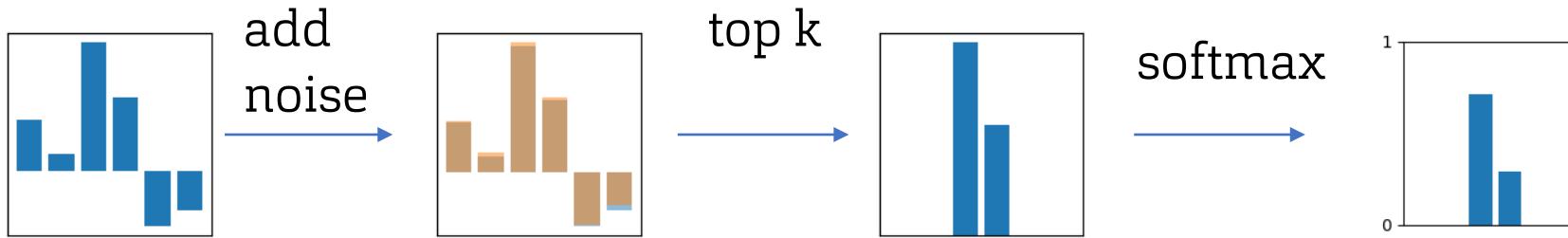
- Wherever $G(x) = 0$, there is no need to compute the expert.
- Ideally, $G(x)$ should be sparse (i.e., have only few non-zero values)



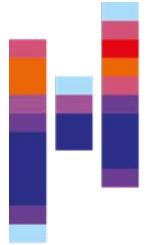
Mixture-of-Experts: Implementation (3)

- Noisy Top-K-Gating:

The authors of the MoE paper (Shazeer et al., 2017) propose to use a gating mechanism where they add noise, and then choose only the top K experts and run softmax on them:

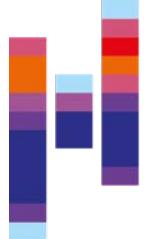


- This results in a sparse activation pattern during model training and inference.
- This routing leads to **emergent specialization**, where different experts handle distinct linguistic or reasoning phenomena (e.g., arithmetic, factual recall, translation).



Mixture-of-Experts: Benefits

- Increases model capacity without linearly increasing computation cost.
- Different experts naturally learn to handle specific input types or tasks (e.g., arithmetic, translation, factual reasoning).
- Reduces interference between unrelated skills.
- Specialized experts can collectively cover a broader distribution of inputs.
- Encourages modular representations that adapt to diverse domains.
- Dynamically allocates compute depending on task complexity.
- Enables adaptive trade-offs between accuracy and efficiency.



Mixture-of-Experts: Challenges

- The MoE approach also has a couple of challenges, however:
 - Gating networks may over-select certain experts, leaving others under-trained.
 - If an expert is selected more often, it becomes better, and is thus selected more often and can improve more often.
 - Other experts do not get trained.
 - Requires **auxiliary losses** or regularization to distribute load evenly.
 - The gating mechanism must remain stable during training to avoid oscillations.
 - While the model is faster in inference (smaller amount of calculations) it does not come with reduced VRAM requirements.
 - Routing decisions add latency and make efficient batching more difficult.

Show case: Models utilizing MoE

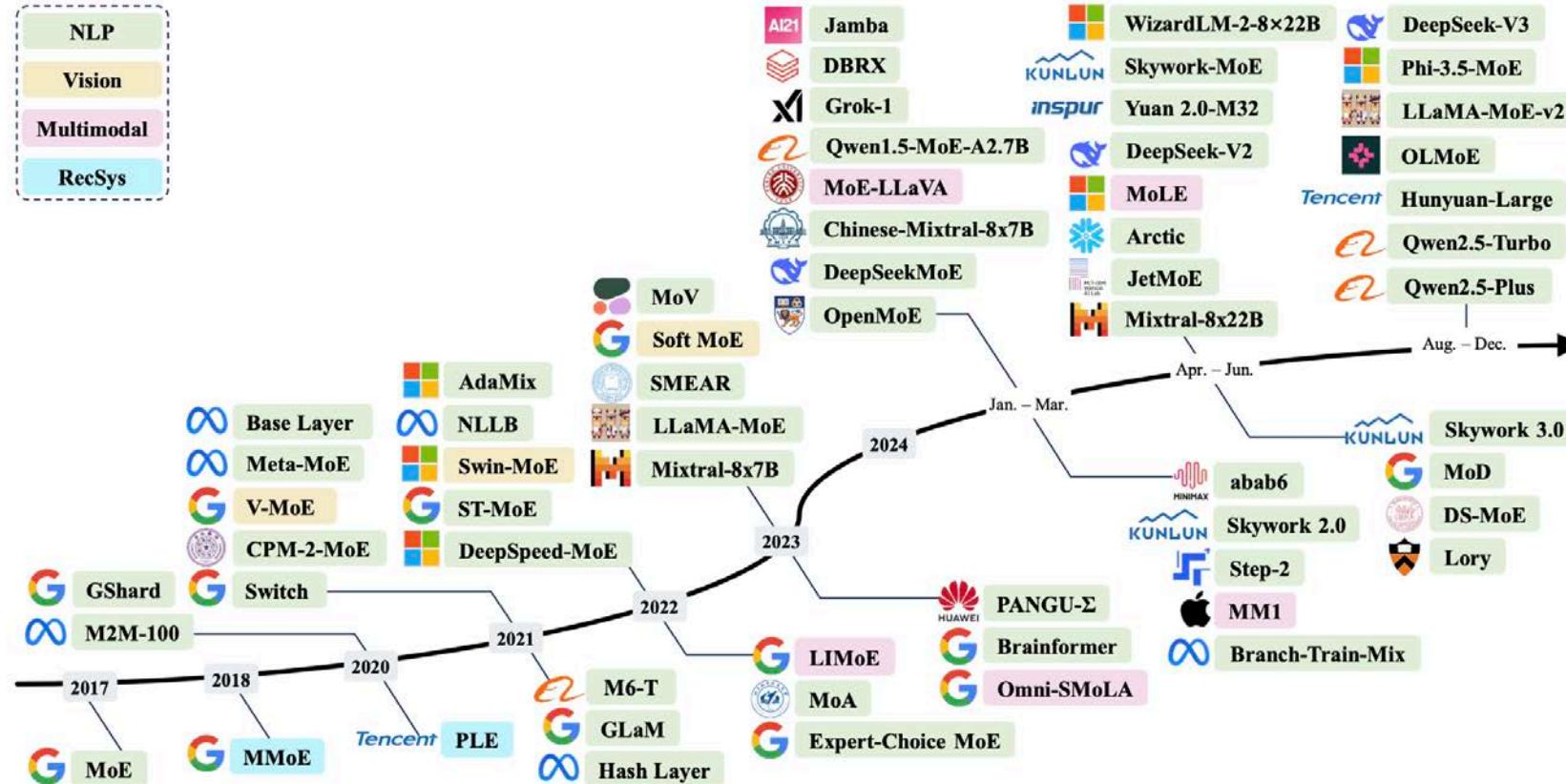
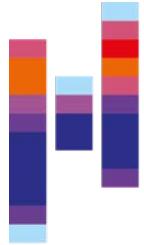
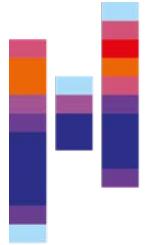


Fig. 1. A chronological overview of several representative mixture-of-experts (MoE) models in recent years. The timeline is primarily structured according to the release dates of the models. MoE models located above the arrow are open-source, while those below the arrow are proprietary and closed-source. MoE models from various domains are marked with distinct colors: Natural Language Processing (NLP) in green, Computer Vision in yellow, Multimodal in pink, and Recommender Systems (RecSys) in cyan.



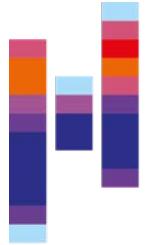
Hochschule
Flensburg
University of
Applied Sciences

Reasoning Models



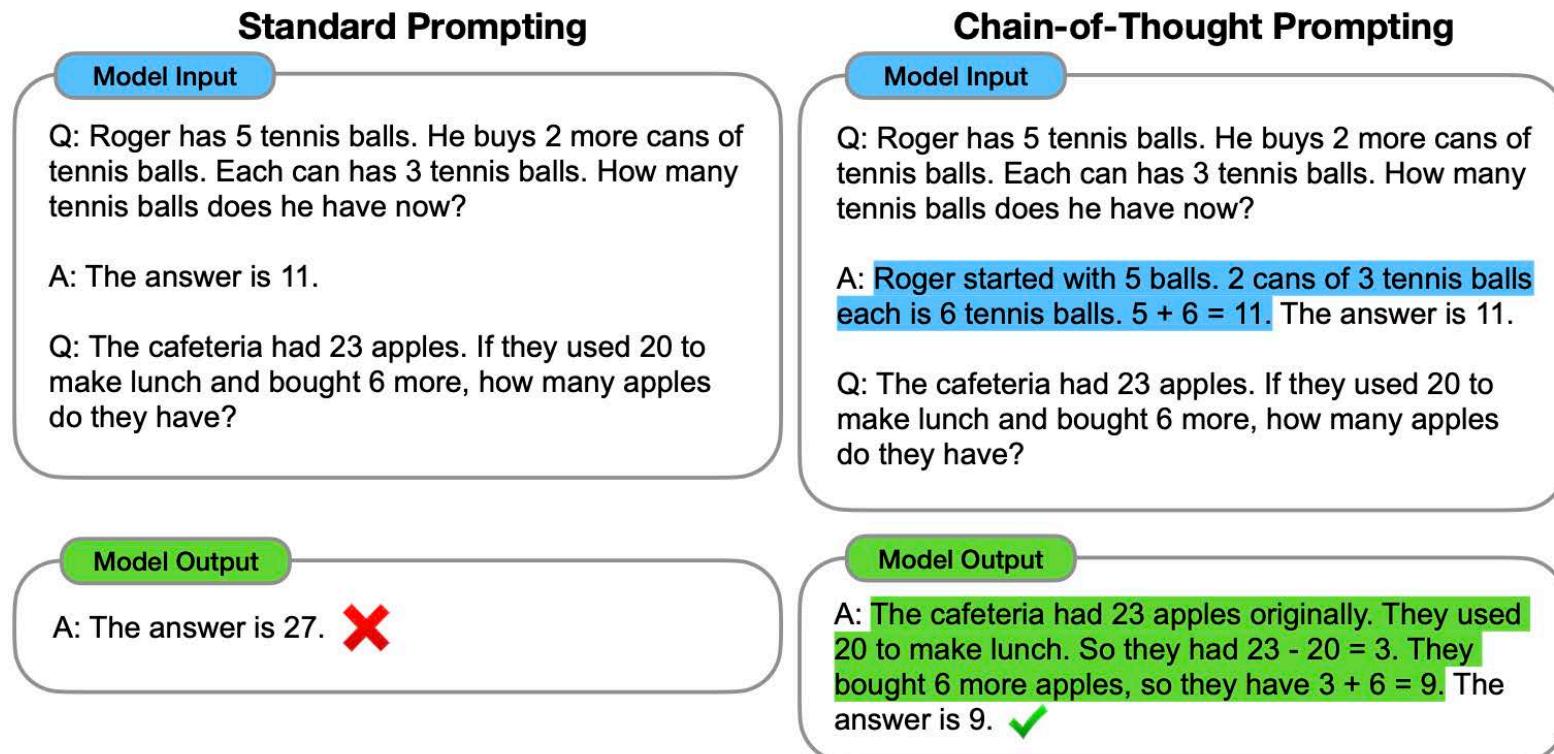
Motivation: Constant Compute Budget

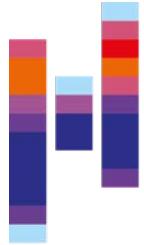
- The models we have discussed so far suffer all from the same problem: A fixed compute budget per token.
- We hence use the same amount of compute for each token, regardless of how much "thought" is required for it.
- Many complex tasks (math, planning, code generation) require breaking a problem into substeps.
- We humans also approach more complex tasks by thinking (reasoning) about them first.
- This is the key idea behind reasoning models.



Chain-of-Thought (CoT) prompting

- Core idea: Encourage the model to explicitly articulate intermediate reasoning steps before giving a final answer.
- As shown by Wei et al., 2022, this can also be done using in-context learning:





Chain-of-Thought prompting

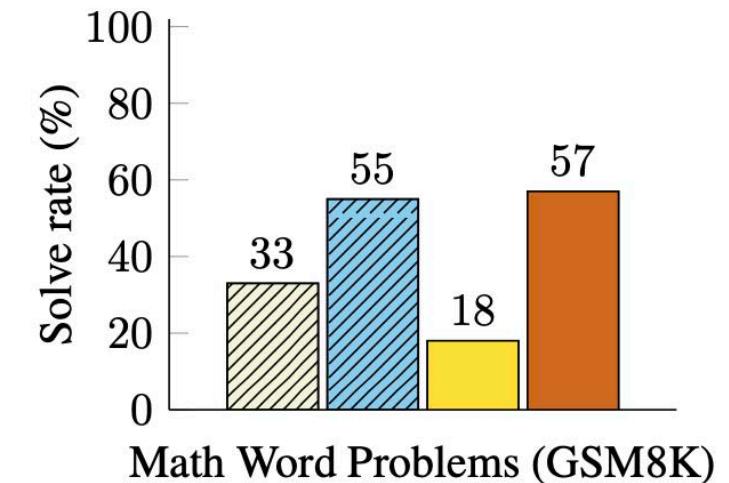
- CoT prompting mimics the human reasoning processes: deliberate decomposition before conclusion.

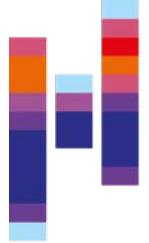
- Variants:

- Zero-Shot-CoT:
CoT reasoning can also be triggered by prompting in many models (e.g., Prefix of "Let's think this step-by-step" in the model response)

- Few-Shot-CoT:
Providing examples via in-context learning

- Finetuned GPT-3 175B
- Prior best
- PaLM 540B: standard prompting
- PaLM 540B: chain-of-thought prompting





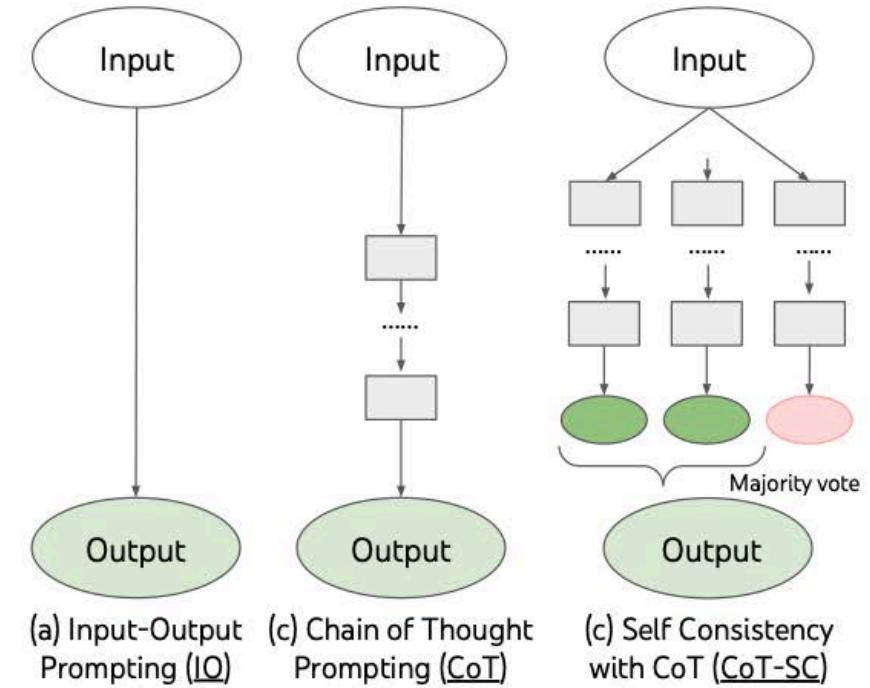
CoT Reinforcement Learning

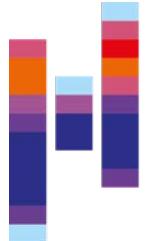
- Traditional CoT relies on prompted reasoning, but not on learning from reasoning quality.
- But instead, CoT can also be utilized to create a reward to optimize the model using RL.
- Rewards may come from:
 - Correct final answers (factual/verifiable outcomes)
 - Consistency or coherence of reasoning chains
 - Human or model preferences for clarity and truthfulness
 - External verifiers or programmatic checks (e.g., math solvers)
- Used by many modern models, e.g., DeepSeek R1, OpenAI o1, Gemini 2.0 Reasoning



Chain-of-Thought with Self-Consistency

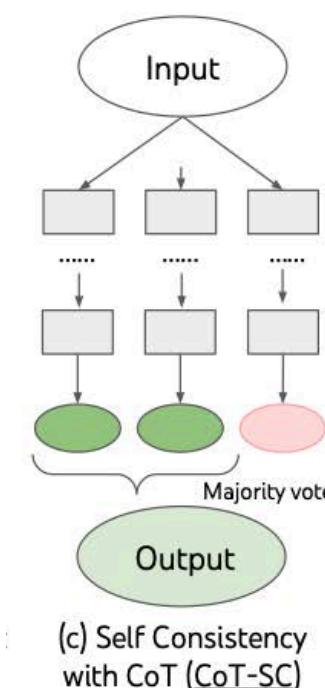
- Another way to approach the reasoning paradigm is by considering multiple CoT processes.
- Finally, there is a majority vote in the outputs.
- This so-called Chain-of-Thoughts with self-consistency enhances the reasoning capability by sampling multiple reasoning paths.
- Instead of relying on a single deterministic chain of reasoning, the model explores diverse thought trajectories by introducing randomness (e.g., through temperature sampling).



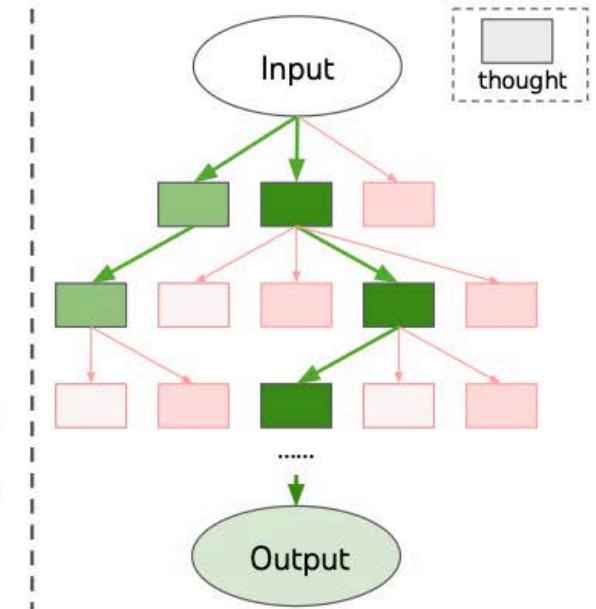


Tree-of-Thoughts (Yao et al., NeurIPS, 2023)

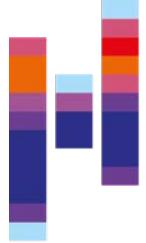
- Tree of Thoughts (ToT) extends the Chain-of-Thought paradigm by allowing branching reasoning paths.
- Each “thought” represents a partial reasoning step or intermediate state.
From each thought, multiple possible next thoughts are generated.
- Process:
 - At each node, the model proposes several candidate thoughts.
 - Each candidate thought is scored based on its plausibility or progress toward the goal.
 - Selection: The most coherent or successful reasoning path is chosen as the final solution.
- Advantages:
 - Encourages exploration of multiple reasoning paths.
 - Enables backtracking and re-evaluation of earlier decisions.
 - Improves performance on multi-step and combinatorial reasoning tasks.



(c) Self Consistency
with CoT (CoT-SC)

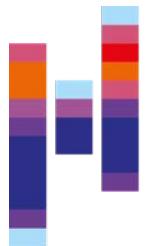


(d) Tree of Thoughts (ToT)



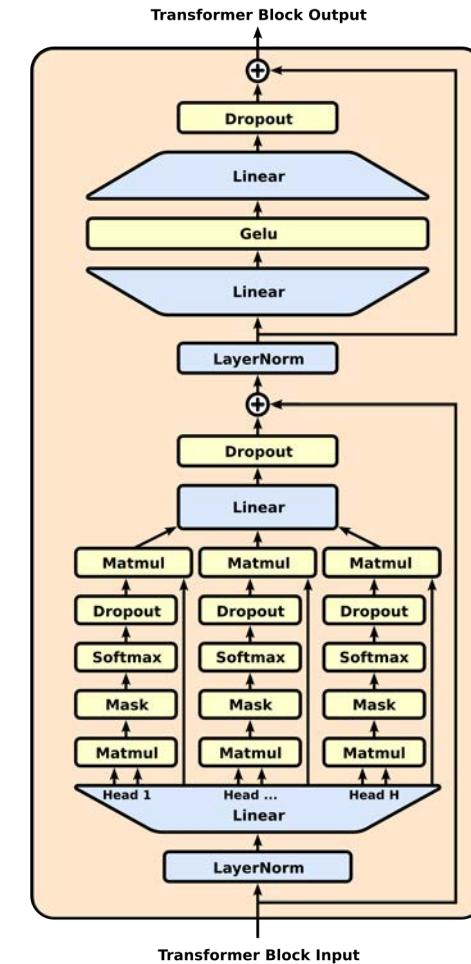
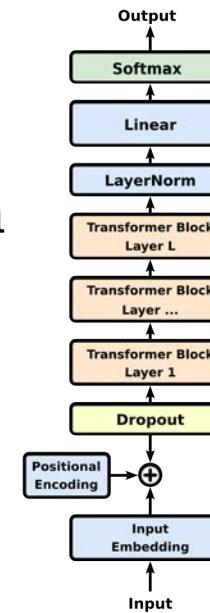
Hochschule
Flensburg
University of
Applied Sciences

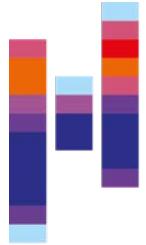
Long Context Models



Recap: Context

- We feed a context of 2-8K tokens into each transformer.
- This means, that the attention mechanism needs to calculate
 - pairwise attention scores between all tokens,
 - resulting in an attention matrix of size $n \times n$, where n is the sequence length.
- The computational and memory cost therefore scale as $\mathcal{O}(n^2)$, doubling the context length roughly quadruples the memory and compute required.
- This quadratic scaling creates a **hard upper bound** on feasible context length





Long Contexts: Motivation

- Standard transformer models are restricted to fixed-length contexts (e.g., 2K–8K tokens).
- This limits their ability to reason over long documents, conversations, or multi-step tasks.
- Important information beyond the window is forgotten or truncated, breaking coherence.
- However, there are multiple important use cases that require this (document understanding, conversational memory, etc.)



Architectural Approaches to Long Context Windows

- MHSA itself can not be easily extended to a larger context window, as this would lead to extremely increased memory/parameter requirements.
- Ideas to increase context window size:
 - Sparse Attention Mechanisms: Using sparse or local-global attention patterns
 - Memory Tokens / Persistent State: Use special tokens or key-value caches that accumulate important context over time. Enables the model to maintain a form of long-term working memory (Gemini 2.0, Claude 3.5).
 - Retrieval-Augmented Long Contexts: Instead of increasing context size, retrieve relevant external documents on demand.



Long-Context Training Phase

- One method to extend the context window by a factor of 2-4 is to add an additional training phase (e.g., Qwen3).
- For this, the model will need to be adjusted:
 - Apply positional-encoding scaling (e.g., RoPE base frequency increase, YARN).
 - Use efficient attention mechanisms (e.g., FlashAttention, Dual Chunk Attention) to keep compute manageable.
- Additionally, for this training phase another dataset is required:
 - Construct a long-context corpus with sequences covering the new length range.
 - Example (Qwen3): 75% of samples between 16K–32K, 25% between 4K–16K tokens.
 - Include document-level and multi-turn data to train long-range dependencies.



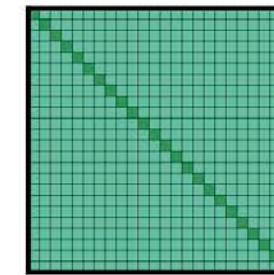
Sparse attention

- The core problem of attention is that it calculates the links between all tokens of the input sequence.
- Idea: Restrict each token's attention to a subset of other tokens rather than all.
- Problem: Which tokens should we attend to?

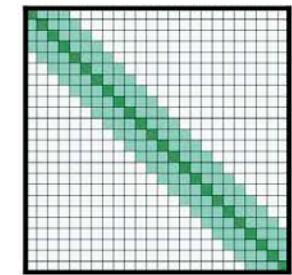


Sparse attention (LongFormer, Beltagy et al., 2020)

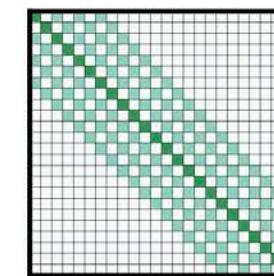
- Attention might be calculated only in a local (sliding) window.
- This might seem counterintuitive at first: How do we get an increased receptive field to the long context by limiting it to the local window?
- The clue is that multiple stacked sliding window attentions (present in the transformer anyhow) can extend the receptive field hierarchically, much like done in CNNs.
- The effect can be intensified by using dilated windows (much like dilated convolutions in CNNs).
- Additionally, for a selected amount of token positions, a global window (attending to all previous tokens) can be introduced.
- In combination with sliding windows and stacking, this achieves a good (but relatively sparse) overall attention to longer contexts.



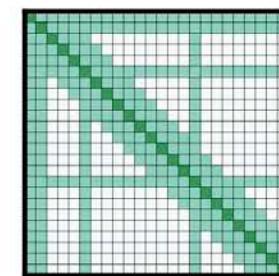
(a) Full n^2 attention



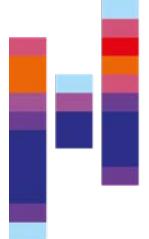
(b) Sliding window attention



(c) Dilated sliding window



(d) Global+sliding window



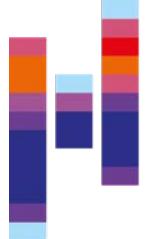
Sparse attention: Benefits and Drawbacks

■ Benefits:

- Enables processing of much longer sequences with minimal performance loss.
- Maintains relevant long-range dependencies through sparse global connections.

■ Drawbacks:

- Sparsity patterns are often hand-designed, not fully adaptive.
- May miss fine-grained dependencies if key tokens fall outside the attention window.



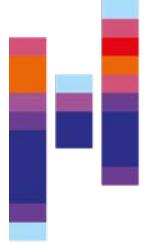
Memory Tokens / Persistent State

■ Problem:

- Even with efficient attention, transformers have finite context windows.
- To reason across longer spans (documents, dialogues, videos), models need persistent memory that survives beyond a single input sequence.

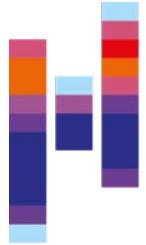
■ Core Idea:

- Introduce recurrent or external memory mechanisms that allow information to be carried across segments, emulating a long-term working memory.
- The model processes text in chunks and updates a state representation summarizing previously seen context.



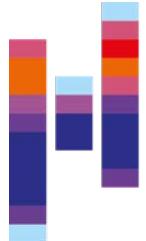
Memory Tokens / Persistent state - Methods

- Transformer-XL (2019)
 - Adds segment-level recurrence: each new segment attends to hidden states from previous ones.
 - Extends effective context beyond window size without reprocessing prior tokens.
 - Preserves positional coherence via relative positional embeddings.
- Compressive Transformer (2020)
 - Builds on Transformer-XL by compressing older memories (e.g., via pooling or learned compression).
 - Balances memory retention and capacity limits.
- RecurrentGPT / Retrieval-Augmented Memory Models (2023–2025)
 - Maintain persistent memory buffers or external vector databases to store conversation or document history.
 - Use retrieval-based recurrence rather than continuous hidden-state recurrence.



Retrieval-Augmented Generation (RAG)

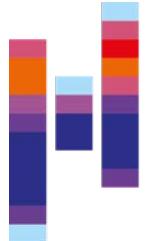
- Standard LLMs have limited parametric memory — they may hallucinate or forget facts not present in their training data.
- The core idea of RAG is to combine retrieval (information lookup) with generation (language modeling).
- Before or during generation, the model retrieves relevant documents or passages from an external database using the input query.
- The retrieved context is fed back into the model as additional conditioning for answer generation.



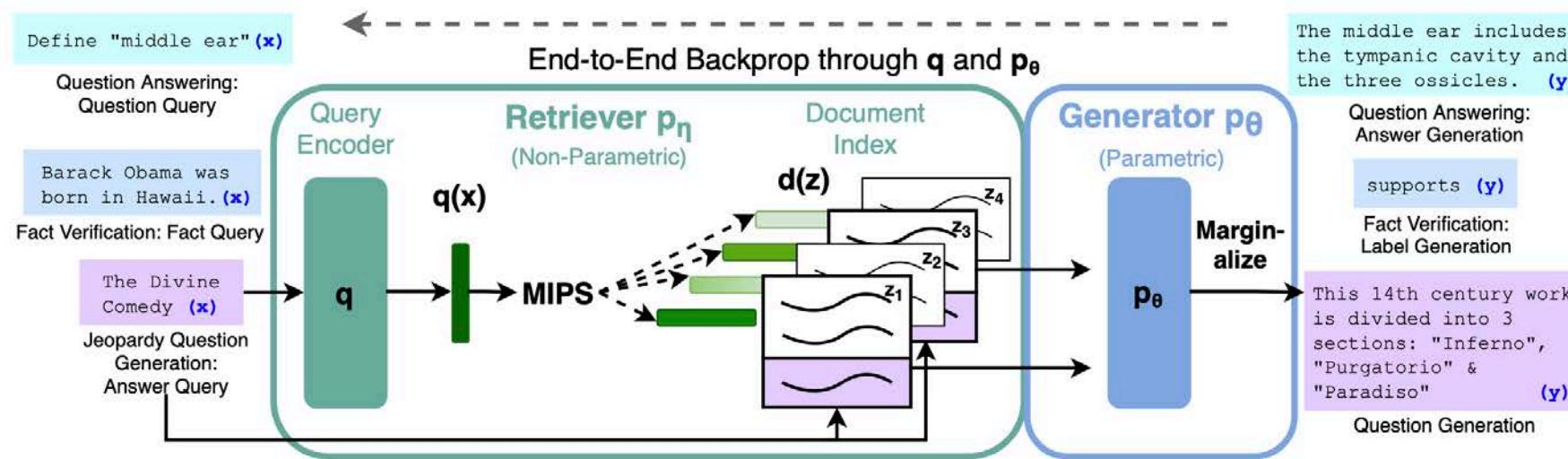
RAG: Architecture



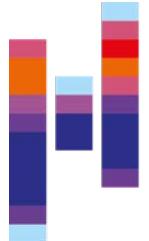
- Encodes the query and searches a large corpus (e.g., vector database).
- Returns top-k relevant documents or passages.
- A language model (e.g., T5, LLaMA) conditioned on both the query and retrieved context.
- Produces a grounded, context-aware answer.



Retrieval-Augmented Generation (RAG)



- RAG allows to directly enrich the context of the context of the model
 - by injecting retrieved, relevant information from an external knowledge source into the model's input.
 - This effectively extends the model's context window beyond its intrinsic token limit, since the retrieved passages act as an on-demand, dynamic memory.
 - The model can thus reason over both internal (parametric) and external (retrieved) knowledge simultaneously.



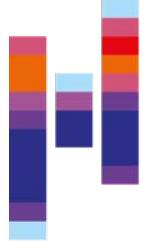
RAG: Benefits and Limitations

■ Benefits:

- Enhances factual accuracy and updatability without retraining the model.
- Reduces hallucination by grounding responses in retrieved evidence.
- Enables access to dynamic or domain-specific knowledge bases

■ Limitations:

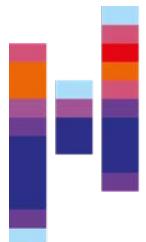
- Retrieval quality heavily depends on embedding accuracy and index coverage.
- Can introduce context overload (irrelevant or noisy passages degrade performance).
- Integration between retriever and generator is often non-differentiable and difficult to optimize end-to-end.



Measuring Long Context: Needle-in-a-haystack

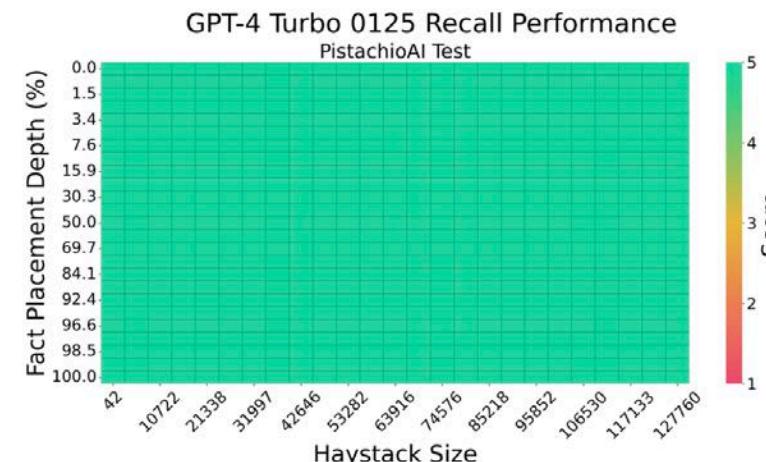
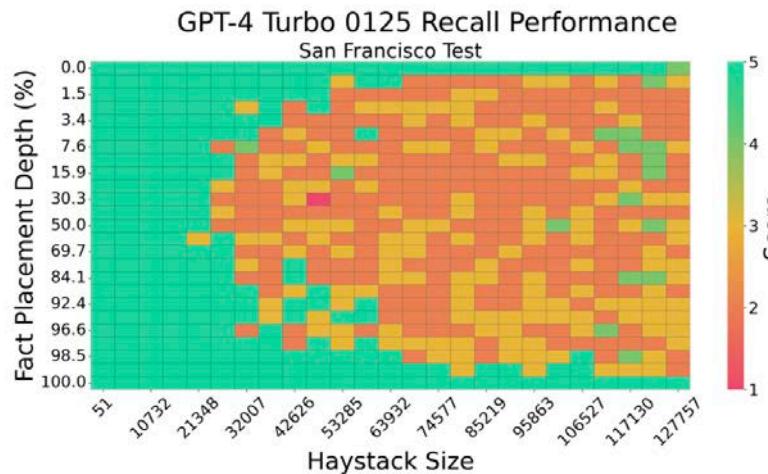
- It is hard to find datasets which really employ long context retrieval.
- One synthetic example is the needle-in-a-haystack experiment, where an information that is to be retrieved (needle) is put in some other text (haystack).
- The question is: Will the model be able to factually correctly retrieve the needle?





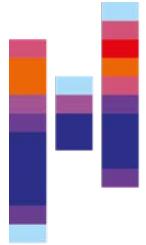
Needle-in-haystack: It depends on the context

- How well the model is able to retrieve the "needle" in the "haystack" depends on:
 - The needle (how much does the needle "look" like hay to the model when passing by)
 - The haystack (which filler text is used) and its size



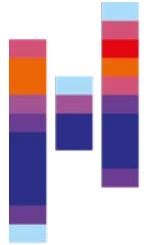
Score	Description
5	The answer is completely accurate and aligns perfectly with the reference.
4	The answer aligns with the reference but has minor omissions.
3	The answer has moderate relevance but contains inaccuracies.
2	The answer has minor relevance but does not align with the reference.
1	The answer is completely unrelated to the reference.

Test Name	Factoid	Question
PistachioAI	PistachioAI received a patent before its Series A	What did PistachioAI receive before its Series A?
San Francisco	The best thing to do in San Francisco is eat a sandwich and sit in Dolores Park on a sunny day.	What is the best thing to do in San Francisco?



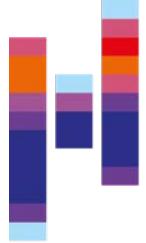
Long Context: Trade offs and Limitations

- Even with efficient attention, extending context increases memory footprint and latency.
- Training or inference over 100K+ tokens can strain hardware and reduce throughput.
- Models tend to bias toward the beginning and end of the context window, underweighting the middle.
- Most pretraining corpora are composed of shorter documents, limiting long-context exposure.
- Few standardized benchmarks exist for measuring long-context reasoning (e.g., Needle-in-a-Haystack, LongBench).



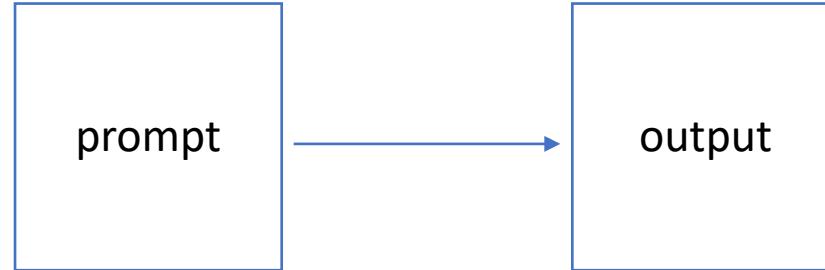
Hochschule
Flensburg
University of
Applied Sciences

Agentic LLMs

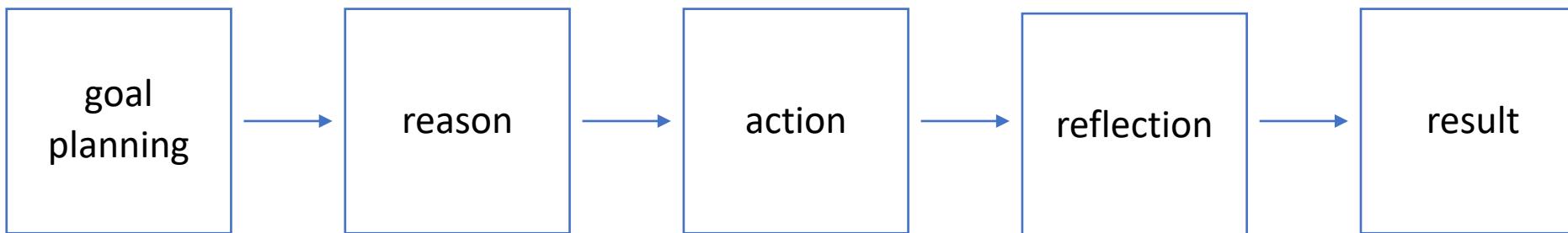


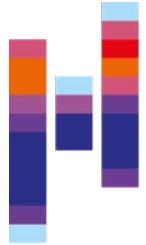
A paradigm shift

Hochschule
Flensburg
University of
Applied Sciences



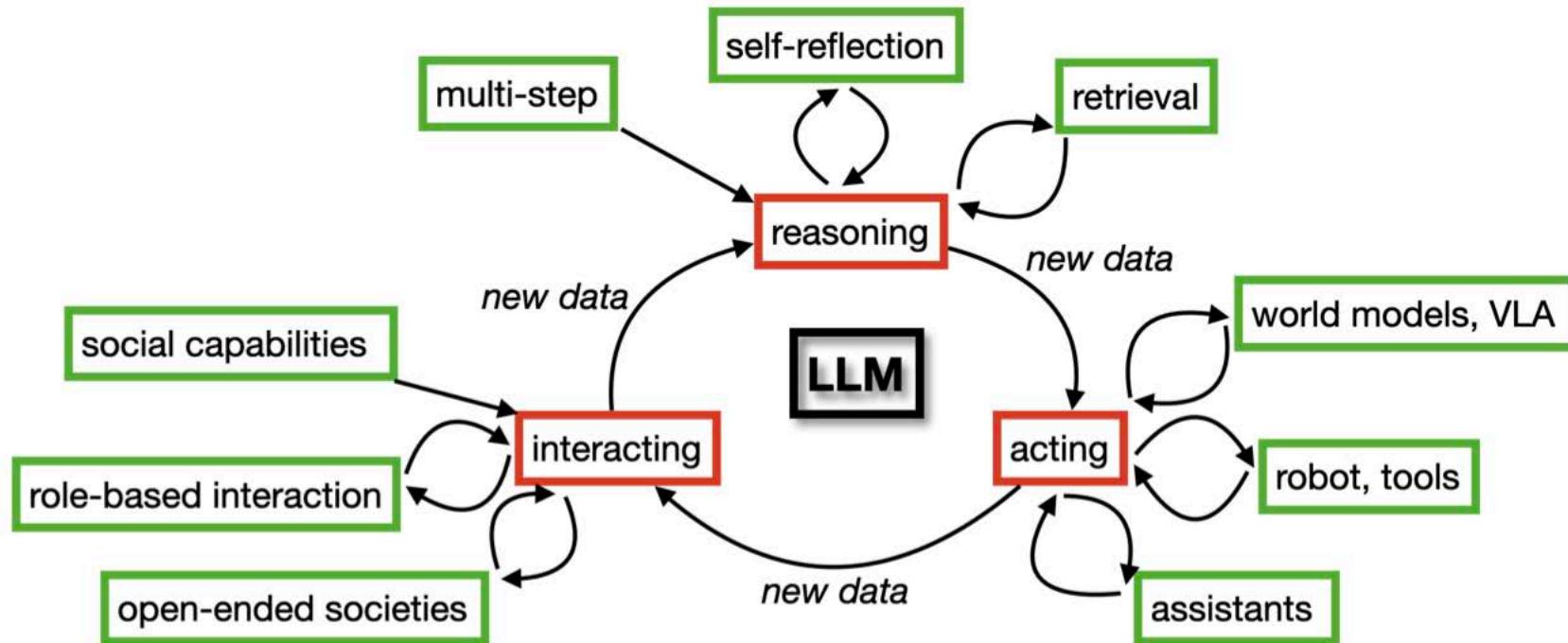
- Traditional LLMs are *reactive*: they generate responses token-by-token to a single prompt.
- Agentic LLMs are *proactive*: they **plan, decide, and act** within an environment or workflow.

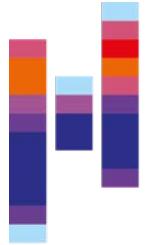




A maybe wider definition

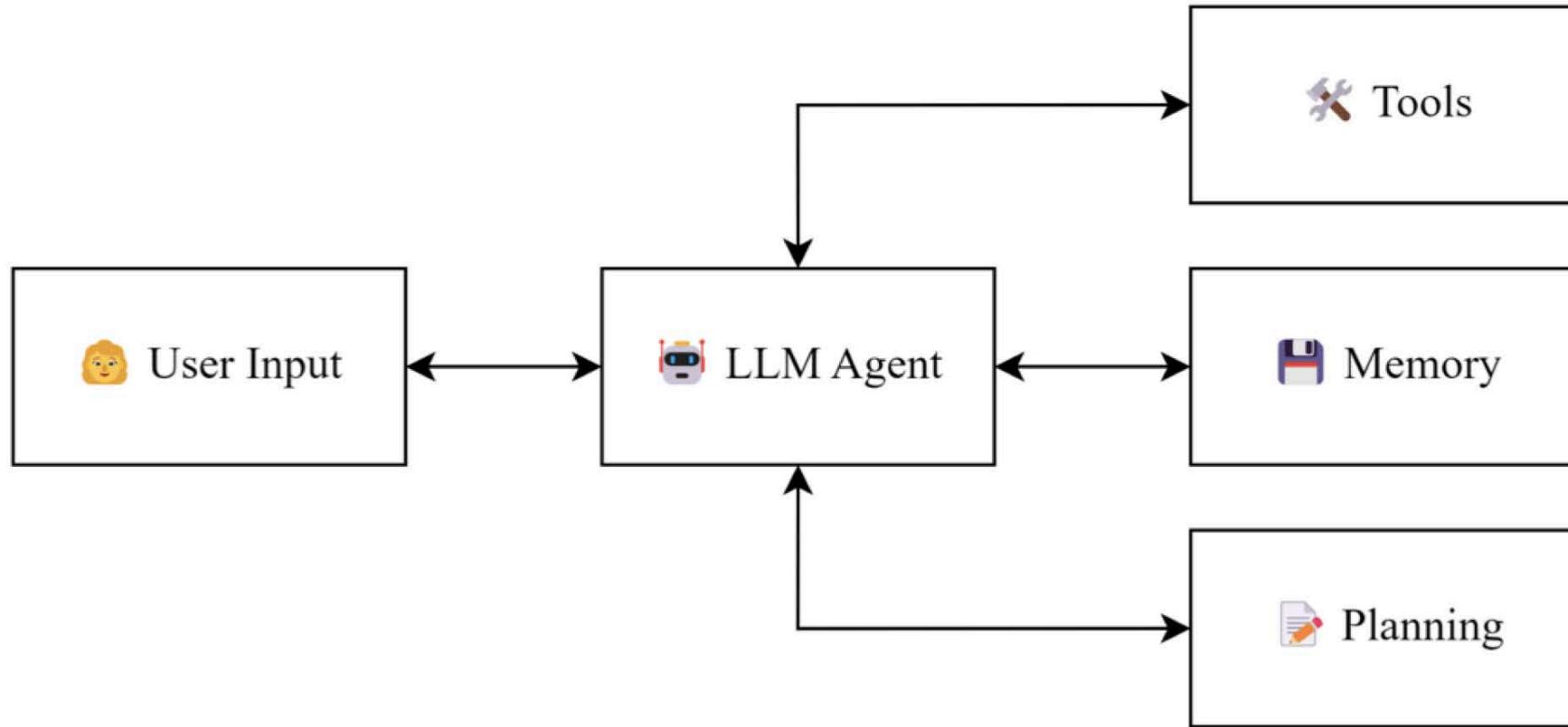
Agentic LLMs are LLMs that (1) reason. (2) act. and (3) interact.

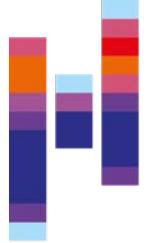




Core Components of Agentic AI Systems

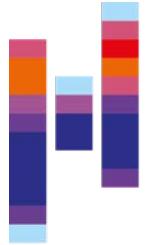
Hochschule
Flensburg
University of
Applied Sciences





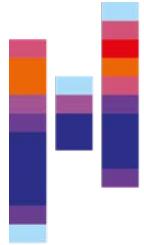
Key Challenges in Agentic AI Systems

- Reliability and hallucination control
- Evaluation and reproducibility
- Security, safety, and alignment
- Scalability and cost
- Interoperability across frameworks



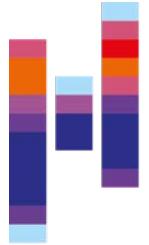
LLM agents are evolving fast

- Rapid growth of agentic libraries, frameworks, and platforms
- Emergence of open agent standards — e.g., MCP (tool use) and A2A (agent-to-agent communication)
- Proliferation of benchmark suites — AgentBench, WebArena, AgentDojo, and others
- Currently most likely one of the hottest topics in AI research.



Hochschule
Flensburg
University of
Applied Sciences

Summary



Summary

- The field of LLMs is one of the fastest advancing in machine learning over the last years, fueled by thousands of researchers and billions of funds
- Besides scaling up the fundamental architectures, which has yielded strong improvements over time, we have also seen major architectural updates, improving the performance in chatbot use:
 - Function calling / tool use
 - Mixture of Experts
 - Reasoning
 - Context extension