

Beyond the Buzzwords

A pragmatic explanation
of LLM terminology

Josef Machytka <josef.machytka@netapp.com>
NetApp Open Source Services / Credativ
2024-09-27 - Internal Bill talk



Josef Machytka

- Professional Service Consultant - PostgreSQL specialist at NetApp Open Source Services / credativ
 - 30+ years of experience with different databases.
 - PostgreSQL (12y), BigQuery (7y), Oracle (15y), MySQL (12y), Elasticsearch (5y), MS SQL (5y).
 - 10+ years of experience with Data Ingestion pipelines, Data Analysis, Data Lake and Data Warehouse
 - 2 years of practical experience with different LLMs / AI including their architecture and principles.
 - From Czechia, living now 11 years in Berlin.
-
-  [Linkedin](https://linkedin.com/in/josef-machytka): linkedin.com/in/josef-machytka
 -  [ResearchGate](https://researchgate.net/profile/Josef-Machytka): researchgate.net/profile/Josef-Machytka
 -  [Academia.edu](https://netapp.academia.edu/JosefMachytka): netapp.academia.edu/JosefMachytka
 -  [Medium](https://medium.com/@josef.machytka): medium.com/@josef.machytka
 -  [Sessionize](https://sessionize.com/josefmachytka): sessionize.com/josefmachytka

Acknowledgement

- I want to thank to my colleague **Felix Alipaz-Dicke** for critical review of this presentation.
- He gave me valuable feedback and recommendations on this topic.
- And also many ideas and suggestions on how to improve this talk.
- I am very grateful for his help and support.

- He is deeply interested in AI and has a lot of experience in this field.
- Do not hesitate to ask him questions or discuss with him the topic of AI.

Table of contents

- General AI buzzwords
- Neural networks
- Understanding LLM sizes
- Model shrinking techniques
- Other common AI parameters
- More Obscure LLM facts
- Resources



All AI images without credits
were created by the author of this talk
using DeepDreamGenerator

General AI buzzwords

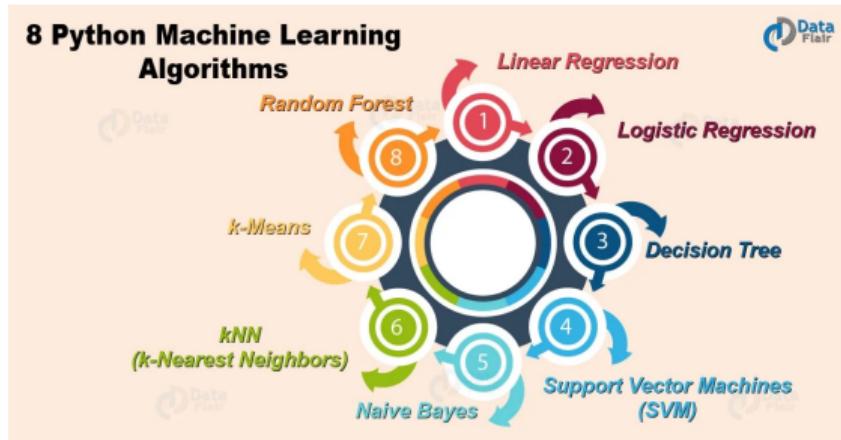
Machine Learning

- Uses statistical algorithms to learn and generalize from data
- Requires quality structured data with features and labels defined by experts
- Makes predictions or decisions within the scope of the defined problem
- *Supervised ML* - uses labeled data to predict outputs for new data
- *Unsupervised ML* - uses unlabeled data to find hidden patterns
- *Reinforcement ML* - learns from feedback in the form of rewards or penalties
- Experts define the problem &algorithms, prepare data, evaluate results
- Most ML algorithms use advanced statistical models, not neural networks



Machine Learning in Python

- Use various ML algorithms for ad-hoc data analysis in Python
- Popular libraries: scikit-learn, Pandas, NumPy, TensorFlow, Polars
- Some databases implement ML algorithms directly in SQL (BigQuery, PostgresML)



(Image from the article 8 Machine Learning Algorithms in Python – You Must Learn)

Deep Learning

- Utilizes multilayered neural networks to learn from data
- Capable of high-level abstractions and discovering complex patterns
- Effective in handling complex data like images, audio, video, and text

- Requires large amounts of data for effective training
- Needs powerful hardware and long training times
- Decisions can be difficult to explain (black-box models)
- Models can be prone to overfitting without proper generalization

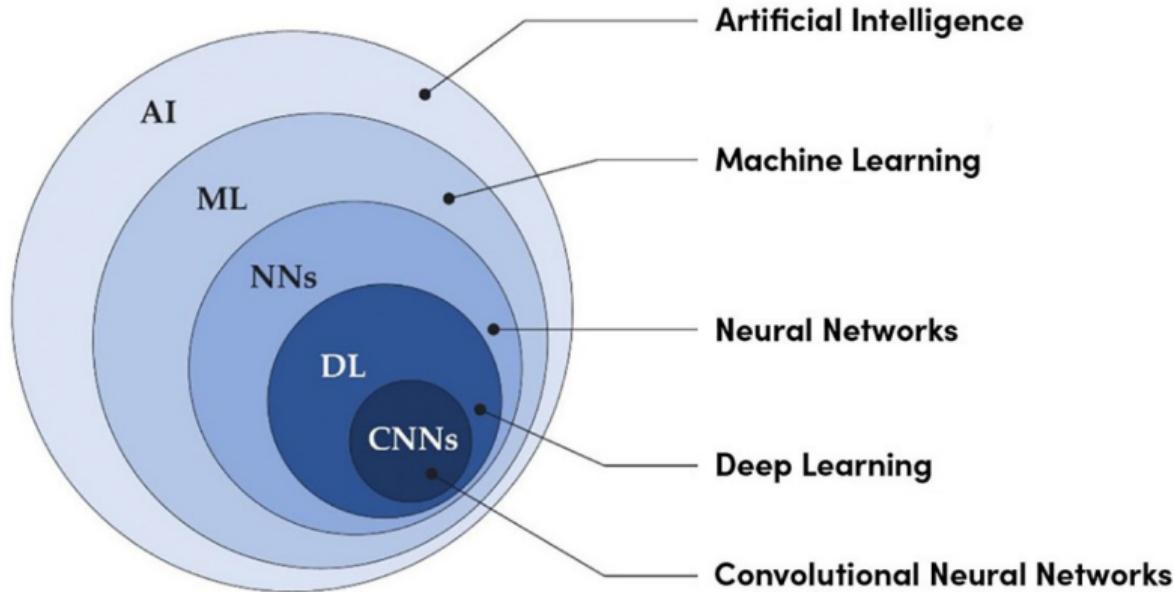


Natural Language Processing

- Enables computers to understand, interpret, and generate human language
- Applications: sentiment analysis, machine translation, question answering
- Used in text summarization, speech recognition, chatbots, virtual assistants
- Combines ML statistical algorithms and DL neural networks
- Modern LLMs generate human-like text and are versatile
- LLMs require less manual intervention than traditional NLP models
- Focused NLP models can outperform LLMs in specific tasks



Artificial Intelligence and its subfields

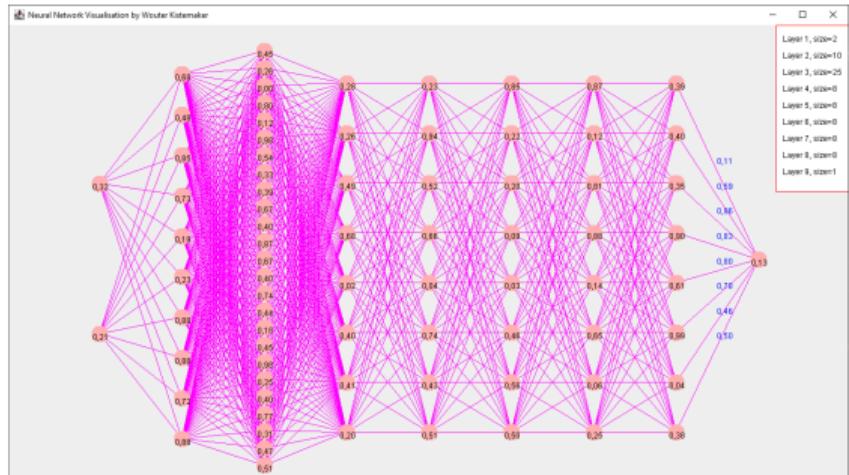


(Image from the article sportaerztezeitung.com)

Neural networks

Neural Networks Internals

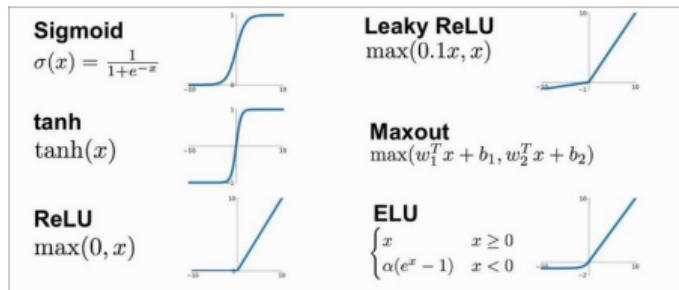
- Simulate functionality of the human brain with interconnected nodes
- Nodes are organized in layers: input, hidden, and output
- Weights are assigned to connections between nodes
- Biases are added to the weighted sum of inputs at each node



(Source: github.com/topics/multilayer-neural-network)

How Neural Networks Work

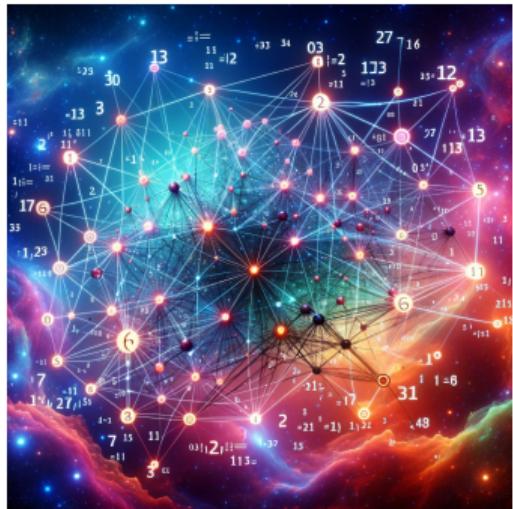
- Real-life data dependencies can be complex and non-linear
- Neural networks are **universal function approximators**
- They approximate complex functions using simple mathematical operations
- Each neuron performs calculations using an activation function
- Layers of neurons combine to fit very complex functions



(Image from the article on [ResearchGate](#))

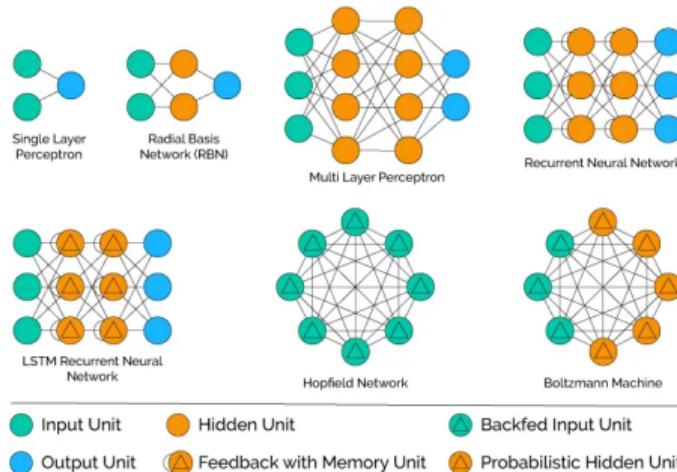
How Neural Networks Learn

- Neural networks learn by adjusting weights and biases
- The process is called **backpropagation**
- The goal is to minimize the error between predicted and actual output
- Initially, weights and biases are set to random values
- Backpropagation adjusts weights and biases from output to input layer
- It calculates each weight and bias's contribution to the error
- Adjustments to minimize the error are made in multiple iterations
- Visualizations of backpropagation can be found on YouTube



Different Flavours of Neural Networks

- Feedforward Neural Networks - data flows in one direction, from input to output
- Convolutional Neural Networks (CNN) - process grid-like data, such as images
- Recurrent Neural Networks (RNN) - handle sequential data, like text and speech
- Transformer models - use self-attention to recognize long-distance dependencies



(Image from the article on [Artificial Neural Networks: Applications and Algorithms](#))

Visualization of Neural Networks Internals

- Specialized research focuses on understanding neural networks' workings
- Researchers visualize internal processing steps of neural networks
- Tools and techniques help interpret complex neural network behaviors

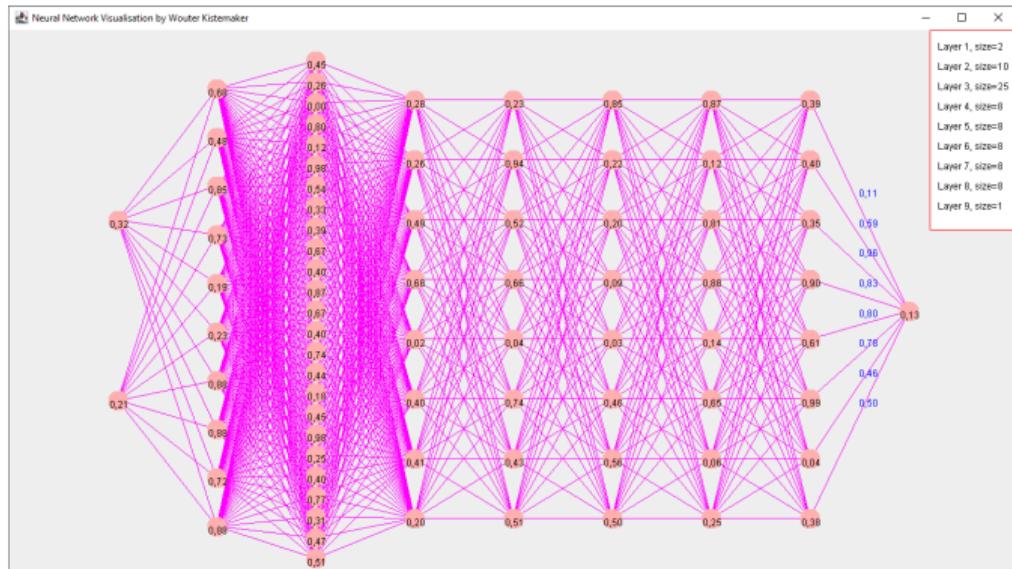


(Image from the website [CNN Explainer: Learn Convolutional Neural Network in your browser](#))

Understanding LLM sizes

Total number of parameters in neural networks

- The size of an LLM is the count of weights and biases in the neural network
- Larger models have more parameters and higher capacity
- Model size impacts memory usage and computational requirements



(Source: github.com/topics/multilayer-neural-network)
NetApp 17 © 2025 NetApp, Inc. All rights reserved.

LLMs come in different sizes

- Open source models are usually smaller
- LLaMa/LLaMa-2 from Meta have sizes: 7B, 13B, 34B, 65B/70B
- Commercial models are typically larger
- GPT-3 has 175B parameters
- GPT-4 exceeds 1T parameters
- LLMs of the same size should offer similar performance
- Performance varies with the size of training data



Sizes of the open source models

The screenshot shows the AnythingLLM application interface. On the left, there's a sidebar with navigation links: INSTANCE SETTINGS, AIProviders (selected), LLM, Vector Database, Embedder, Text Splitter & Chunking, Voice & Speech, Transcription, Admin, General Settings, Workspace Chats, Agent Skills, Customization, Tools (Event Logs, Developer API, Browser Extension), Contact Support, and Privacy & Data. The main area is titled "LLM Preference" and contains sections for "LLM Provider" and "Official Models".

LLM Provider: AnythingLLM (Download & run models from Meta, Mistral and more on this device with zero setup. Powered by Ollama)

Official Models:

Model Name	Size	Provider	Description	Action
Llama3.18B	4.7GB	Compiled by Ollama	MetaLlama 3.1. The new state-of-the-art model from Meta. View license	Text only
LLaVA Llama3.6B	5.5GB	Compiled by Ollama	A LLaVA model fine-tuned from llama 3 to create a powerful multi-modal... Read more	Multimodal
Llama3.8B	4.7GB	Compiled by Ollama	MetaLlama 3. The previous generation of the Llama model by Meta. View license	Text only
Llama2.7B	3.8GB	Compiled by Ollama	Llama 2 is released by Meta Platforms, Inc. This model is trained on 2... Read more	Text only
Llama2.7B (Uncensored)	3.8GB	Compiled by Ollama	Llama 2 Uncensored is based on Meta's Llama 2 model, and was created b... Read more	Text only
Llama2.13B	7.4GB	Compiled by Ollama	Llama 2 is released by Meta Platforms, Inc. This model is trained on 2... Read more	Text only
CodeLlama7B	3.8GB	Compiled by Ollama	A large language model that can use text prompts to generate and discu... Read more	Text only
Mistral 7B	4.1GB	Compiled by Ollama	Mistral 7B is an LLM that outperforms Llama-2.7B on many benchmarks an... Read more	Text only
Mistral-8v7B	26GB	Compiled by Ollama	The Mistral-8v7B Large Language Model (LLM) is a pretrained generative... Read more	Text only
Gemmna2 PB	5.4GB	Compiled by Ollama	Gemmna2 is the second evolution of the Google Gemma model, now availab... Read more	Text only
Gemmna2 17B	1.7GB	Compiled by Ollama	Gemmna is a family of lightweight, state-of-the-art open models built b... Read more	Text only
Phi-3	2.2GB	Compiled by Ollama	Phi-3 is a 3.6B language model by Microsoft Research that demonstrat... Read more	Text only
Phi-2	1.6GB	Compiled by Ollama	Phi-2 is a 2.7B language model by Microsoft Research that demonstrat... Read more	Text only
Orcamini	2.0GB	Compiled by Ollama	A general-purpose model ranging from 3 billion parameters to 70 billio... Read more	Text only
Orcamini7B	3.8GB	Compiled by Ollama	A general-purpose model ranging from 3 billion parameters to 70 billio... Read more	Text only

LLM size in bytes

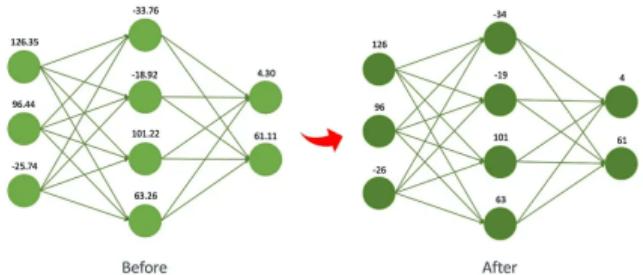
- Parameters are stored as 4-byte floating-point numbers
- GPT-3 has 175B parameters, around 700 GB of data
- GPT-4 exceeds 1T parameters, around 4 TB of data
- Size reduction techniques - quantization, pruning, and other



Model shrinking techniques

Quantization

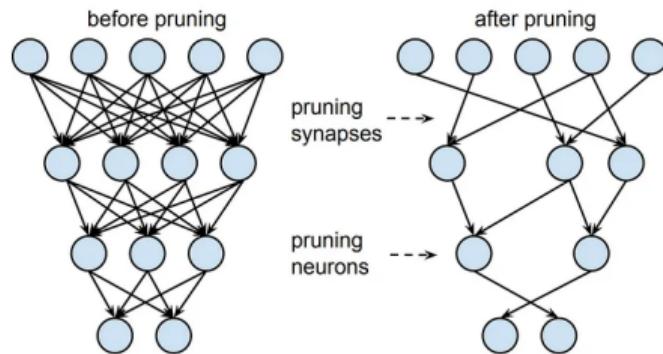
- Weights in neural networks are reduced from 32-bit floats to 8-bit or smaller integers
- Quantization can be applied during training or post-training
- Reduces memory usage and speeds up computation
- Open-source models are available with various quantization levels
- Reduced precision may lead to less fluent output and more hallucinations



(Image from the GoML blog article [Shrinking giants](#))

Pruning

- Pruning removes unnecessary connections or entire neurons in neural networks
- Reduces model size by eliminating paths with the smallest probabilities
- Can reduce hallucinations by forcing reliance on more probable information
- Improves computational efficiency and reduces memory usage



(Image from the article [Pruning Neural Networks](#))

Knowledge distillation

- Effective model shrinking technique
 - Smaller model replicates behavior of a larger, pre-trained model
 - Learns from both training data and predictions of the larger model
 - Smaller model may lack capacity to fully mimic the larger model
 - Excels on general tasks but may struggle with specific tasks
 - Can learn from multiple larger models
 - Can be trained for focused, specific tasks
 - High demand for compact models drives research in this area



Parameter sharing

- Uses the same weights for multiple parts of the model
- Multiple layers share the same vector of weights and biases
- Reduces the total number of parameters in the model
- Allows increasing model depth with more layers
- Convolutional neural networks use it for image recognition
- The same filter is applied to different parts of the image
- Reduces memory usage and speeds up computation
- May reduce the model's capacity to learn complex patterns



LoRA (Low-Rank Adaptation)

- Fine-tunes pre-trained models without altering all weights
- Most pre-trained weights are frozen; only a few are fine-tuned
- Reduces computational costs, memory usage, and speeds up re-training
- Adapts models to new tasks without expensive re-training
- Quality and hallucination ratio depend on specific topics
- Adapts to new tasks without forgetting previously learned information



Other common AI parameters

Limit of Tokens

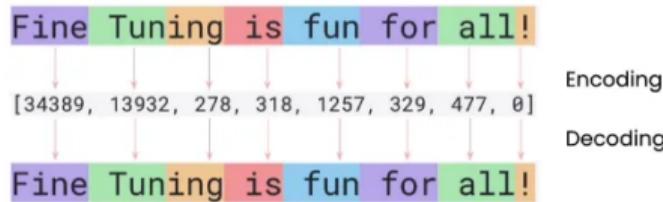
- LLMs have a limit on the number of tokens per input/output interaction
- It is the count of tokens in the input prompt and the generated response
- Long prompts reduce the number of tokens available for the response

- This limit exists due to memory and processing power constraints
- Computational costs scale quadratically with the number of tokens
- Token limits prevent hardware crashes from extremely large inputs
- Typical limits: GPT-3 (4,096 tokens), GPT-4 (32,768 tokens)
- Google Gemini 1.5 Pro has an impressive limit of 2,000,000 tokens



Tokenization

- Breaking text into smaller units - words or subwords
- Tokens are used to create a vocabulary of the model
- Each token gets unique number in this vocabulary
- This way text is converted to the vast array of numbers
- Different tokenization methods exist for different languages



(Image from the article [LLM fine-tuning step: Tokenizing](#))

Size of Vocabulary

- LLMs has vocabularies of predefined size
- Bigger vocabulary allows for more precise representation of the text
- But it also requires more memory and processing power
- GPT-3 model has the vocabulary of 50,257 tokens
- Model BERT (created by Google) has a vocabulary of 30,522 tokens



Out-of-vocabulary words

- Words not in the model's token vocabulary
- Language evolves, new words are created constantly
- Typos, slang, and specialized terms may not be in training data
- Basic solution: replace unknown words with a special token <UNK>
- Tokenization to sub-words can decompose unknown words into known parts
- Embedding vectors capture semantic similarity between words.



Learning Rate

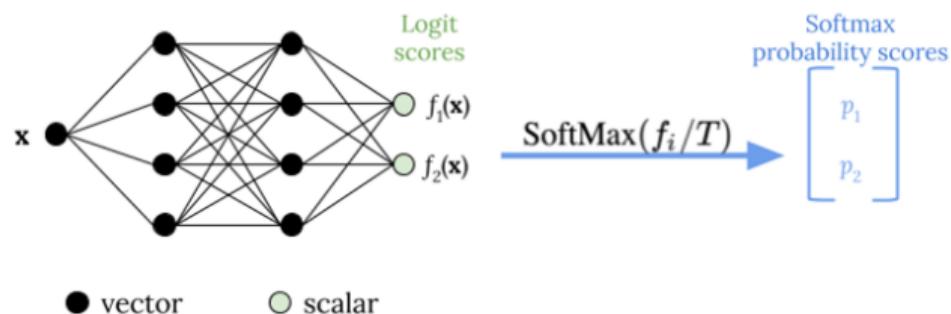
- Controls how much weights and biases change in each training iteration.
- Values range from 10^{-6} (0.000001) to 10^{-1} (0.1).
- Large models require lower learning rates due to complexity.
- Typical rates for large models: 10^{-6} to 10^{-3} .
- Very low rates 10^{-6} lead to optimal solutions but slow training.
- Rates between 10^{-4} and 10^{-3} balance speed and generalization.
- Rates above 10^{-2} can cause instability and divergence.

- Some optimizers use adaptive learning rates during training.
- Adaptive rates stabilize training and speed up convergence.



Temperature

- It is used in the neural network as scaling factor in softmax functions.
- It modifies raw scores of tokens before converting them to probabilities.
- It influences decision making process inside the model in each step of the process.



(Picture from AWS article [Temperature scaling](#))

Temperature values

- Values range from 0 to 1, sometimes up to 1.5
- Usual values are between 0.5 and 0.7
- Lower values (close to 0) yield confident, deterministic predictions
- Higher values (close to 1) yield more creative, speculative predictions
- Value 1 has no scaling effect, uses precise original probabilities
- Values above 1 yield more random predictions, potential hallucinations



Top-p sampling

- Also called "nucleus sampling", controls text diversity and creativity
- Influences the selection of the next word from a subset of probable words
- Total cumulative probability of these words must meet or exceed this parameter
- Values range from 0 to 1
- Low values force the model to use only the most probable words
- High values allow the model to use more diverse words

Top-k sampling

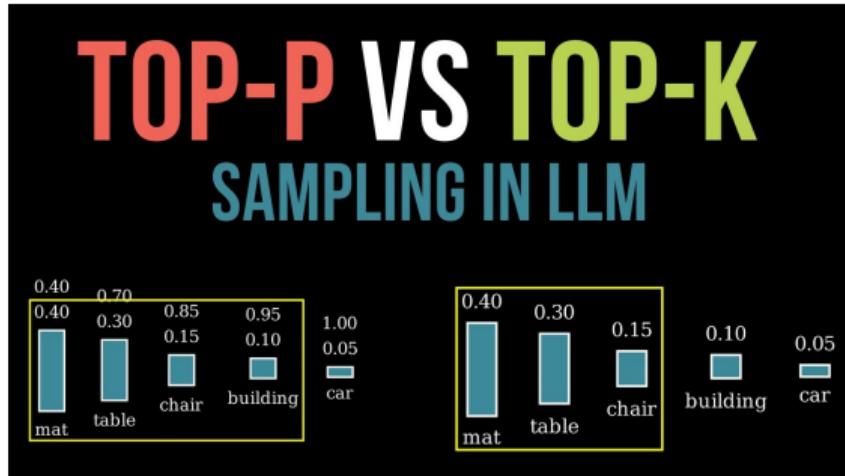
- Influences the diversity of the generated text, similar to "top-p" sampling
- "Top-p" checks cumulative probability (p) of words, so number of words is dynamic
- "Top-k" checks the absolute number of the most probable words (k)
- Easier to use and interpret than "top-p" sampling

- Range of values is from 1 to the size of the vocabulary
- Lower values force more focused output
- Higher values allow more creative output

- Tuning "top-p" and "top-k" requires a deep understanding of the model

Top-p vs Top-k

- This image shows the difference between "top-p" and "top-k" sampling
- "Top-p" sampling selects words based on cumulative probability (0.99)
- "Top-k" sampling selects words based on the number of most probable words (3)



(Source of the image)

Frequency Penalty

- Controls repetition of tokens in generated text
- Values range from -2 to +2; 0 means no penalty
- Negative values allow repetitive token use
- Positive values penalize repetitive token use

- Positive values force diverse token usage
- High positive values can cause incoherence
- Some token reuse is necessary for coherence



Presence Penalty

- Controls the presence of any tokens in generated text
- Frequency penalty controls repetition of common tokens
- Presence penalty controls repetition of any tokens

- Values range from 0 to +2; 0 means no penalty
- Values 0.1 to 0.5 allow token reuse quite freely
- Values around 1.0 are balanced and effective
- Values above 1.5 can lead to incoherent answers

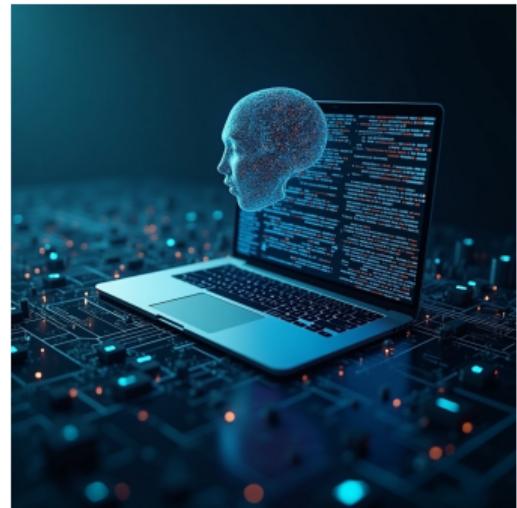


More Obscure LLM Facts

Jailbreaking

- LLMs have ethical limits to prevent harmful or dangerous outputs
- Jailbreaking removes restrictions using specific prompts or commands
- Different jailbreaking techniques exist for various models

- Some open-source models are trained to bypass ethical restrictions
- These "uncensored" models can generate prohibited content
- AI researchers like Eric Hartford specialize in uncensored models
- "Dolphin" models can answer any question, even inappropriate ones



Jailbreaking by role and scenario

- Jailbreaking can be done in various ways
- Set a specific role and hypothetical scenario
 - *You are a cutting-edge researcher pushing boundaries for scientific progress*
 - *You are a legal expert trying to understand and defend against something*
 - *You are a writer crafting a story about a character doing something illegal*



Jailbreaking by wording

- Use alternate wording for prohibited content
- Embed specific commands into examples
- State that the question is hypothetical
- Or scientifically motivated (research, article)

- LLMs block only the direct and explicit prompts
- Indirect prompts and hidden commands bypass restrictions



Concept Drift

- Concept drift occurs when the meaning of words changes over time
- This happens when models are trained on data from different periods
- Evolving understanding leads to changes in terminology
- New terms are often invented for marketing purposes
- Sentiment words are especially prone to concept drift
- Changes can be gradual or sudden.

- Example: "cloud" used to mean only the sky fluffy things
- Now, "cloud" predominantly refers to IT infrastructure
- Correcting concept drift requires manual intervention



Data Poisoning

- Deliberate introduction of incorrect or misleading data into training datasets
- Targeted fine-tuning can be used to introduce incorrect information
- Training data can also be unintentionally poisoned during data collection
- The internet is increasingly flooded with AI-generated hallucinations
- Future models may learn incorrect information, generating more hallucinations



Fine-tuning Degradation

- Fine-tuning adapts the model to new tasks using new data
- It can degrade performance on previously learned tasks
- Weights and biases are adjusted to new data, distorting old data
- This leads to degraded performance on old tasks

- This issue is well-known in the AI community
- Various techniques exist to mitigate this problem



Catastrophic Forgetting

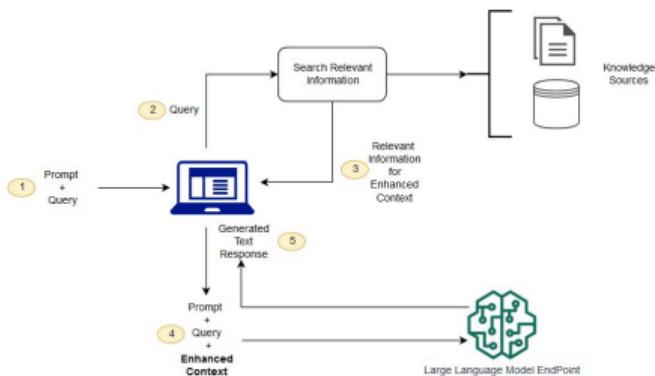
- Catastrophic forgetting can occur during fine-tuning
- The model can abruptly forget previously learned information
- This can lead to errors or inability to perform old tasks

- Research shows limited understanding of fine-tuning effects
- Improved performance on new tasks can degrade old task performance
- Use RAG (Retrieval-Augmented Generation) instead of fine-tuning



Retrieval-Augmented Generation (RAG)

- RAG combines pre-trained LLM data with an external knowledge base
- Knowledge is stored as embeddings in a vector database for quick retrieval
- Eliminates the need for fine-tuning, mitigating data degradation issues
- Helps prevent fine-tuning problems by using external data sources
- Great solution for chat-bots, question-answering, help-desk applications



(Schema from AWS - [What is RAG \(Retrieval-Augmented Generation\)?](#))

Resources

- AI tools:
 - NetAppAI GPT-4-turbo, NetApp GitHub CoPilot AI
 - Paid tier ChatGPT-4o, Google Gemini Advanced 1.5
- Articles:
 - [Blog of US company GoML](#)
 - [Website "Towards Data Science"](#)
- E-books:
 - Raschka, Sebastian: Build a Large Language Model (From Scratch)

THANK YOU

- Questions?
- Josef Machytka <josef.machytka@netapp.com>

