

DAV/DGVFM
Jahrestagung
2026

Zoran Nikolić & Sara Kristin Schmidt, Deloitte

Mathematische Grundlagen generativer KI

Fachgruppe ADS/AI, 30.04.2026

Fragen an MS Copilot (GPT-5)



Wie wird man Mitglied bei der DAV?



Um Mitglied beim **Deutschen Alpenverein (DAV)** zu werden, gehst du in der Regel in vier einfachen Schritten vor.[...]



Wie wird **eine Mathematikerin** Mitglied bei der DAV?



Um Mitglied der **Deutschen Aktuarvereinigung (DAV)** zu werden, folgt **eine Mathematikerin** in der Regel diesem Weg: [...]

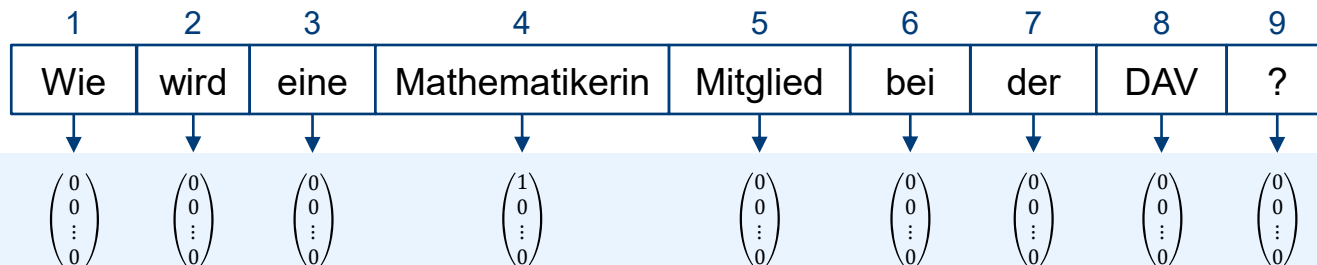


Wenn du mit DAV die **Deutsche Aktuarvereinigung e. V.** meinst (was im Kontext Mathematikerin am besten passt), dann läuft der Weg zur Mitgliedschaft über die Ausbildung zur Aktuarin DAV. [...]



Tokenisierung

Vokabular ca. vom Umfang
 $|V| \approx 50 \text{ Tsd.}$

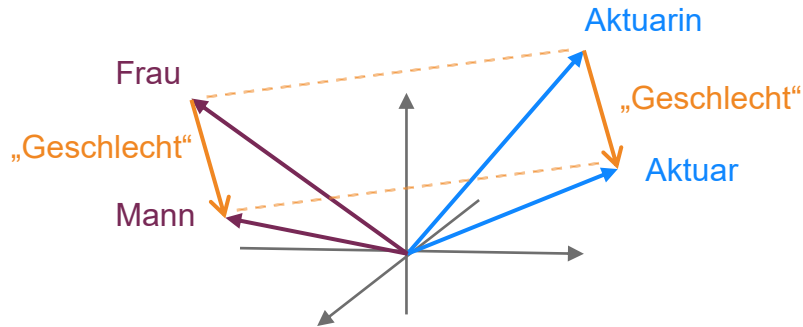


Repräsentation im $\mathbb{R}^{|V|}$?

Embeddings

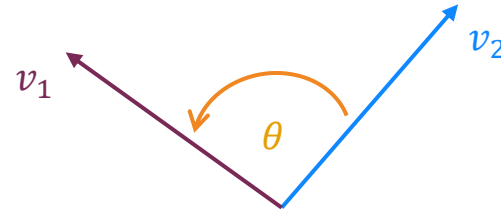
Repräsentation der Tokens als numerische Vektoren in „kleinerem“ Raum \mathbb{R}^{4096} mit Hilfe von einer trainierbaren Embedding-Matrix $E \in \mathbb{R}^{4096 \times 50.000}$, die das Folgende mathematisch erfasst:

Bedeutung von Tokens



$$\text{Aktuar} \approx \text{Aktuarin} + (\text{Mann} - \text{Frau})$$

Ähnlichkeit von Tokens



$$\text{Ähnlichkeit}(v_1, v_2) := \frac{1}{|v_1||v_2|} v_1 \cdot v_2 = \cos(\theta)$$

Tokenisierung

Vokabular ca. vom Umfang
 $|V| = 50 \text{ Tsd.}$

1	2	3	4	5	6	7	8	9
Wie	wird	eine	Mathematikerin	Mitglied	bei	der	DAV	?

Embedding

Finden geeigneter
Repräsentationen im \mathbb{R}^{4096}

$\begin{pmatrix} 0,0107 \\ 0,0043 \\ \vdots \\ -0,0312 \end{pmatrix}$	$\begin{pmatrix} -0,0020 \\ 0,0237 \\ \vdots \\ -0,0194 \end{pmatrix}$	$\begin{pmatrix} 0,0483 \\ 0,0456 \\ \vdots \\ 0,0240 \end{pmatrix}$	$\begin{pmatrix} -0,0067 \\ -0,0348 \\ \vdots \\ -0,0082 \end{pmatrix}$	$\begin{pmatrix} 0,0236 \\ -0,0185 \\ \vdots \\ 0,0025 \end{pmatrix}$	$\begin{pmatrix} 0,0002 \\ 0,0159 \\ \vdots \\ -0,0109 \end{pmatrix}$	$\begin{pmatrix} 0,0489 \\ 0,0213 \\ \vdots \\ 0,0238 \end{pmatrix}$	$\begin{pmatrix} 0,0190 \\ 0,004 \\ \vdots \\ -0,0243 \end{pmatrix}$	$\begin{pmatrix} 0,0081 \\ -0,0033 \\ \vdots \\ 0,0194 \end{pmatrix}$
---	--	--	---	---	---	--	--	---

Transformerblöcke

Sukzessive Transformation der
Embeddings im \mathbb{R}^{4096}



Transformerblöcke

Idee: Schrittweise Verbesserung der Embeddings im \mathbb{R}^{4096} durch sukzessives Durchlaufen der Transformerblöcke, wobei jeder Block alle Tokens parallel verarbeitet.

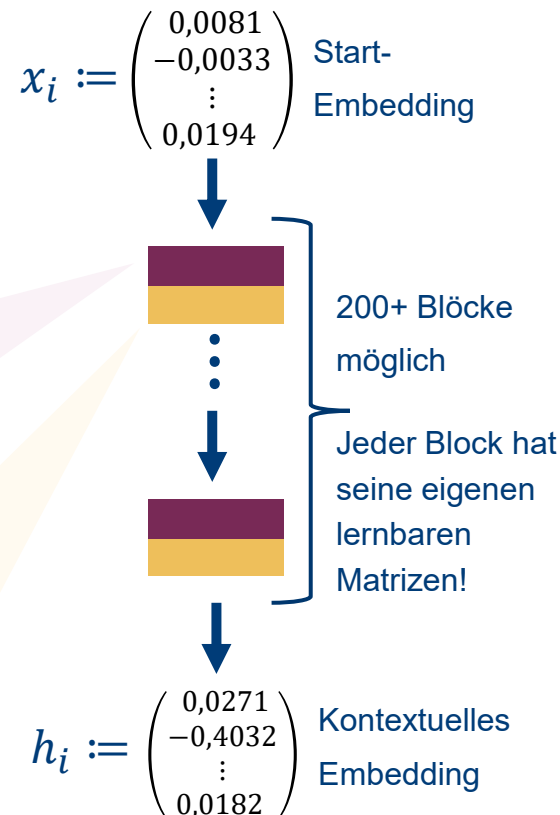
Multi-Head Attention: Jedes Token bezieht kontextuelle Informationen von anderen Tokens der Sequenz (gleich mehr!).

Feed-Forward Layer: Neuronales Netz für nicht-lineare Transformation, das von jedem Token i einzeln durchlaufen wird:

$$FFN(x_i) = W_2 \cdot \sigma(W_1 \cdot x_i)$$

Nicht-lineare Funktion

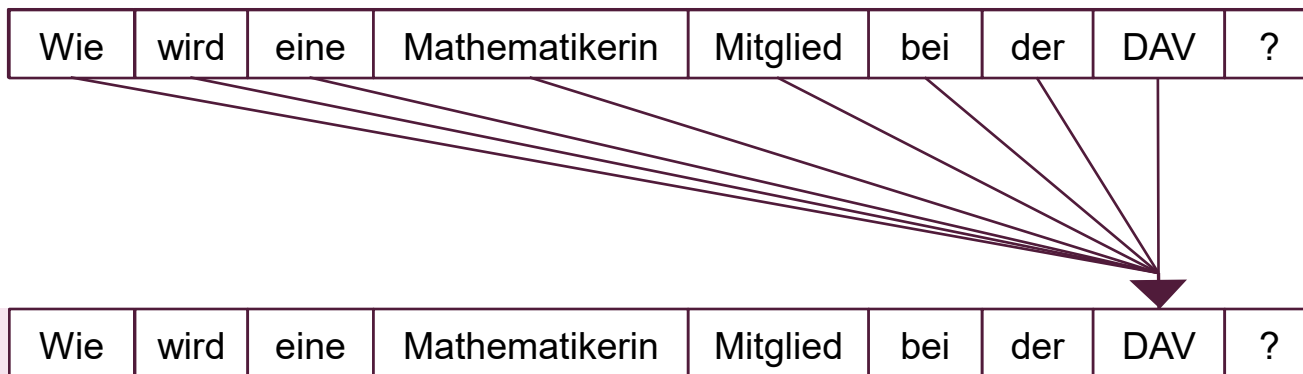
Lernbare Matrizen



Attention: Idee

Jedes Token bezieht kontextuelle Informationen von den anderen Tokens, die in der Sequenz davor liegen.

Beispielhaft für das Token **DAV** gilt:



Embeddings aus
vorherigem Block

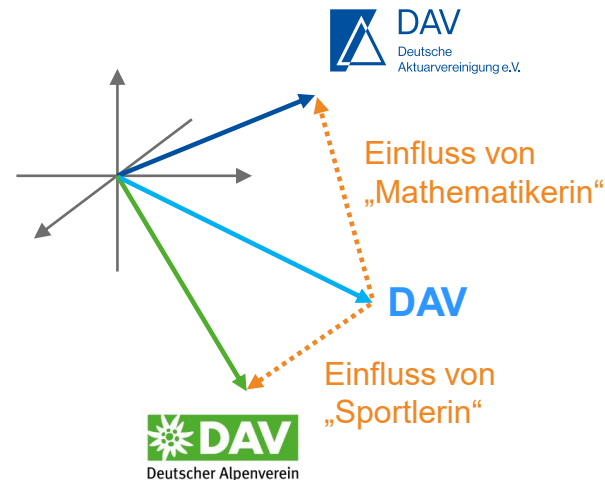
Attention im
aktuellen Block



Multi-Head Attention: Mathematische Intuition

$$Head_i \approx \sum_{j=1}^i \underbrace{\left(\underbrace{(W^Q x_i)}_{\text{DAV Vektor „Ich suche“}} \cdot \underbrace{(W^K x_j)}_{\text{Mathematikerin Vektor „Ich biete“}} \right)}_{\text{Ähnlichkeit } \in \mathbb{R}} \cdot \underbrace{(W^V x_j)}_{\text{Anpassungsrichtung } \in \mathbb{R}^{4096}}$$

→ W^Q , W^K und W^V sind lernbare Matrizen des Attention Heads



Multi-Head Attention: Lineare Kombination mehrerer Heads (mit jeweils eigenen Matrizen), um der Komplexität der natürlichen Sprache gerecht zu werden.

Tokenisierung

Vokabular ca. vom Umfang
 $|V| = 50 \text{ Tsd.}$

1	2	3	4	5	6	7	8	9
Wie	wird	eine	Mathematikerin	Mitglied	bei	der	DAV	?

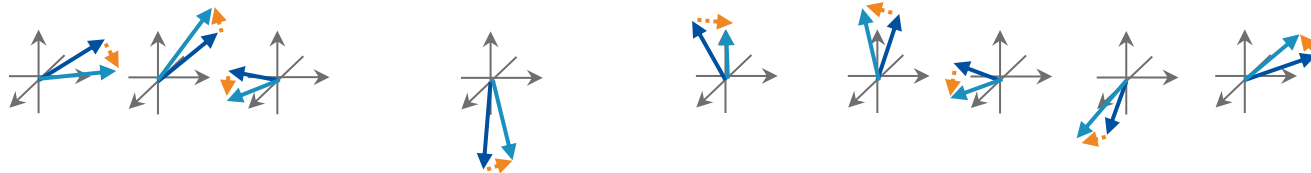
Embedding

Finden geeigneter
Repräsentationen im \mathbb{R}^{4096}

$\begin{pmatrix} 0,0107 \\ 0,0043 \\ \vdots \\ -0,0312 \end{pmatrix}$	$\begin{pmatrix} -0,0020 \\ 0,0237 \\ \vdots \\ -0,0194 \end{pmatrix}$	$\begin{pmatrix} 0,0483 \\ 0,0456 \\ \vdots \\ 0,0240 \end{pmatrix}$	$\begin{pmatrix} -0,0067 \\ -0,0348 \\ \vdots \\ -0,0082 \end{pmatrix}$	$\begin{pmatrix} 0,0236 \\ -0,0185 \\ \vdots \\ 0,0025 \end{pmatrix}$	$\begin{pmatrix} 0,0002 \\ 0,0159 \\ \vdots \\ -0,0109 \end{pmatrix}$	$\begin{pmatrix} 0,0489 \\ 0,0213 \\ \vdots \\ 0,0238 \end{pmatrix}$	$\begin{pmatrix} 0,0190 \\ 0,004 \\ \vdots \\ -0,0243 \end{pmatrix}$	$\begin{pmatrix} 0,0081 \\ -0,0033 \\ \vdots \\ 0,0194 \end{pmatrix}$
---	--	--	---	---	---	--	--	---

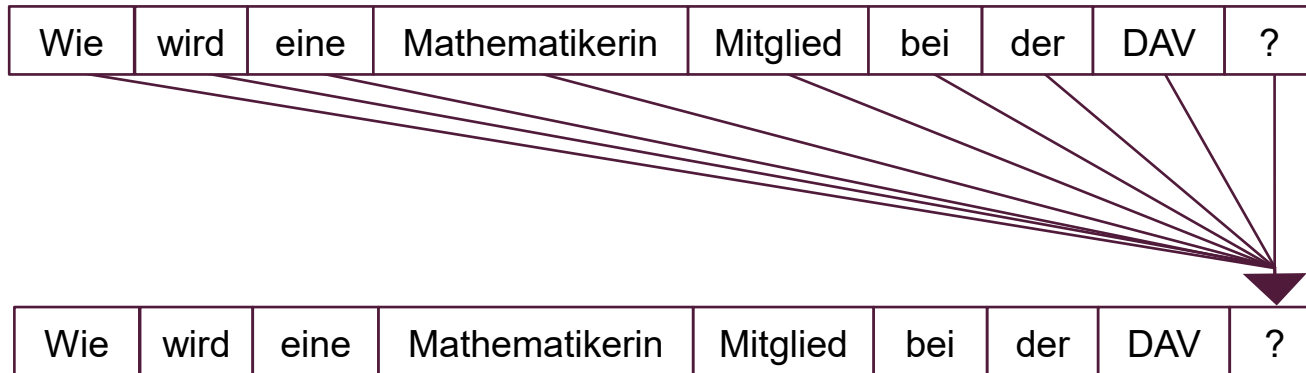
Transformerblöcke

Sukzessive Transformation der
Embeddings im \mathbb{R}^{4096}



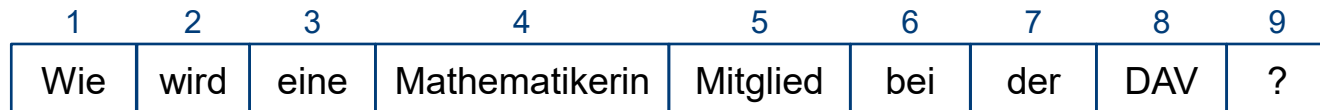
Transformerblöcke: Output

Das letzte Token bezieht kontextuelle Informationen von **allen** anderen Tokens. Es genügt daher, aus dem Output des letzten Transformerblocks nur das kontextuelle Embedding des letzten Tokens weiterzuverarbeiten.



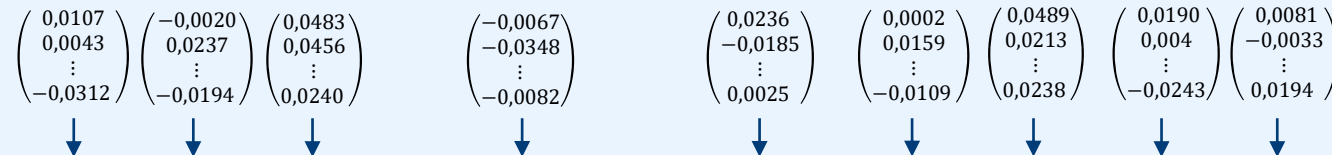
Tokenisierung

Vokabular ca. vom Umfang
 $|V| = 50 \text{ Tsd.}$



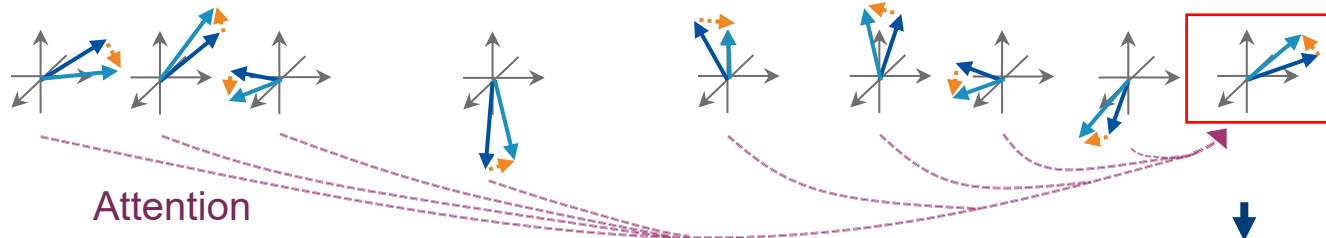
Embedding

Finden geeigneter
Repräsentationen im \mathbb{R}^{4096}



Transformerblöcke

Sukzessive Transformation der
Embeddings im \mathbb{R}^{4096}



Language Modeling Head

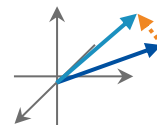
Unembedding

Projektion des „komprimierten“ Vektors im \mathbb{R}^{4096} zurück auf den größeren Raum $\mathbb{R}^{|V|} = \mathbb{R}^{50.000}$ mittels trainierbarer Unembedding-Matrix $U \in \mathbb{R}^{50.000 \times 4096}$.

Softmax-Funktion

Ermittlung der Wahrscheinlichkeitsverteilung über alle Tokens für das nächste zu wählende Token bedingt auf bisherige Token-Sequenz als

$$p_i = \frac{\exp(u_i)}{\sum_{j=1}^{|V|} \exp(u_j)} \quad \forall i \in 1, \dots, |V|.$$



$$h_9 := \begin{pmatrix} 0,0271 \\ -0,4032 \\ \vdots \\ 0,0182 \end{pmatrix} \in \mathbb{R}^{4096}$$

$$u := \begin{pmatrix} -0,3578 \\ 4,1241 \\ \vdots \\ -2,0765 \end{pmatrix} \in \mathbb{R}^{|V|}$$

$$p := \begin{pmatrix} 0,0001 \\ 0,1240 \\ \vdots \\ 0,0000 \end{pmatrix} \begin{array}{l} \text{Mathematiker} \\ \text{Um} \\ \text{Sonnig} \end{array}$$

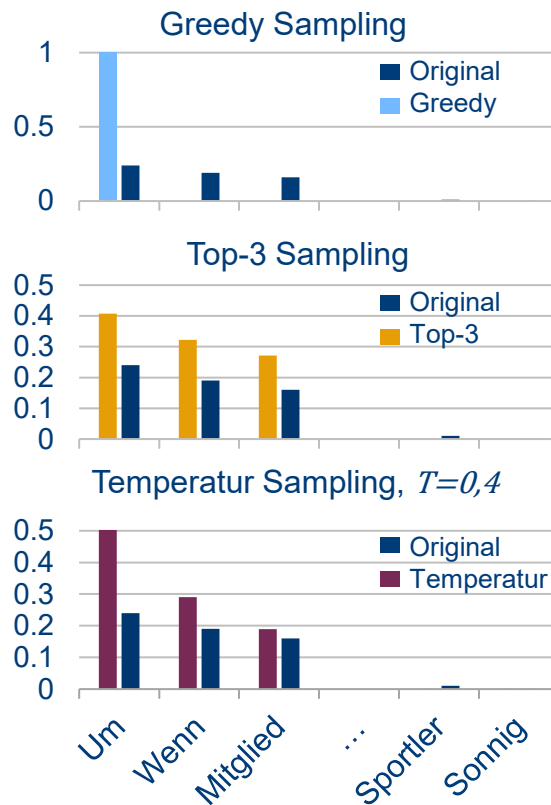
Sampling

Greedy Sampling: Wähle das Token mit der höchsten Wahrscheinlichkeit.

Top-k Sampling ($k \in \mathbb{N}$): Ziehe aus der Menge der k Tokens mit den k höchsten Wahrscheinlichkeiten gemäß ihrer (reskalierten) Wahrscheinlichkeiten.

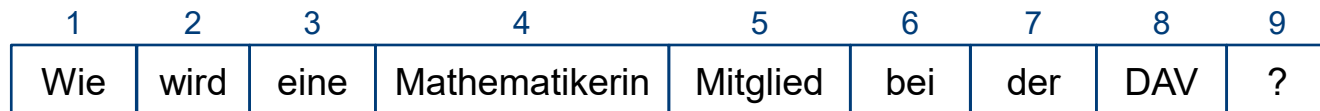
Temperatur Sampling: Modifiziere die Softmax-Funktion mittels des Temperatur-Parameters $T \in \mathbb{R}^+$ als

$$p_i = \frac{\exp\left(\frac{u_i}{T}\right)}{\sum_{j=1}^{|V|} \exp\left(\frac{u_j}{T}\right)} \quad \forall i \in 1, \dots, |V|.$$



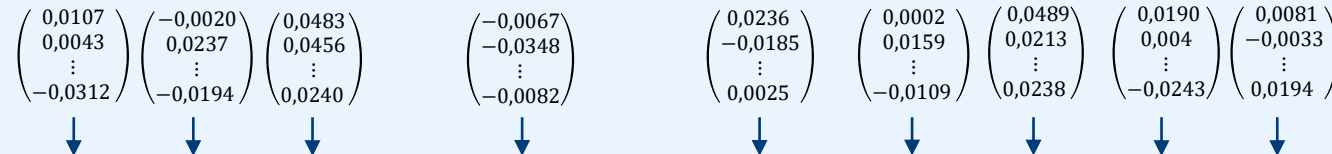
Tokenisierung

Vokabular ca. vom Umfang
 $|V| = 50 \text{ Tsd.}$



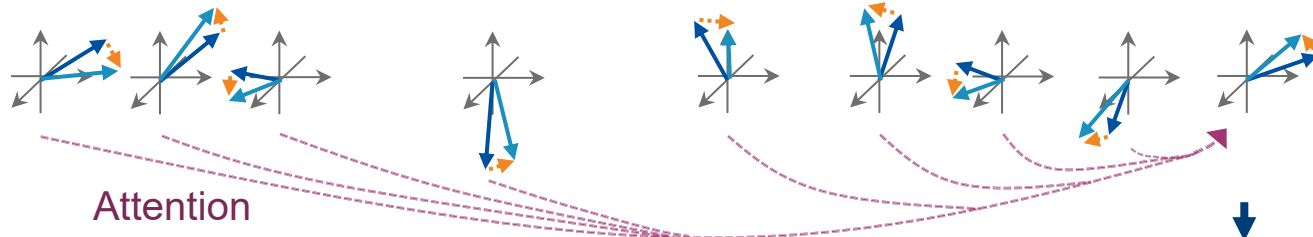
Embedding

Finden geeigneter
Repräsentationen im \mathbb{R}^{4096}



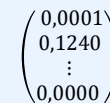
Transformerblöcke

Sukzessive Transformation der
Embeddings im \mathbb{R}^{4096}



Language Modeling Head

Projektion in $\mathbb{R}^{|V|}$, Umwandlung
in Wahrscheinlichkeiten



Sampling



Take Home Messages

Architektur der Sprachmodelle:

- Basiert auf überraschend einfachen mathematischen Methoden der Linearen Algebra (Vektoren, Matrixmultiplikation), Stochastik (Sampling) und Analysis (Gradientenabstieg im Training).
- Durchbruch im letzten Jahrzehnt durch den Attention-Mechanismus: „Mathematikerin“ im Beispielsatz beeinflusst die Bedeutung von „DAV“.
- Neueste Modelle besitzen mehr als 1 Billion Parameter: Extrem eingeschränkte Interpretierbarkeit.

Hinweise zur Anwendung im Aktuarat:

- **Don't:** Anwendungen, bei denen exakte, reproduzierbare und nachvollziehbare Antworten nötig sind.
- **Do:** Text-Zusammenfassungen, Modell-Dokumentationen, Code-Erklärungen und Code-Erstellung.

Zeit für Fragen!

Für Interessierte:

Fit4AI kompakt

Nächste Termine:

21.-22.09.2026

13.-14.10.2026

03.-04.11.2026

07.-08.12.2026

[Veranstaltungen für Aktuare -
Deutsche Aktuarvereinigung e.V.](#)

Speech and Language Processing

An Introduction to Natural Language Processing,
Computational Linguistics, and Speech Recognition
with Language Models

Third Edition draft

Daniel Jurafsky
Stanford University

James H. Martin
University of Colorado at Boulder

Copyright ©2026. All rights reserved.

Draft of January 6, 2026. Comments and typos welcome!

<https://web.stanford.edu/~jurafsky/slp3>

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez*[†]
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin*[†]
illia.polosukhin@gmail.com

<https://arxiv.org/pdf/1706.03762>

**Vielen Dank für Ihre
Aufmerksamkeit.**
