

# Word Embeddings und genderlinguistische Forschung

**Dominic Schmitz** 

Heinrich-Heine-Universität Düsseldorf

GENELLI Projektgruppe 20. Oktober 2025

## **Fahrplan**

- 1. Generelles zu Word Embeddings
- 2. Anwendungsbeispiele

## **Word Embeddings**

Harris (1954)

"Words that occur in similar contexts tend to have similar meanings."

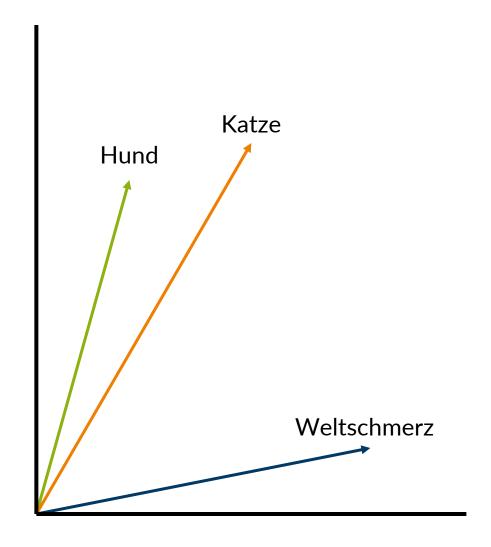
- bedeutungsgleiche oder -ähnliche Wörter haben ähnliche syntaktische Umgebungen
- Firth (1957)

"You shall know a word by the company it keeps."

- Bedeutung ergibt sich aus regelmäßiger Kookurrenz
- Harris lieferte die distributionelle Hypothese Firth die prägnanteste
   Formulierung und die korpuslinguistische Brücke dazu

## **Vektorraum-Modell**

- Wortbedeutung wird als Vektor repräsentiert
- dabei kodieren Vektoren die statistische Verteilung eines Wortes über relevante sprachliche Kontexte
- Vektoren = Wegweiser zu
   Punkten im semantischen Raum
- semantische Ähnlichkeit =Distanz zwischen Vektoren



#### Kookurrenz

- über welchen Kontext hinweg wird Kookurrenz erfasst?
- das ist von Modell zu Modell und Implementierung zu Implementierung verschieden
- Beispiele
  - Wörter im Satz, Absatz, Dokument
  - Wörter in einem festen Wortfenster
  - Wörter in bestimmten syntaktischen Beziehungen

• etc.

Die Äpfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Äpfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Apfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h2ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Die Apfel (Malus) bilden eine Pflanzengattung der Kernobstgewächse (Pyrinae) aus der Familie der Rosengewächse (Rosaceae). Nicht zu verwechseln mit den Äpfeln sind die nicht näher verwandten Granatäpfel (Punica granatum). Das Wort Apfel wird auf die indogermanische Grundform \*h₂ébōl zurückgeführt. Der wissenschaftliche Gattungsname Malus ist abgeleitet von dem lateinischen Wort malum, was auf Deutsch so viel wie Apfel oder apfelförmige Frucht bedeutet. Die Arten der Gattung Äpfel (Malus) sind sommergrüne Bäume oder Sträucher. Die Frucht des Apfelbaumes wird Apfel genannt. Der Blühbeginn des Apfels markiert im phänologischen Kalender den Beginn des Vollfrühlings.

Spalten: Eigenschaften des Vokabulars

	Malus	verwandt	Granatapfel	Frucht	Baum	•••
Apfel	2	1	1	2	1	•••

Zeilen: Wörter, an denen wir interessiert sind

Spalten: Eigenschaften des Vokabulars

	Malus	verwandt	Granatapfel	Frucht	Baum	•••
Apfel	2	1	1	2	1	•••
Baum	1	0	0	1	0	•••

Zeilen: Wörter, an denen wir interessiert sind

## Apfel, Baum und Birne

	Malus	verwandt	Granatapfel	Frucht	Baum	•••
Apfel	2	1	1	2	1	•••
Baum	1	0	0	1	0	•••

• 
$$\vec{v}_{Apfel} = \langle 2,1,1,2,1,... \rangle$$

• 
$$\vec{v}_{Baum} = \langle 1, 0, 0, 1, 0, ... \rangle$$

	Malus	verwandt	Granatapfel	Frucht	Baum	•••
Birne	1	1	0	2	1	•••

•  $\vec{v}_{Birne} = \langle 0, 1, 0, 2, 1, ... \rangle$ 

#### Apfel, Baum und Birne

- ein Standardmaß für die Ähnlichkeit zweier Vektoren ist der Kosinus des Winkels zwischen den Vektoren, die Kosinus-Ähnlichkeit
  - Kosinus = 1 → Vektoren sind identisch
  - Kosinus = 0 → Vektoren sind orthogonal (rechtwinkling)
- weitere Maße
  - euklidische Distanz: kleiner Wert = große Ähnlichkeit
  - Korrelation: hoher Wert = große Ähnlichkeit

#### Kosinus-Ähnlichkeit

	Baum	Birne
Apfel	0.85	0.91

#### euklidische Distanz

	Baum	Birne
Apfel	2.24	1.41

## Entscheidungen und Varianten

- Vorverarbeitung
  - Wortformen vs. Lemmata
  - mit oder ohne Funktionswörter
  - etc.
- Kontextart: Wortfenster, Textteil, etc.
- Gewichte: Häufigkeiten, Wahrscheinlichkeiten, etc.
- ...und damit man nicht alles selbst kodieren muss: welches Modell?

#### word2vec

- entwickelt von Mikolov et al. (2013)
- modelliert die Wahrscheinlichkeit, dass ein Wort in einem bestimmten Kontext vorkommt
- d.h. Bedeutung wird aus wiederkehrenden Wortnachbarschaften erlernt
- Vektoren können in "zwei Richtungen" trainiert werden
  - CBOW (continuous bag of words): Kontext → Zielwort vorhersagen
     Der \_\_\_\_ wächst am Baum. das Modell lernt, dass Apfel hier am besten passt
  - Skip-Gram: Zielwort → Kontextwörter vorhersagen
     Apfel das Modell lernt, dass passende Kontexte Baum, Frucht sind
- CBOW ist bei großen Korpora eindeutig schneller

#### fastText

- entwickelt von Bojanowski et al. (2016)
- baut auf word2vec auf und erweitert es um "subword"-Informationen
  - hierdurch kann das Modell auch Vektoren für seltene und unbekannte Wortformen erstellen
- Wörter bestehen aus Buchstabenfolgen, deren Vektoren kombiniert werden
  - Apfel = <Ap, pf, fe, el>
  - *Birne* = <Bi, ir, rn, ne>
  - Baum = <Ba, au, um>

## GloVe

- entwickelt von Pennington et al. (2014)
- kombiniert lokale Kontexte mit globalen Häufigkeiten
  - erstellt eine Kookurrenzmatrix (Wörter x Kontexte) mit log-transformierten
     Häufigkeiten
  - das Modell lernt, diese Matrix mit möglichst geringem Fehler zu rekonstruieren
- Apfel und Birne treten ähnlich oft mit Baum, Frucht, süß, Markt auf
  - → große Nähe
- Baum teilt globale Kontexte wie Natur, Obst, Wald
  - → moderate Nähe

## **Naive Discriminative Learning (NDL)**

- entwickelt von Baayen et al. (2011), basiert auf psychologischer
   Grundlage von Rescorla & Wagner (1972)
- Lernen = Aufbau und Anpassung von assoziativen Gewichten zwischen
   Cues (Spalten) und Outcomes (Zeilen)
- kein Vorhersagemodell, sondern Lernsimulation
  - wenn Apfel mit Baum auftritt → Gewichtungen verstärkt
  - wenn Baum ohne Apfel vorkommt  $\rightarrow$  Gewichtung abgeschwächt

## Anwendungsbeispiele

aus der genderlinguistischen Forschung

#### **Semantischer Wandel**

Sökefeld, C., & Amaral, P. (2025). Semantic change of female-denoting nouns in diachronic German corpora. Proceedings of the Linguistic Society of America, 10(1), 5905. https://doi.org/10.3765/plsa.v10i1.5905

#### Fragestellung

Wie verändert sich die Semantik weiblich bezeichnender Nomen im Deutschen über die Jahrhunderte?

#### Methode

- Trainieren mehrerer word2vec-Modelle basierend auf Texten von 1350 bis 1899, je ein Modell pro Zeitraum
- In jedem Zeitraum werden dann für die Zielwörter die "nächsten Nachbarn" gefunden und verglichen

#### **Semantischer Wandel**

Sökefeld, C., & Amaral, P. (2025). Semantic change of female-denoting nouns in diachronic German corpora. *Proceedings of the Linguistic Society of America*, 10(1), 5905. doi.org/10.3765/plsa.v10i1.5905

#### Ergebnisse

Rang	frühes 16. Jhd	spätes 17. Jhd
1	reyne	knecht
2	gebärmutter	gesinde
3	latona (Name)	dirne
4	hertsens	weib
5	jungfraw	rizarize (Name)
6	färber	frau

Magd als ,Dienstmädchen'

Magd als ,Mädchen, Jungfrau'

Schmitz, D., Schneider, V., & Esser, J. (2023). No genericity in sight: An exploration of the semantics of masculine generics in German. *Glossa Psycholinguistics*, 2(1): 12, pp. 1–33. doi.org/10.5070/G6011192

#### Fragestellung

Ist die Semantik s.g. generischer Maskulina trotz ihrer Formgleichheit mit spezifischen Maskulina unbiased?

#### Methode

- Wortform-zu-Wortform-Matrix wird mit NDL basierend auf einem Korpus deutscher Nachrichtentexte trainiert
- zusätzlich zu Wortformen sind auch GENUS, GENERIZITÄT und NUMERUS als Cues (Spalten) eingeführt

Schmitz, D., Schneider, V., & Esser, J. (2023). No genericity in sight: An exploration of the semantics of masculine generics in German. *Glossa Psycholinguistics*, 2(1): 12, pp. 1–33. doi.org/10.5070/G6011192

#### Methode

- Zielformen, d.h. generische und spezifische Maskulina sowie Feminina, lassen sich anhand der *NDL*-Vektoren berechnen
- $\vec{v}_{W\_generisch} = \vec{v}_W + \vec{v}_{generisch}$
- $\vec{v}_{W\_spezifisch} = \vec{v}_W + \vec{v}_{spezifisch}$
- $\vec{v}_{W\_generisch\_mask} = \vec{v}_W + \vec{v}_{generisch} + \vec{v}_{mask}$
- usw.

Schmitz, D., Schneider, V., & Esser, J. (2023). No genericity in sight: An exploration of the semantics of masculine generics in German. *Glossa Psycholinguistics*, 2(1): 12, pp. 1–33. doi.org/10.5070/G6011192

#### Ergebnisse

	Numerus	spezifisches Maskulinum	Femininum
generisches	Singular	0.996	0.934
Maskulinum	Plural	0.991	0.822
spezifisches	Singular		0.939
Maskulinum	Plural		0.835

ähnlicher dem Mask als dem Fem

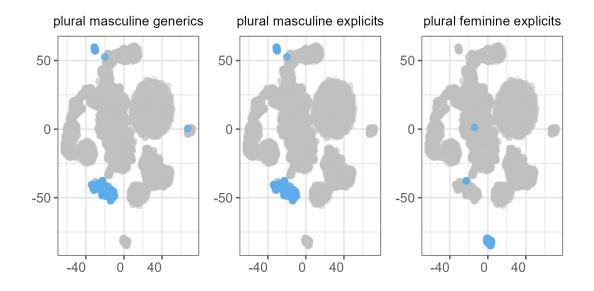
höchste Ähnlichkeit

#### **Einschub: t-SNE**

- t-Distributed Stochastic Neighbor Embedding (van der Maaten & Hinton, 2008)
- Ziel: Hochdimensionale Daten auf 2 oder 3 Dimensionen reduzieren, ohne die lokale Struktur zu verlieren
- Prinzip
  - für jedes Wort wird berechnet, welche anderen Wörter in seiner Nähe liegen
  - t-SNE versucht, eine 2D-Darstellung zu finden, in der diese Nachbarschaften ähnlich bleiben
  - Ergebnis: ähnliche Wörter bilden Cluster, unähnliche liegen weiter auseinander

Schmitz, D., Schneider, V., & Esser, J. (2023). No genericity in sight: An exploration of the semantics of masculine generics in German. *Glossa Psycholinguistics*, 2(1): 12, pp. 1–33. doi.org/10.5070/G6011192

#### Ergebnisse



Maskulina "leben" im gleichen Bereich, Feminina woanders

## Einschub: Homophonie, Polysemie

Die in einer **Batterie** gespeicherte elektrische Ladung...

Die Batterie ist bei der Artillerie der Bundeswehr...

90 Prozent der Hühner werden in **Batterien** gehalten...

Er sitzt im Park auf einer **Bank**...

Sie holt eine große Summe Geld in ihrer **Bank** ab...

Alle **Lehrer** sind schon im Raum...

Alle **Lehrer** haben ihren Wehrpflichtdienst bereits hinter sich...

Alle **Lehrer** sind schon im Raum; die Lehrerinnen stehen noch davor...

Problem: Generell erzeugen Algorithmen pro Wortform ein Embedding

## Einschub: Homophonie, Polysemie

- Lösungsansätze
  - Clustering von Kontextvektoren (Huang et al., 2012)
  - Sense-specific embeddings (Neelakantan et al., 2014)
  - Contextual clustering (Arora et al., 2018)
  - lexikalisch informierte Modelle (lacobacci et al., 2015)
  - Instance Vectors (Schmitz, 2024)
  - Contextualised embedding (Schmitz et al., in prep)

Schmitz, D. (2024). Instances of bias: The gendered semantics of generic masculines in German revealed by instance vectors. Zeitschrift für Sprachwissenschaft, 43(2). doi.org/10.1515/zfs-2024-2010

#### Fragestellung

Ist die Semantik s.g. generischer Maskulina trotz ihrer Formgleichheit mit spezifischen Maskulina unbiased?

#### Methode

- mit fastText generierte Vektoren für alle Wortformen eines Korpus
- Vektoren für Zielwörter berechnet als Durchschnitt der Vektoren der 2/5/8
   Wörter vor und nach dem Zielwort

Schmitz, D. (2024). Instances of bias: The gendered semantics of generic masculines in German revealed by instance vectors. Zeitschrift für Sprachwissenschaft, 43(2). doi.org/10.1515/zfs-2024-2010

#### Ergebnisse

	+/- 2	+/- 5	+/- 8
generisches Mask spezifisches Mask	0.415	0.590	0.776
generisches Mask spezifisches Fem	0.377	0.566	0.632
spezifisches Mask spezifisches Fem	0.525	0.726	0.583

gen vs spez
Mask
immer ähnlicher
als
gen Mask vs Fem

## Einschub: Homophonie, Polysemie

- Lösungsansätze
  - Clustering von Kontextvektoren (Huang et al., 2012)
  - Sense-specific embeddings (Neelakantan et al., 2014)
  - Contextual clustering (Arora et al., 2018)
  - lexikalisch informierte Modelle (lacobacci et al., 2015)
  - Instance Vectors (Schmitz, 2024)
  - Contextualised embedding (Schmitz et al., in prep)

Schmitz, D., Ochs, S., Müller-Spitzer, C., & Rüdiger, J. O. (in prep). The male bias of generic masculines remains stable independent of the amount of provided context.

#### Fragestellung

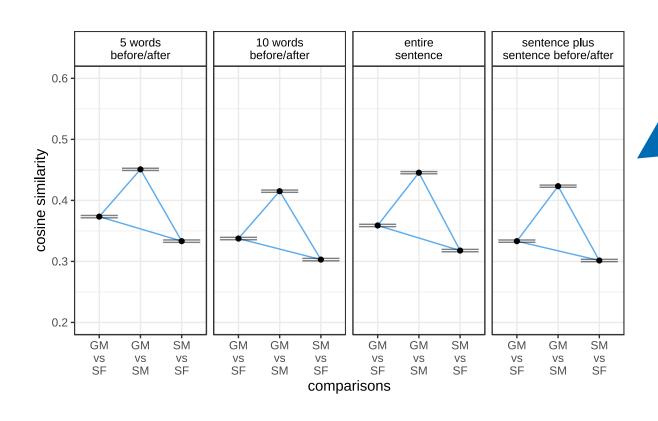
Ist die Semantik s.g. generischer Maskulina trotz ihrer Formgleichheit mit spezifischen Maskulina unbiased?

#### Methode

- mit BERT (Bidirectional Encoder Representations from Transformers; Devlin et al., 2018) generierte Vektoren
- BERT generiert pro Token einen Vektor
- zur Generierung wird Kontext bereitgestellt; die Menge kann frei variiert werden (ganzer Satz, ganzer Absatz, n Wörter, etc.)

Schmitz, D., Ochs, S., Müller-Spitzer, C., & Rüdiger, J. O. (in prep). The male bias of generic masculines remains stable independent of the amount of provided context.

#### Ergebnisse



gen vs spez
Mask
immer ähnlicher
als
gen Mask vs Fem

## Zusammenfassung

- Word Embeddings
  - sind Vektorrepräsentationen von Wortbedeutungen, positionieren Wörter nach ihrer semantischen Ähnlichkeit im Raum
  - können mit verschiedenen Algorithmen und verschiedenen Ausgangsdaten generiert werden, dabei sollte die Wahl der Ausgangsdaten und des jeweiligen Algorithmus theoretisch motiviert geschehen
  - bilden Gebrauch, nicht "objektive" Bedeutung ab

Danke fürs Zuhören!

#### Literatur 1/2

- Arora, S., Li, Y., Liang, Y., Ma, T., & Risteski, A. (2018). Linear Algebraic Structure of Word Senses, with Applications to Polysemy. *Transactions of the Association for Computational Linguistics*, 6, 483–495.
- Baayen, R. H., Milin, P., Đurđević, D. F., Hendrix, P., & Marelli, M. (2011). An amorphous model for morphological processing in visual comprehension based on naive discriminative learning. *Psychological Review*, 118(3), 438–481.
- **Bojanowski**, P., **Grave**, E., **Joulin**, A., & **Mikolov**, T. (2016). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, *5*, 135–146.
- **Devlin**, J., **Chang**, M. W., **Lee**, K., & **Toutanova**, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL HLT 2019 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies Proceedings of the Conference, 1, 4171–4186.
- Firth, J.R. (1957). A synopsis of linguistic theory 1930-1955. Studies in Linguistic Analysis: 1–32. London: Longman.
- Harris, Z. S. (1954). Distributional structure. WORD, 10(2-3), 146-162.
- Huang, E., Socher, R., Manning, C., & Ng, A. (2012). Improving Word Representations via Global Context and Multiple Word Prototypes. In H. Li, C.-Y. Lin, M. Osborne, G. G. Lee, & J. C. Park (Eds.), *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers) (pp. 873–882). Association for Computational Linguistics.
- lacobacci, I., Pilehvar, M. T., & Navigli, R. (2015). SensEmbed: Learning Sense Embeddings for Word and Relational Similarity. In C. Zong & M. Strube (Eds.), Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers) (pp. 95–105). Association for Computational Linguistics.

Maaten, L. van der, & Hinton, G. (2008). Visualizing Data using t-SNE. Journal of Machine Learning Research, 9(86), 2579–2605.

#### Literatur 2/2

- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. 1st International Conference on Learning Representations, ICLR 2013 Workshop Track Proceedings.
- **Neelakantan**, A., **Shankar**, J., **Passos**, A., & **McCallum**, A. (2015). *Efficient Non-parametric Estimation of Multiple Embeddings per Word in Vector Space* (No. arXiv:1504.06654). arXiv.
- Pennington, J., Socher, R., & Manning, C. D. (2014). GloVe: Global vectors for word representation.
- **Rescorla**, R. A., & **Wagner**, A. R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.), *Classical conditioning II: Current research and theory* (pp. 64–99). Appleton-Century-Crofts.
- **Schmitz**, D. (2024). Instances of bias: The gendered semantics of generic masculines in German revealed by instance vectors. *Zeitschrift für Sprachwissenschaft*, 43(2).
- **Schmitz**, D., **Ochs**, S., **Müller-Spitzer**, C., & **Rüdiger**, J. O. (in prep). The male bias of generic masculines remains stable independent of the amount of provided context.
- Schmitz, D., Schneider, V., & Esser, J. (2023). No genericity in sight: An exploration of the semantics of masculine generics in German. *Glossa Psycholinguistics*, 2(1): 12, pp. 1–33.
- Sökefeld, C., & Amaral, P. (2025). Semantic change of female-denoting nouns in diachronic German corpora. *Proceedings of the Linguistic Society of America*, 10(1), 5905.