

# KI generierte Texte erkennen

Dürrmeier,Wagner

19. Januar 2024

1 Eingrenzungen

2 Funktion

3 Merkmale

4 Detektion

5 Probleme

6 Lösung

- Behandeln von flüssigen Texten
- Konzentration auf den Lehrkontext
- nur generative Modelle bis Mitte 2023, da für spätere Modelle noch nicht genügend Studienlage vorhanden ist
- bei allen Erkennungsmethoden und Merkmalen gilt grundlegend: je länger/größer der Text desto genauer

- **Natural Language Processing:** komplette Maschine
- **Natural Language Understanding:** zieht sich die benötigten Informationen aus eingegebenen Text
- **Natural Language Generation:** nutzt ein Language Model (LM) um Antwort in natürlicher Sprache zu generieren

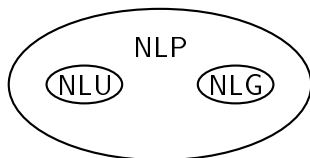


Abbildung: nach Maddala; 2022; *Medium* „NLP, NLU, and NLG“

- **neuronale LMs:**<sup>1</sup>neuronale Netze und Maschine Learning
- **wissensbasierte LMs:**<sup>1</sup>linguistisches Wissen
- **statistische LMs:**<sup>1</sup>Auftrittswahrscheinlichkeiten und Wahrscheinlichkeitsverteilungen

---

<sup>1</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

- $P(w)$  beschreibt Wahrscheinlichkeit, dass auf ein Wort Token  $w$  folgt
- $V$  ist die Menge aller möglich folgenden Tokens
- **Top- $k$ -Sampling:**<sup>2</sup> Auswahl von  $k$  Tokens, die zusammenaddiert eine möglichst hohe Wahrscheinlichkeit ergeben ( $\sum w \in V | P(w)$ ) und zufällige Auswahl aus der erstellten Menge

---

<sup>2</sup>Rieke; 2023; [10.13140/RG.2.2.35567.00163.]

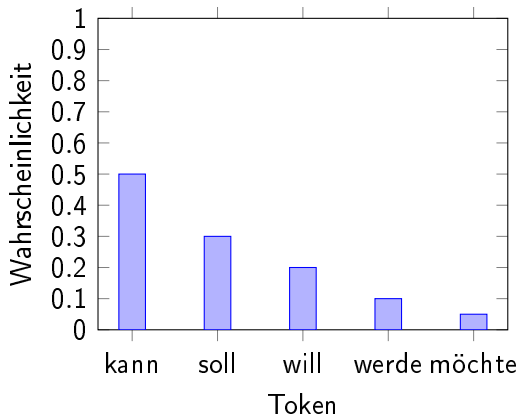
- $P(w)$  beschreibt Wahrscheinlichkeit, dass auf ein Wort Token  $w$  folgt
- $V$  ist die Menge aller möglich folgenden Tokens
- **Top- $k$ -Sampling:**<sup>2</sup> Auswahl von  $k$  Tokens, die zusammenaddiert eine möglichst hohe Wahrscheinlichkeit ergeben ( $\sum w \in V | P(w)$ ) und zufällige Auswahl aus der erstellten Menge
- **Top- $p$ -Sampling:**<sup>2</sup> möglichst kleine Menge von Tokens, die den Wert  $p$  nicht überschreiten ( $\sum w \in V | P(w) \geq p$ ) und zufällige Auswahl aus der erstellten Menge

---

<sup>2</sup>Rieke; 2023; [10.13140/RG.2.2.35567.00163.]

# Beispiel Top- $k$ -Sampling<sup>3</sup>

- Top- $k$ -Sampling; Auswahl des nächsten Tokens nach „Ich“
- mit  $k := 3$

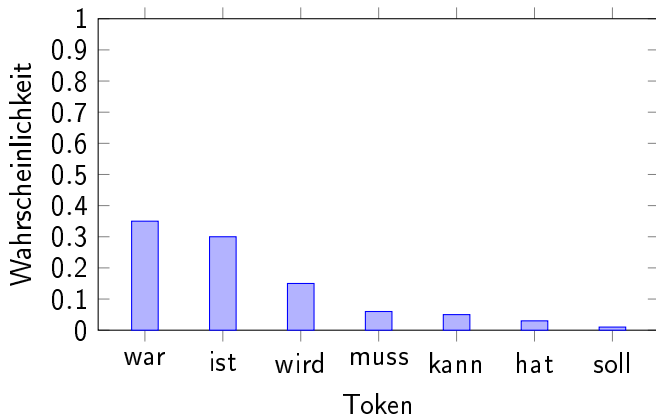


<sup>3</sup>nach Rieke; 2023; [10.13140/RG.2.2.35567.00163.]



# Beispiel Top- $p$ -Sampling<sup>4</sup>

- Top- $p$ -Sampling; Auswahl des nächsten Tokens nach „Es“
- mit  $p := 0.9$



<sup>4</sup>nach Rieke; 2023; [10.13140/RG.2.2.35567.00163.]

- **Häufigkeitsmerkmale:**<sup>5</sup> bei Menschen eine Verteilung nach dem Zipfschen Gesetz

$$f \approx \frac{\frac{1}{k^s}}{\sum_{n=1}^N \frac{1}{n^s}}$$

(mit  $s \in \mathbb{R} | s \geq 1$  und mit der Häufigkeit  $f$  eines Tokens des Ranges  $k$  aus  $N$  unterschiedlichen Tokens)<sup>6</sup>

---

<sup>5</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>6</sup>Zipf; 1949; „Human behavior and the principle of least effort.“

<sup>7</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

- **Häufigkeitsmerkmale:**<sup>5</sup> bei Menschen eine Verteilung nach dem Zipfschen Gesetz

$$f \approx \frac{\frac{1}{k^s}}{\sum_{n=1}^N \frac{1}{n^s}}$$

(mit  $s \in \mathbb{R} | s \geq 1$  und mit der Häufigkeit  $f$  eines Tokens des Ranges  $k$  aus  $N$  unterschiedlichen Tokens)<sup>6</sup>

- **Geläufigkeitsmerkmale:**<sup>7</sup> Konsistenz und Kohärenz des Textes; Maß dafür ist z. B. der Flesch-Index

---

<sup>5</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>6</sup>Zipf; 1949; „Human behavior and the principle of least effort.“

<sup>7</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

- **linguistische Merkmale:**<sup>8</sup>

- Messen von Phrasalverben und vergleichen mit Auflösung der Koreferenzen
- Messen der Verteilung von Part-of-Speech-Tags und Named-Entity-Tags

---

<sup>8</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>9</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>10</sup>Vogelgesang, et.al.;Uni Hohenheim; 2023

- **linguistische Merkmale:**<sup>8</sup>
  - Messen von Phrasalverben und vergleichen mit Auflösung der Koreferenzen
  - Messen der Verteilung von Part-of-Speech-Tags und Named-Entity-Tags
- **komplexe phrasale Merkmale:**<sup>9</sup> idiomatische Phrasen bei Menschen häufiger; neuere NLPs verwenden diese inzwischen aber auch

---

<sup>8</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>9</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>10</sup>Vogelgesang, et.al.; Uni Hohenheim; 2023

- **linguistische Merkmale:**<sup>8</sup>
  - Messen von Phrasalverben und vergleichen mit Auflösung der Koreferenzen
  - Messen der Verteilung von Part-of-Speech-Tags und Named-Entity-Tags
- **komplexe phrasale Merkmale:**<sup>9</sup> idiomatische Phrasen bei Menschen häufiger; neuere NLPs verwenden diese inzwischen aber auch
- **grundlegende Textmerkmale:**<sup>9</sup> Häufigkeit von Satzzeichen, Länge von Texten und Absätzen sowie Fehleranfälligkeit

---

<sup>8</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>9</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>10</sup>Vogelgesang, et.al.; Uni Hohenheim; 2023

- **linguistische Merkmale:**<sup>8</sup>
  - Messen von Phrasalverben und vergleichen mit Auflösung der Koreferenzen
  - Messen der Verteilung von Part-of-Speech-Tags und Named-Entity-Tags
- **komplexe phrasale Merkmale:**<sup>9</sup> idiomatische Phrasen bei Menschen häufiger; neuere NLPs verwenden diese inzwischen aber auch
- **grundlegende Textmerkmale:**<sup>9</sup> Häufigkeit von Satzzeichen, Länge von Texten und Absätzen sowie Fehleranfälligkeit
- **zeitabhängige Merkmale:**<sup>10</sup> Aussagen über den Wissensstand einer KI hinaus sind vage und ungenau bzw. falsch

---

<sup>8</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>9</sup>Crothers, Japkowicz, Viktor; 2023; [arXiv:2210.07321v4]

<sup>10</sup>Vogelgesang, et.al.; Uni Hohenheim; 2023

- Detektion auch in Kombination von Mensch und Maschine möglich, aber bisher zu wenig Studienlage
- Vorstellen der zwei erfolgreicherer maschinellen Möglichkeiten



- **Zero-Shot-Ansatz:**

- kleiner NLG wird auf Erkennungsmerkmale trainiert <sup>11</sup>
- Beispiel: GROVER für Fake-News-Erkennung
- vergleichsweise ineffizient<sup>12</sup>

---

<sup>11</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>12</sup>Solaiman, Brundage, et al.; 2019

<sup>13</sup>Liu, Ott, et al.; 2019; [arXiv:1907.11692]

- **Zero-Shot-Ansatz:**

- kleiner NLG wird auf Erkennungsmerkmale trainiert <sup>11</sup>
- Beispiel: GROVER für Fake-News-Erkennung
- vergleichsweise ineffizient<sup>12</sup>

- **Fine-Tuning-Ansatz:**<sup>13</sup>

- NLP mit Texten von Menschen und KIs als Trainingsdatensatz
- Beispiel: RoBERTa für ChatGPT-2 Erkennung
- im Allgemeinen sehr effizient (ja nach Trainingsatz)

---

<sup>11</sup>Crothers, Japkowicz, et al.; 2022; [arXiv:2203.07983 ]

<sup>12</sup>Solaiman, Brundage, et al.; 2019

<sup>13</sup>Liu, Ott, et al.; 2019; [arXiv:1907.11692]

- Entscheidung wird auf Basis von semantischen Ungenauigkeiten getroffen<sup>14</sup>
- mit Experten erhält man eine Erkennungsgenauigkeit von ca. 74%<sup>14</sup>
- nicht trainierte Menschen erkennen im Regelfall nicht KI-Texte<sup>15</sup>

---

<sup>14</sup>Ippolito, Duckworth, et al.; *ACL 2020 Annual Conference*

<sup>15</sup>Albrecht; 2023; *API*

- **Überarbeitung:**<sup>16</sup>
  - Überarbeiten von KI geschriebener Texte durch Menschen im Nachhinein
  - linguistisch basierte Merkmale fallen dadurch weg
  - KIs fehlen so wichtige Merkmale

---

<sup>16</sup>Ippolito, Duckworth, et al.; *ACL 2020 Annual Conference*

<sup>17</sup>Vogelgesang, et al.; Uni Hohenheim; 2023

<sup>18</sup>Satzung Universität Tübingen; 2023

<sup>19</sup>Albrecht; 2023; *API*

- **Überarbeitung:**<sup>16</sup>
  - Überarbeiten von KI geschriebener Texte durch Menschen im Nachhinein
  - linguistisch basierte Merkmale fallen dadurch weg
  - KIs fehlen so wichtige Merkmale
- **Rechtssicherheit:**<sup>17</sup>
  - Detektionsmethoden liefern keine zuverlässige Aussagen
  - also selbst bei vermuteter KI-Autorenschaft kein Grund zum Nicht-Bestehen
  - Universität Tübingen verbietet in ihrer Satzung<sup>18</sup> die Nutzung von KIs im Lehrkontext → können durch fehlende Beweisbarkeit nicht durchsetzen
- Verbessern der Situation in Zukunft nicht zu erwarten<sup>19</sup>

---

<sup>16</sup>Ippolito, Duckworth, et al.; *ACL 2020 Annual Conference*

<sup>17</sup>Vogelgesang, et al.; Uni Hohenheim; 2023

<sup>18</sup>Satzung Universität Tübingen; 2023

<sup>19</sup>Albrecht; 2023; *API*

- KIs zulassen und Aufgaben komplex genug machen (z.B. ausdrückliche kritische Auseinandersetzung mit einem Thema fordern)<sup>20</sup>
- KIs ausschließen, indem man neue Themen abfragt, und erörtern lässt<sup>21</sup>

---

<sup>20</sup>Albrecht; 2023; *API*

<sup>21</sup>Vogelgesang, et al.; Universität Hohenheim; 2023

Abschließend möchte ich betonen, dass die Welt der AI-generierten Texte manchmal so raffiniert ist, dass selbst mein Kaffeekocher neidisch wäre. Aber im Ernst, die Erkennung solcher Texte ist entscheidend für die digitale Hygiene unserer Informationen. Vielen Dank für eure Aufmerksamkeit! Wenn ihr jetzt denkt, dass euer Laptop heimlich Romane schreibt, stehe ich bereit, eure Fragen zu entschlüsseln.

-ChatGPT 3.5