

### Klexikon: A German Dataset for Joint Summarization and Simplification

#### **Dennis Aumiller** and Michael Gertz

Heidelberg University Institute of Computer Science Database Systems Research Group aumiller@informatik.uni-heidelberg.de





- **1. Reading a Simplified Article**
- 2. Constructing the Klexikon Dataset
- **3. Exploratory Analysis**
- 4. Future Work



### **Reading a Simplified Article**

## Reading a Simplified Article



#### UNIVERSITÄT **HEIDEI BERG ZUKUNFT** SEIT 1386

#### Expectation

#### Adler (Biologie)

Der Begriff Adler (von mittelhochdeutsch adelarie) für "Edel-Aar", wobei Aar Adler und Adlerartige (unedle Adler) bezeichnete) ist im weiteren Sinn eine nicht genau definierte Sammelbezeichnung für große und kräftige Arten von Greifvögeln in der Familie der Habichtartigen, wie beispielsweise für die Gattung der Seeadler, den Schlangenadler oder den Philippinenadler. Im engeren Sinn steht die Bezeichnung Adler für die Gattung Echte Adler (Aguila) mit den ihr zugehörigen Arten. Deren größte lebende Vertreter sind der Steinadler und der Keilschwanzadler. Außerhalb der Habichtartigen wird im Deutschen auch der in seine eigene Familie gestellte Fischadler unter die Adler gezählt.



Mit einer Spannweite bis zu 3 Meter gilt der um 1400 n. Chr. ausgestorbene Haastadler (Harpagornis moorei) als größter bekannter Adler. Sein Lebensraum befand sich auf Neuseeland.

#### Jagd und Fortpflanzung [Bearbeiten | Quelltext bearbeiten ]

Die außerordentliche Sehschärfe ihrer Augen gestattet es allen Adlern, ihre Beute auch aus großer Höhe zu erspähen und zu erjagen. Ihre Hauptbeute sind bodenlebende Wirbeltiere, in erster Linie Säuger. Nur selten ernähren sie sich von Aas. Adler sind sehr gut daran angepasst, zum Fliegen Aufwinde zu nutzen. Der Adler packt die Beute beim Niederstoßen aus dem Flug mit den Fängen. Dabei weisen drei seiner Zehen nach vorne, die Hinterkralle durchsticht den Beutekörper in einer Zangenbewegung. Durch den Griff wird die Beute erstickt

In das aus Zweigen gebaute Nest (Adlerhorst) legt das Weibchen im Regelfall zwei bis drei Eier. Häufig wird jedoch nur eines der Jungen aufgezogen. Alle Adlerarten sind geschützt, ihr Bestand ist durch menschlichen Einfluss gefährdet.

#### Kunst und Symbolik [Bearbeiten | Quelitext bearbeiten ]

Die Kunst des Alten Orients verwendet den Adler in vielfachen Abwandlungen als Symboltier. Bereits 3000 v. Chr. findet sich das Mischwesen Löwenadler in mesopotamischen Darstellungen, Doppeladler sind aus dem Babylonien des 23. Jahrhunderts v. Chr. bekannt. Die Antike kennt unzählige Abbildungen des Adlers auf Vasen, Gegenständen des täglichen Gebrauchs, Münzen und Schmuckstücken sowie in der Architektur auf Reliefs, Akroterien und Giebeln.

Seit dem 5. Jahrhundert nach Christus fertigten die Goten große Adlerfibeln und Schnallen mit Adlerköpfen, Aus der romanischen Periode des Mittelalters sind mit Adlern verzierte figürliche Darstellungen von Kapitellen bekannt. Ebenfalls hervorzuheben sind die Adlerpulte christlicher Altarräume. In der Malerei erscheint der Adler häufig als Symbol für den Evangelisten Johannes.

In den Überlieferungen, dem Mythos und dem Volksglauben der frühen Hochkulturen (Ägypten, Mesopotamien) sowie der Völker des Altertums gilt der Adler allgemein als Symbol der Herrschaft und des Göttlichen.<sup>[1]</sup> Verstärkt wurde die Symbolik durch die Darstellung von Mischwesen wie dem Greif oder den Harpyien, welche die Kraft, Stärke oder auch Gesinnung durch unterschiedliche Tierkörperteile hervorhoben. Im altindischen Rigveda überbringt der Garuda als schlangentötender Bote die Nachrichten der Götter, insbesondere von Vishnu. Der Sage nach wurde der sumerische König Gilgamesch von einem Adler gerettet. Für die antiken Griechen war der Adler das Symbol des obersten olympischen Gottes Zeus, im alten Rom war er das Zeichen der obersten römischen Gottheit lupiter sowie der kaiserlichen Macht und in der Apotheose Zeichen der Göttlichkeit des Kaisers.



Zeichen der Tapferkeit, aus den Knochen fertigten sie Pfeifen.

#### Bei vielen nordamerikanischen Indianerstämmen galten ihre Federn - meist als Kopfschmuck - als



Adler sind große Greifvögel. Es gibt mehrere Arten, wie zum Beispiel Steinadler, Seeadler oder Fischadler, Sie ernähren sich von kleinen und größeren Tieren. Sie greifen ihre Beute mit ihren starken Krallen im Flug, am Boden oder im Wasser



Zum Text für

Hier streiten sich Seeadler um Beute

Adler bauen ihre Nester. die

man Horste nennt, meist auf Felsen oder hohen Bäumen. Dort hinein legt das Weibchen ein bis vier Eier. Die Brutzeit beträgt ie nach Art 30 bis 45 Tage. Die Küken sind anfangs weiß, ihr dunkles Federkleid wächst erst später. Nach ungefähr 10 bis 11 Wochen können die lungen fliegen.

Die bekannteste Adlerart in Mitteleuropa ist der Steinadler. Seine Federn sind braun und seine ausgestreckten Flügel sind etwa zwei Meter breit. Er lebt vor allem in den Alpen und rund ums Mittelmeer, aber auch in Nordamerika und Asien. Der Steinadler ist sehr kräftig und kann Säugetiere iagen, die schwerer sind als er selbst. Meist fängt er Hasen und Murmeltiere, aber auch junge Rehe und Hirsche, manchmal auch Reptilien und Vögel.

Aumiller and Gertz: "Klexikon: A German Dataset for Joint Summarization and Simplification"

### • Reading a Simplified Article



#### UNIVERSITÄT HEIDELBERG ZUKUNFT SEIT 1386

#### • Reality

#### Adler (Biologie)

Der Begriff Adler (von mittelhochdeutsch adelar[e] für "Edel-Aar", wobei Aar Adler und Adlerartige (unedle Adler) bezeichnete) ist im weiteren Sinn eine nicht genau definierte Sammelbezeichnung für große und kräftige Arten von Greifvögen in der Familie der Habichattigen, wie beispielsweise für die Gattung der Seeadler, den Schlangenadler oder den Philippinenadler. Im engeren Sinn steht die Bezeichnung Adler für die Gattung Echte Adler (Aquila) mit den ihr zugehörigen Arten. Deren größte lebende Vertreter sind der Steinadler und der Kellschwanzadler. Außerhalb der Habichtartigen wird im Deutschen auch der in seine eigene Familie gestellte Fischadler unter die Adler gezählt.



Mit einer Spannweite bis zu 3 Meter gilt der um 1400 n. Chr. ausgestorbene Haastadler (Harpagornis moorei) als größter bekannter Adler. Sein Lebensraum befand sich auf Neuseeland.

#### Jagd und Fortpflanzung [Bearbeiten | Quelitext bearbeiten ]

Die außerordentliche Sehschärfe ihrer Augen gestattet es allen Adlern, ihre Beute auch aus großer Höhe zu erspähen und zu erjagen. Ihre Hauptbeute sind bodenbende Wirbeitliere, in erster Linie Säuger. Nur seiten ernähren sie sich von Aas. Adler sind sehr gut daran angepasst, zum Fliegen Aufwinde zu nutzen. Der Adler packt die Beute beim Niederstoßen aus dem Flug mit den Fängen. Dabei weisen drei seiner Zehen nach vorne, die Hinterkralle durchsticht den Beutekörper in einer Zangenbewegung. Durch den Griff wird die Beute erstickt.

In das aus Zweigen gebaute Nest (Adlerhorst) legt das Weibchen im Regelfall zwei bis drei Eier. Häufig wird jedoch nur eines der Jungen aufgezogen. Alle Adlerarten sind geschützt, ihr Bestand ist durch menschlichen Einfluss gefährdet.

#### Kunst und Symbolik [Bearbeiten | Quelltext bearbeiten ]

Die Kunst des Alten Orients verwendet den Adler in vielfachen Abwandlungen als Symboltier. Bereits 3000 v. Chr. findet sich das Mischwesen Löwenadler in mesopotamischen Darstellungen. Doppeladler sind aus dem Babylonien des 23. Jahrhunderts v. Chr. bekannt. Die Antike kennt unzählige Abbildungen des Adlers auf Vasen, Gegenständen des täglichen Gebrauchs, Münzen und Schmuckstücken sowie in der Architektur auf Reliefs, Akroterien und Glebeln.

Seit dem 5. Jahrhundert nach Christus fertigten die Goten große Adlerfibeln und Schnallen mit Adlerköpfen. Aus der romanischen Periode des Mittelalters sind mit Adlern verzierte figürliche Darstellungen von Kapitellen bekannt. Ebenfalls hervorzuheben sind die Adlerpulte christlicher Altraräume. In der Malerei erscheint der Adler häufig als Symbol für die Evangelisten Johannes.

In den Überlieferungen, dem Mythos und dem Volksglauben der frühen Hochkulturen (Ågypten, Mesopotamien) sowie der Völker des Altertums gilt der Adler allgemein als Symbol der Herrschaft und des Göttlichen.<sup>[11]</sup> Verstänkt wurde die Symbolik durch die Darstellung von Mischwesen wie dem Greif oder den Harpylen, welche die Kraft, Stärke oder auch Gesinnung durch unterschiedliche Tierkörperteile hervonhoben. Im altindischen Rijveda überbringt der Garuda als schlangenöttender Bote die Nachrichten der Götter, insbesondere von Vishnu. Der Sage nach wurde der sumerische König Gilgamesch von einem Adler gerettet. Für die antiken Griechen war der Adler das Symbol des obersten olympischen Gottes Zeus, im alten Rom war er das Zeichen der obersten römischen Gottheit Jupiter sowie der kalserlichen Macht und in der Apotheose Zeichen der Göttlichkeit des Kaisers.



Bei vielen nordamerikanischen Indianerstämmen galten ihre Federn – meist als Kopfschmuck – als Zeichen der Tapferkeit, aus den Knochen fertigten sie Pfeifen.

#### Aumiller and Gertz: "Klexikon: A German Dataset for Joint Summarization and Simplification"

#### Jagdweise [Bearbeiten | Quelltext bearbeiten ]

Steinadler jagen meist in offenen oder halboffenen Landschaften im bodennahen Flug unter optimaler Ausnutzung jeglicher Deckung. Sie gleiten dabei dicht an Hängen entlang, über Kuppen und kleine Hügel und versuchen ihre Beute auf kruze Distanz zu überraschen. Häufig jagen sie auch von einem Ansitz aus. Die Beute greifen die Adler meist auf dem Boden oder im bodennahen Luftraum und töten sie mit den außerordentlich kräftigen Zehen und Krallen. Sehr große Beutetiere wie Kitze des Steinbocks oder junge Gämsen greifen sie am Kopf. Der Steinadler schlägt dabei seine Krallen durch die Schädeldecke in das Gehirn. In den wenigen beobachteten Fällen wurden diese großen Beutetiere innerhalb von Sekunden etötett.

Weniger häuftig ist die Jagd im freien Luftraum; die Erbeatung von ziehenden Kormoranen ist jedoch zum Beispiel schon mehrfach beobachtet worden. In Anbetracht ihrer Größe bewegen sich Steinadler in der Luft außerordentlich wendig und schnell, so wurde mehrfach beobachtet, wie sich ein Steinadler im Flug auf den Rücken drehte und so zum Beispiel einen verfolgenden Kolkraben erbeutete. Steinadler können keine Kadaver im Flug tragen, deren Gewicht das eigene Körpergewicht überrifft. Schwere Beutetiere zerteilen sie daher und deponieren portionsweise, oder sie fliegen den Kadaver über mehrere Tage an.

#### Nahrung [Bearbeiten | Quelltext bearbeiten ]

Steinadier sind außerordentlich kräftig und sehr geschickt. Sie erbeuten regelmäßig Tiere, die erheblich schwerer sind als sie selbst. Das maximale Beutegewicht liegt bei etwa 15 Kilogramm. Es gibt nur einen dokumentierten Fall, in dem ein noch schwererer Sikahirsch erlegt wurde.<sup>[4]</sup> Angriffsversuche auf annähernd ausgewachsene Gämsen sind dokumentiert.<sup>[5116]</sup> Überteidigungsstrategie der Gämsen besteht darin, hangabwärts zu springen und sich überschlagend zu rollen, was eine erhebliche Verletzungsgefahr für beide bedeutet.

Im Beutespektrum dominieren meist bodenbewohnende, kleine bis mittelgroße Säugetiere von Ziesel- bis Steinbockkitz-Größe, Vögel spielen meist nur eine kleinere Rolle. Meist bilden wenige Säugerarten den Hauptteil der Nahrung. Daneben erbeutet der Steinadier jedoch fast alle kleinen und mittelgroßen Säuger und Vögel, die im jeweiligen Gebiet vorkommen. Insbesondere im Süden des Verbreitungsgebietes frisst er auch regelmäßig Reptillen, dort lassen Steinadler ähnlich wie Bartgeier auch Landschildkröten auf Felsen fallen, um so deren harten Panzer zu zerbrechen. Insbesondere im Winter, regional aber auch im Sommer, spielt Aas eine wichtige Rolle in der Ernährung.

Im Schweizer Kanton Graubünden dominierten zur Brutzeit im Beutespektrum Alpenmurmeltiere mit 60,2 % aller Beutetiere, an zweiter Stelle folgten junge Gämsen mit 8,0 %. Danach folgten Schwechase, Alpenschneehulm und Birkhuhm mit Jeweils 5,2  $(3^{17})$  Im schweizerschen Alpenvorland bestand die Nestlings-Nahrung in 4 Revieren vor allem aus Ecklassen (36,2 % aller Beutetiere), danach folgten Hauskatzen (27,5 %), Rehkitze (14,1 %) und Haushühner (8,1 %),<sup>[8]</sup> Populationen im Zentralmassiv Frankreichs jagen hauptsächlich Wildkaninchen. In Schottland wurden je nach Region Hasenattige in 10,7 % bis 46,9 % aller im Sommer gefundenen Gewölle nachgewiesen. Weitere wichtige Beutetiere waren dort Schafe und Ziegen (in 0,6 bis 26,8 % aller Gewölle), Raufußhühner (5,4 %) und Rothirsche (als Aas) (1,2 bis 22,3 %).<sup>[9]</sup>

#### Raumnutzung und Siedlungsdichte [Bearbeiten | Quelitext bearbeiten ]

Trotz des großen Verbreitungsgebietes liegen bisher nur wenige Daten zur Größe des Aktionsraumes, also zu der von einem Brutpaar genutzten Fläche vor. Die festgesteilten Werte schwanken je nach Lebensraum und Nahrungsangebot erheblich. Im Schweizer Kanton Graubünden betrug die Größe des Aktionsraumes nach Sichtbeobachtungen in 26 Revieren zwischen 29 und 88 km<sup>-</sup>, im Mittel 33 km<sup>-</sup>, in I daho beflogen vier Steinadierpaare ebenfalls nach Sichtbeobachtungen Flächen zwischen 11.6 und 49.0 km<sup>-</sup>, im Mittel 32,8 km<sup>-</sup>; <sup>[10]</sup> Methodisch bedingt stellen die hier dargesteilten Werte vermutlich eher die Untergrenze der tatsschlichen Aktionsraumgrößen dar. Steinadier verteidigen ihren gesamten Aktionsraum ganzjährig vehement gegen Artgenossen, eine früher gelegentlich vorgenommene Trennung zwischen dem verteidigten Revier und der zur Nahrungssuche genutzten Fläche ist also nicht gerechtfertigt.

## Limitations of Current Text Simplification Systems



- Simplification:
  - " "Translate" text into a simplified version (sentence-by-sentence)
  - <sup>2</sup> Length of input corresponds to length of output!

- Instead, we should consider simplification as document-level summarization!
  - Simultaneously compress article length and simplify texts

## Limitations of Current Text Simplification Systems



• But what about training data?

Corpora mostly provide sentence-level alignments

• Document-level resources barely exist!

	-		Aligned	Avg. #Sentences	
		Resource	$\mathbf{Articles}$	Source	$\mathbf{Simple}$
~ Fo	For English:	(Kauchak, 2013)	59,775	64.52	8.46
-		(Xu et al., 2015)	$1,\!130$	49.59	51.27
~ ⊑	For German:	(Battisti et al., 2020)	378	45.29	55.75
ΓC		(Hewett and Stede, 20	21) 978	10.12	43.54

Table 1: Existing document-aligned resources for simplification.

### Contributions



• Frame document-level simplification as a joint summarization and simplification problem

 Contribute a large-scale (document-aligned) resource for German based on alignments between Wikipedia and a German children's encyclopedia



### **Constructing the Klexikon Dataset**



## Klexikon



- Started in 2014 as a resource specifically for children (age 6-13)
- Each article is internally reviewed before release
- Almost 3,300 articles on diverse topics
  - Crawled during April 2021

# **Corpus Alignment**



- **Basic idea:** Klexikon articles are (summarized) simplifications of the Wikipedia article on the same topic
- Find Wikipedia article for each Klexikon entry (alignment)!
- **Hypothesis:** Aligned articles represent a suitable dataset for training a joint summarization/simplification objective



# **Corpus Alignment: Disambiguation**



• Not always that easy: Only 90% match directly



• Requires manual resolution of conflicts

# **Corpus Alignment: Multi-Pages**

UNIVERSITÄT HEIDELBERG **ZUKUNFT** SEIT 1386

- Content on Wikipedia can spread across pages
  - Manual review of disambiguated articles
  - Keep if more than 66% of the Klexikon article are corresponding to a *single* Wikipedia article



Warum sind Adler oft in Wappen?

Ein Wappen ist ein Bild, das für

erreichten.

Adler sind große Greifvögel, Es

gibt mehrere Arten, wie zum Beispiel Steinadler, Seeadler

ein Land, eine Stadt oder

Familie steht, Schon seit dem Altertum sind Mensch großen Vögeln, die am Himmel gleiten, Forscher vermuten sogar, dass der Name Adler von dem Wort "edel" kommt. Bei den alten Griechen galt der Adler als Zeichen für den Göttervater Zeus, bei den Römern für Jupiter.

Auch im Mittelalter war der Adler ein Zeichen für königliche Macht und Vornehmheit. Deshalb durften nur Könige und Kaiser den Adler als Wappentier führen. So kam er in die Wappen vieler Länder, zum Beispiel Deutschland, Österreich, Polen oder Russland. Sogar die USA haben ein Adler-Wappen, obwohl sie nie einen König hatten. Der amerikanische Adler ist ein Weißkopfadler, der deutsche ein Steinadler.

Adler (Eagle)







Adler (Heraldry)



### **Exploratory Analysis**

## **Resulting Corpus**



- Around 2,900 articles aligned successfully!
- 3x more articles than largest previous corpus
- Significantly longer source documents

A	ligned	Avg. #Sentences		Compr.				
Resource A	$\mathbf{rticles}$	Source	$\mathbf{Simple}$	Ratio				
(Battisti et al., 2020)	378	45.29	55.75	0.81				
(Hewett and Stede, 2021)	978	10.12	43.54	0.23				
Klexikon (Ours)	2,898	242.09	32.51	7.45				
Table 2: German resources in comparison.								

Aumiller and Gertz: "Klexikon: A German Dataset for Joint Summarization and Simplification"

### **Length Distributions**





Distribution of the dataset including median (red line), mean (dotted black) and standard deviation (thin black line) for both Wikipedia and Klexikon.

Aumiller and Gertz: "Klexikon: A German Dataset for Joint Summarization and Simplification"

# Suitability for Summarization



UNIVERSITÄT HEIDELBERG ZUKUNFT SEIT 1386

- How do we know the dataset is appropriate?
  - Test against (extractive) baselines!
- Compared methods:
  - Lead-3 (first three sentences)
  - Lead-k (opening paragraph in Wikipedia)
  - Full article
  - Luhn's Algorithm (Luhn, 1958)
  - LexRank (Erkan and Radev, 2004)
    (with sentence-transformers embeddings)
  - Extractive Oracle
- Evaluation with ROUGE

	<b>R-1</b>	<b>R-2</b>	R-L
Lead-3	16.95	3.77	9.81
Lead- $k$	24.87	5.10	12.01
Full article	16.81	4.23	6.95
Luhn	31.86	5.55	11.57
LexRank S-T	<b>33.90</b>	6.11	12.86
Oracle	41.85	10.68	16.00

Table 3: ROUGE F1 scores of baselines.

# **Suitability for Simplification**



- Use modified Flesch score for German (Amstad, 1978)
- Average sentence length in tokens
- Average word length in characters
- Relative usage of the top 1000 lemmas of each corpus (Wikipedia/Klexikon) in relation to the total number of lemmas
  - <sup>~</sup> Restrict to nouns, adjectives, verbs or adverbs
  - Less biased against longer texts

# **Suitability for Simplification**



WikipediaKlexikonMean Flesch score $40.1 \pm 7.3$  $66.7 \pm 6.0$ Mean sentence length $22.7 \pm 2.6$  $13.5 \pm 1.5$ Mean word length $8.7 \pm 4.0$  $6.9 \pm 3.0$ Share of top 1000 lemmas68.8%82.3%

Table 4: Indicators of simplified texts across all metrics.



### **Open Challenges and Future Work**

### **Open Challenges**



- No "sophisticated" system that performs both simplification and summarization exists yet
  - How to integrate simplification as part of summarizers?
  - Abstractive systems are limited by length
- How to align documents on a fine-grained context (sentence/paragraph-level) without manual annotation?
  - Existing solutions assume linear alignments and English

### **Future Work**



- Provide additional sentence-level alignments
- Hybrid retrieval systems to cut down texts to suitable lengths for abstractive systems
- Experiment with regularization to incorporate simplification into neural networks

### Resources



- Amstad, T. (1978). Wie verständlich sind unsere Zeitungen? Studenten-Schreib-Service.
- Battisti, A., Pfütze, D., Säuberli, A., Kostrzewa, M., and Ebling, S. (2020). A corpus for automatic readability assessment and text simplification of German. In Proceedings of the 12th Language Resources and Evaluation *Conference*, pages 3302–3311, Marseille, France, May. European Language Resources Association.
- Erkan, G. and Radev, D. R. (2004). Lexrank: Graph-based lexical centrality as salience in text summarization. Journal of Artificial Intelligence Research, 22:457–479.
- Hewett, F. and Stede, M. (2021). Automatically evaluating the conceptual complexity of German texts. In Proceedings of the 17th Conference on Natural Language Processing (KONVENS 2021), pages 228–234, Düsseldorf, Germany, 6–9 September. KONVENS 2021 Organizers.
- Kauchak, D. (2013). Improving text simplification language modeling using unsimplified text data. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1537–1546, Sofia, Bulgaria, August. Association for Computational Linguistics.
- Luhn, H. P. (1958). The automatic creation of literature abstracts. IBM Journal of Research and Development, 2(2):159-165.
- Xu, W., Callison-Burch, C., and Napoles, C. (2015). Problems in current text simplification research: New data can help. Transactions of the Association for Computational Linguistics, 3:283–297. Aumiller and Gertz: "Klexikon: A German Dataset for Joint Summarization and Simplification" 24

# Thank you for your attention!



Check out the dataset on Huggingface: <u>https://huggingface.co/datasets/dennlinger/klexikon</u>

The experimental code is also on Github: <u>https://github.com/dennlinger/klexikon</u>

### **Any questions?**

Or send a mail to <u>aumiller@informatik.uni-heidelberg.de</u>