

Zwischen Interpretation und Berechnung

Didaktische Exploration computerlinguistischer
Verfahren
mit augmentierten Transkripten von
Verkaufsgesprächen

Paul Koop

Lehr- und Lernmaterial 2026

Zusammenfassung

Dieses Lehr- und Lernmaterial dient der didaktischen Exploration computerlinguistischer Verfahren auf der Grundlage der acht Transkripte von Verkaufsgesprächen. Im Unterschied zu den vorangegangenen ARS-Versionen 2.0 und 3.0, die auf interpretativ gebildeten Terminalzeichen basierten, wird hier der Schritt in Richtung automatischer Sprachverarbeitung vollzogen. Die Verfahren werden zu Demonstrationszwecken auf augmentierten Daten trainiert, um ihre Funktionsweise transparent zu machen. Der Fokus liegt auf dem didaktischen Erkenntnisgewinn, nicht auf empirischer Validität. Die Szenarien C (Computerlinguistische Integration) und D (Hybride Modellierung) werden schrittweise entwickelt und miteinander verglichen.

Inhaltsverzeichnis

1 Einleitung: Didaktische Ziele und methodologische Reflexion	2
2 Die acht Transkripte: Rohdaten und Terminalzeichen	3
2.1 Die Rohdaten	3
2.1.1 Transkript 1 - Metzgerei	3
2.1.2 Transkript 2 - Marktplatz (Kirschen)	3
2.1.3 Transkript 3 - Fischstand	4
2.1.4 Transkript 4 - Gemüsestand (ausführlich)	4
2.1.5 Transkript 5 - Gemüsestand (mit KAV zu Beginn)	4
2.1.6 Transkript 6 - Käseverkaufsstand	4
2.1.7 Transkript 7 - Bonbonstand	5
2.1.8 Transkript 8 - Bäckerei	5
2.2 Die Terminalzeichenketten (ARS 3.0)	6
3 Szenario C: Computerlinguistische Integration	7
3.1 Didaktische Augmentierung	7
3.2 Speech Act Recognition mit Transformer-Modellen	8
3.3 Word Embeddings und semantische Ähnlichkeit	12
3.4 Topic Modeling mit BERTopic	15
3.5 Rhetorical Structure Theory (RST)	19
3.6 Integration der Komponenten in Szenario C	24
4 Szenario D: Hybride Modellierung	29
4.1 CRF für sequenzielle Abhängigkeiten	29
4.2 Transformer-Embeddings als Ergänzung	33
4.3 Graph Neural Networks für die Nonterminal-Hierarchie	36
4.4 Attention-Mechanismen für relevante Vorgänger	40
4.5 Integration der Komponenten in Szenario D	44
5 Vergleich der Szenarien und methodologische Reflexion	48
5.1 Gegenüberstellung der Ansätze	48
5.2 Didaktische Erkenntnisse aus Szenario C	49
5.3 Didaktische Erkenntnisse aus Szenario D	49
5.4 Fazit für die Lehrpraxis	50
6 Ausblick	50

1 Einleitung: Didaktische Ziele und methodologische Reflexion

Die vorangegangenen Versionen der Algorithmisch Rekursiven Sequenzanalyse (ARS 2.0 und 3.0) haben gezeigt, wie aus interpretativ gewonnenen Terminalzeichenketten formale Grammatiken induziert werden können. Diese Verfahren bleiben methodologisch kontrolliert: Die Kategorienbildung erfolgt durch qualitative Interpretation, die formalen Modelle explizieren lediglich die beobachtbaren Regularitäten.

Die nun folgenden Szenarien C und D wagen einen Schritt über diese methodologische Grenze hinaus. Sie explorieren, wie computerlinguistische Verfahren – insbesondere neuronale Netze, Word Embeddings und Topic Models – auf die acht Transkripte angewendet werden könnten, wenn man sie zu Demonstrationszwecken augmentiert.

Dieses Dokument ist als Lehr- und Lernmaterial konzipiert. Es verfolgt folgende didaktische Ziele:

1. **Verständnis neuronaler Architekturen:** Wie funktionieren Transformer, LSTM-Netze und Attention-Mechanismen auf Sequenzdaten?
2. **Data Augmentation als Technik:** Wie kann man mit kleinen Datensätzen umgehen, um die Funktionsweise von Verfahren zu demonstrieren?
3. **Vergleich verschiedener Modellierungsebenen:** Welche Unterschiede bestehen zwischen rein computerlinguistischen (C) und hybriden (D) Ansätzen?
4. **Methodologische Reflexion:** Wo liegen die Grenzen automatischer Verfahren im Vergleich zur interpretativen Kategorienbildung?

Alle hier vorgestellten Implementierungen arbeiten mit augmentierten Daten – die acht Originaltranskripte wurden künstlich vervielfacht, um das Training neuronaler Netze zu ermöglichen. Die Ergebnisse sind daher nicht empirisch valide, sondern dienen ausschließlich der didaktischen Veranschaulichung.

2 Die acht Transkripte: Rohdaten und Terminalzeichen

2.1 Die Rohdaten

Die folgenden acht Transkripte dokumentieren Verkaufsgespräche auf dem Aachener Marktplatz im Juni/Juli 1994. Sie bilden die empirische Grundlage aller folgenden Analysen.

2.1.1 Transkript 1 - Metzgerei

Datum: 28. Juni 1994, Ort: Metzgerei, Aachen, 11:00 Uhr

```
1 Kunde: Guten Tag!
2 Verk uferin: Guten Tag!
3 Kunde: Einmal von der groben Leberwurst, bitte.
4 Verk uferin: Wie viel darf's denn sein?
5 Kunde: Zwei hundert Gramm.
6 Verk uferin: Sonst noch etwas?
7 Kunde: Ja, dann noch ein St ck von dem Schwarzw lder Schinken.
8 Verk uferin: Wie gro soll das St ck sein?
9 Kunde: So um die dreihundert Gramm.
10 Verk uferin: Das macht dann acht Mark zwanzig.
11 Kunde: Bitte.
12 Verk uferin: Danke und einen sch nen Tag noch!
13 Kunde: Danke, ebenfalls!
```

Listing 1: Transkript 1 - Rohdaten

2.1.2 Transkript 2 - Marktplatz (Kirschen)

Datum: 28. Juni 1994, Ort: Marktplatz, Aachen

```
1 Verk ufer: Kirschen kann jeder probieren hier!
2 Kunde 1: Ein halbes Kilo Kirschen, bitte.
3 Verk ufer: Ein halbes Kilo? Oder ein Kilo?
4 Verk ufer: Drei Mark, bitte.
5 Kunde 1: Danke sch n!
6 Verk ufer: Kirschen kann jeder probieren hier!
7 Kunde 2: Ein halbes Kilo, bitte.
8 Verk ufer: Drei Mark, bitte.
9 Kunde 2: Danke sch n!
```

Listing 2: Transkript 2 - Rohdaten

2.1.3 Transkript 3 - Fischstand

Datum: 28. Juni 1994, Ort: Fischstand, Marktplatz, Aachen

```
1 Kunde: Ein Pfund Seelachs, bitte.  
2 Verk ufer: Seelachs, alles klar.  
3 Verk ufer: Vier Mark neunzehn, bitte.  
4 Kunde: Danke sch n!
```

Listing 3: Transkript 3 - Rohdaten

2.1.4 Transkript 4 - Gemüsestand (ausführlich)

Datum: 28. Juni 1994, Ort: Gemüsestand, Aachen, Marktplatz, 11:00 Uhr

```
1 Kunde: H ren Sie, ich nehme ein paar Champignons mit.  
2 Verk ufer: Braune oder helle?  
3 Kunde: Nehmen wir die hellen.  
4 Verk ufer: Die sind beide frisch, keine Sorge.  
5 Kunde: Wie ist es mit Pfifferlingen?  
6 Verk ufer: Ah, die sind super!  
7 Kunde: Kann ich die in Reissalat tun?  
8 Verk ufer: Eher kurz anbraten in der Pfanne.  
9 Kunde: Okay, mache ich.  
10 Verk ufer: Sch nen Tag noch!  
11 Kunde: Gleichfalls!
```

Listing 4: Transkript 4 - Rohdaten

2.1.5 Transkript 5 - Gemüsestand (mit KAV zu Beginn)

Datum: 26. Juni 1994, Ort: Gemüsestand, Aachen, Marktplatz, 11:00 Uhr

```
1 Kunde 1: Auf Wiedersehen!  
2 Kunde 2: Ich h tte gern ein Kilo von den Granny Smith pfeln hier  
. .  
3 Verk ufer: Sonst noch etwas?  
4 Kunde 2: Ja, noch ein Kilo Zwiebeln.  
5 Verk ufer: Sechs Mark f nfundzwanzig, bitte.  
6 Kunde 2: Auf Wiedersehen!
```

Listing 5: Transkript 5 - Rohdaten

2.1.6 Transkript 6 - Käseverkaufsstand

Datum: 28. Juni 1994, Ort: Käseverkaufsstand, Aachen, Marktplatz

```

1 Kunde 1: Guten Morgen!
2 Verk ufer: Guten Morgen!
3 Kunde 1: Ich h tte gerne f nfhundert Gramm holl ndischen Gouda.
4 Verk ufer: Am St ck?
5 Kunde 1: Ja, am St ck, bitte.

```

Listing 6: Transkript 6 - Rohdaten

2.1.7 Transkript 7 - Bonbonstand

Datum: 28. Juni 1994, **Ort:** Bonbonstand, Aachen, Marktplatz, 11:30 Uhr

```

1 Kunde: Von den gemischten h tte ich gerne hundert Gramm.
2 Verk ufer: F r zu Hause oder zum Mitnehmen?
3 Kunde: Zum Mitnehmen, bitte.
4 Verk ufer: F nfzig Pfennig, bitte.
5 Kunde: Danke!

```

Listing 7: Transkript 7 - Rohdaten

2.1.8 Transkript 8 - Bäckerei

Datum: 9. Juli 1994, **Ort:** Bäckerei, Aachen, 12:00 Uhr

```

1 (Schritte h rbar, Hintergrundger usche, teilweise unverst ndlich
)
2 Kunde: Guten Tag!
3 (Unverst ndliche Begr ung im Hintergrund)
4 Verk uferin: Einmal unser bester Kaffee, frisch gemahlen, bitte.
5 (Ger usche der Kaffeem hle, Verpackungsger usche)
6 Verk uferin: Sonst noch etwas?
7 Kunde: Ja, noch zwei St ck Obstsalat und ein Sch lchen Sahne.
8 Verk uferin: In Ordnung!
9 (Ger usche der Kaffeem hle, Papierger usche)
10 Verk uferin: Ein kleines Sch lchen Sahne, ja?
11 Kunde: Ja, danke.
12 (T rger usch, Lachen, Papierger usche)
13 Verk uferin: Keiner k mmert sich darum, die T ren zu len .
14 Kunde: Ja, das ist immer so.
15 (Lachen, Ger usche von M nzen und Verpackung)
16 Verk uferin: Das macht vierzehn Mark und neunzehn Pfennig, bitte.
17 Kunde: Ich zahle in Kleingeld.
18 (Lachen und Ger usche von M nzen)
19 Verk uferin: Vielen Dank, sch nen Sonntag noch!
20 Kunde: Danke, Ihnen auch!

```

Listing 8: Transkript 8 - Rohdaten

2.2 Die Terminalzeichenketten (ARS 3.0)

Für die ARS 3.0 wurden diese Rohdaten in Terminalzeichenketten überführt, die als Grundlage für die hierarchische Grammatikinduktion dienten:

Tabelle 1: Terminalzeichenketten der acht Transkripte

Transkript	Terminalzeichenkette
1 (Metzgerei)	KBG, VBG, KBBd, VBBd, KBA, VBA, KBBd, VBBd, KBA, VAA, KAA, VAV
2 (Kirschen)	VBG, KBBd, VBBd, VAA, KAA, VBG, KBBd, VAA, KAA
3 (Fischstand)	KBBd, VBBd, VAA, KAA
4 (Gemüse)	KBBd, VBBd, KBA, VBA, KBBd, VBA, KAE, VAE, KAA, VAV, KAV
5 (Gemüse KAV)	KAV, KBBd, VBBd, KBBd, VAA, KAV
6 (Käse)	KBG, VBG, KBBd, VBBd, KAA
7 (Bonbon)	KBBd, VBBd, KBA, VAA, KAA
8 (Bäckerei)	KBG, VBBd, KBBd, VBA, VAA, KAA, VAV, KAV

Die Bedeutung der Terminalzeichen:

- **KBG:** Kunden-Gruß
- **VBG:** Verkäufer-Gruß
- **KBBd:** Kunden-Bedarf (konkret)
- **VBBd:** Verkäufer-Nachfrage
- **KBA:** Kunden-Antwort
- **VBA:** Verkäufer-Reaktion
- **KAE:** Kunden-Erkundigung
- **VAE:** Verkäufer-Auskunft
- **KAA:** Kunden-Abschluss
- **VAA:** Verkäufer-Abschluss
- **KAV:** Kunden-Verabschiedung
- **VAV:** Verkäufer-Verabschiedung

3 Szenario C: Computerlinguistische Integration

Szenario C realisiert eine vollständig computerlinguistische Modellierung der acht Transkripte. Es umfasst vier Komponenten:

1. **Speech Act Recognition**: Automatische Erkennung der Sprechakte aus den Rohdaten
2. **Word Embeddings**: Vektorielle Repräsentation der Äußerungen
3. **Topic Modeling**: Identifikation thematischer Verschiebungen
4. **Rhetorical Structure Theory (RST)**: Analyse der argumentativen Struktur

3.1 Didaktische Augmentierung

Da neuronale Netze für ihr Training große Datenmengen benötigen, werden die acht Transkripte zu Demonstrationszwecken augmentiert:

```
1 def augment_transcripts_for_teaching(transcripts, factor=20):  
2     """  
3         Augmentiert die acht Transkripte für didaktische Zwecke.  
4  
5         Didaktischer Hinweis: Diese Augmentierung dient  
6             ausschließlich der  
7             Veranschaulichung der Methodik. Die resultierenden Daten  
8             sind nicht  
9             empirisch valide, sondern ermöglichen lediglich die  
10                Demonstration  
11                der Funktionsweise neuronaler Verfahren.  
12  
13    augmented = []  
14  
15    # 1. Basis-Augmentierung: einfaches Kopieren  
16    for _ in range(factor):  
17        augmented.extend(transcripts)  
18  
19    # 2. Syntaktische Variationen (didaktisch kontrolliert)  
20    import copy  
21    import random  
22  
23    for transcript in transcripts:  
24        augmented.append(copy.deepcopy(transcript))  
25        if random.random() < 0.5:  
26            augmented.append(transcript.replace(" ", " " * random.randint(1, 5)))  
27        else:  
28            augmented.append(transcript.replace(" ", " " * random.randint(1, 5)))  
29  
30    return augmented
```

```

21     for _ in range(factor // 4):
22         var = copy.deepcopy(transcript)
23         # Vertausche zwei benachbarte Wörterungen (selten)
24         if len(var) > 3 and random.random() < 0.1:
25             idx = random.randint(0, len(var)-2)
26             var[idx], var[idx+1] = var[idx+1], var[idx]
27             augmented.append(var)
28
29     # 3. Lexikalische Variationen (Synonyme)
30     synonyms = {
31         'Guten Tag': ['Guten Morgen', 'Hallo', 'Guten Abend'],
32         'Danke': ['Vielen Dank', 'Danke schön', 'Merci'],
33         'Bitte': ['Bitte sehr', 'Gern geschehen']
34     }
35
36     # Hier könnten weitere Variationen implementiert werden
37
38     return augmented

```

Listing 9: Data Augmentation für Lehrzwecke

3.2 Speech Act Recognition mit Transformer-Modellen

Die automatische Erkennung der Sprechakte erfolgt mit einem feinabgestimmten BERT-Modell:

```

1 """
2 Speech Act Recognition mit transformer-basierten Modellen
3 Didaktische Implementierung für Lehrzwecke
4 """
5
6 import torch
7 import torch.nn as nn
8 from transformers import BertTokenizer, BertModel
9 import numpy as np
10 from sklearn.preprocessing import LabelEncoder
11 from torch.utils.data import Dataset, DataLoader
12
13 class SpeechActDataset(Dataset):

```

```

14     """Dataset f r Speech Act Recognition"""
15     def __init__(self, utterances, labels, tokenizer,
16                  max_length=128):
17         self.utterances = utterances
18         self.labels = labels
19         self.tokenizer = tokenizer
20         self.max_length = max_length
21
22     def __len__(self):
23         return len(self.utterances)
24
25     def __getitem__(self, idx):
26         utterance = self.utterances[idx]
27         label = self.labels[idx]
28
29         encoding = self.tokenizer(
30             utterance,
31             truncation=True,
32             padding='max_length',
33             max_length=self.max_length,
34             return_tensors='pt'
35         )
36
37         return {
38             'input_ids': encoding['input_ids'].flatten(),
39             'attention_mask': encoding['attention_mask'].
40                             flatten(),
41             'label': torch.tensor(label, dtype=torch.long)
42         }
43
44     class BertSpeechActClassifier(nn.Module):
45         """
46             BERT-basierter Klassifikator f r Sprechakte
47             Didaktisch vereinfachte Architektur
48         """
49
50         def __init__(self, num_classes=12, dropout=0.3):
51             super().__init__()
52             self.bert = BertModel.from_pretrained('bert-base-
53                                                 german-cased')
54             self.dropout = nn.Dropout(dropout)

```

```

51         self.classifier = nn.Linear(768, num_classes)
52
53     # Freeze BERT layers f r didaktische Zwecke (
54     # schnelleres Training)
55     for param in self.bert.parameters():
56         param.requires_grad = False
57
58     def forward(self, input_ids, attention_mask):
59         outputs = self.bert(input_ids=input_ids,
60                             attention_mask=attention_mask)
61         pooled_output = outputs.pooler_output
62         dropped = self.dropout(pooled_output)
63         logits = self.classifier(dropped)
64         return logits
65
66     def prepare_speech_act_data(transcripts, terminal_chains):
67         """
68         Bereitet die Daten f r das Speech Act Training vor
69         """
70
71         utterances = []
72         labels = []
73
74         # Extrahiere alle uerungen aus den Rohdaten
75         # Hier vereinfacht: Mapping der Terminalzeichen auf
76         # Sprechakte
77         for trans, chain in zip(transcripts, terminal_chains):
78             # In einer vollst ndigen Implementierung m sssten
79             # die Rohdaten geparst werden
80             # F r didaktische Zwecke verwenden wir die
81             # Terminalzeichen direkt
82             for symbol in chain:
83                 utterances.append(f"Beispiel u erung f r {symbol}")
84                 labels.append(symbol)
85
86         # Label-Encoding
87         label_encoder = LabelEncoder()
88         y_encoded = label_encoder.fit_transform(labels)
89
90         return utterances, y_encoded, label_encoder

```

```

85
86 def train_speech_act_model(utterances, labels, epochs=10):
87     """
88     Trainiert das Speech Act Recognition Modell
89     """
90
91     tokenizer = BertTokenizer.from_pretrained('bert-base-
92         german-cased')
93     dataset = SpeechActDataset(utterances, labels, tokenizer)
94     dataloader = DataLoader(dataset, batch_size=8, shuffle=
95         True)
96
97
98     print("\n==== Speech Act Recognition Training (didaktisch)
99         ===")
100
101    for epoch in range(epochs):
102        total_loss = 0
103
104        for batch in dataloader:
105            optimizer.zero_grad()
106            outputs = model(batch['input_ids'], batch['
107                attention_mask'])
108
109            loss = criterion(outputs, batch['label'])
110
111            loss.backward()
112            optimizer.step()
113            total_loss += loss.item()
114
115
116            print(f"Epoch {epoch+1}: Loss = {total_loss/len(
117                dataloader):.4f}")
118
119
120    return model, tokenizer, label_encoder
121
122
123 # Didaktischer Hinweis
124 print("\n" + "="*70)
125 print("DIDAKTISCHER HINWEIS ZUR SPEECH ACT RECOGNITION")
126 print("=*70")
127 print("Die hier gezeigte Implementierung verwendet
128     augmentierte")

```

```

118 print("Daten und dient ausschließlich Lehrzwecken. Die
119     automatische")
120 print("Erkennung von Sprechakten wurde in der Praxis:")
121 print("    Millionen von annotierten Trainingsdaten
122     benötigen")
123 print("    Auf spezifische Domänen (Verkaufsgespräche)
124     feinabgestimmt werden")
125 print("    Mit erheblichen Unsicherheiten behaftet sein")

```

Listing 10: Speech Act Recognition mit BERT

3.3 Word Embeddings und semantische Ähnlichkeit

Für die Quantifizierung semantischer Ähnlichkeit werden vortrainierte Word Embeddings verwendet:

```

1 """
2 Word Embeddings für semantische Ähnlichkeitsanalysen
3 Didaktische Implementierung mit vortrainierten Modellen
4 """
5
6 from sentence_transformers import SentenceTransformer
7 import numpy as np
8 from sklearn.metrics.pairwise import cosine_similarity
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 class SemanticAnalyzer:
13 """
14     Analysiert semantische Ähnlichkeiten zwischen
15     Utterungen
16 """
17     def __init__(self, model_name='paraphrase-multilingual-
18         MiniLM-L12-v2'):
19         print(f'Lade vortrainiertes Modell: {model_name}')
20         self.model = SentenceTransformer(model_name)
21         self.embeddings = []
22
23     def encode_utterances(self, utterances):
24 """
25     Erzeugt Embeddings für eine Liste von Utterungen

```

```

24
25     """
26     embeddings = self.model.encode(utterances)
27     for utt, emb in zip(utterances, embeddings):
28         self.embeddings[utt] = emb
29     return embeddings
30
31
32     def similarity_matrix(self, utterances):
33         """
34             Berechnet die    hnlichkeitsmatrix      f   r alle
35             uerungen
36         """
37
38         embeddings = self.encode_utterances(utterances)
39         sim_matrix = cosine_similarity(embeddings)
40         return sim_matrix
41
42
43     def find_similar(self, query, utterances, top_k=5):
44         """
45             Findet die    hnlichsten      uerungen      zu einer Query
46         """
47
48         query_emb = self.model.encode([query])[0]
49         utt_embs = self.encode_utterances(utterances)
50
51         similarities = cosine_similarity([query_emb],
52                                         utt_embs)[0]
53         top_indices = np.argsort(similarities)[-top_k:][::-1]
54
55         results = []
56         for idx in top_indices:
57             results.append({
58                 'utterance': utterances[idx],
59                 'similarity': similarities[idx]
60             })
61
62
63         return results
64
65
66     def visualize_similarity(self, utterances, labels=None):
67         """
68             Visualisiert die    hnlichkeitsmatrix      als Heatmap
69         """
70
71         sim_matrix = self.similarity_matrix(utterances)

```

```

62
63     plt.figure(figsize=(12, 10))
64     sns.heatmap(sim_matrix,
65                 xticklabels=labels if labels else range(
66                     len(utterances)),
67                 yticklabels=labels if labels else range(
68                     len(utterances)),
69                 cmap='viridis', vmin=0, vmax=1)
70     plt.title('Semantische Ähnlichkeit zwischen
71                 Ueरungen')
72     plt.tight_layout()
73     plt.savefig('semantic_similarity.png', dpi=150)
74     plt.show()

75
76 # Didaktisches Beispiel
77 def demonstrate_semantic_analysis():
78     """
79         Demonstriert die semantische Analyse an Beispielen
80     """
81     analyzer = SemanticAnalyzer()

82     # Beispieldauerungen aus den Transkripten
83     utterances = [
84         "Guten Tag!",
85         "Guten Morgen!",
86         "Einmal Leberwurst, bitte.",
87         "Ich hätte gerne Wurst.",
88         "Danke schön!",
89         "Vielen Dank!",
90         "Auf Wiedersehen!",
91         "Tschüss!"
92     ]
93
94     print("\n==== Semantische Ähnlichkeitsanalyse ====")
95
96     # Ähnlichkeitsmatrix berechnen
97     sim_matrix = analyzer.similarity_matrix(utterances)
98
99     # Ähnlichste Ueरungen für "Guten Tag!"
100    similar = analyzer.find_similar("Guten Tag!", utterances,

```

```

        top_k=3)
99     print("\n hnlichste zu 'Guten Tag!':")
100    for r in similar:
101        print(f" {r['utterance']}: {r['similarity']:.3f}")

102
103    # Visualisierung
104    analyzer.visualize_similarity(utterances, utterances)

105
106    return analyzer

107
108 # Didaktischer Hinweis
109 print("\n" + "="*70)
110 print("DIDAKTISCHER HINWEIS ZU WORD EMBEDDINGS")
111 print("="*70)
112 print("Die verwendeten Embeddings wurden auf gro en Korpora"
      )
113 print("(Wikipedia, Nachrichten, Webtexte) vorgenutzt. Sie")
114 print("erfassen allgemeinsprachliche hnlichkeiten , nicht"
      die")
115 print("spezifischen Kategorien der Verkaufsgespr che.")

```

Listing 11: Semantische Ähnlichkeit mit Word Embeddings

3.4 Topic Modeling mit BERTopic

Für die Identifikation thematischer Verschiebungen wird BERTopic verwendet:

```

1 """
2 Topic Modeling zur Identifikation thematischer Verschiebungen
3 Didaktische Implementierung mit BERTopic
4 """
5
6 from bertopic import BERTopic
7 from sklearn.feature_extraction.text import CountVectorizer
8 import pandas as pd
9 import matplotlib.pyplot as plt
10
11 class TranscriptTopicModeler:
12     """
13         F hrt Topic Modeling auf den Transkripten durch
14     """

```

```

15     def __init__(self):
16         self.model = None
17         self.topics = None
18         self.probs = None
19
20     def prepare_documents(self, transcripts):
21         """
22             Bereitet die Transkripte als Dokumente f r Topic
23             Modeling vor
24         """
25
26         documents = []
27         metadata = []
28
29         for i, transcript in enumerate(transcripts, 1):
30             # Jedes Transkript als ein Dokument
31             doc = ' '.join(transcript)
32             documents.append(doc)
33             metadata.append(f'Transkript {i}')
34
35             # Alternativ: Jede uerung als Dokument
36             # for j, utterance in enumerate(transcript):
37             #     documents.append(utterance)
38             #     metadata.append(f'T{i}_U{j}')
39
40         return documents, metadata
41
42     def fit_model(self, documents):
43         """
44             Trainiert das Topic Model
45         """
46
47         # Benutzerdefinierte Stopw rter
48         stopwords = ['bitte', 'danke', 'gern', 'mal', 'noch',
49                     'ja', 'nein']
50
51         vectorizer = CountVectorizer(stop_words=stopwords)

52         self.model = BERTopic(
53             embedding_model="paraphrase-multilingual-MiniLM-
54             L12-v2",
55             vectorizer_model=vectorizer,
56             verbose=True,

```

```

52     nr_topics='auto'
53 )
54
55     self.topics, self.probs = self.model.fit_transform(
56         documents)
57
58     return self.topics, self.probs
59
60
61 def visualize_topics(self):
62 """
63     Visualisiert die gefundenen Themen
64 """
65
66     if self.model is None:
67
68         return
69
70
71     fig = self.model.visualize_topics()
72     fig.write_html("topic_visualization.html")
73
74
75     # Statistik
76     topic_counts = pd.Series(self.topics).value_counts()
77     print("\n==== Themenverteilung ===")
78
79     for topic, count in topic_counts.items():
80
81         if topic == -1:
82
83             print(f"Outlier: {count} Dokumente")
84
85         else:
86
87             words = self.model.get_topic(topic)[:5]
88             words_str = ', '.join([w for w, _ in words])
89             print(f"Thema {topic}: {count} Dokumente - {words_str}")
90
91
92     def visualize_topics_over_time(self, documents,
93         timestamps):
94
95         """
96         Visualisiert Themenentwicklung über die Zeit
97         """
98
99         if self.model is None:
100
101             return
102
103
104         topics_over_time = self.model.topics_over_time(
105             documents,
106             timestamps,
107             )

```

```

89         nr_bins=8
90     )
91
92     fig = self.model.visualize_topics_over_time(
93         topics_over_time)
94     fig.write_html("topics_over_time.html")
95
96 def demonstrate_topic_modeling(transcripts):
97     """
98     Demonstriert Topic Modeling an den Transkripten
99     """
100    modeler = TranscriptTopicModeler()
101    documents, metadata = modeler.prepare_documents(
102        transcripts)
103
104    print("\n==== Topic Modeling der acht Transkripte ===")
105    topics, probs = modeler.fit_model(documents)
106
107    for i, (doc, topic, prob, meta) in enumerate(zip(
108        documents, topics, probs, metadata)):
109        if topic != -1:
110            words = modeler.model.get_topic(topic)[:3]
111            words_str = ', '.join([w for w, _ in words])
112            print(f"{meta}: Thema {topic} (Confidence: {prob
113                :.2f}) - {words_str}")
114        else:
115            print(f"{meta}: Kein klares Thema (Outlier)")
116
117    modeler.visualize_topics()
118    return modeler
119
120
121 # Didaktischer Hinweis
122 print("\n" + "="*70)
123 print("DIDAKTISCHER HINWEIS ZUM TOPIC MODELING")
124 print("=". * 70)
125 print("Topic Modeling identifiziert latente Themen in
126      Textkorpora.")
127 print("Bei nur acht Dokumenten ist die Themenfindung instabil
128      .")
129 print("Die Ergebnisse dienen daher nur der Veranschaulichung")

```

```
der")  
123 print("Methodik, nicht der inhaltlichen Analyse.")
```

Listing 12: Topic Modeling mit BERTopic

3.5 Rhetorical Structure Theory (RST)

Für die Analyse der argumentativen Struktur wird ein RST-Parser implementiert:

```

27     'Contrast': ['aber', 'jedoch', 'hingegen', 'dagegen'],
28     'Cause': ['weil', 'da', 'denn', 'deshalb', 'daher'],
29     'Condition': ['wenn', 'falls', 'sofern'],
30     'Purpose': ['um zu', 'damit'],
31     'Sequence': ['dann', 'danach', 'zuerst', 'schließlich']
32 }
33
34 def __init__(self):
35     self.relations = []
36     self.graph = nx.DiGraph()
37
38 def segment_transcript(self, transcript):
39     """
40     Segmentiert ein Transkript in elementare
41     Diskuseinheiten (EDUs)
42     Vereinfacht: Jede Wörterung ist eine EDU
43     """
44
45     return transcript
46
47 def identify_relations(self, segments):
48     """
49     Identifiziert RST-Relationen zwischen Segmenten
50     Didaktisch vereinfachte Implementierung
51     """
52
53     relations = []
54
55     for i in range(len(segments)-1):
56         current = segments[i]
57         next_seg = segments[i+1]
58
59         # Prüfe auf Cue Phrases
60         for rel_type, cues in self.cue_phrases.items():
61             for cue in cues:
62                 if cue in current.lower() or cue in
next_seg.lower():
63                     relations.append(RSTRelation(
64                         type_name=rel_type,
65                         nucleus=i,

```

```

63                     satellite=i+1
64             ))
65             break
66
67     # Standard: Sequenz-Relation
68     if i < len(segments)-1:
69         relations.append(RSTRelation(
70             type_name='Sequence',
71             nucleus=i,
72             satellite=i+1
73         ))
74
75     return relations
76
77 def build_tree(self, segments, relations):
78 """
79     Baut einen RST-Baum aus den identifizierten
80     Relationen
81 """
82
83     # Knoten hinzuf gen
84     for i, seg in enumerate(segments):
85         self.graph.add_node(i, text=seg[:30] + '...' if
86                             len(seg) > 30 else seg)
87
88     # Kanten hinzuf gen
89     for rel in relations:
90         self.graph.add_edge(rel.nucleus, rel.satellite,
91                             relation=rel.type)
92
93     return self.graph
94
95 def parse(self, transcript):
96 """
97     Vollst ndige RST-Analyse eines Transkripts
98 """
99
100    segments = self.segment_transcript(transcript)
101    relations = self.identify_relations(segments)
102    tree = self.build_tree(segments, relations)

```

```

101
102     return {
103         'segments': segments,
104         'relations': relations,
105         'tree': tree
106     }
107
108     def visualize(self, title="RST-Struktur"):
109         """
110             Visualisiert den RST-Baum
111         """
112         pos = nx.spring_layout(self.graph)
113         plt.figure(figsize=(12, 8))
114
115         # Knoten zeichnen
116         nx.draw_networkx_nodes(self.graph, pos, node_color='lightblue',
117                               node_size=500)
118
119         # Kanten zeichnen mit Relationstyp als Label
120         for edge in self.graph.edges(data=True):
121             nx.draw_networkx_edges(self.graph, pos, [(edge
122                 [0], edge[1])])
123             nx.draw_networkx_edge_labels(
124                 self.graph, pos,
125                 {(edge[0], edge[1]): edge[2]['relation']}
126             )
127
128         # Knotenlabels
129         labels = {node: f"{node}: {self.graph.nodes[node]['text']}"
130                   for node in self.graph.nodes()}
131         nx.draw_networkx_labels(self.graph, pos, labels,
132                               font_size=8)
133
134         plt.title(title)
135         plt.axis('off')
136         plt.tight_layout()
137         plt.savefig('rst_structure.png', dpi=150)
138         plt.show()

```

```

137
138 def demonstrate_rst_analysis(transcripts):
139     """
140     Demonstriert RST-Analyse an den Transkripten
141     """
142     parser = SimpleRSTParser()
143
144     print("\n==== RST-Analyse der Transkripte ===")
145
146     for i, transcript in enumerate(transcripts, 1):
147         print(f"\nTranskript {i}:")
148         result = parser.parse(transcript)
149
150         # Zeige identifizierte Relationen
151         for rel in result['relations'][:5]: # Nur erste 5
152             seg1 = result['segments'][rel.nucleus][:20] + ,
153                         ...
154             seg2 = result['segments'][rel.satellite][:20] + ,
155                         ...
156             print(f" {rel.type}: {seg1} {seg2}")
157
158         if i == 1: # Nur erstes Transkript visualisieren
159             parser.visualize(f"RST-Struktur Transkript {i}")
160
161     return parser
162
163 # Didaktischer Hinweis
164 print("\n" + "="*70)
165 print("DIDAKTISCHER HINWEIS ZUR RST-ANALYSE")
166 print("="*70)
167 print("Die hier implementierte RST-Analyse ist stark")
168     vereinfacht.)
169 print("Ein vollständiger RST-Parser wäre:")
170 print("    Aufwendige manuelle Annotation erfordern")
171 print("    Mit trainierten neuronalen Modellen arbeiten")
172 print("    Mehrere Hierarchieebenen von Diskursrelationen")
173     berücksichtigen")

```

Listing 13: Rhetorical Structure Theory Parser

3.6 Integration der Komponenten in Szenario C

Die vollständige Integration aller Komponenten in Szenario C:

```

1 """
2 Szenario C: Vollständige computerlinguistische Integration
3 Didaktische Implementierung für Lehrzwecke
4 """
5
6 import os
7 import json
8 from datetime import datetime
9
10 class ScenarioC:
11     """
12         Integriert alle computerlinguistischen Komponenten:
13         - Speech Act Recognition
14         - Word Embeddings / Semantische Analyse
15         - Topic Modeling
16         - RST-Analyse
17     """
18
19     def __init__(self, transcripts, terminal_chains):
20         self.transcripts = transcripts
21         self.terminal_chains = terminal_chains
22         self.results = {}
23
24         print("\n" + "="*70)
25         print("SZENARIO C: COMPUTERLINGUISTISCHE INTEGRATION")
26             )
27         print("="*70)
28         print("\nDieses Szenario demonstriert die Anwendung")
29         print("computerlinguistischer Verfahren auf die acht")
30             )
31         print("Transkripte. Alle Ergebnisse dienen")
32             didaktischen")
33         print("Zwecken und sind nicht empirisch valide.\n")
34
35     def run_speech_act_recognition(self):
36         """
37             Führt die Speech Act Recognition aus

```

```

35
36     """
37     print("\n--- Speech Act Recognition ---")
38     utterances, labels, encoder = prepare_speech_act_data
39     (
40         self.transcripts, self.terminal_chains
41     )
42
43     model, tokenizer, label_encoder =
44         train_speech_act_model(
45             utterances, labels, epochs=5
46         )
47
48     self.results['speech_act'] = {
49         'model': model,
50         'tokenizer': tokenizer,
51         'label_encoder': label_encoder,
52         'num_classes': len(label_encoder.classes_)
53     }
54
55     return self.results['speech_act']
56
57
58     def run_semantic_analysis(self):
59         """
60         F hrt die semantische hnlichkeitsanalyse aus
61         """
62         print("\n--- Semantische hnlichkeitsanalyse ---")
63         analyzer = SemanticAnalyzer()
64
65         # Sammle alle uerungen
66         all_utterances = []
67         for transcript in self.transcripts:
68             all_utterances.extend(transcript)
69
70         # hnlichkeitsmatrix
71         sim_matrix = analyzer.similarity_matrix(
72             all_utterances[:20]) # Nur erste 20
73
74         self.results['semantic'] = {
75             'analyzer': analyzer,
76             'utterances': all_utterances,
77         }

```

```

72         'similarity_matrix': sim_matrix
73     }
74
75     return self.results['semantic']
76
77 def run_topic_modeling(self):
78     """
79     F hrt das Topic Modeling aus
80     """
81     print("\n--- Topic Modeling ---")
82     modeler = TranscriptTopicModeler()
83     documents, metadata = modeler.prepare_documents(self.
84         transcripts)
85     topics, probs = modeler.fit_model(documents)
86     modeler.visualize_topics()
87
88     self.results['topic'] = {
89         'modeler': modeler,
90         'topics': topics,
91         'probabilities': probs,
92         'documents': documents,
93         'metadata': metadata
94     }
95
96     return self.results['topic']
97
98 def run_rst_analysis(self):
99     """
100     F hrt die RST-Analyse aus
101     """
102     print("\n--- RST-Analyse ---")
103     parser = SimpleRSTParser()
104
105     rst_results = []
106     for i, transcript in enumerate(self.transcripts, 1):
107         result = parser.parse(transcript)
108         rst_results.append({
109             'transcript_id': i,
110             'segments': result['segments'],
111             'relations': [(r.type, r.nucleus, r.satellite]

```

```

                ) for r in result['relations']]
111        })
112
113        if i == 1:
114            parser.visualize(f"RTS-Struktur Transkript {i}
115                           }")
116
117        self.results['rst'] = rst_results
118
119    def run_all(self):
120        """
121            F hrt alle Analysen aus
122        """
123
124        self.run_speech_act_recognition()
125        self.run_semantic_analysis()
126        self.run_topic_modeling()
127        self.run_rst_analysis()
128
129        # Zusammenfassung
130        print("\n" + "="*70)
131        print("ZUSAMMENFASSUNG SZENARIO C")
132        print("="*70)
133        print(f"    Speech Act Recognition: {self.results['
134                         speech_act']['num_classes']} Klassen")
135        print(f"    Semantische Analyse: {len(self.results['
136                         semantic']['utterances'])} uerungen ")
137        print(f"    Topic Modeling: {len(set(self.results['
138                         topic']['topics']))} Themen")
139        print(f"    RST-Analyse: {len(self.results['rst'])}
140                           Transkripte analysiert")
141
142        return self.results
143
144    # Didaktische Ausf hrung
145    def run_scenario_c_demonstration():
146        """
147            F hrt die vollst ndige Demonstration von Szenario C aus
148        """
149
150        # Lade die Transkripte

```

```

145     from ars_data import transcripts, terminal_chains
146
147     # Augmentiere die Daten f r didaktische Zwecke
148     augmented_transcripts = augment_transcripts_for_teaching(
149         transcripts, factor=10)
150     augmented_chains = augment_transcripts_for_teaching(
151         terminal_chains, factor=10)
152
153     print("\n" + "="*70)
154     print("DIDAKTISCHE AUGMENTIERUNG")
155     print("="*70)
156     print(f"Original: {len(transcripts)} Transkripte")
157     print(f"Augmentiert: {len(augmented_transcripts)}"
158           " Transkripte")
159
160
161     # F hre Szenario C aus
162     scenario = ScenarioC(augmented_transcripts,
163                           augmented_chains)
164     results = scenario.run_all()
165
166
167     # Speichere Ergebnisse
168     with open('scenario_c_results.json', 'w') as f:
169         # Konvertiere nicht-serialisierbare Objekte
170         serializable = {
171             'speech_act': {'num_classes': results['speech_act']
172                            ['num_classes']},
173             'semantic': {'num_utterances': len(results['semantic']
174                                         ['utterances'])},
175             'topic': {'num_topics': len(set(results['topic']
176                                            ['topics']))},
177             'rst': results['rst']
178         }
179         json.dump(serializable, f, indent=2)
180
181
182         print("\nErgebnisse gespeichert in 'scenario_c_results."
183               "json'")
184
185
186     return results
187
188
189 if __name__ == "__main__":

```

```
177     run_scenario_c_demonstration()
```

Listing 14: Szenario C - Vollständige Integration

4 Szenario D: Hybride Modellierung

Szenario D integriert computerlinguistische Verfahren mit den interpretativ gebildeten Kategorien der ARS 3.0. Es überspringt die vollständige Automatisierung der Kategorienbildung (Szenario C) und nutzt die neuen Verfahren komplementär.

4.1 CRF für sequenzielle Abhängigkeiten

Conditional Random Fields modellieren Abhängigkeiten der Sprechakte vom weiteren Kontext:

```
1 """
2 Conditional Random Fields (CRF) f r sequenzielle
3      Abh ngigkeiten
4 Didaktische Implementierung mit sklearn-crfsuite
5 """
6
7 import sklearn_crfsuite
8 from sklearn_crfsuite import metrics
9 import numpy as np
10
11 class CRFSequenceModel:
12     """
13         CRF-Modell f r die Sequenzmodellierung der
14             Terminalzeichen
15     """
16
17     def __init__(self):
18         self.crf = sklearn_crfsuite.CRF(
19             algorithm='lbfgs',
20             c1=0.1, # L1-Regularisierung
21             c2=0.1, # L2-Regularisierung
22             max_iterations=100,
23             all_possible_transitions=True
24         )
25         self.label_encoder = None
```

```

24
25     def word2features(self, tokens, i):
26         """
27             Erzeugt Features f r Position i in der Sequenz
28         """
29         word = tokens[i]
30
31         features = {
32             'bias': 1.0,
33             'word': word,
34             'word.is_first': i == 0,
35             'word.is_last': i == len(tokens) - 1,
36             'word.prefix_K': word.startswith('K'),
37             'word.prefix_V': word.startswith('V'),
38             'word.suffix_A': word.endswith('A'),
39             'word.suffix_B': word.endswith('B'),
40             'word.suffix_E': word.endswith('E'),
41             'word.suffix_G': word.endswith('G'),
42             'word.suffix_V': word.endswith('V'),
43         }
44
45         # Kontext-Features
46         if i > 0:
47             word_prev = tokens[i-1]
48             features.update({
49                 '-1:word': word_prev,
50                 '-1:word.prefix_K': word_prev.startswith('K')
51                         ,
52                 '-1:word.prefix_V': word_prev.startswith('V')
53                         ,
54                 '-1:word.suffix_A': word_prev.endswith('A'),
55             })
56         else:
57             features['BOS'] = True
58
59         if i < len(tokens) - 1:
60             word_next = tokens[i+1]
61             features.update({
62                 '+1:word': word_next,
63                 '+1:word.prefix_K': word_next.startswith('K')
64             })

```

```

        ,
        '+1:word.prefix_V': word_next.startswith('V')
        ,
        '+1:word.suffix_A': word_next.endswith('A'),
    })
else:
    features['EOS'] = True

return features

def extract_features(self, sequences):
    """
    Extrahiert Features f r alle Sequenzen
    """
    X = []
    for seq in sequences:
        X.append([self.word2features(seq, i) for i in
                  range(len(seq))])
    return X

def fit(self, sequences, labels):
    """
    Trainiert das CRF-Modell
    """
    X = self.extract_features(sequences)
    self.crf.fit(X, labels)
    return self

def predict(self, sequences):
    """
    Sagt Labels f r neue Sequenzen vorher
    """
    X = self.extract_features(sequences)
    return self.crf.predict(X)

def evaluate(self, test_sequences, test_labels):
    """
    Evaluiert das Modell
    """
    pred = self.predict(test_sequences)

```

```

99
100     # Flatten f r Metriken
101     y_true = [label for seq in test_labels for label in
102               seq]
103     y_pred = [label for seq in pred for label in seq]
104
105     return {
106         'accuracy': np.mean(np.array(y_true) == np.array(
107             y_pred)),
108         'classification_report': metrics.
109             flat_classification_report(
110                 test_labels, pred, labels=sorted(set(y_true)))
111     }
112
113
114
115 def demonstrate_crf(terminal_chains):
116     """
117     Demonstriert CRF-Modellierung auf den Terminalzeichen
118     """
119
120     print("\n==== CRF-Modellierung der Terminalzeichen ===")
121
122     # Train-Test-Split (didaktisch)
123     train_size = int(len(terminal_chains) * 0.7)
124     train_chains = terminal_chains[:train_size]
125     test_chains = terminal_chains[train_size:]
126
127     # Features extrahieren
128     model = CRFSequenceModel()
129     X_train = model.extract_features(train_chains)
130
131     # Training
132     print(f"Trainiere CRF mit {len(train_chains)} Sequenzen
133         ...")
134     model.fit(train_chains, train_chains) # Labels sind die
135         Sequenzen selbst
136
137     # Evaluation
138     results = model.evaluate(test_chains, test_chains)
139     print(f"\nGenauigkeit: {results['accuracy']:.3f}")
140
141
142
143

```

```
134     return model
```

Listing 15: CRF für sequenzielle Abhängigkeiten

4.2 Transformer-Embeddings als Ergänzung

Transformer-Embeddings werden zusätzlich zu den kategorialen Terminalzeichen verwendet:

```
1 """
2 Transformer-Embeddings als Ergänzung zu kategorialen
3   Terminalzeichen
4
5 import torch
6 import numpy as np
7 from sentence_transformers import SentenceTransformer
8
9 class TerminalEmbeddingEnricher:
10   """
11     Ergänzt Terminalzeichen um semantische Embeddings der
12       zugrundeliegenden Tolerungen
13   """
14
15   def __init__(self, model_name='paraphrase-multilingual-
16     MiniLM-L12-v2'):
17     self.model = SentenceTransformer(model_name)
18     self.symbol_to_embedding = {}
19     self.symbol_to_text = self._create_symbol_mapping()
20
21   def _create_symbol_mapping(self):
22     """
23       Erstellt eine Mapping von Terminalzeichen zu
24         Beispieltextrnen
25     """
26
27     return {
28       'KBG': ['Guten Tag', 'Guten Morgen', 'Hallo'],
29       'VBG': ['Guten Tag', 'Guten Morgen', 'Hallo
30             zur ck'],
31       'KBBd': ['Einmal Leberwurst', 'Ich hätte gerne
32             K se', 'Ein Kilo pfel bitte'],
```

```

27     'VBBd': ['Wie viel darf es sein?', 'Welche Sorte?',
28                 ', 'Sonst noch etwas?'],
29     'KBA': ['Zweihundert Gramm', 'Die hellen bitte',
30               ', Ja, gerne'],
31     'VBA': ['In Ordnung', 'Kommt sofort', 'Alles klar
32               '],
33     'KAE': ['Kann ich das in Salat tun?', 'Woher
34               kommen die?', 'Ist das frisch?'],
35     'VAE': ['Eher anbraten', 'Aus der Region', 'Ja,
36               ganz frisch'],
37     'KAA': ['Bitte', 'Danke', 'Ja, danke'],
38     'VAA': ['Das macht 8 Mark 20', '3 Mark bitte', '
39               14 Mark 19'],
40     'KAV': ['Auf Wiedersehen', 'Tsch ss', 'Einen
41               sch nen Tag'],
42     'VAV': ['Vielen Dank', 'Sch nen Tag noch', 'Auf
43               Wiedersehen']
44   }
45
46
47   def get_embedding(self, symbol):
48     """
49       Gibt das Embedding f r ein Terminalzeichen zur ck
50     """
51
52     if symbol in self.symbol_to_embedding:
53       return self.symbol_to_embedding[symbol]
54
55     # Durchschnitt der Beispieltext-Embeddings
56     texts = self.symbol_to_text.get(symbol, [symbol])
57     embeddings = self.model.encode(texts)
58     avg_embedding = np.mean(embeddings, axis=0)
59
60     self.symbol_to_embedding[symbol] = avg_embedding
61     return avg_embedding
62
63
64   def enrich_sequence(self, sequence):
65     """
66       Erweitert eine Sequenz von Terminalzeichen um
67       Embeddings
68     """
69
70     symbols = sequence

```

```

58     embeddings = np.array([self.get_embedding(sym) for
59                           sym in symbols])
60
61     return {
62         'symbols': symbols,
63         'embeddings': embeddings,
64         'combined': np.column_stack([
65             self._one_hot_encode(symbols),
66             embeddings
67         ]) if len(symbols) > 0 else np.array([])
68     }
69
70     def _one_hot_encode(self, symbols):
71         """
72             One-Hot-Encoding der Terminalzeichen
73         """
74
75         unique_symbols = sorted(set(self.symbol_to_text.keys()
76                               ()))
77
78         symbol_to_idx = {sym: i for i, sym in enumerate(
79                         unique_symbols)}
80
81         one_hot = np.zeros((len(symbols), len(unique_symbols))
82                           )
83         for i, sym in enumerate(symbols):
84             if sym in symbol_to_idx:
85                 one_hot[i, symbol_to_idx[sym]] = 1
86
87         return one_hot
88
89     def demonstrate_embedding_enrichment():
90         """
91             Demonstriert die Anreicherung von Terminalzeichen mit
92             Embeddings
93         """
94
95         enricher = TerminalEmbeddingEnricher()
96
97
98         print("\n==== Anreicherung der Terminalzeichen mit
99             Embeddings ===")
100
101        # Beispieldsequenz

```

```

92     sequence = [ 'KBG' , 'VBG' , 'KBBd' , 'VBBd' , 'KBA' ]
93
94     enriched = enricher.enrich_sequence(sequence)
95
96     print(f"\nSequenz: {''''.join(sequence)}")
97     print(f"Embedding-Dimension: {enriched['embeddings'].
98           shape[1]}")
99     print(f"One-Hot-Dimension: {enriched['combined'].shape[1]
100        - enriched['embeddings'].shape[1]}")
101    print(f"Kombinierte Dimension: {enriched['combined'].
102          shape[1]}")
103
104    return enricher

```

Listing 16: Transformer-Embeddings für Terminalzeichen

4.3 Graph Neural Networks für die Nonterminal-Hierarchie

Die Nonterminal-Hierarchie wird als Graph Neural Network modelliert:

```

1 """
2 Graph Neural Network f r die Nonterminal-Hierarchie
3 Didaktische Implementierung mit PyTorch Geometric
4 """
5
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9 from torch_geometric.nn import GCNConv, GATConv
10 from torch_geometric.data import Data
11 import networkx as nx
12
13 class GrammarGNN(nn.Module):
14 """
15     Graph Neural Network f r die Grammatik-Hierarchie
16 """
17
18     def __init__(self, input_dim, hidden_dim=64, num_classes
19                  =12):
20         super().__init__()
21         self.conv1 = GCNConv(input_dim, hidden_dim)

```

```

21         self.conv2 = GCNConv(hidden_dim, hidden_dim)
22         self.classifier = nn.Linear(hidden_dim, num_classes)
23
24     def forward(self, x, edge_index):
25         x = self.conv1(x, edge_index)
26         x = F.relu(x)
27         x = F.dropout(x, training=self.training)
28         x = self.conv2(x, edge_index)
29         x = F.relu(x)
30         x = self.classifier(x)
31         return F.log_softmax(x, dim=1)
32
33 class GrammarHierarchyGNN:
34     """
35     Verwaltet das GNN f r die Nonterminal-Hierarchie
36     """
37
38     def __init__(self, grammar_rules):
39         self.grammar = grammar_rules
40         self.graph = self._build_graph()
41         self.model = None
42
43     def _build_graph(self):
44         """
45         Baut einen Graphen aus der Grammatik-Hierarchie
46         """
47         G = nx.DiGraph()
48
49         # Knoten: Terminale und Nonterminale
50         all_symbols = set()
51
52         # Nonterminale als Knoten
53         for nt, productions in self.grammar.items():
54             all_symbols.add(nt)
55             for prod, _ in productions:
56                 for sym in prod:
57                     all_symbols.add(sym)
58
59         # Kanten: Ableitungsrelationen
60         for nt, productions in self.grammar.items():

```

```

61         for prod, prob in productions:
62             for sym in prod:
63                 G.add_edge(nt, sym, weight=prob)
64
65     return G
66
67 def prepare_data(self):
68     """
69     Bereitet die Daten f r das GNN vor
70     """
71
72     # Knoten-Indizes
73     nodes = list(self.graph.nodes())
74     node_to_idx = {node: i for i, node in enumerate(nodes)}
75
76     # Feature-Matrix (vereinfacht: One-Hot)
77     x = torch.eye(len(nodes))
78
79     # Kanten-Index
80     edge_index = []
81     for u, v, data in self.graph.edges(data=True):
82         edge_index.append([node_to_idx[u], node_to_idx[v]])
83
84     edge_index = torch.tensor(edge_index, dtype=torch.long).t().contiguous()
85
86     return Data(x=x, edge_index=edge_index)
87
88 def train(self, epochs=100):
89     """
90     Trainiert das GNN
91     """
92
93     data = self.prepare_data()
94     self.model = GrammarGNN(input_dim=data.x.shape[1])
95
96     optimizer = torch.optim.Adam(self.model.parameters(), lr=0.01)
97
98     print("\n==== Training des Grammar GNN ====")

```

```

97     for epoch in range(epochs):
98         self.model.train()
99         optimizer.zero_grad()
100        out = self.model(data.x, data.edge_index)
101
102        # Selbst berwachtes Lernen: Rekonstruktion des
103        # Graphen
104        # Vereinfacht: Vorhersage der Nachbarn
105        loss = F.nll_loss(out[data.edge_index[0]], data.
106                           edge_index[1])
107
108        loss.backward()
109        optimizer.step()
110
111
112        if epoch % 20 == 0:
113            print(f"Epoch {epoch}: Loss = {loss.item():.4
114                           f}")
115
116
117        return self.model
118
119
120    def demonstrate_gnn(grammar_rules):
121        """
122        Demonstriert das GNN f r die Grammatik-Hierarchie
123        """
124
125        print("\n==== Graph Neural Network f r Nonterminal-
126             Hierarchie ===")
127
128        gnn = GrammarHierarchyGNN(grammar_rules)
129        print(f"Graph: {gnn.graph.number_of_nodes()} Knoten, "
130              f"{gnn.graph.number_of_edges()} Kanten")
131
132
133        model = gnn.train(epochs=100)
134
135
136        return gnn, model

```

Listing 17: Graph Neural Network für Nonterminal-Hierarchie

4.4 Attention-Mechanismen für relevante Vorgänger

Attention-Mechanismen identifizieren besonders relevante Vorgänger für aktuelle Entscheidungen:

```
1 """
2 Attention-Mechanismen f r die Identifikation relevanter
3 Vorg nger
4 """
5 import torch
6 import torch.nn as nn
7 import torch.nn.functional as F
8 import numpy as np
9
10 class SequenceAttention(nn.Module):
11     """
12         Attention-Mechanismus f r Sequenzmodellierung
13     """
14
15     def __init__(self, embedding_dim, hidden_dim=64):
16         super().__init__()
17         self.embedding_dim = embedding_dim
18         self.hidden_dim = hidden_dim
19
20         # Attention-Parameter
21         self.W_q = nn.Linear(embedding_dim, hidden_dim, bias=
22                             False)
23         self.W_k = nn.Linear(embedding_dim, hidden_dim, bias=
24                             False)
25         self.W_v = nn.Linear(embedding_dim, hidden_dim, bias=
26                             False)
27         self.scale = hidden_dim ** 0.5
28
29     def forward(self, x, mask=None):
30         """
31             x: (seq_len, batch, embedding_dim)
32             """
33
34             # Query, Key, Value berechnen
35             Q = self.W_q(x) # (seq_len, batch, hidden_dim)
36             K = self.W_k(x) # (seq_len, batch, hidden_dim)
```

```

33     V = self.W_v(x)  # (seq_len, batch, hidden_dim)
34
35     # Attention-Scores
36     scores = torch.matmul(Q.transpose(0, 1), K.transpose
37         (0, 1).transpose(1, 2))
38     scores = scores / self.scale
39
40     if mask is not None:
41         scores = scores.masked_fill(mask == 0, -1e9)
42
43     # Attention-Gewichte
44     attention_weights = F.softmax(scores, dim=-1)
45
46     # Gewichtete Summe
47     context = torch.matmul(attention_weights, V.transpose
48         (0, 1))
49
50     return context, attention_weights
51
52
53 class SymbolPredictorWithAttention(nn.Module):
54     """
55     Sagt das n chste Symbol mit Attention auf die Vorg nger
56     vorher
57     """
58
59     def __init__(self, num_symbols, embedding_dim=50,
60                  hidden_dim=64):
61         super().__init__()
62         self.embedding = nn.Embedding(num_symbols,
63                                     embedding_dim)
64         self.attention = SequenceAttention(embedding_dim,
65                                           hidden_dim)
66         self.lstm = nn.LSTM(embedding_dim, hidden_dim,
67                            batch_first=True)
68         self.classifier = nn.Linear(hidden_dim +
69                                     embedding_dim, num_symbols)
70
71     def forward(self, x):
72         """
73         x: (batch, seq_len) mit Symbol-Indizes

```

```

65 """
66 # Embeddings
67 embedded = self.embedding(x) # (batch, seq_len,
68 # embedding_dim)
69
70 # LSTM f r sequenzielle Abh ngigkeiten
71 lstm_out, (hidden, cell) = self.lstm(embedded)
72
73 # Attention ber die Sequenz
74 # Transponieren f r Attention (seq_len, batch,
75 # embedding_dim)
76 context, attention_weights = self.attention(embedded.
77 transpose(0, 1))
78
79 # Kombiniere letzten LSTM-State mit Attention-Context
80 last_hidden = hidden[-1] # (batch, hidden_dim)
81 last_context = context[-1] # (batch, hidden_dim)
82
83 # Vorhersage
84 combined = torch.cat([last_hidden, last_context], dim
85 =-1)
86 logits = self.classifier(combined)
87
88 return logits, attention_weights
89
90
91 def demonstrate_attention(terminal_chains, symbol_to_idx):
92 """
93 Demonstriert Attention-Mechanismen auf den Sequenzen
94 """
95 print("\n==== Attention-Mechanismen f r relevante
96 Vorg nger ===")
97
98 # Daten vorbereiten
99 sequences = []
100 for chain in terminal_chains:
101     seq = [symbol_to_idx[sym] for sym in chain]
102     sequences.append(seq)
103
104 # Padding f r Batch-Verarbeitung
105 from torch.nn.utils.rnn import pad_sequence

```

```

100     sequences_padded = pad_sequence([torch.tensor(seq) for
101         seq in sequences],
102                                         batch_first=True,
103                                         padding_value=0)
104
105
106 # Modell initialisieren
107 model = SymbolPredictorWithAttention(num_symbols=len(
108     symbol_to_idx))
109
110 # Forward-Pass
111 logits, attention_weights = model(sequences_padded[:2])
112     # Nur erste 2 Sequenzen
113
114 print(f"\nInput-Shape: {sequences_padded[:2].shape}")
115 print(f"Attention-Weights-Shape: {attention_weights.shape
116     }")
117 print(f"Logits-Shape: {logits.shape}")
118
119 # Visualisierung der Attention-Gewichte
120 plot_attention_weights(attention_weights[0].detach().
121     numpy(),
122     sequences[0], sequences[0])
123
124 return model
125
126
127 def plot_attention_weights(attention, source_tokens,
128     target_tokens):
129     """
130     Visualisiert Attention-Gewichte als Heatmap
131     """
132
133     import matplotlib.pyplot as plt
134     import seaborn as sns
135
136     plt.figure(figsize=(10, 8))
137     sns.heatmap(attention[:len(target_tokens), :len(
138         source_tokens)],
139                 xticklabels=source_tokens,
140                 yticklabels=target_tokens,
141                 cmap='viridis', annot=True, fmt='.2f')
142
143     plt.title('Attention-Gewichte zwischen Vorgängern und

```

```

    Vorhersage')
132 plt.xlabel('Vorg nger-Symbole')
133 plt.ylabel('Vorhersage-Position')
134 plt.tight_layout()
135 plt.savefig('attention_weights.png', dpi=150)
136 plt.show()

```

Listing 18: Attention-Mechanismen für Sequenzmodellierung

4.5 Integration der Komponenten in Szenario D

Die vollständige Integration aller Komponenten in Szenario D:

```

1 """
2 Szenario D: Hybride Modellierung
3 Integration computerlinguistischer Verfahren mit
4     interpretativen Kategorien
5 """
6
7 import json
8 import numpy as np
9
10
11 class ScenarioD:
12     """
13         Integriert computerlinguistische Verfahren komplement r
14             zu den
15         interpretativ gebildeten Kategorien der ARS 3.0
16     """
17
18     def __init__(self, terminal_chains, grammar_rules,
19                  reflection_log):
20         self.terminal_chains = terminal_chains
21         self.grammar_rules = grammar_rules
22         self.reflection_log = reflection_log
23         self.results = []
24
25         print("\n" + "="*70)
26         print("SZENARIO D: HYBRIDE MODELLIERUNG")
27         print("="*70)
28         print("\nDieses Szenario integriert
29             computerlinguistische")

```

```

25     print("Verfahren KOMPLEMENT R zu den interpretativen
26         ")
27     print("Kategorien der ARS 3.0. Die interpretative
28         Basis")
29     print("bleibt erhalten, wird aber durch neue
30         Verfahren")
31     print("angereichert.\n")

32     def run_crf_modeling(self):
33         """
34             F hrt CRF-Modellierung auf den Terminalzeichen durch
35         """
36         print("\n--- CRF-Modellierung ---")
37         crf_model = demonstrate_crf(self.terminal_chains)
38         self.results['crf'] = {'model': crf_model}
39         return crf_model

40     def run_embedding_enrichment(self):
41         """
42             Reichert Terminalzeichen mit Transformer-Embeddings
43             an
44         """
45         print("\n--- Embedding-Anreicherung ---")
46         enricher = demonstrate_embedding_enrichment()

47         # Beispiel f r angereicherte Sequenz
48         example_seq = self.terminal_chains[0][:5]
49         enriched = enricher.enrich_sequence(example_seq)

50         self.results['embeddings'] = {
51             'enricher': enricher,
52             'example': enriched
53         }

54         return enricher

55     def run_gnn_hierarchy(self):
56         """
57             Modelliert die Nonterminal-Hierarchie als GNN
58         """

```

```

61     print("\n--- GNN f r Nonterminal-Hierarchie ---")
62     gnn, model = demonstrate_gnn(self.grammar_rules)
63     self.results['gnn'] = {'gnn': gnn, 'model': model}
64     return gnn, model
65
66 def run_attention_analysis(self):
67     """
68     Analysiert Attention-Mechanismen auf den Sequenzen
69     """
70     print("\n--- Attention-Analyse ---")
71
72     # Symbol zu Index Mapping
73     all_symbols = set()
74     for chain in self.terminal_chains:
75         all_symbols.update(chain)
76     symbol_to_idx = {sym: i for i, sym in enumerate(
77                     sorted(all_symbols))}

78     model = demonstrate_attention(self.terminal_chains,
79                                   symbol_to_idx)
80     self.results['attention'] = {'model': model}
81
82     return model
83
84
85 def run_all(self):
86     """
87     F hrt alle Analysen aus (komplement r, nicht
88     ersetzend)
89     """
90     self.run_crf_modeling()
91     self.run_embedding_enrichment()
92     self.run_gnn_hierarchy()
93     self.run_attention_analysis()
94
95     # Zusammenfassung
96     print("\n" + "="*70)
97     print("ZUSAMMENFASSUNG SZENARIO D")
98     print("="*70)
99     print("    CRF-Modellierung: Sequenzielle
100        Abh ngigkeiten modelliert")

```

```

97     print("      Embedding-Anreicherung: Terminalzeichen
98         semantisch angereichert")
99     print("      GNN-Hierarchie: Nonterminal-Struktur als
100        Graph modelliert")
101    print("      Attention-Analyse: Relevante Vorg nger
102        identifiziert")
103    print("\nDie interpretativen Kategorien der ARS 3.0
104        bleiben")
105    print("die Grundlage aller Analysen. Die
106        computerlinguistischen")
107    print("Verfahren dienen der komplement ren
108        Erkenntnisgewinnung.")

109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126

```

`print(" Embedding-Anreicherung: Terminalzeichen
semantisch angereichert")
print(" GNN-Hierarchie: Nonterminal-Struktur als
Graph modelliert")
print(" Attention-Analyse: Relevante Vorg nger
identifiziert")
print("\nDie interpretativen Kategorien der ARS 3.0
bleiben")
print("die Grundlage aller Analysen. Die
computerlinguistischen")
print("Verfahren dienen der komplement ren
Erkenntnisgewinnung.")

def run_scenario_d_demonstration(terminal_chains,
grammar_rules, reflection_log):
 """
 F hrt die vollst ndige Demonstration von Szenario D aus
 """
 scenario = ScenarioD(terminal_chains, grammar_rules,
 reflection_log)
 results = scenario.run_all()

 # Speichere Ergebnisse
 with open('scenario_d_results.json', 'w') as f:
 # Vereinfachte serialisierbare Version
 serializable = {
 'crf': {'status': 'completed'},
 'embeddings': {'status': 'completed'},
 'gnn': {'num_nodes': results['gnn'][0].graph.
 number_of_nodes()},
 'attention': {'status': 'completed'}
 }
 json.dump(serializable, f, indent=2)

 print("\nErgebnisse gespeichert in 'scenario_d_results.
 json'")

return results`

```

127
128 # Didaktischer Hinweis
129 print("\n" + "="*70)
130 print("METHODOLOGISCHER HINWEIS ZU SZENARIO D")
131 print("="*70)
132 print("Szenario D behält die interpretative Basis der ARS
      3.0 bei.")
133 print("Die computerlinguistischen Verfahren werden
      KOMPLEMENT R")
134 print("eingesetzt, nicht als Ersatz der manuellen
      Kategorienbildung.")
135 print("Dies entspricht der methodologischen Forderung nach")
136 print("Kontrolle und Transparenz im Sinne der XAI-Kriterien."
      )

```

Listing 19: Szenario D - Vollständige hybride Integration

5 Vergleich der Szenarien und methodologische Reflexion

5.1 Gegenüberstellung der Ansätze

Tabelle 2: Vergleich der Szenarien C und D

Kriterium	Szenario C	Szenario D
Kategorienbildung	Automatisch (Speech Act Recognition)	Interpretativ (ARS 3.0)
Datenbasis	Augmentierte Rohdaten	Terminalzeichenketten
Repräsentation	Vektorielle Embeddings	Diskrete Kategorien + Embeddings
Hierarchie	Automatisch gelernt	Explizit induziert (ARS 3.0)
Transparenz	Gering (Black Box)	Hoch (dokumentierte Entscheidungen)
Didaktischer Wert	Funktionsweise neuronaler Verfahren	Integration alter und neuer Methoden
Empirische Validität	Nicht gegeben	Eingeschränkt (basierend auf Interpretation)
Methodologische Kontrolle	Verloren	Erhalten

5.2 Didaktische Erkenntnisse aus Szenario C

Die Implementierung von Szenario C hat gezeigt:

1. **Notwendigkeit großer Datenmengen:** Neuronale Verfahren benötigen für valide Ergebnisse Datenmengen, die weit über die acht Transkripte hinausgehen. Die Augmentierung ermöglicht zwar die Demonstration der Funktionsweise, ersetzt aber keine echten Daten.
2. **Opazität der Entscheidungen:** Die automatisch gelernten Kategorien und Attention-Gewichte sind für Dritte nicht ohne weiteres nachvollziehbar. Die XAI-Kriterien der Verständlichkeit und Transparenz werden verletzt.
3. **Verlust der interpretativen Basis:** Die automatische Speech Act Recognition bildet nicht die qualitativ gehaltvollen Unterscheidungen der ARS ab (z.B. zwischen KBA und KAA), sondern lernt statistische Korrelationen im Vektorraum.

5.3 Didaktische Erkenntnisse aus Szenario D

Die Implementierung von Szenario D hat gezeigt:

1. **Komplementarität statt Substitution:** Die computerlinguistischen Verfahren können wertvolle Zusatzinformationen liefern (z.B. semantische Ähnlichkeiten zwischen verschiedenen Äußerungen), ohne die interpretative Basis zu ersetzen.
2. **Validierungsmöglichkeiten:** Die Embedding-Ähnlichkeiten können zur Validierung der interpretativen Kategorienbildung genutzt werden: Ähnliche Äußerungen sollten ähnliche Terminalzeichen erhalten.
3. **Visualisierung von Abhängigkeiten:** Attention-Mechanismen und CRF-Modelle visualisieren, welche Vorgänger für aktuelle Entscheidungen besonders relevant sind – dies kann die sequenzielle Struktur der Gespräche veranschaulichen.
4. **Methodologische Kontrolle bleibt erhalten:** Da die interpretativen Kategorien die Grundlage bilden, bleiben alle Ergebnisse an die qualitativen Entscheidungen rückgebunden und intersubjektiv prüfbar.

5.4 Fazit für die Lehrpraxis

Die didaktische Exploration der Szenarien C und D führt zu folgenden Schlussfolgerungen:

1. **Szenario C eignet sich zur Demonstration der Funktionsweise neuroaler Verfahren**, sollte aber mit explizitem Hinweis auf die fehlende empirische Validität und die methodologischen Probleme eingesetzt werden.
2. **Szenario D ist methodologisch vorzuziehen**, da es die interpretative Basis erhält und computerlinguistische Verfahren komplementär nutzt. Es vermittelt, wie alte und neue Methoden produktiv verbunden werden können.
3. **Data Augmentation ist ein wertvolles didaktisches Werkzeug**, um mit kleinen Datensätzen die Funktionsweise von Verfahren zu demonstrieren. Der augmentierte Charakter der Daten muss dabei stets transparent gemacht werden.
4. **Die XAI-Kriterien** (Verständlichkeit, Genauigkeit, Wissensgrenzen) bieten einen geeigneten Rahmen, um die verschiedenen Ansätze zu bewerten und ihre Stärken und Schwächen zu reflektieren.

6 Ausblick

Die hier vorgestellten didaktischen Implementierungen können in mehreren Richtungen weiterentwickelt werden:

1. **Erweiterung der Augmentierungsstrategien**: Über einfaches Kopieren hinaus könnten komplexere Augmentierungen (Paraphrasierung, kontrollierte Variation) implementiert werden.
2. **Integration weiterer Verfahren**: z.B. PETRI-Netze für Nebenläufigkeit, Bayessche Netze für Inferenz, oder formale Verifikationsmethoden.
3. **Entwicklung von Vergleichsmetriken**: Wie kann man die Ergebnisse der verschiedenen Szenarien quantitativ vergleichen, ohne die qualitative Basis zu verlieren?
4. **Übertragung auf andere Datensätze**: Die Methodik lässt sich auf andere Interaktionstypen übertragen (Arzt-Patient-Gespräche, Unterrichtsinteraktionen, etc.).

Entscheidend bleibt dabei stets die methodologische Kontrolle: Die formalen Verfahren müssen den interpretativen Charakter der Analyse respektieren und dürfen nicht zu dessen Automatisierung führen.

Literatur

- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbadó, A., García, S., Gil-López, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of NAACL-HLT 2019*, 4171-4186.
- Flick, U. (2019). *Qualitative Sozialforschung: Eine Einführung* (9. Aufl.). Rowohlt.
- Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. *Proceedings of ICML 2001*, 282-289.
- Mann, W. C., & Thompson, S. A. (1988). Rhetorical Structure Theory: Toward a functional theory of text organization. *Text*, 8(3), 243-281.
- Przyborski, A., & Wohlrab-Sahr, M. (2021). *Qualitative Sozialforschung: Ein Arbeitsbuch* (5. Aufl.). De Gruyter Oldenbourg.
- Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Proceedings of EMNLP-IJCNLP 2019*, 3982-3992.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems 30*, 5998-6008.